*2016 Operations Stay Treat*

# Improving System Reliability

*A Case Study of Accelerator UPSs*
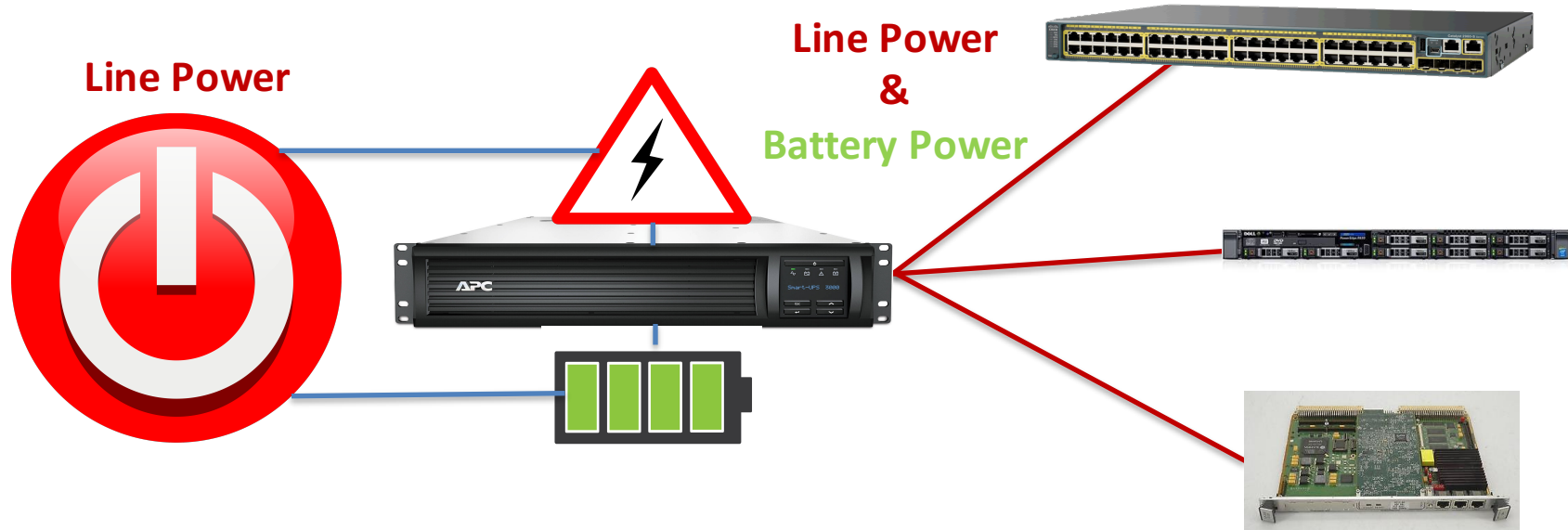
## Anthony Cuffe

**Jefferson Lab**

# Summary     *A Case Study of Improving UPS Reliability*

- ⊙ Typical UPS Usage
- ⊙ Backstory
- ⊙ Initial Status and Analysis
- ⊙ Reliability Study
- ⊙ Maximizing Availability, Reliability and Manageability
- ⊙ Automation, Monitoring and Logging Tools
- ⊙ Results
- ⊙ Key Points
- ⊙ Q & A

Jefferson Lab

# Typical UPS Usage in A.C.E.



**Line Power**

**Line Power & Battery Power**

## Protection from Power Loss and Poor Line Quality for:

- ⊙ Critical Workstations and Servers
- ⊙ Network Infrastructure
- ⊙ RF, Vacuum and other critical IOCs

Jefferson Lab

# Backstory    *What just happened?*

◉ Historically our group maintained around 10-15 rack-mounted, enterprise grade UPSs.

*... This was manageable without much effort.*

◉ After a re-org, our group became the proud owners of 100+ UPSs.

*... And most of them weren't shiny and new!*

◉ After staging through denial, anger, bargaining, and depression, we accepted it.

*... This is how we got over most our grief.*

Jefferson Lab

# Initial Status   *The Reality of Square One*

- Responsible for 100+ UPSs *(and growing!)*
- ~10 were outright failed
- ~10 were in a faulted state
- 40+ were found to be 8-16 years old *(5 year lifetime)*
- Configuration was manual, onerous and inconsistent
- Deployment took a very long time *(1++ hour)*
- Monitoring was impossible
  - *Sheer volume of error messages (X,000+/per day)*
  - *System age*
- Our attitudes towards the system were awful
  - *Dread*
  - *Apathy*
  - *A few resumes were tuned up*

Jefferson Lab

# Initial Analysis    *Is the system worth improving?*

- ⊙ Is the system's functionality still necessary?
  - ○ *Servers, network gear and critical IOCs need to survive power interruptions.*
  - ○ *UPSs provide the only low-cost solution to fulfill this need.*

- ⊙ Is the benefit worth the cost and staff time?
  - ○ *~ $50K to make the overall system reliable and ~$10-20K/year to maintain*
  - ○ *The recovery costs and reduced beam availability associated with frequent power-outages greatly exceeds the cost of maintaining UPSs for critical systems.*

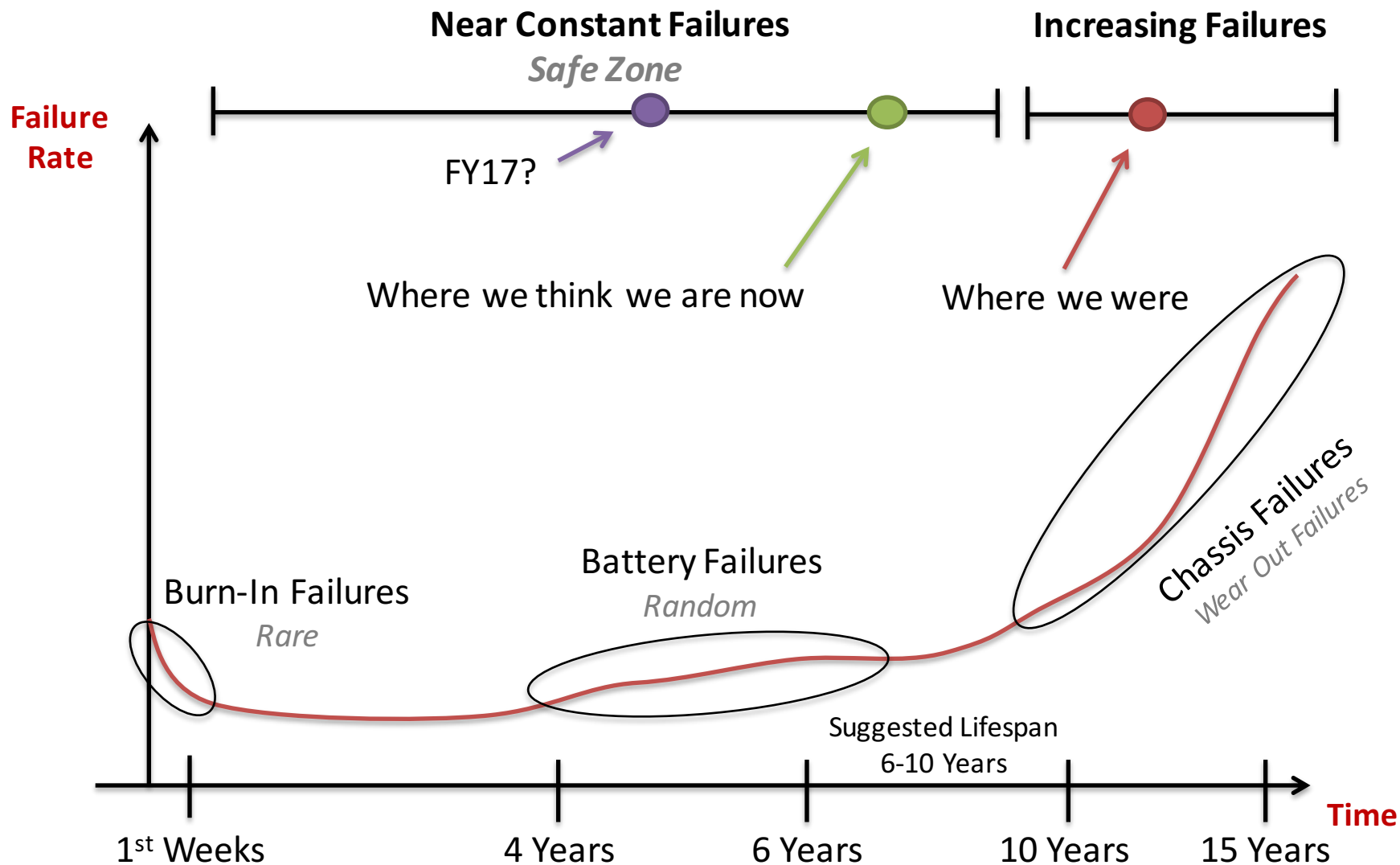- ⊙ Is it possible to improve the cost/benefit ratio?
  - ○ *We felt we could decrease the time spent maintaining these systems by applying practices common to our other systems.*
  - ○ *The benefit could be improved with better reliability and selective application.*
  - ○ *Costs could be spread over years with a well maintained system.*

- ⊙ Do you have management buy in?
  - ○ *We had to convince management that the benefit would be worth the time, effort and procurement dollars.*
  - ○ *It is impossible to improve or maintain a system without resources!*

Jefferson Lab

*Cost Benefit Analysis*

# Reliability Study  *UPS Failure Rate*



**Near Constant Failures**
*Safe Zone*

**Increasing Failures**

**Failure Rate**

FY17?

Where we think we are now

Where we were

Chassis Failures
*Wear Out Failures*

Battery Failures
*Random*

Burn-In Failures
*Rare*

Suggested Lifespan
6-10 Years

Time

1st Weeks      4 Years      6 Years      10 Years    15 Years

* Curve estimated based on experience, dependent on operating environment

**Jefferson Lab**

# Improving Availability

**Availability Depends on Reliability and Manageability**

$$AVAILABILITY \approx \frac{MTBF}{(MTBF + MTTR)}$$

⊙ **Maximize Reliability**
- o *Make **MTBF** as big as possible*
- o *Mitigate unreliability where **MTBF** cannot be reduced*

⊙ **Maximize Manageability**
- o *Make **MTTR** as small as possible*
- o *Proactive replacement/maintenance where **MTTR** cannot be reduced*

**MTTR** - *Mean Time to Repair.*
**MTBF** - *Mean Time Between Failures.*

Jefferson Lab          *Availability, Reliability and Manageability*

# Maximize Reliability _Increase MTBF_

⊙ Get the System to a Manageable State (stability point).
  - o _Inventory audit – You have to know what is out there._
  - o _Replace failed units._
  - o _Repair/replace units in error/faulted state._

⊙ Reduce Complexity
  - o _Get rid of unnecessary units._
  - o _Remove unnecessary components (PDUs, Transfer Switches …)_
  - o _Reduce unwanted loads on the system._

⊙ Implement Redundancy
  - o _Design/Deploy systems with redundant power supplies._
  - o _Implement redundant systems so failures are actually tolerable._

⊙ Monitor for Faults, Failures and Predicted Failures
  - o _Email Alerts – from units themselves and active polling_
  - o _Alerts for low Runtimes, Faults and Failures_
  - o _Performance Graphing for historical perspective and statistics_
  - o _Failure Prediction_
  - o _Centralized Logging for history and analytics_

Jefferson Lab

# Maximize Maintainability *Decrease MTTR*

- ◉ *Corrective* rather than *Preventative* maintenance
  - o *Predicted failures are replaced before problems arise*
  - o *Actual failures are replaced quickly*
  - o *Batteries are replaced when there are faults not on a schedule*
  - o *Configuration changes and updates are done globally*

- ◉ On the shelf spares

- ◉ Implement redundancy wherever possible
  - o *MTBF and MTTR are zero if system is redundant!*
  - o *Critical systems always have dual power supplies and UPSs*

- ◉ Automate as much as possible
  - o *Research available tools (open source and vendor provided!)*
  - o *Automation produces consistent and complete configuration*
  - o *Custom scripts can automate tasks and deployment (~5 min)*
  - o *"If you do something more than once, you should write a script to do it!"*

Jefferson Lab

## **APC UPS Network Management Tools**

⊙ *Monitoring, Control and Remote Access*

- o *Web browser*
- o *Command line interface*
- o *SNMP, SSH, FTP …*

⊙ Event and Data Logging

⊙ *Scheduled Self-Tests*

⊙ *Email Notification*

- o *Faults, Failures & Warnings*
- o *Self-Test Results*
- o *Configuration Changes*

**Jefferson Lab**

# Custom Tools  *Automation*

## ⊙ Scripting Framework

- ○ *Global management tasks via custom scripts*
- ○ *Scripts leverage vendor tools*
- ○ *Simple text file database of UPSs*
- ○ *Centralized configuration management*
- ○ *Configuration consistency checks*
- ○ *Information gathering tools*

```
[root@opsfs scripts]# ./find-old-firmware.pl
======================================
          BAD FIRMWARE REPORT
======================================

Problems found: 0
```

```
[root@opsfs scripts]# ./get-batt-replacement-date.pl mccups10
Community [public]:
mccups10        09/01/05
[root@opsfs scripts]# ./get-batt-replacement-date.pl
Community [public]:
aceups01        05/27/2014
ams01b03ups01   04/17/2013
ams01b03ups02   03/22/2013
bs04b07ups01    11/15/2015
cl02c06ups01    06/01/13
cl02c10ups01    09/28/2015
cl03ups01       09/17/2014
cl03ups01a      03/25/2015
cl03ups10       07/11/2012
```

Jefferson Lab

# Custom Tools   *Automation*

## ◉ **Rapid Deployment System (~5 mins)**

- o *Automated deployment with custom scripts*

- o *Process leverages vendor provided tools*

- o *Easy and well documented procedure*

- o *Consistent and complete configuration*

- o *Rapid and automated reconfiguration*

```
[upsadm@node scripts]#  ./deploy_new_ups.pl mccups13
Username [apc]:
Password [apc]:
Community [public]:

Ready to run on the following UPSes.
mccups13
Continue (y/N):
```

Jefferson Lab

# Monitoring Tools  *Graphing with Cacti*

**Cacti** – Open Source Polling and Graphing Framework

- ⦿ *Performance Monitoring via SNMP*
- ⦿ *Adaptable to Graph Virtually Anything*
- ⦿ *Historical Statistics*
- ⦿ *Problem Identification*
- ⦿ *Failure Prediction*
- ⦿ *Global Overview*
- ⦿ *Free*



mccups11a (mc02ups11a) - APC Battery Details

| | | | |
|---|---|---|---|
| Battery Capacity % | Current: 100.00 | Average: 99.97 | Minimun: 99.18 |
| Battery Load % | Current: 18.05 | Average: 18.50 | Maximum: 38.06 |
| Battery Temp | Current: 15.00 | Average: 16.14 | Maximum: 21.41 |
| Battery Status | Current: 2.00 | Max: 2.00 | |

Unable to sustain current load if battery status is over 3

Last Updated: Fri 17 Jun 12:55:03 EDT 2016



mccups11a (mc02ups11a) - Battery Runtime Remaining

Runtime Remaining    Current: 35.04    Average: 42.38    Max: 75.00

Jefferson Lab

# Monitoring Tools *Cacti*

## Custom Global Overviews

# Monitoring Tools *Cacti*

## Adaptable to Graph Virtually Anything

*… even mis-behaving users*



Jefferson Lab

# Monitoring Tools  *Status Polling with Nagios*

## **Nagios** – Open Source Monitoring Framework

⊙ *System and Service Status Polling via SSH, SNMP, …*
   o *Runtime Warnings*
   o *Fault and Failure Detection*
   o *Active and Negative Service Monitoring (want SSH but not FTP running)*
⊙ *Easily adapted to almost any need and for multiple groups*
   o *Already monitoring services for others groups like AHLA*
   o *Group based service/system notification*

⊙ *Email/SMS Alerts*
⊙ *Web based H.U.D*
⊙ *Monitor for Monitoring*
⊙ *Free*

```
** NAGIOS PROBLEM Alert: hareboot1.acc.jlab.org/FTP is CRITICAL **
** NAGIOS PROBLEM Alert: eosups1.acc.jlab.org/RUNTIME is CRITICAL **

From System Account ACE Group <accadm@jlab.org>
Subject ** NAGIOS PROBLEM Alert: eosups1.acc.jlab.org/RUNTIME is CRITICAL **
   To accadm@jlab.org

***** Nagios *****

Notification Type: PROBLEM

Service: RUNTIME
Host: eosups1.acc.jlab.org
Address: eosups1.acc.jlab.org
State: CRITICAL

Date/Time: Mon Jun 6 10:19:11 EDT 2016

Additional Info:

CRIT: 2 minute(s) remaining. Waring:10. Crit:5

Misc Info:
```

Jefferson Lab

# Monitoring Tools *Status Polling with Nagios*

## Heads-Up-Display

# Logging Tools   *Centralized Logging w/Analytics*

◉ **Swatch** (*Syslog Watch*)

- *Filter and Summary Tool for centralized logging – filter out normal, noisy logs*
- *Daily Reporting via Email*
- *Extensible Filtering – easily add custom filters*
- *Open Source*

◉ **Splunk**

*Splunk is a fully featured, powerful platform that collects and indexes any machine data from virtually any source in real time.   We feed it our centralized logs.*

- *Index any Text Data*
- *Interactive Search Results*
- *Monitoring and Alerting*
- *Reporting and Analysis*
- *Event Correlation*
- *Custom Dashboards*
- *Add-in apps*
- *Already extended to other groups*
- *Not Free but DOE contract in works!*



**Jefferson Lab**

# Logging Tools   *Splunk*

## Deep-Dive Your Data

# Results *Higher Availability with Lower TCO?*

⊙ **Before (10–15 units):**

- *Configuration was manual, onerous and inconsistent*
- *Deployment took a long time (+1 hour)*
- *Attitudes towards system were apathetic or dreadful*
- *Monitoring was nonexistent or a deluge.*
- *Failures were commonplace and faults were ignored.*

⊙ **After (100+ units):**

- *Configuration is automatic, fast and consistent.*
- *Deployment is very quick (~5 minutes).*
- *Replacement is now simply a routine task .*
- *Monitoring is proactive, accurate and very predictive.*
- *Failures are rare and faults are dealt with quickly.*

Jefferson Lab

# Key Points to Take Away

◉ **Buy in by management a must.**

 o *You cannot maintain a system  without resources!!*

◉ **Redundancy wherever possible**

 o *MTTR and MTBF are virtually zero with redundancy!*

◉ **Automate as much as possible**

 o *Ensures standardization.*

 o *Reduces staff time.*

 o *Improves the culture.*

◉ **Monitor as much as possible**

 o *Failure notification and prediction is critical!*

 o *Fixing things before they break means zero downtime!*

 o *What you don't know usually hits you hardest!*

 o *Pretty graphs and real data translate into better budgets!*

◉ **Find or create tools to make things faster and easier**

 o *Research vendor and open source solutions*

 o *Simple solutions like shell scripts can be very powerful*

 o *Leverage the expertise of other groups*

Jefferson Lab

# Q & A