**NORTHROP GRUMMAN**

# Six Years of Controlled Peer Reviews

Practical Lessons Learned
25 June 2009

Richard L. W. Welch, PhD

Associate Technical Fellow
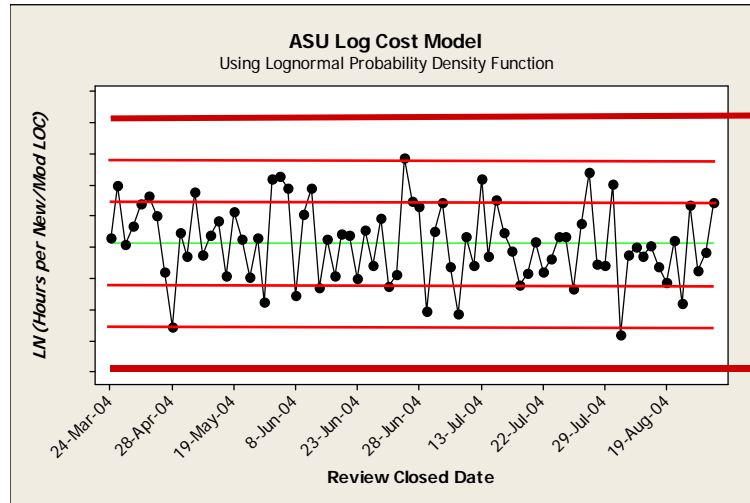
Steve D. Tennant

SEPG Lead

# Topics

- Background

- Getting started (2004)
  - Overcoming technical issues

- Handling human issues and institutionalizing the process (2005)

- Growing the benefits (2006)
  - More processes, projects, disciplines

- Increasing our effectiveness (2007-2008)
  - Exploring new techniques

- Future pathways

# Why Peer Reviews?

- Ubiquity
  - Many work products reviewed throughout software development life cycle
    - System & software design artifacts
    - Source code
    - Test plan, procedures & reports

- Frequency
  - High data rates

- Influence
  - Approximately 10% of the software development effort is spent on peer reviews and inspections
  - Code walkthroughs represent biggest opportunity & most advantageous starting point

- Serendipity
  - All engineering disciplines peer review their design products
  - Techniques & lessons learned have demonstrated extensibility

3

# Why Statistical Process Control?
## Successful Quantitative Project Management

**ASU Log Cost Model**
Using Lognormal Probability Density Function

LN (Hours per New/Mod LOC)

Review Closed Date

24-Mar-04, 28-Apr-04, 19-May-04, 8-Jun-04, 23-Jun-04, 28-Jun-04, 13-Jul-04, 22-Jul-04, 29-Jul-04, 19-Aug-04

**Upper Control Limit**
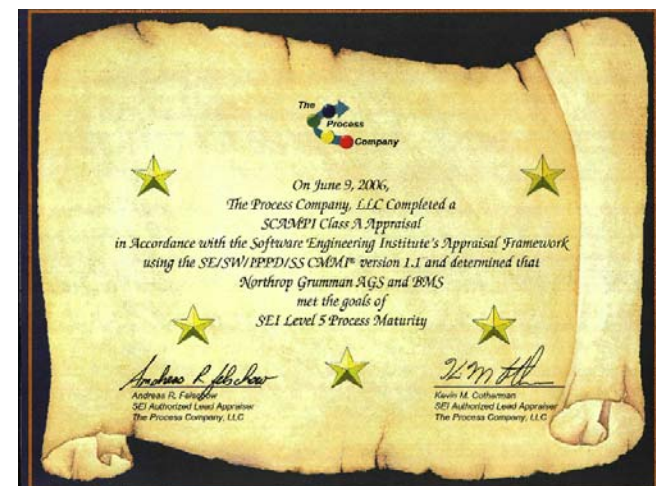
**Average performance**

**Lower Control Limit**

**A stable process** operates within the control limits 99.7% of the time

- Analysis of special cause variation focuses on recognizing & *preventing* deviations from this pattern
- Analysis of common cause variation focuses on *improving* the average and tightening the control limits
- SPC offers opportunities for systematic process improvement that NGC & industry benchmarks indicate will yield an *ROI averaging between 4:1 & 6:1*

# Case Study Essentials

- Data represent software-related peer reviews conducted at Northrop Grumman's Aerospace Systems facility in Melbourne, Florida facility between March 2004 and October 2008

- One peer review process (now standard for Aerospace Systems)
  - Covers the entire system life cycle from system requirements analysis & architecture through maintenance
  - Requires the peer review of all major systems, software & test artifacts
  - Uses an automated data base tool that integrates data quality & process control features

- All peer review records captured in the data base
  - > 5,700 source code peer reviews
  - > 1,100 other software-related peer reviews

- CMMI Level 5 appraisals
  - CMMI-SE/SW (V1.1) in 2005
  - CMMI-SE/SW/IPPD/SS (V1.1) in 2006
  - CMMI-DEV+IPPD (V1.2) scheduled for 2009



The Process Company

On June 9, 2006,
The Process Company, LLC Completed a
SCAMPI Class A Appraisal
in Accordance with the Software Engineering Institute's Appraisal Framework
using the SE/SW/IPPD/SS CMMI® version 1.1 and determined that
Northrop Grumman AGS and BMS
met the goals of
SEI Level 5 Process Maturity

Andreas R. Felschow
SEI Authorized Lead Appraiser
The Process Company, LLC

Kevin M. Cotharman
SEI Authorized Lead Appraiser
The Process Company, LLC
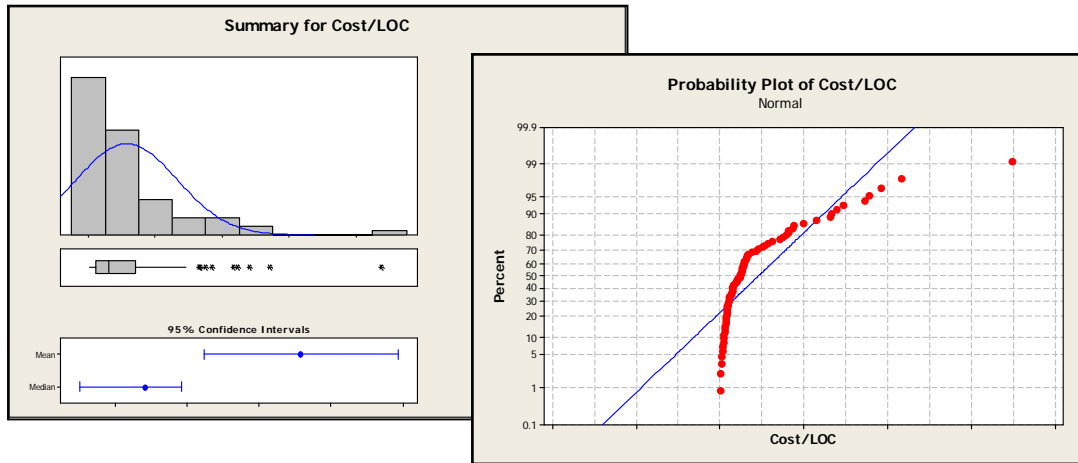
# Getting Started - 2004

Overcoming Technical Difficulties &
Learning To Love Logarithms

# SW-CMM Level 4 Prior State (1998-2003)

- Software development baseline characterized by life cycle phase
  - SW Requirements-Design-Code & Verification-SW Integration-Software Test
  - 10+ year process improvement record resulted in costs reduced by over 67%

- But we had no CMMI "Gestalt"
  - No insight into the statistical behavior of lower level elements
  - No "above the shop floor" experience with statistical sub-process control
  - No insight into downstream behavior

- We wanted to control product quality, but were thwarted by issues with our process quality
  - Inconsistent data
  - Superficial results

- Root cause analysis traced this to indifferent attention paid to *managing* peer reviews
  - We realized we had to control the *efficiency* of our peer reviews in terms of the effort spent (peer review cost), based on classic industry guidelines that efficient reviewers operate in a "sweet spot" of about 200 lines of code per review hour

# Our Problem

Summary for Cost/LOC



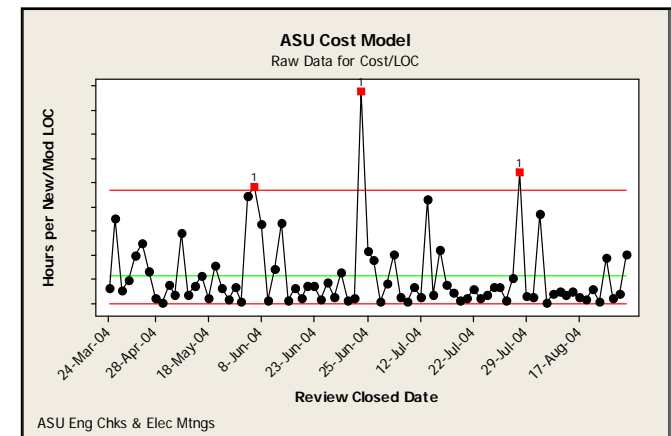Probability Plot of Cost/LOC
Normal

## Data Characteristics
- Anderson-Darling test p < 0.005
- Data non-normality & asymmetry violated probability model assumptions
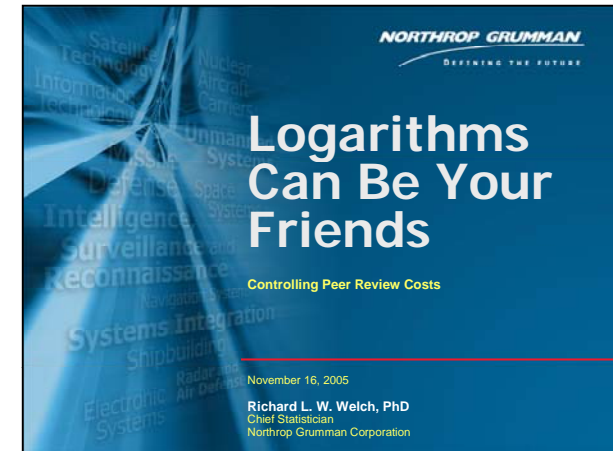
## Control Chart Difficulties

- 11% false alarm rate (Chebyshev's inequality)
  - *Penalized due diligence in reviewing code*

- No meaningful lower control limit
  - *Did not flag superficial reviews*

- Arithmetic mean distorted the central tendency
  - *Apparent cost did not meet budget*



ASU Cost Model
Raw Data for Cost/LOC

## Could we control our peer reviews?

# Stabilizing the Data

- Senior author's presentation at 2005 CMMI[SM] Technology Conference demonstrated how a log-cost model can successfully control software code inspections
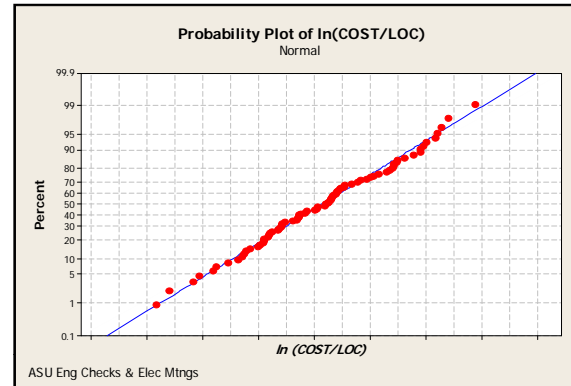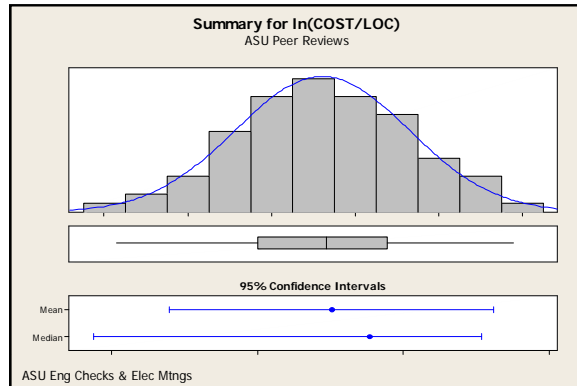


  - Peer review unit costs (hours per line of code) behave like commodity prices in the short term
  - Short term commodity price fluctuations follow a lognormal distribution
  - As a consequence, commodity prices follow a lognormal distribution
  - Therefore, taking the natural logarithm of a sequence of peer review costs transforms the sequence to a normally distributed series
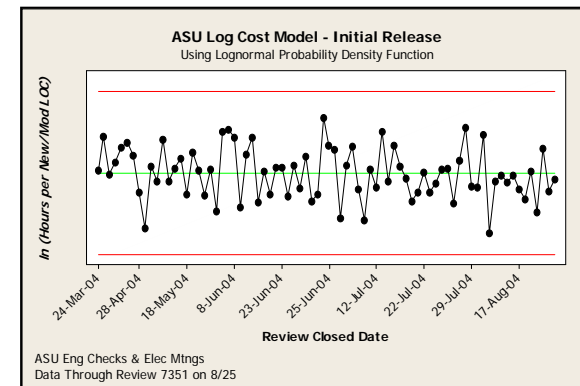
Notes:
- Details on the log-cost model, "one of the most ubiquitous models in finance," can be found at riskglossary.com (http://www.riskglossary.com/articles/lognormal_distribution.htm)
- Prior CMMI Technology Conference & User Group papers are published on-line at: http://www.dtic.mil/ndia/

# Our Data on Logs

**Summary for ln(COST/LOC)**
ASU Peer Reviews

95% Confidence Intervals

ASU Eng Checks & Elec Mtngs



**Probability Plot of ln(COST/LOC)**
Normal

ASU Eng Checks & Elec Mtngs

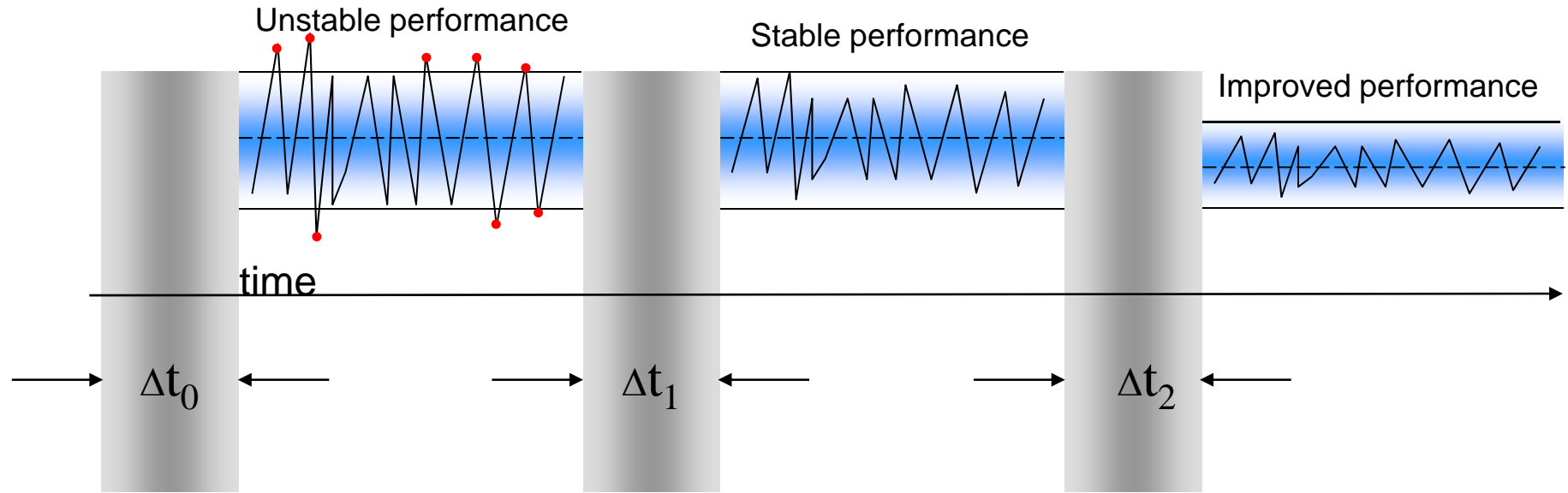Anderson-Darling
test p < 0.759

- Impacts
  - False alarms minimized
  - Meaningful lower control limit
  - Geometric mean preserves the budget
    - *OK, you still have to find the antilog*

- Demonstrated utility & applicability
  - > 6,800 peer reviews over 5 years provide large sample validation



**ASU Log Cost Model - Initial Release**
Using Lognormal Probability Density Function

ln (Hours per New/Mod LOC)

Review Closed Date

ASU Eng Checks & Elec Mtngs
Data Through Review 7351 on 8/25

**A textbook demonstration of an in-control, stable process**

# The High Maturity Data Dilemma
## Why Management Can't Have It Tomorrow

Unstable performance

Stable performance

Improved performance

time

$\Delta t_0$

$\Delta t_1$

$\Delta t_2$

$\Delta t_0$ :
- Process selection
- Analysis of suitability
  for SPC

$\Delta t_1$ & $\Delta t_2$ :
- Identify improvement proposals
- Evaluate & prioritize proposals
- Select improvement
- Pilot improvement
- Deploy improvement

We can minimize $\Delta t_0$, $\Delta t_1$ & $\Delta t_2$ by careful
management, but the length of the data runs will
depend on the periodicity of the process itself

# Handling Human Issues - 2005

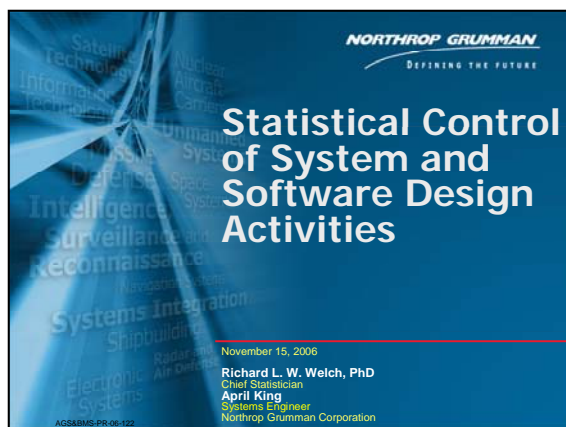## and Institutionalizing the Process

# 2005 Challenges

- First demonstration of CMMI$^{SM}$ Level 4 and 5 capabilities focused on code inspections and parallel effort to control peer reviews of software test plans, procedures and reports
  - High data rates inherent in these back-end processes helped us to understand and overcome statistical difficulties
  - We gained practical lessons learned on the obstacles that had to be overcome

- Desire to introduce successful SPC techniques for quantitative project management in the front-end system and software design phases

- When coding starts
  - Product development is one-half over
  - Opportunities to recognize and correct special & common cause variation in the design process are gone

> First-year decisions determine up to 70% of total life cycle cost on DoD programs.
> Early, effective statistical control offers great practical benefit

# Practical Difficulties at Level 4

- Getting started
  - Selecting good candidates for statistical management

- Statistical innumeracy
  - Discipline needs to own the right skill set

- Little historical data & inherently low data rates
  - Personnel need familiarity with robust statistical procedures

- Cautionary note: you must also take care of the basics (CMMI$^{SM}$ Level 3)
  - Budget and charter
    - Project impacts
  - Metrics infrastructure across engineering
    - Metric definitions
    - Data collection mechanisms
    - Consistency of processes across projects



NORTHROP GRUMMAN
DEFINING THE FUTURE

**Statistical Control of System and Software Design Activities**

November 15, 2006

Richard L. W. Welch, PhD
Chief Statistician
April King
Systems Engineer
Northrop Grumman Corporation

AGS&BMS-PR-06-122

– "Outstanding Presentation for High Maturity"
– "Conference Winner"

Note: Prior CMMI Technology Conference & User Group papers are published on-line at: http://www.dtic.mil/ndia/

14

# Getting Started
## Process Selection for Statistical Management

- Statistical control is imposed on sub-processes at an elemental level in the process architecture

- Processes are selected based on their
  - Business significance – "sufficient conditions"
  - Statistical suitability – "necessary conditions"

- Business checklist
  - Is the candidate sub-process a component of a project's defined key process?
    - Is it significant to success of a business plan goal?
    - Is it a significant contributor to an important estimating metric in the discipline?
  - Is there an identified business need for predictable performance?
    - Cost, schedule or quality
  - How does it impact the business?
    - Need to map sub-process ↔ process ↔ business goal

- Statistical checklist (table)

| PRINCIPAL FACTORS INVOLVED IN SUB-PROCESS SELECTION FOR STATISTICAL PROCESS CONTROL | | |
|---|---|---|
| **PRIMARY QUESTION** | **SUPPORTING QUESTION** | **DEMONSTRATED INDICATOR** |
| Are data collected? | Does the data collection system require update or redeployment? | Data collection system ready for deployment |
| What is the data rate? That is, how often will the process be repeated on the project? | Will the process be repeated frequently enough to develop control limits if such limits do not exist from baseline historical performance? | At least 20-30 data points exist or will be produced |
| Are there historical performance data? | Stable performance: Will the process be performed in roughly the same manner as on previous projects? | A documented procedure or training materials are used by those performing the process |
| Has a control metric been identified? | Has statistical analysis of past project performance identified measures that are indicative of overall process performance? | Control metric can be computed from the collected data |
| Does a baseline exist for the control metric? | What are the average and statistical variation of previous performance? | Performance excursions outside the control limits can be identified & attributed to their root causes |
| Do specification limits exist for process performance? (Optional) | Do limits exist beyond which process performance is deemed unacceptable? | Specification limits are documented |

15

# Overcoming Statistical Innumeracy
## Success Factors

- Minitab

- "Dark green belt" training
  - Curriculum tailored to focus on applied statistical techniques and Minitab familiarity
  - Deming principle applied in the class room
    - In God we trust, *all others bring data*
  - Lean and process management training covered in other courses

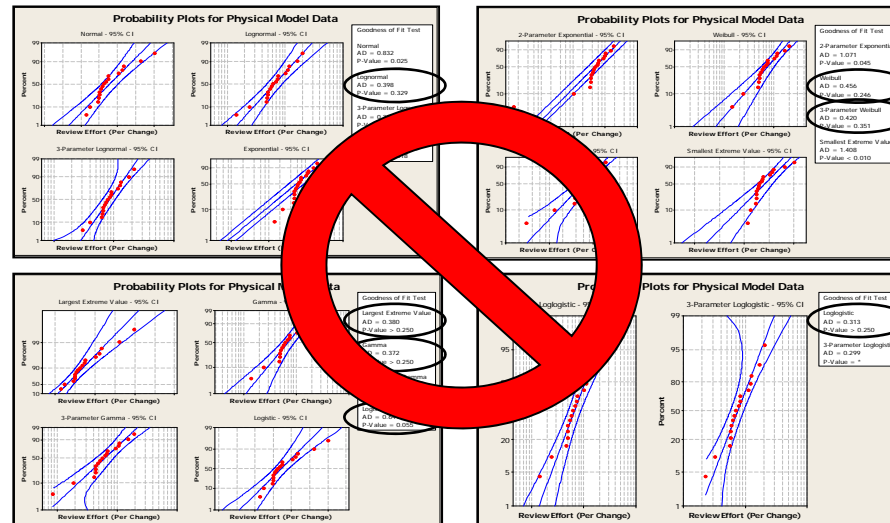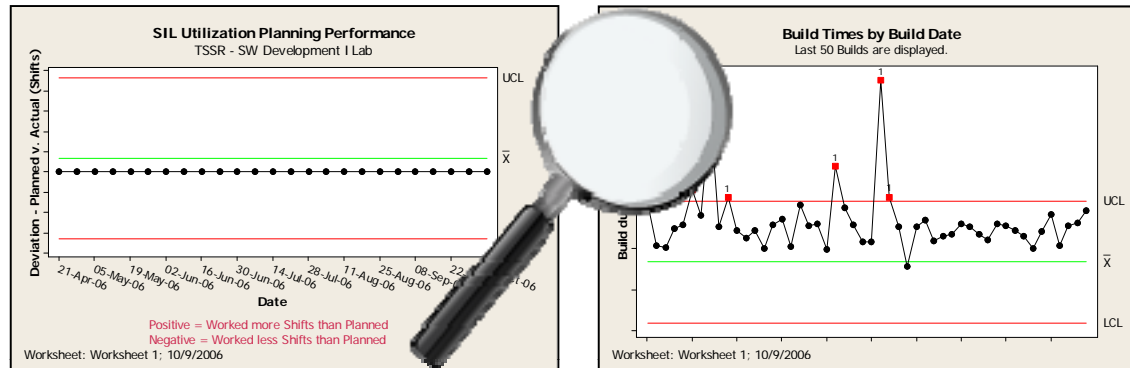- Green belt community of practice

- Chief statistician

**Key success factors:**
**Management recognition & support for the investment**

# Post 2005 Follow-up

- Sector standards for certifying Green Belts, Black Belts & Master Black Belts (2006)
  - Training
  - Project portfolio

- Green Belt certification (2006)

- Black Belt cadre (2007-2008)

- Future Master Black Belt cadre (2009+)

**Success creates a continuing need to grow the infrastructure**

# Pitfalls



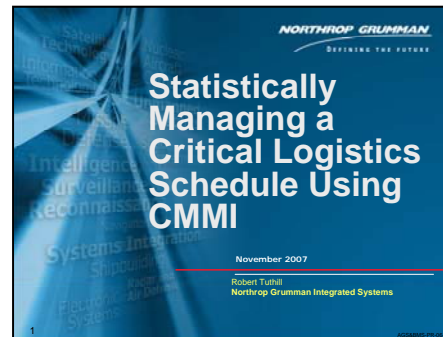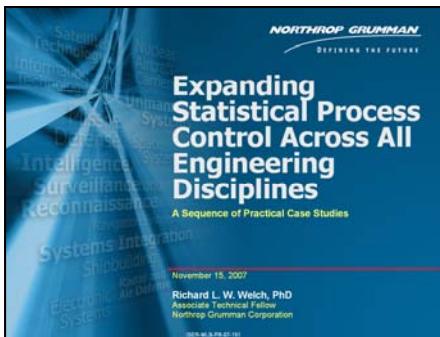**Dealing with what was taught – but not learned – in green belt class**

# Growing the Benefits – 2006

## More Processes, Projects, Disciplines

# Growth

- After our initial success, we aggressively expanded the use of SPC techniques in all engineering projects
  - Led by senior management
  - Clear expectations of significant benefit to the business
  - Particular focus on our hardware and logistics disciplines

- By year-end 2006, we had gone from the original 4 sub-processes under control in 2 Engineering homerooms to 30 sub-processes under control in 6 Engineering homerooms
  - Expect ~45 sub-processes that are significant to our business under active control by year-end 2008





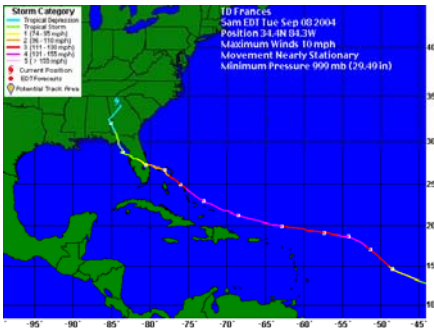- "Outstanding Presentation for High Maturity"
- "Conference Winner"

Note: Prior CMMI Technology Conference & User Group papers are published on-line at: http://www.dtic.mil/ndia/
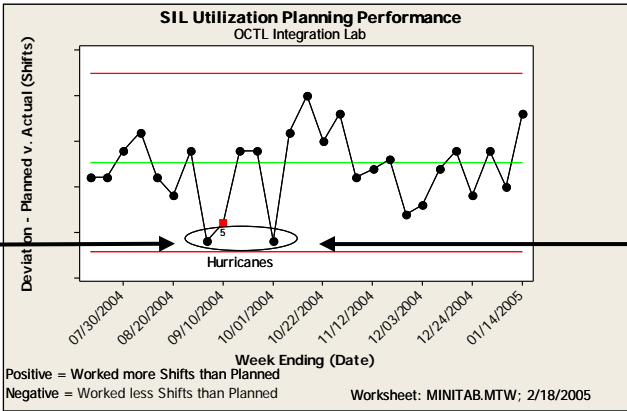
# A Humorous Sidebar
## Identifying Special Causes

- As part of this effort, our Test & Evaluation personnel analyzed some 2004-2005 baseline data, and asked what has become one of our favorite statistical questions:

  Can isolated points be considered as special cause points, and be deleted from a data set as outliers, even though they don't fall outside of the 3-sigma control limits?



**Francis**



**Jeanne**

**NORTHROP GRUMMAN**

# Increasing Our Effectiveness – 2007-2008

## Exploring New Techniques
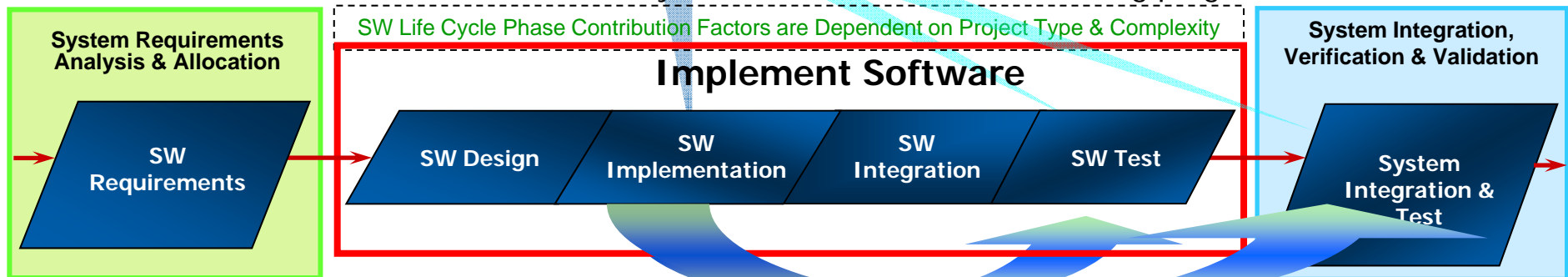
# What Is the Benefit?

- Organizational Impact
  - 100% of delivered code is peer reviewed
  - Average of 105 reviews completed per month
  - This activity affects all major development & test activities after software design
    - SW Implementation
    - Software Test
    - System Test
    - Ground & Flight Test Support
  - Code review effort constitutes a significant portion of the earned value credit in these phases
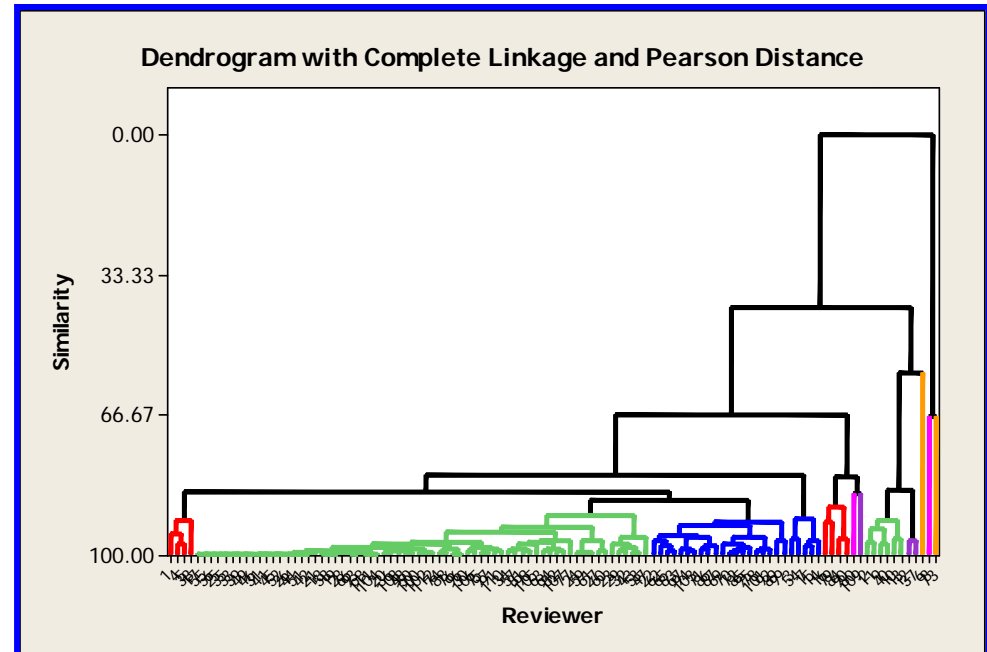
- Benefits
  - Increased early defect detection
  - Fewer delivered defects
  - Increased code maintainability, reduced cost on future sustaining programs

| System Requirements Analysis & Allocation | SW Life Cycle Phase Contribution Factors are Dependent on Project Type & Complexity | System Integration, Verification & Validation |
|---|---|---|
| SW Requirements | **Implement Software** — SW Design \| SW Implementation \| SW Integration \| SW Test | System Integration & Test |

**Peer reviews have a significant impact on downstream product quality and development costs**

# Analytical Approach

- Use accumulated data to explore factors related to reviewer performance and experience

- Use a multivariate clustering procedure (agglomerative hierarchical method) to identify groups of reviewers with similar performance characteristics (initially not known)

- Decide how many groups are logical for the data and classify accordingly

- Three reviewer performance categories support the needed level of insight

  - Group 1 reviewers have lots of review experience, review at the best rates, & identify the most defects

  - Group 2 reviewers are newer & less experienced (reflected by the number of reviews they have completed), with a wide range of rates and discovered defects

  - Group 3 reviewers have lots of experience, review at fast rates, but identify significantly fewer defects



Dendrogram with Complete Linkage and Pearson Distance

24

# A Serious Sidebar
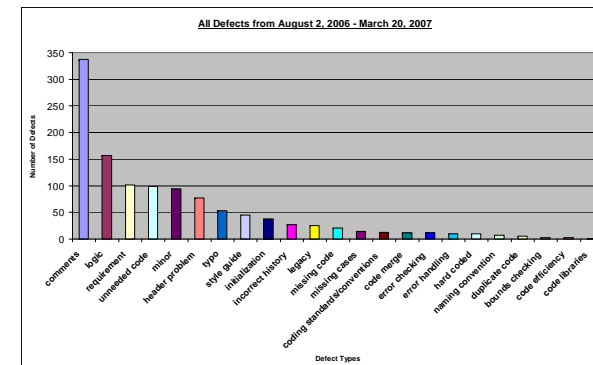## Measuring Individual Performance

- It violates the peer review process to use numbers of defects to measure the Author's performance ("killing the goose that lays the golden egg")
  - Need reviewers free to report any issues they find, even if they are not totally sure that the item is truly a defect
  - Even the very best and most conscientious engineers create defects – the primary objective of the review is to find and remove any defects

- Peer review database design enables study of individual reviewer performance – *with the express goal of increasing skills through vital, focused training*
  - Good reviewers provide an essential contribution both for the author and the company – reviewer diligence should be encouraged, recognized and rewarded
  - Cumulative data on reviewer performance provides valuable insight - similar to measurements applied in sports

> Reviewer knowledge and skill are key–
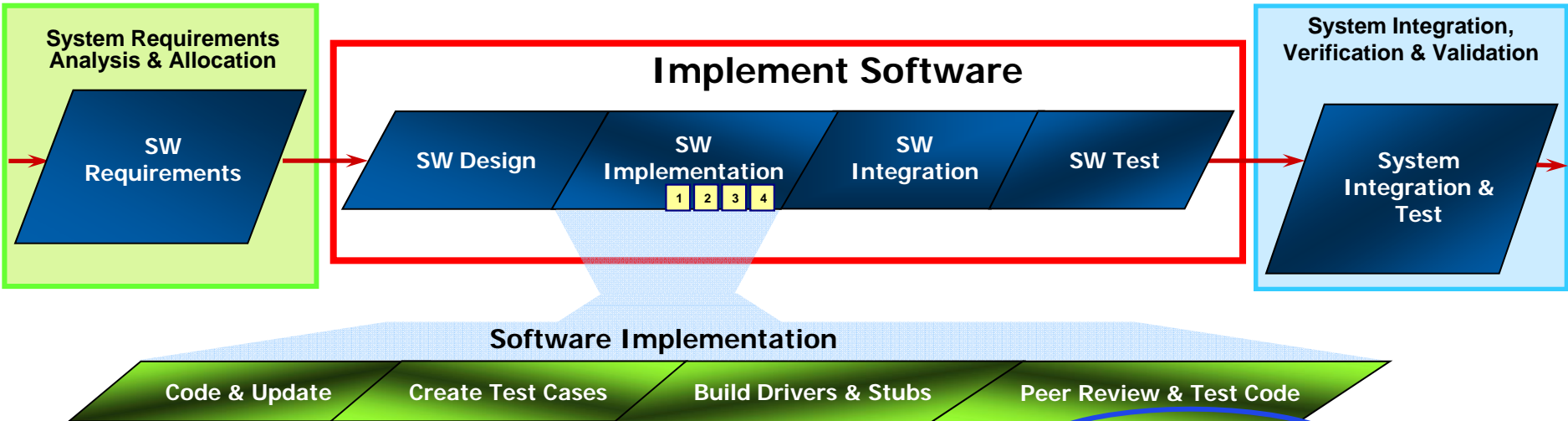> Knowing what to look for & how to find it...

# Causal Analysis & Resolution

- After the fact analysis by Group 1 reviewers indicated Group 2 and Group 3 reviewers consistently miss defects

- A retrospective study focused on the common types of defects being discovered

- An improvement team identified ways to increase the skills of Group 2/3 reviewers
  - Pair Group 2/3 reviewers with Group 1 mentors
  - Review and update coding standards to clarify descriptions or address missing elements
  - Develop and deliver a technical-level review training course to provide refreshed or deeper insight into 'problematic' programming issues
    - 'Problematic' programming issues were identified based on team member experience and results from the retrospective study
  - Enhance checklists

Common issues emerge
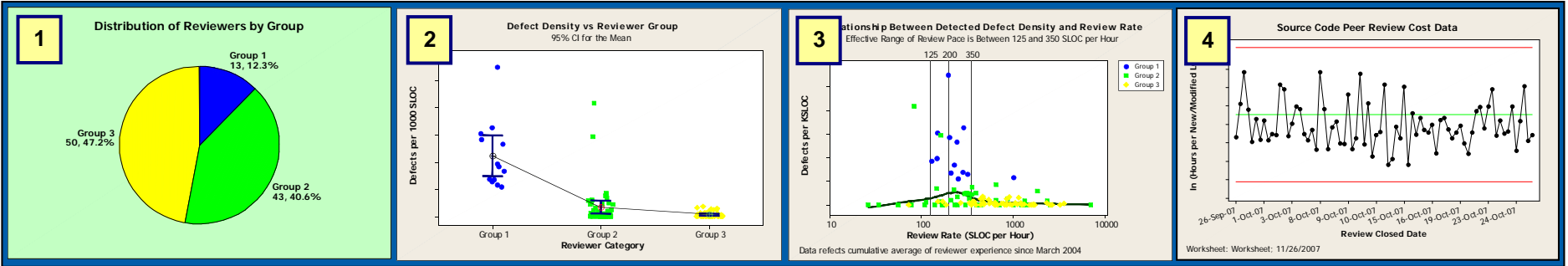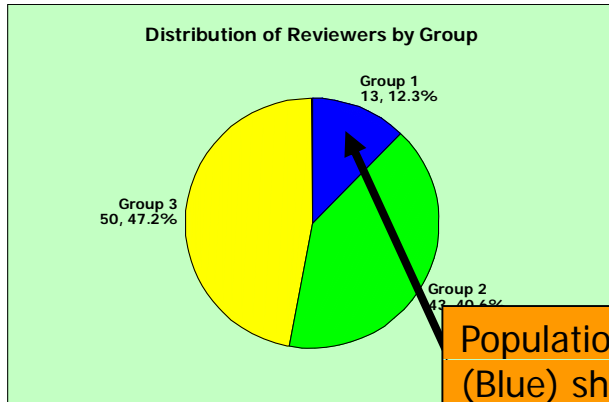when we examine
defect types & frequencies

**All Defects from August 2, 2006 - March 20, 2007**

# Peer Review Effectiveness Metrics

NORTHROP GRUMMAN

**System Requirements Analysis & Allocation**

SW Requirements

**Implement Software**

SW Design | SW Implementation `1` `2` `3` `4` | SW Integration | SW Test

**System Integration, Verification & Validation**

System Integration & Test

**Software Implementation**

Code & Update | Create Test Cases | Build Drivers & Stubs | Peer Review & Test Code

`1` `2` `3` `4`

## Metrics:

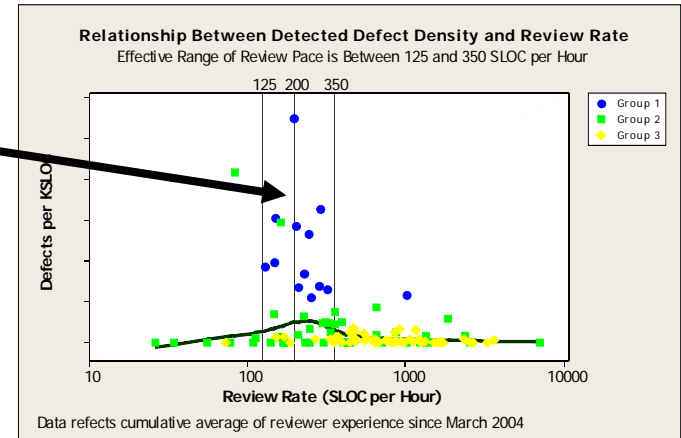| Peer Review Effectiveness Metrics | Units |
|---|---|
| 1 Percentage of Group 1 Reviewers | Percent of Total Reviewers |
| 2 Average Detected Defect Density | Defects per Thousand SLOC |
| 3 Average Review Rate | SLOC per Review Hour |
| 4 Log Cost Model | Log(Total Hours per SLOC) |

Note: SW Test & System Test are decomposed similarly

**1 Distribution of Reviewers by Group**
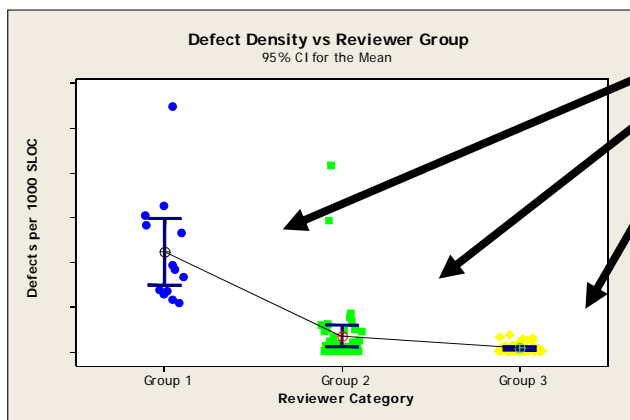- Group 1: 13, 12.3%
- Group 2: 43, 40.6%
- Group 3: 50, 47.2%

**2 Defect Density vs Reviewer Group**
95% CI for the Mean
Defects per 1000 SLOC
Reviewer Category (Group 1, Group 2, Group 3)

**3 Relationship Between Detected Defect Density and Review Rate**
Effective Range of Review Pace is Between 125 and 350 SLOC per Hour
125 200 350
Defects per KSLOC
Review Rate (SLOC per Hour)
Group 1, Group 2, Group 3
Data refects cumulative average of reviewer experience since March 2004

**4 Source Code Peer Review Cost Data**
In (Hours per New/Modified L...
Review Closed Date
26-Sep-07 ... 24-Oct-07
Worksheet: Worksheet; 11/26/2007

27

CLEARED FOR PUBLIC RELEASE BM&ES-MLB-PR-09-56

# Forecasting the Outcome

**Distribution of Reviewers by Group**

Group 1
13, 12.3%

Group 3
50, 47.2%

Group 2
43, 40.6%

**Relationship Between Detected Defect Density and Review Rate**
Effective Range of Review Pace is Between 125 and 350 SLOC per Hour

125 200 350

Group 1
Group 2
Group 3

Defects per KSLOC

10    100    1000    10000

**Review Rate (SLOC per Hour)**

Data refects cumulative average of reviewer experience since March 2004

Population should cluster around ideal review rate of 200 LOC/Hr (industry std)

Population of Group 1 (Blue) should increase

We used four ways to measure changes from the initial March 2007 performance baseline. Demonstrating skill development in review effectiveness does not lend to routine control chart monitoring

**Defect Density vs Reviewer Group**
95% CI for the Mean

Defects per 1000 SLOC

Group 1    Group 2    Group 3

**Reviewer Category**

Detected defect density should increase or remain the same

**Source Code Peer Review Cost Data**

Hours per New/Modified LOC)

7-Oct-07  8-Oct-07  9-Oct-07  10-Oct-07  15-Oct-07  16-Oct-07  19-Oct-07  23-Oct-07  24-Oct-07
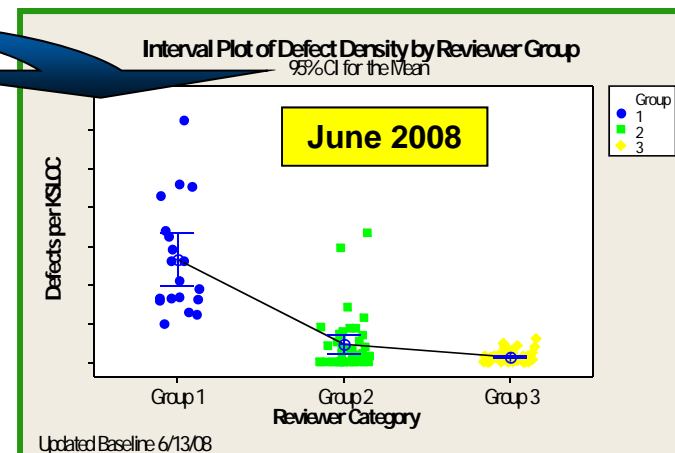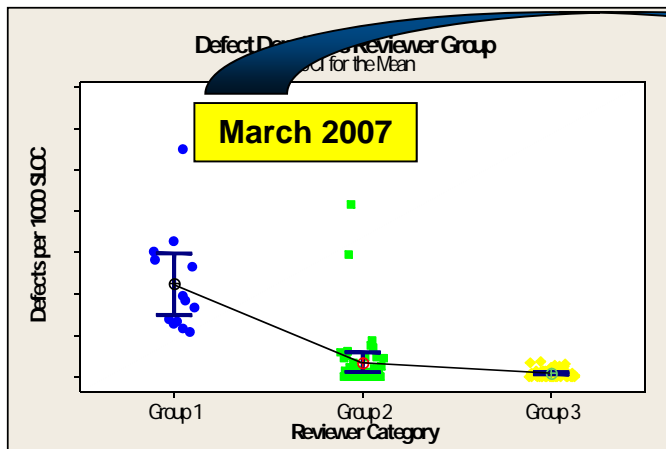
**Review Closed Date**

Worksheet: Worksheet; 11/26/2007

Process cost performance should remain stable
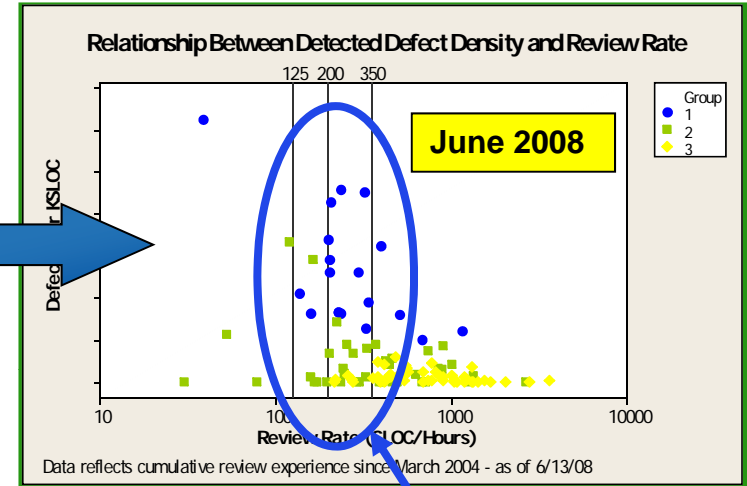
28

# Verifying the Outcome

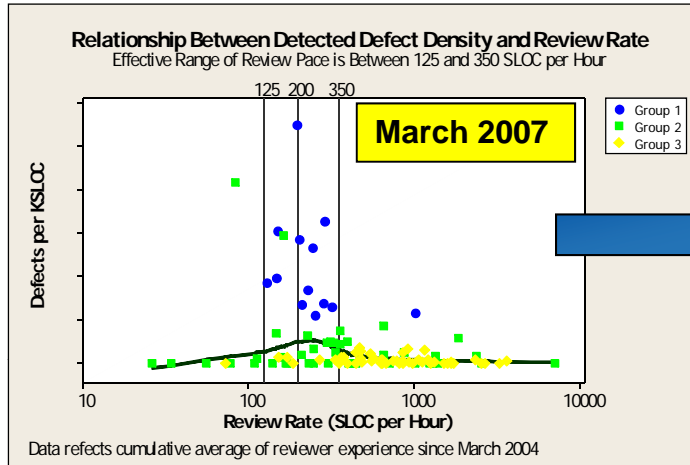NORTHROP GRUMMAN



**Group 1 reviewers increased by 53%**

**Overall discovered defect density increased by 56%**

29

CLEARED FOR PUBLIC RELEASE BM&ES-MLB-PR-09-56

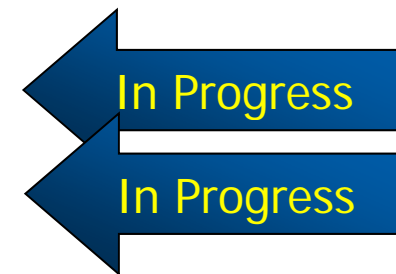# Checking the Control Variables

**Stable Review rate**

**Predictable process**

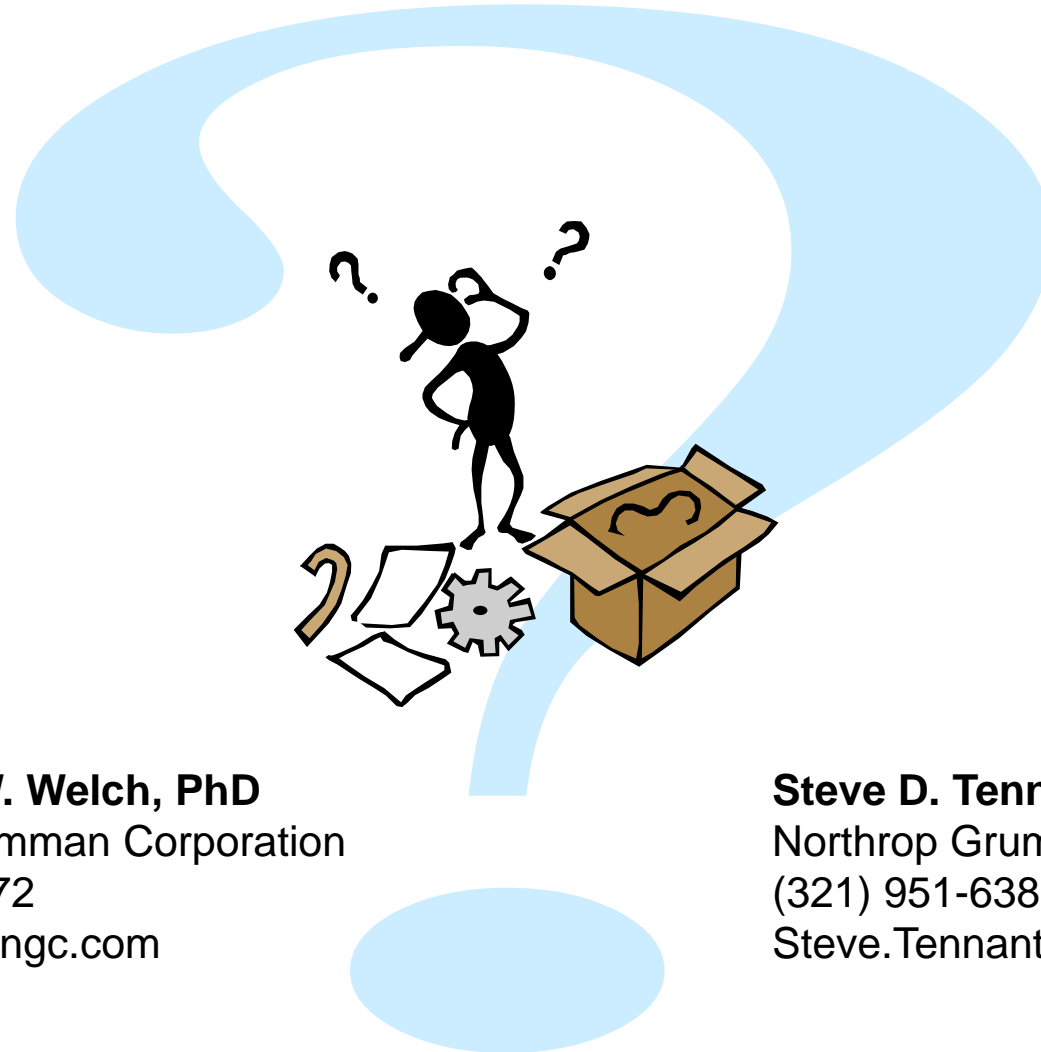**Number of reviewers performing in the ideal range increased**

# Future Pathways

✓ Maintain strategic focus to sharpen skills

    ✓ Continue to support inexperienced developers with mentoring

    ✓ Maintain periodic skill enhancement training

- Continue the quest to remove impediments

    • Better integrated toolsets        **In Progress**

    • Improved coding standards    **In Progress**

---

**Bottom line motivator:**
The 2007-2008 initiative has resulted in a 12% reduction
in the number of software bugs per release

# Questions

**Richard L. W. Welch, PhD**
Northrop Grumman Corporation
(321) 951-5072
Rick.Welch@ngc.com

**Steve D. Tennant**
Northrop Grumman Corporation
(321) 951-6387
Steve.Tennant@ngc.com