

# Problem Management Lite

## *Quarterly Deep Dive*

### Agenda

➤ **Objectives:**

<input type="checkbox"/> RCA Review	5 minutes	1:05 – 1:10
<input type="checkbox"/> Mega-Problem Review	5 minutes	1:10 – 1:15
<input type="checkbox"/> Priority 1 & 2 Review	50 minutes	1:15 – 2:05
<input type="checkbox"/> Next Steps	20 minutes	2:10 – 2:25

➤ **Goals:**

- Keep leadership apprised of technical risks
- Identify next steps

**December 19, 2012**

**Stuart Kendrick**

**Problem** A cause, or potential cause, of an incident that has already, or may in the future, interfere with a defined IT service

## RCA: A Year in Review

### Five RCAs

1. Q4 2011 Rhino RCA Phase I
2. Q1 2012 Tungsten RCA Phase I
3. Q2 2012 Tungsten RCA Phase II
4. Q2-Q3 2012 Rhino RCA Phase II

### Results

1. Found and fixed two issues *Success*
2. Identified handful of issues *Success*
3. Identified next steps *Iced*
4. Mitigated third issue *Risk Accepted*

### Inhibitors

Typically, each RCA stumbles during its first month:

1. Free up staff
2. Expand permissions to include RCA techs
3. Hack together monitoring & analysis tools

And thereafter picks up steam.

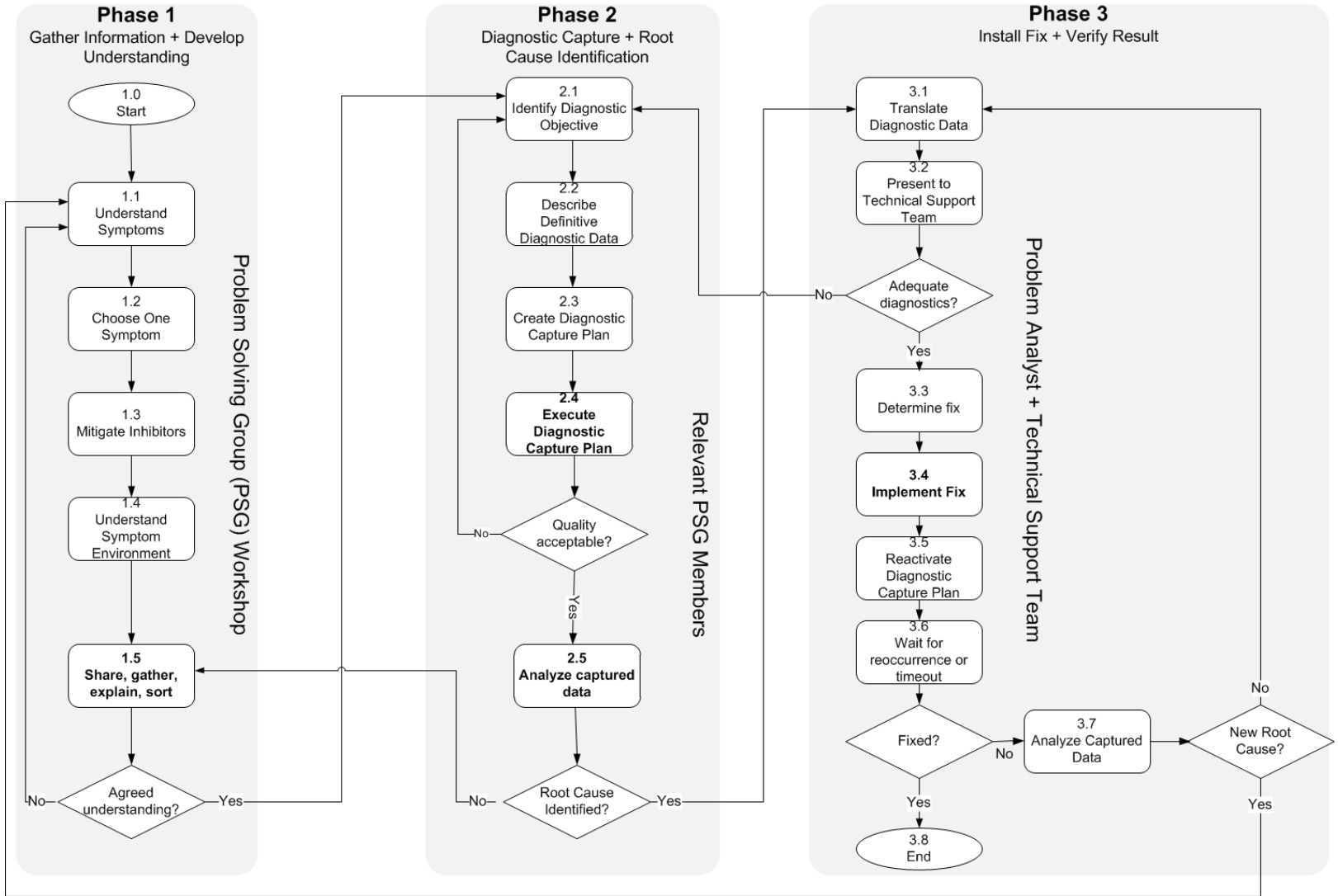
### Future

Kick-off each RCA with a *Mitigate Inhibitors* meeting with leadership, to clarify resource, permissions, and tools constraints

Develop a standardized methodology which adapts to Inhibitors

Refresh toolkit

# RCA Methodology



## Transition

Next we review Mega-Problems

Questions regarding RCAs?

## Criteria for a Mega-Problem

### The Bucket

*The Bucket* is structured for techs and their Supervisors, chunked into mouthfulls which make sense at this level, typically:

- Contained within a single technology group and ideally
- Doable by a single tech

*More work needed to implement this fully*

### Mega-Problems

For leadership, I maintain a separate *Mega-Problem* queue, abstracted from *The Bucket* -- these are the items which I believe are worth leadership attention.

Historically, *Mega-Problems*, and only *Mega-Problems*, are what you saw during a Deep Dive.

*Mega-Problems* typically span technology groups and/or require Project-level intervention to address.

### Criteria

To become a *Mega-Problem*, the issue must offer a noticeable ... my judgment ... chance of inflicting at least one of the following:

1. Disrupt services Center-wide
2. Disrupt our core business: *Grant/contract funded research*
3. Impact Center leadership (reputational issues)

### Synergy

Most of these *Mega-Problems* encapsulate multiple *Bucket Entries*: typically when those entries interact synergistically: domino effect

One day, we may see a one-to-one mapping between a *Mega-Problem* and a *Bucket Entry* (aka a *Problem*)

... but this hasn't happened yet.

One-third of the issues in the Queue belong to a Mega-Problem

# Mega-Problem Status

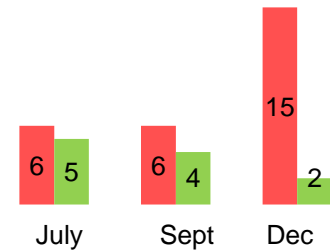
## Tungsten Stumbles

Owner + Implementer = InfraOps + InfraOps  
 Roadmap = Data & Storage



## Data Center Redundancy Degrading

Owner + Implementer = InfraOps + InfraOps  
 Roadmaps = Data & Storage, Server, Core Transport



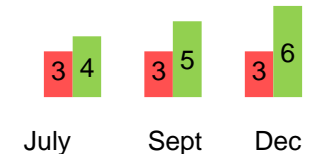
## Servers Running Unsupported OS and Applications

Owner + Implementer = Security + Numerous  
 Roadmaps = Application, Security, Server



## Security Patching

Owner + Implementer = Security + InfraOps  
 Roadmaps = Application, Security, Server



Green = Assigned  
 Red = Open

## Transition

Next we review the Priority 1 and 2 Problems

Questions regarding Mega-Problems?

# {Problem Title}

{Mega-Problem Icon}

<b>Priority</b> 1-5	<b>Status</b> Open   Assigned	<b>Owner</b> CIT Dept	<b>Service Impacted</b> Service Catalogue item	<b>Customer Impacted</b> End-Users Affected	<b>Start Date</b>
<b>Incident Likelihood</b>	High   Medium   Low				{How likely is this to occur?}
<b>Incident Impact</b>	Widespread   Significant   Limited   Local				{Description of effect on end-users}
<b>Incident Recovery</b>	High   Medium   Low				{How much effort in terms of staff hours & \$\$ to recover from an Incident?}
<b>Root Cause</b>	Known   Unknown				{Description}
<b>Summary</b>					
{Description}					
<b>Status - Degrading   Static   Improving</b>					
{Is the Problem worsening, staying the same, or getting better?}					
<b>{Analysis   Mitigation   Resolution} Plan</b>			<b>{Analysis   Mitigation   Resolution} Effort - High   Medium   Low</b>		
{What might we do to analyze, mitigate, and/or fix this Problem?}			{Cost in terms of staff hours & \$\$}		

- Do you understand the risk? If not, what do you need to learn?
- Next steps



# Tungsten-Cobalt Fumble



<b>Priority</b> 2	<b>Status</b> Open	<b>Owner</b> InfraOps	<b>Service Impacted</b> Data Storage	<b>Customer Impacted</b> FHCRC	<b>Start</b> June 2010
<b>Incident Likelihood</b>	Unknown Three Incidents across three years: March 2011, January 2012, February 2012				
<b>Incident Impact</b>	Widespread Possible application disruption for hours; perhaps a few applications unavailable for days				
<b>Incident Recovery</b>	High ~200 CIT + NAG hours across several weeks, gradual restoration of end-user services*				
<b>Root Cause</b>	Unknown				
<b>Summary</b>					
<ul style="list-style-type: none"> <li>• Weekly: Tungsten logs errors about delays in reaching Cobalt</li> <li>• Monthly: Fragile servers crash and require a reboot</li> <li>• Three Incidents: Wide-spread disruption</li> <li>• NetApp and HP tech support warn that these error messages indicate a serious problem</li> </ul>					
<b>Status - Static</b>					
<b>Analysis Plan</b>			<b>Analysis Effort - High</b>		
Spin up an RCA			30 – 200 hours 2-3 Staff		

- Do you understand the risk? If not, what do you need to learn?
- Next steps: RCA | Risk Accept | Leave Open ...

\*Based on past experience

# Tungsten Performance Degraded



<b>Priority</b> 2	<b>Status</b> Open	<b>Owner</b> Architecture	<b>Service Impacted</b> Data Storage	<b>Customer Impacted</b> FHCRC & SCCA	<b>Start</b> Oct 2011
<b>Incident Likelihood</b>	Ongoing				
<b>Incident Impact</b>	Widespread			Slow response from applications	
<b>Incident Recovery</b>	Low			Users restart applications if they time out	
<b>Root Cause</b>	Known				Load
<b>Summary</b>					
Load has increased to the point where application performance suffers noticeably					
<b>Status – Degrading</b>					
Historically, we have deployed various measures to mitigate Tungsten OS overload, including:					
<ul style="list-style-type: none"> <li>• Giving up high-availability (migrating from Active/Standby to Active/Active to employ both Heads)</li> <li>• Buying additional NetApps (Thorium for SCHARP, Copper for SciComp)</li> <li>• Application optimizations</li> </ul>					
<b>Mitigation Plan</b>			<b>Mitigation Effort - Low</b>		
Spin up Project to identify options			15 hours 2 staff		

- Do you understand the risk? If not, what do you need to learn?
- Next steps: Project | Risk Accept | Leave Open ...

# SQL Server Crumps During Storage Events



<b>Priority</b> 2	<b>Status</b> Open	<b>Owner</b> InfraOps	<b>Service Impacted</b> Database Hosting Enterprise App Mgmt	<b>Customer Impacted</b> FHCRC + SCCA	<b>Start</b> 2008
<b>Incident Likelihood</b>	Unknown Unplanned events interrupting Carbon, Tungsten, or the J4-401 Network would trigger this				
<b>Incident Impact</b>	Widespread Enterprise SQL and MIS Applications unavailable, possible data loss				
<b>Incident Recovery</b>	High ~50-100 CIT hours across a week, gradual restoration of end-user services*				
<b>Root Cause</b>	Known Software bugs + configuration errors				
<b>Summary</b>					
<ul style="list-style-type: none"> <li>• Our various SQL Servers do not tolerate minor bumps in access to Storage, unmounting their databases abruptly, sometimes suffering database corruption</li> <li>• As a work-around, we always shut down SQL Server whenever we anticipate Storage or Network maintenance</li> <li>• Naturally, this work-around is only effective for planned events</li> </ul>					
<b>Status - Static</b>					
<b>Resolution Plan</b>			<b>Resolution Effort - Medium</b>		
<ol style="list-style-type: none"> <li>1. Use the Development environment to test fixes</li> <li>2. Use a maintenance window to verify that the fixes work</li> </ol>			<ol style="list-style-type: none"> <li>1. 30 hours</li> <li>2. 30 hours</li> </ol> <p>1 staff</p>		

- Do you understand the risk? If not, what do you need to learn?
- Next steps: Assign | Leave Open | Risk Accept ...

\*Based on past experience

# Data Centers Intermittently Isolated During Router Failure



<b>Priority</b> 2	<b>Status</b> Open	<b>Owner</b> InfraOps	<b>Service Impacted</b> Connectivity	<b>Customer Impacted</b> FHCRC + SCCA	<b>Start</b> March 2010
<b>Incident Likelihood</b>	Low				
<b>Incident Impact</b>	Widespread <a href="#">Brief</a> isolation of a Building plus its Data Centers; <a href="#">Possible</a> data loss and application corruption				
<b>Incident Recovery</b>	High ~20-200 CIT + NAG hours across one week, gradual restoration of end-user services				
<b>Root Cause</b>	Unknown				
<b>Summary</b>					
<ul style="list-style-type: none"> <li>• By design, our Buildings can survive the loss of either of their Redundant Routers</li> <li>• However, during routine validation, we discovered that portions of a Building became isolated during the loss of a Router</li> <li>• The highest profile example has been the DF-120 Data Center (Data Protection)</li> </ul>					
<b>Status - Static</b>					
<b>Analysis Plan</b>			<b>Analysis Effort - Low</b>		
During a maintenance window, replicate the problem, analyze			20 hours 2 staff		

- Do you understand the risk? If not, what do you need to learn?
- Next steps: Leave Open | Risk Accept ...

# Uncertain Router Failover for NetApps



<b>Priority</b> 2	<b>Status</b> Open	<b>Owner</b> InfraOps	<b>Service Impacted</b> Data Storage	<b>Customer Impacted</b> FHCRC	<b>Start</b> Dec 2009
<b>Incident Likelihood</b>	Low Loss of a Redundant Router or Switch might trigger this				
<b>Incident Impact</b>	Widespread Isolation from Storage for Clients located outside a given Building				
<b>Incident Recovery</b>	Medium ~2-50 CIT + NAG hours across one week, gradual restoration of end-user services				
<b>Root Cause</b>	Known Configuration error				
<b>Summary</b>					
<ul style="list-style-type: none"> <li>By default, our NetApps are <u>not</u> configured to take advantage of our Redundant Router scheme</li> <li>This is surprising, as we use one of several popular and standard Redundant Router schemes</li> <li>In contrast, out-of-the-box Windows, Linux, and Mac OS hosts are configured suitably</li> <li>Nevertheless, our analysis says that NetApps ship by default with a configuration choice ('FastPath') which dodges our scheme</li> <li>This issue only affects Clients located outside a given Building (e.g. home drive users of Tungsten located outside Yale)</li> </ul>					
<b>Status - Static</b>					
<b>Analysis Plan</b>			<b>Analysis Effort – Low</b>		
<ul style="list-style-type: none"> <li>Make the configuration change on a guinea pig NetApp</li> <li>Reboot the associated Redundant Router</li> <li>Analyze results</li> </ul>			10 hours 2-3 staff		

- Do you understand the risk? If not, what do you need to learn?
- Next steps: Leave Open | Risk Accept ...

# vColo Does Not Reliably Recover from Storage Bumps



<b>Priority</b> 2	<b>Status</b> Open	<b>Owner</b> InfraOps	<b>Service Impacted</b> Hosting	<b>Customer Impacted</b> FHCRC	<b>Start</b> June 2010
<b>Incident Likelihood</b>	High				
<b>Incident Impact</b>	Widespread		vColo disrupted, various applications unavailable		
<b>Incident Recovery</b>	Low ~2 - 20 CIT + NAG hours across one day, gradual restoration of end-user services				
<b>Root Cause</b>	Unknown				
<b>Summary</b>					
<ul style="list-style-type: none"> <li>• Sometimes when Storage stumbles, some vColo Guests (and vColo mgmt functions) require reboots to recover</li> <li>• During pathological Storage Incidents, this is unsurprising</li> <li>• But it happens even during clean + fast Storage head fail-overs, when other Storage clients recover without issue</li> </ul>					
<b>Status - Static</b>					
<b>Analysis Plan</b>			<b>Analysis Effort – Medium</b>		
Spin up RCA			200 hours 3 staff		

- Do you understand the risk? If not, what do you need to learn?
- Next steps: Leave Open | Risk Accept ...

# Uncertain Router High-Availability



<b>Priority</b> 2	<b>Status</b> Open	<b>Owner</b> InfraOps	<b>Service Impacted</b> Connectivity	<b>Customer Impacted</b> FHCRC + SCCA	<b>Start</b> April 2011
<b>Incident Likelihood</b>	Unknown <span style="float: right;">Loss of a Redundant Router or Switch might trigger this</span>				
<b>Incident Impact</b>	High <u>Isolation</u> of a Building plus its Data Centers from the rest of Campus; <u>Possible</u> data loss and application corruption				
<b>Incident Recovery</b>	High ~20-200 CIT + NAG hours across one week, gradual restoration of end-user services				
<b>Root Cause</b>	Known Software bugs + configuration errors				
<b>Summary</b>					
We are not confident that the Data Center, Building, and Campus Core switches/routers are still highly-available (HA)					
<b>Status - Degrading</b>					
We are stalled on patching Data Center switches in part because we fear triggering this issue					
<b>Resolution Plan</b>			<b>Resolution Effort – Low to Medium</b>		
<ol style="list-style-type: none"> <li>1. Characterize current state, Resolve current issues</li> <li>2. Schedule quarterly Building + Campus validation</li> <li>3. Pick one Data Center per quarter, validate behavior</li> </ol>			<ol style="list-style-type: none"> <li>1. 12 hours / 1 staff</li> <li>2. 8 hours / quarter</li> <li>3. 100 hours / quarter</li> </ol>		

- Do you understand the risk? If not, what do you need to learn?
- Next steps: Assign | Leave Open | Risk Accept ...

# Intermittent Data Protection Gaps



<b>Priority</b> 3	<b>Status</b> Open	<b>Owner</b> Architecture	<b>Service Impacted</b> Data Protection	<b>Customer Impacted</b> FHCRC	<b>Start</b> March 2010
<b>Incident Likelihood</b>	Low Two Incidents in recent memory: June 2011 Fred and November 2011 Silo				
<b>Incident Impact</b>	High Loss of data and/or servers				
<b>Incident Recovery</b>	High ~~50-500 CIT + NAG hours across days to months, gradual restoration of end-user services				
<b>Root Cause</b>	Known Insufficient resources in data protection systems / lack of archiving				
<b>Summary</b>					
<ul style="list-style-type: none"> <li>• Our data protection systems no longer backup all our data during nightly incrementals and weekly fulls</li> <li>• Depending on the timing of a failure, we might lose one day of data (best case), multiple days (likely), or even weeks (worst case)</li> <li>• The gaps vary day-from-day and by system (vColo, Tungsten ...)</li> <li>• Backups fail intermittently (not clear why), requiring manual intervention and widening the gaps</li> </ul>					
<b>Status - Degrading</b>					
<ul style="list-style-type: none"> <li>• As data volume grows and systems become slower, our gaps become wider</li> <li>• Long-term resolution will involve a substantial re-engineering effort</li> </ul>					
<b>Analysis Plan</b>			<b>Analysis Effort - Low</b>		
Spin up an RCA to quantify the gaps			30 hours 4 staff		

- Do you understand the risk? If not, what do you need to learn?
- Next steps: Leave Open | RCA | Risk Accept ...



# ACD not Highly Available



<b>Priority</b> 3	<b>Status</b> Open	<b>Owner</b> InfraOps	<b>Service Impacted</b> Voice	<b>Customer Impacted</b> FHCRC	<b>Start</b> 2007
<b>Incident Likelihood</b>	Low				
<b>Incident Impact</b>	Medium 'Fast Busy' for these numbers for 2 hours – 1 week				
<b>Incident Recovery</b>	Low ~1 hour – 1 day of CIT staff time				
<b>Root Cause</b>	Known Design Choice				

## Summary

- We provide an Automated Call Distribution (ACD) service for two numbers 667-5000 (FHCRC's main number) and 667-5700 (CIT's HelpDesk)
- An ACD service allows multiple handsets to ring simultaneously when a call arrives, allowing multiple people to service the number
- We have only one server (*Audhumla*) behind this capability
- During the recent J4-401 Ethernet switch chassis upgrade, *Audhumla* was isolated (it has only one NIC) and Security noticed (calls to x5000 were receiving 'Fast Busy')

## Status - Static

- This is a design decision we made years ago – we weren't willing to spend the \$\$ to buy a second one
- The proposal to add a second *Audhumla* did not make the FY'13 funding cut
- Loss of an Ethernet switch, *Audhumla*'s single NIC, or all of *Audhumla* itself would result in 'Fast Busy' to x5000 and x5700
- I'm using Problem Management to verify that this audience knows and accepts this risk

## Resolution Plan

Purchase a second *Audhumla*  
Configure it as a hot-standby

## Resolution Effort - Low

\$xx capital  
1-2 staff / yy hours

- Do you understand the risk? If not, what do you need to learn?
- Next steps: Leave Open | Project | Risk Accept ...

# Intermittent Failed Network Connections



<b>Priority</b> 2	<b>Status</b> Open	<b>Owner</b> InfraOps	<b>Service Impacted</b> Connectivity	<b>Customer Impacted</b> FHCRC + SCCA	<b>Start</b> March 2010
<b>Incident Likelihood</b>	High				
<b>Incident Impact</b>	Significant InfrastructureTasks <u>interrupted</u> ; <u>Possible</u> data loss and application corruption				
<b>Incident Recovery</b>	Low - High ~2-200 CIT + NAG hours across one week, gradual restoration of end-user services				
<b>Root Cause</b>	Unknown				
<b>Summary</b>					
<ul style="list-style-type: none"> <li>• Hosts inside J4-401 are briefly unable to talk with one another</li> <li>• This leads to failed Infrastructure Tasks, e.g. Backups, Database Maintenance ... and occasional Server Freezes</li> <li>• Hard to tell how frequent the issue is ... Infrastructure Tasks tend to be Patient (retry many times) and Automated (no human involved)</li> <li>• <u>Documented</u> cases occur rarely (~once/month)</li> </ul>					
<b>Status - Degrading</b>					
<b>Analysis Plan</b>			<b>Analysis Effort - High</b>		
Spin up RCA			200 hours 4 staff		

- Do you understand the risk? If not, what do you need to learn?
- Next steps: Leave Open | RCA | Risk Accept ...

## Transition

Next we review Metrics

Questions regarding P1s & P2s?

## Quarterly Snapshots

### March 30

**New** 34

**Open** 81

**Assigned** 9

**Project** 0

**RCA** 1

Escalated 0

**Total Active** 91

### July 5

**New** 9

**Open** 72

**Assigned** 39

**Project** 0

**RCA** 0

Escalated 9

**Total Active** 120

### Sept 25

**New** 5

**Open** 68

**Assigned** 29

**Project** 0

**RCA** 0

Escalated 0

**Total Active** 97

### Dec 20

**New** 9

**Open** 79

**Assigned** 27

**Project** 0

**RCA** 0

Escalated 0

**Total** 106

**Risk Accepted** 0

**Rejected** 0

Resolved 11

**Total Closed** 11

**Risk Accepted** 0

**Rejected** 0

Resolved 4

**Total Closed** 4

**Risk Accepted** 11

**Rejected** 5

Resolved 7

**Total Closed** 23

**Risk Accepted** 4

**Rejected** 0

Resolved 6

**Total Closed** 10

#### Problem Status

**Assigned** Tech is working on the Problem

**Escalated** Disagree on Status, seeking leadership direction

**Open** Inert: not working on it

**Project** Project instantiated

**RCA** Root Cause Analysis team instantiated

**Risk Accepted** We intend to live with this

**Rejected** Not a Problem

**Resolved** Fixed + Closed

## Next Steps

What action items have we generated?

Who owns them?

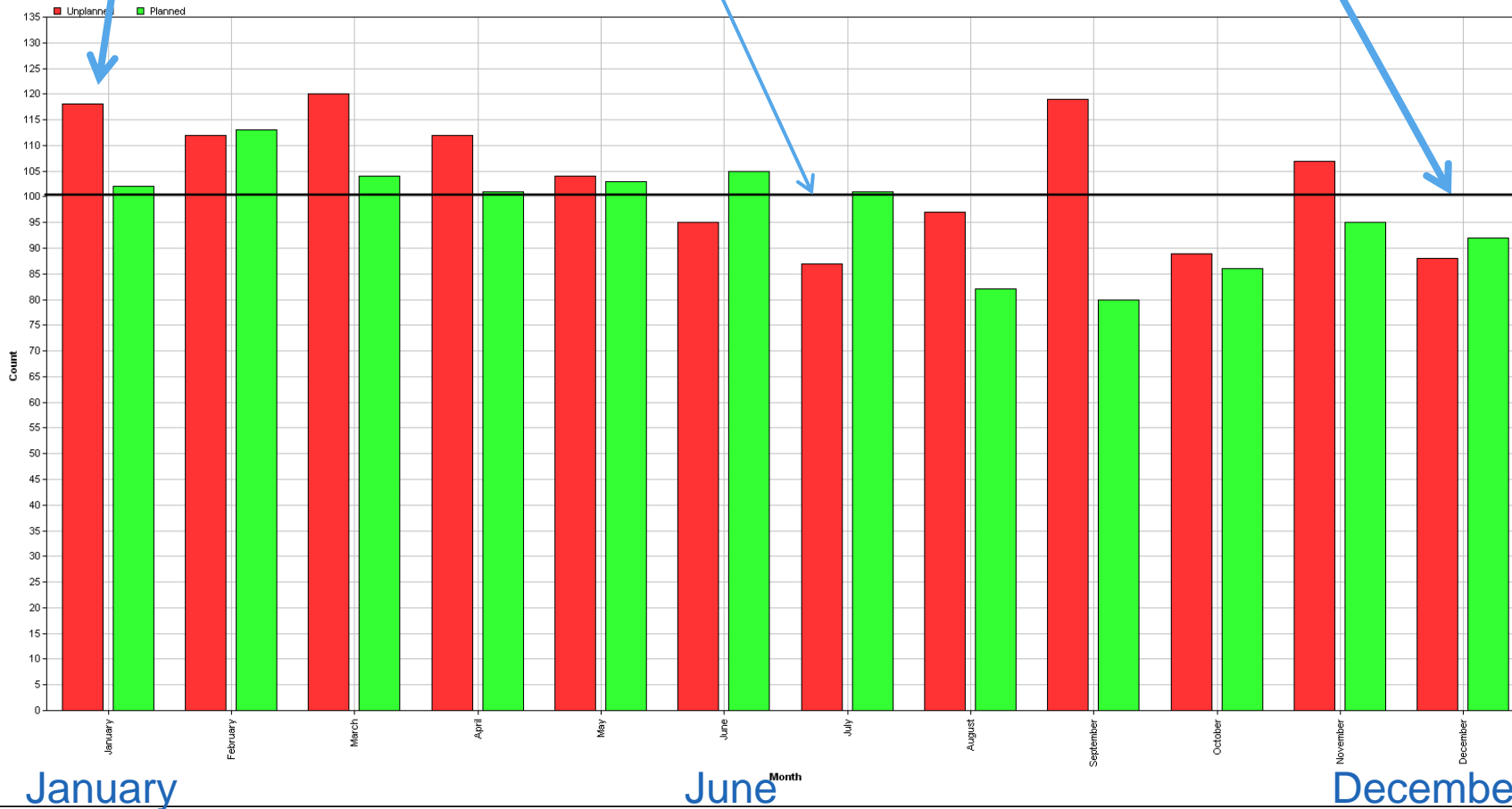
# Planned & Unplanned Outages by Month, October 2000 – June 2012

Coming Soon

Line marking '100'

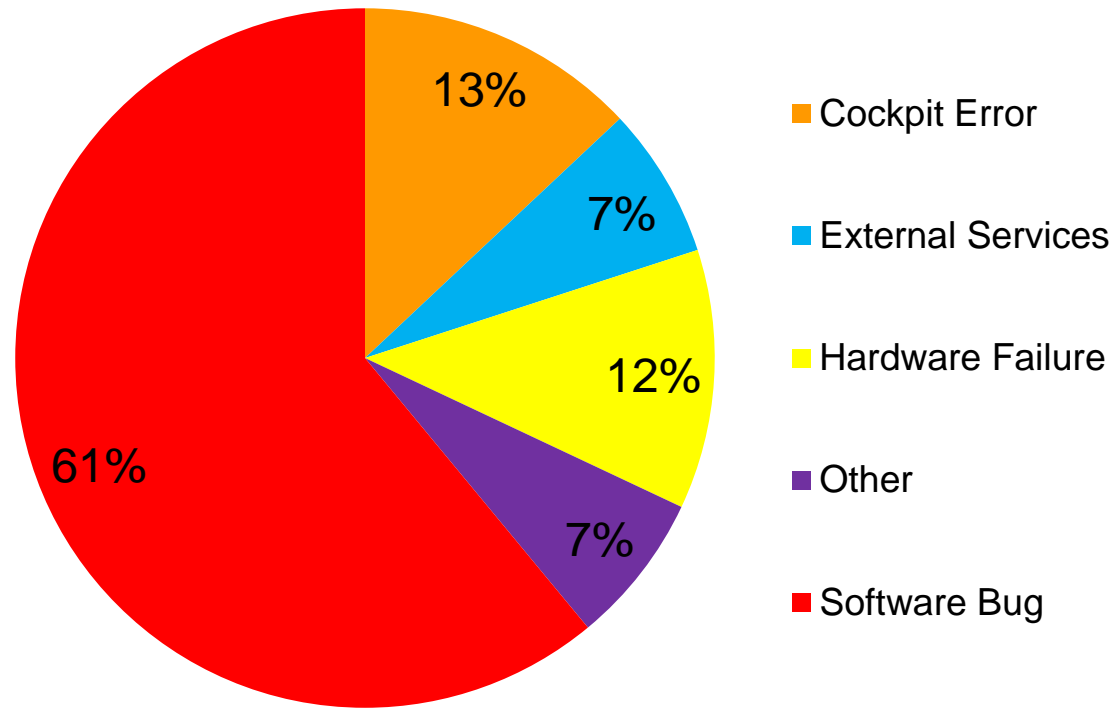
You are Here

Planned & Unplanned Outages by Month, October 2000 - June 2012



Past performance does not guarantee future returns ...

## Causes of Unplanned Outages: October 2000 – June 2012



*Would regular patching shrink **Software Bug**?*

### Is this what you want from a Prob Mgmt Deep Dive?

Goals:

- Keep leadership apprised of current technical risks
- Identify next steps

### What would you like to see in the March Deep Dive?

1. Review P1s & P2s
2. Metrics
3. Process Engineering



# Appendix

<u>Topic</u>	<u>Page</u>
Stuart's Summary	25
Process Engineering	26
Life Cycle of a Problem	27
Priority Matrix	28

## Stuart's Summary

- (A) Consolidated Storage experiences on-going issues
- (B) Portions of the deep infrastructure are degrading
  - Highly-Available becoming Highly-Unavailable
  - Data Centers vulnerable to complex disruptions
  - Interferes with capacity upgrades, **bug fixes**, security patching
- (C) Aging applications living on aging servers
  - Hard to identify owner
  - Exposed to security vulnerabilities
  - General ignorance and age suggests *unknown unknowns*
- (D) Limited visibility into our flock of Problems: *Unknown unknowns*

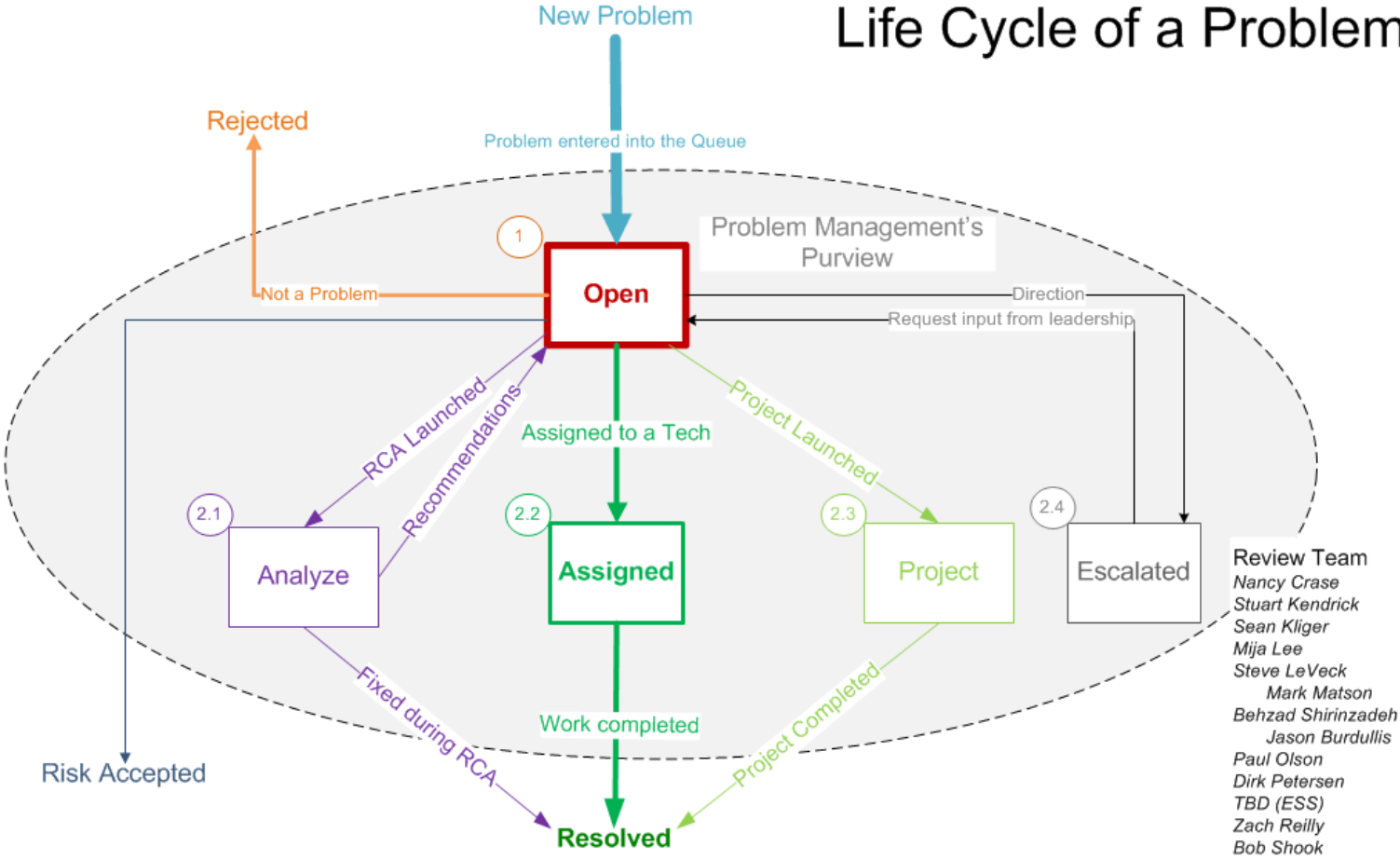
*Predicting is hard, particularly the future* –Niels Bohr

## Process Engineering

1. Do we want to track time spent fixing Problems? If so, how?
2. How do we prioritize fixing Problems into our Maintain / Build / Improve work?
3. How do we track progress toward resolution?
4. What value are we delivering with Problem Mgmt?

Action item: Joan to pursue 'pooled' resource model

# Life Cycle of a Problem



Problem Management States		Problem Management Sources/Destinations	
Open	Inert: Not yet reviewed or stalled on resources	New Problem	From Incident Mgmt, Techs, Project Managers ...
Analyze	Root Cause Analysis team instantiated	Rejected	Not a Problem
Assigned	Tech is working on the Problem	Risk Accepted	We intend to live with this
Project	Project instantiated	Resolved	Fixed + Closed
Escalated	Disagree, seeking leadership direction		

skendric 2012-08-16

## Priority Matrix

**Priority is the relative urgency of the problem** , calculated using the look-up table below

Rating	Description
1	Most important
2	
3	Average importance
4	
5	Less important

		Impact				
		1	2	3	4	5
Likelihood	1	1	1	2	3	4
	2	1	2	2	3	4
	3	2	2	3	4	4
	4	3	3	4	5	5
	5	4	4	4	5	5