# coverity

# Static Analysis Use Case

# Samba and Coverity

David Maxwell for
Scale 9x

coverity

- Launched, March 2006

- DHS sponsored "Open Source Hardening Project"
  - 2006-2009

- Using Coverity's commercial static analysis product to identify bugs at the source code level

- 35 open source projects on day one

- Since grown to 300+ projects

- Over 15,000 bugs fixed

There is no single measure of the effectiveness of a tool on the software development process.

coverity

Since we can never run the same development effort twice, with identical teams, portions of this evaluation are highly subjective.

coverity

- Objective measures
  - Static Analysis produced defect counts
  - Numbers of Bug Reports
  - Defects confirmed as 'real' by the developers

- Subjective measures
  - Anecdotal comments by developers
  - Community feedback
  - 'Support Load' reduction

coverity

- Static Analysis produced defect counts
  - Good objective measure
    - Reproducible
    - Consistent
    - Low effort to collect
    - Automatable
    - "Static Analysis Tools as Early Indicators of Pre-Release Defect Density"  - Microsoft Research Paper

coverity

- Numbers of Bug Reports
  - Potentially useful if all other factors are controlled
  - Not the case in our example
    - Multiple development branches
    - Concurrent new development during defect resolution
    - Userbase changes over time
    - Platform support changes over time

coverity

- Defects confirmed as 'real' by the developers
  - A high False Positive rate would bring the defect count metric into question
  - Would also affect future developer trust in the analysis tool

coverity

- Anecdotal comments by developers
  - Informative, but not comparable between projects

- Community feedback
  - Dependent on the nature of each project's community

- 'Support Load' reduction
  - Difficult to quantify in an open source environment, due to the variety of support channels

coverity

As in most engineering problems…

What do you want to minimize?

- Immediate Cost
- Long Term Cost
- Time
- Manpower
- Ongoing Support

coverity

- Samba
  - Open source networking suite
  - Provides Microsoft protocol compatibility
  - International team, started in Australia
  - Project founded in 1992
  - ~300KLOC -> 850KLOC   2006-now

coverity

- Started regular scanning March 2006
- 14 Developers accessing the results
- Database available 24/7, SAAS
- New analysis every 2 days on average
  - (797 builds in database)

# Use Case – Samba & Coverity

- ## Static Analysis defect counts, 310KLOC



**Defect Count**

# Use Case – Samba & Coverity

**Defect Count**

Day 1: Fixed

4 NULL Pointer derefs
10 Resource leaks
1 Uninitialized data
31 Use after free

But – other changes
that day introduced
new defects

**118** if (!brl_lock) {

**119** return False;

**120** }

Event **func_conv**: Suspicious implicit conversion to function pointer:

"&brl_lock == 0";

did you intend to call the function?

**118**    if (!brl_lock) {

**119**            return False;

**120**    }

```
688 /********************************************************************
689 Lock a range of bytes.
690 ********************************************************************/
691
692 NTSTATUS brl_lock(struct byte_range_lock *br_lck,
693                     uint16 smbpid,
694                     struct process_id pid,
695                     br_off start,
696                     br_off size,
697                     enum brl_type lock_type,
698                     enum brl_flavour lock_flav,
699                     BOOL *my_lock_ctx)
700 {
701         NTSTATUS ret;
702         struct lock_struct lock;
703
704         *my_lock_ctx = False;
```

**Defect Count**

**Defect Count**

coverity

**Defect Count**



Would this graph be solid blue?

coverity

**Defect Count**

Would this graph be solid blue?

■ Defect Count

coverity

- Defects confirmed as 'real' by the developers

13 defects marked False Positive

216 total defects

13 / 216  = 6%

coverity

- Subjective measures
  - Anecdotal comments by developers

"This tool has become part of our process"

coverity

"Using […] source code analysis technology is like having a developer on the team with an inhuman attention to detail, who points out all the corner cases and boundary conditions developers didn't consider when they first wrote the code."

"I code more carefully, because I know my laziness will be caught and embarrass me."

- Community feedback

- Community feedback
  – Invited to give opening keynote at annual Samba conference in 2009

- Whitepaper series  - http://scan.coverity.com/report/
  - Open Source Report 2008
  - Open Source Report 2009
  - Open Source Report 2010 (Android & Supply Chain)

Stop.