

PLDI 2005 -- Ben's slides!

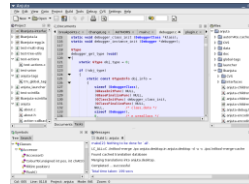


Scalable Statistical Bug Isolation

*Ben Liblit, Mayur Naik, Alice Zheng,
Alex Aiken, and Michael Jordan*

*University of Wisconsin, Stanford
University, and UC Berkeley*

Post-Deployment Monitoring



Goal: Measure Reality

- Where is the black box for software?
 - Crash reporting systems are a start
- Actual runs are a vast resource
 - Number of real runs >> number of testing runs
 - Real-world executions are most important
- This talk: post-deployment bug hunting
 - Mining feedback data for causes of failure

What Should We Measure?

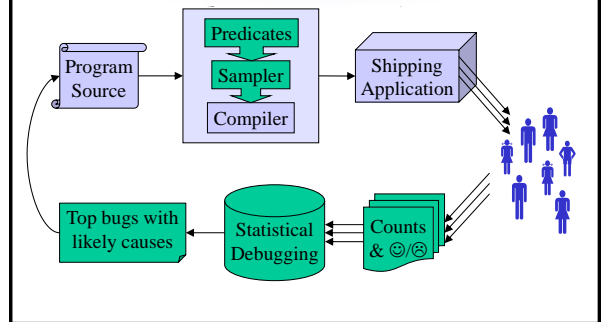
- **Function return values** `err = fetch(file, &obj);`
 - **Control flow decisions** `if (!err && count < size)`
 - **Minima & maxima** `list[count++] = obj;`
 - **Value relationships** `else`
 - **Pointer regions** `unref(obj);`
 - **Reference counts**
 - **Temporal relationships**
- In other words,
lots of things

Our Model of Behavior

Any interesting behavior is expressible as a predicate P on program state at a particular program point.

Count how often “ P observed true” and “ P observed” using sparse but fair random samples of complete behavior.

Bug Isolation Architecture



Find Causes of Bugs

- Gather information about *many* predicates
 - 298,482 predicates in bc
 - 857,384 predicates in Rhythmbox
- Most are not predictive of anything
- How do we find the useful **bug predictors**?
 - Data is incomplete, noisy, irreproducible, ...

Look For Statistical Trends

How likely is failure when P happens?

$F(P)$ = # of failures where P observed true

$S(P)$ = # of successes where P observed true

$$\text{Failure}(P) = \frac{F(P)}{F(P) + S(P)}$$

Good Start, But Not Enough

```
if (f == NULL) {
  x = 0;
  *f;
}
```

Failure(f == NULL) = 1.0
Failure(x == 0) = 1.0

- Predicate `x == 0` is an innocent bystander
 - Program is already doomed

Context

What is the background chance of failure regardless of `P`'s truth or falsehood?

$F(\text{P observed}) = \# \text{ of failures observing P}$

$S(\text{P observed}) = \# \text{ of successes observing P}$

$$\text{Context(P)} = \frac{F(\text{P observed})}{F(\text{P observed}) + S(\text{P observed})}$$

Isolate the Predictive Value of `P`

Does `P` being true *increase* the chance of failure over the background rate?

$$\text{Increase(P)} = \text{Failure(P)} - \text{Context(P)}$$

(a form of likelihood ratio testing)

Increase() Isolates the Predictor

```
if (f == NULL) {
  x = 0;
  *f;
}
```

Increase(f == NULL) = 1.0
Increase(x == 0) = 0.0

Isolating a Single Bug in bc

```

void more_arrays ()
{
  ...

  /* Copy the old arrays. */
  for (indx = 1; indx < old_count;
       arrays[indx] = old_ary[indx];
       #1: indx > scale
       #2: indx > use_math
       #3: indx > opterr
       #4: indx > next_func
       #5: indx > i_base

  /* Initialize the new elements. */
  for (; indx < v_count; indx++)
    arrays[indx] = NULL;
  ...
}

```

It Works!

...for programs with just one bug.

- Need to deal with multiple, unknown bugs
- Redundant predictors are a major problem

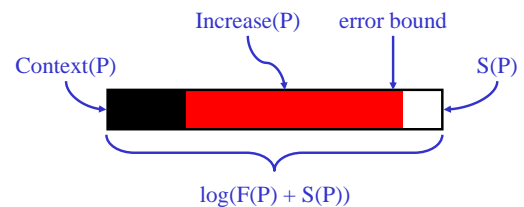
Goal: Isolate the best predictor for each bug, with no prior knowledge of the number of bugs.

Multiple Bugs: Some Issues

- A bug may have many redundant predictors
 - Only need one, provided it is a good one
- Bugs occur on vastly different scales
 - Predictors for common bugs may dominate, hiding predictors of less common problems

Guide to Visualization

- Multiple interesting & useful predicate metrics
- Graphical representation helps reveal trends



Bad Idea #1: Rank by Increase(P)

Thermometer	Context	Increase	S	F	F+S	Predicate
	0.065	0.935±0.019	0	23	23	((!(f1 + i)))>this.last_token < filename
	0.065	0.935±0.020	0	10	10	((!(f1 + i)))>other.last_line == last
	0.071	0.929±0.020	0	18	18	((!(f1 + i)))>other.last_line == filename
	0.073	0.927±0.020	0	10	10	((!(f1 + i)))>other.last_line == yy_n_chars
	0.071	0.929±0.028	0	19	19	bytes == filename
	0.075	0.925±0.022	0	14	14	((!(f1 + i)))>other.first_line == 2
	0.076	0.924±0.022	0	12	12	((!(f1 + i)))>this.first_line < mid
	0.077	0.923±0.023	0	10	10	((!(f1 + i)))>other.last_line == yy_init
..... 2732 additional predictors follow						

- High Increase() but very few failing runs!
- These are all *sub-bug predictors*
 - Each covers one special case of a larger bug
- Redundancy is clearly a problem

Bad Idea #2: Rank by F(P)

Thermometer	Context	Increase	S	F	F+S	Predicate
	0.176	0.007±0.012	22554	5045	27599	files[filesindex].language != 15
	0.176	0.007±0.012	22566	5045	27611	tmp == 0 is FALSE
	0.176	0.007±0.012	22571	5045	27616	strcmp != 0
	0.176	0.007±0.013	18094	4251	23145	tmp == 0 is FALSE
	0.176	0.007±0.013	18885	4240	23125	files[filesindex].language != 14
	0.176	0.008±0.013	17557	4007	21764	filesindex == 25
	0.177	0.008±0.014	16453	3731	20184	new value of M < old value of M
	0.176	0.261±0.023	4800	3716	8516	config.windowing_window_size != arg
..... 2732 additional predictors follow						

- Many failing runs but low Increase()!
- Tend to be *super-bug predictors*
 - Each covers several bugs, plus lots of junk

A Helpful Analogy

- In the language of information retrieval
 - Increase(P) has high precision, low recall
 - F(P) has high recall, low precision
- Standard solution:
 - Take the harmonic mean of both
 - Rewards high scores in both dimensions

Rank by Harmonic Mean

Thermometer	Context	Increase	S	F	F+S	Predicate
	0.176	0.824±0.009	0	1585	1585	files[filesindex].language > 16
	0.176	0.824±0.009	0	1584	1584	strcmp == 0
	0.176	0.824±0.009	0	1580	1580	strcmp == 0
	0.176	0.824±0.009	0	1577	1577	files[filesindex].language == 17
	0.176	0.824±0.009	0	1576	1576	tmp == 0 is TRUE
	0.176	0.824±0.009	0	1573	1573	strcmp == 0
	0.116	0.883±0.012	1	724	725	((!(f1 + i)))>this.last_line == 1
	0.116	0.883±0.012	1	776	777	((!(f1 + i)))>other.last_line == yyieng
..... 2732 additional predictors follow						

- It works!
 - Large increase, many failures, few or no successes
- But redundancy is still a problem

Redundancy Elimination

- One predictor for a bug is interesting
 - Additional predictors are a distraction
 - Want to explain each failure once
- Similar to minimum set-cover problem
 - Cover all failed runs with subset of predicates
 - Greedy selection using harmonic ranking





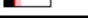
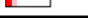
Simulated Iterative Bug Fixing

1. Rank all predicates under consideration
2. Select the top-ranked predicate **P**
3. Add **P** to bug predictor list
4. Discard **P** and all runs where **P** was true
 - Simulates fixing the bug predicted by **P**
 - Reduces rank of similar predicates
5. Repeat until out of failures or predicates

Simulated Iterative Bug Fixing

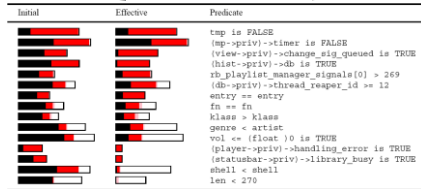
1. Rank all predicates under consideration
2. Select the top-ranked predicate **P**
3. Add **P** to bug predictor list
4. Discard **P** and all runs where **P** was true
 - Simulates fixing the bug predicted by **P**
 - Reduces rank of similar predicates
5. Repeat until out of failures or predicates

Experimental Results: `exif`

Initial	Effective	Predicate
		<code>i < 0</code>
		<code>maxlen > 1900</code>
		<code>o + s > buf_size is TRUE</code>

- 3 bug predictors from 156,476 initial predicates
- Each predicate identifies a distinct crashing bug
- All bugs found quickly using analysis results

Experimental Results: Rhythmbox

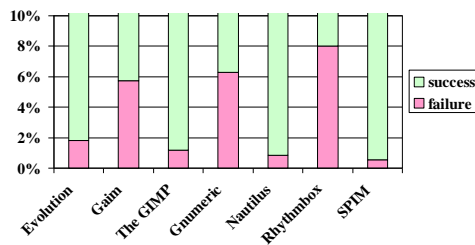


- 15 bug predictors from 857,384 initial predicates
- Found and fixed several crashing bugs

Lessons Learned

- Can learn a lot from actual executions
 - Users are running buggy code anyway
 - We should capture some of that information
- Crash reporting is a good start, but...
 - Pre-crash behavior can be important
 - Successful runs reveal correct behavior
 - Stack alone is not enough for 50% of bugs

Public Deployment in Progress



Join the Cause!

The Cooperative Bug Isolation Project

<http://www.cs.wisc.edu/cbi/>



"Borrowed" with only minor reduction – thank you, Amitabh and Jay!

Effectively Prioritizing Tests in Development Environment

ISSTA 2002

Amitabh Srivastava
Jay Thiagarajan
PPRC, Microsoft Research

Using program changes

- Source code differencing
 - S. Elbaum, A. Malishevsky & G. Rothermel “Test case prioritization: A family of empirical studies”, Feb. 2002
 - S. Elbaum, A. Malishevsky & G. Rothermel “Prioritizing test cases for regression testing” Aug. 2000
 - F. Vokolos & P. Frankl, “Pythia: a regression test selection tool based on text differencing”, May 1997

Using program changes

- Data and control flow analysis
 - T. Ball, “On the limit of control flow analysis for regression test selection” Mar. 1998
 - G. Rothermel and M.J. Harrold, “A Safe, Efficient Regression Test Selection Technique” Apr. 1997
- Code entities
 - Y. F. Chen, D.S. Rosenblum and K.P. Vo “TestTube: A System for Selective Regression Testing” May 1994

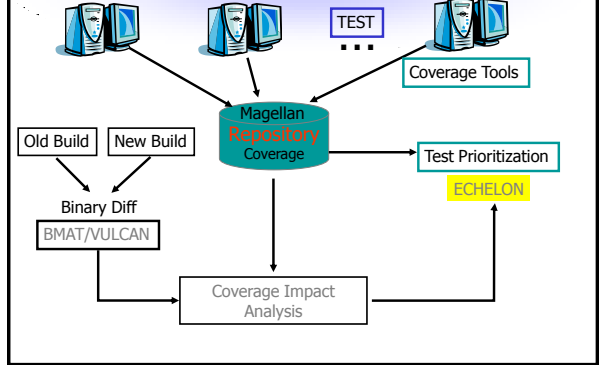
Analysis of various techniques

- Source code differencing
 - Simple and fast
 - Can be built using commonly available tools like “diff”
 - Simple renaming of variable will trip off
 - Will fail when macro definition changes
 - To avoid these pitfalls, static analysis is needed
- Data and control flow analysis
 - Flow analysis is difficult in languages like C/C++ with pointers, casts and aliasing
 - Interprocedural data flow techniques are extremely expensive and difficult to implement in complex environment

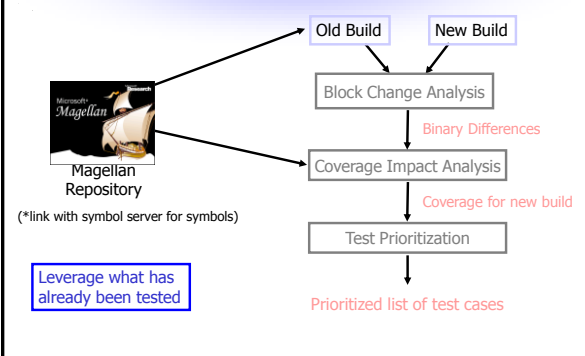
Our Solution

- Focus on change from previous version
 - Determine change at very fine granularity – basic block/instruction
- Operates on binary code
 - Easier to integrate in production environment
 - Scales well to compute results in minutes
- Simple heuristic algorithm to predict which part of code is impacted by the change

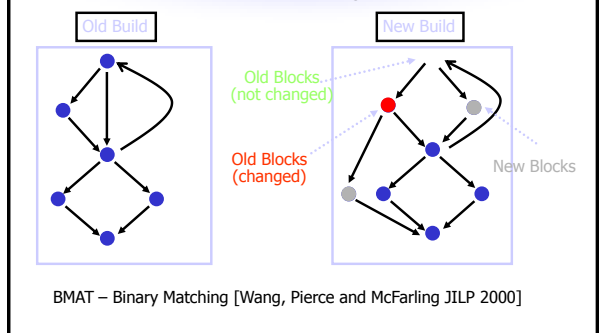
Test Effectiveness Infrastructure



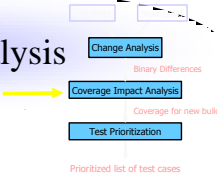
Echelon : Test Prioritization



Block Change Analysis: Binary Matching

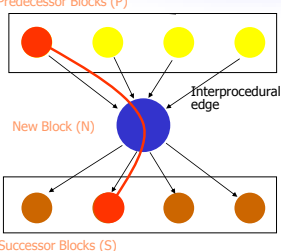


Coverage Impact Analysis



- Terminology
 - Trace: collection of one or more test cases
 - Impacted Blocks: old modified and new blocks
- Compute the coverage of traces for the new build
 - Coverage for old (unchanged and modified) blocks are same as the coverage for the old build
 - Coverage for new nodes requires more analysis

Coverage Impact Analysis



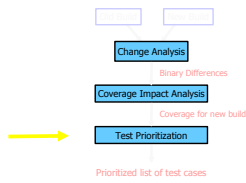
- A Trace may cover a new block N if it covers at least one Predecessor block and at least one Successor Block
- If P or S is a new block, then its Predecessors or successors are used (iterative process)

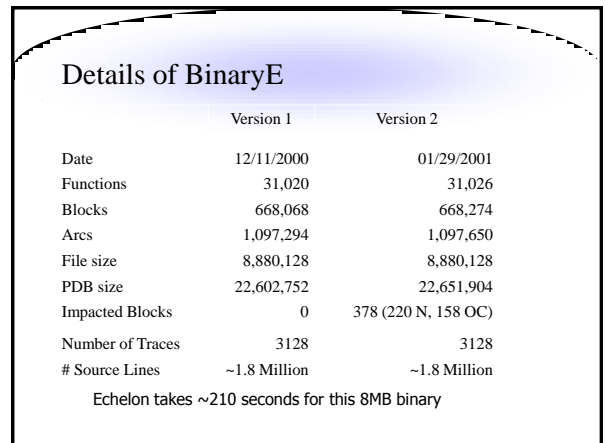
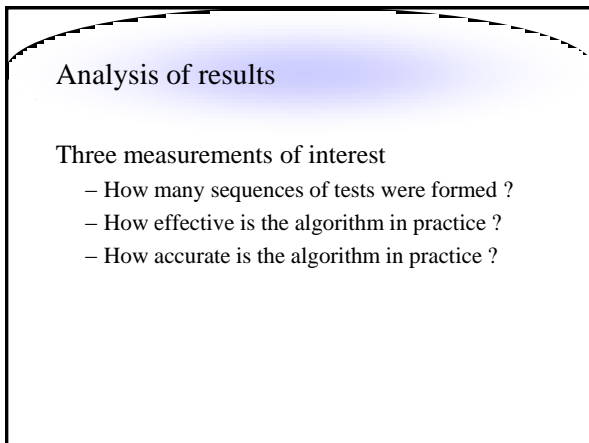
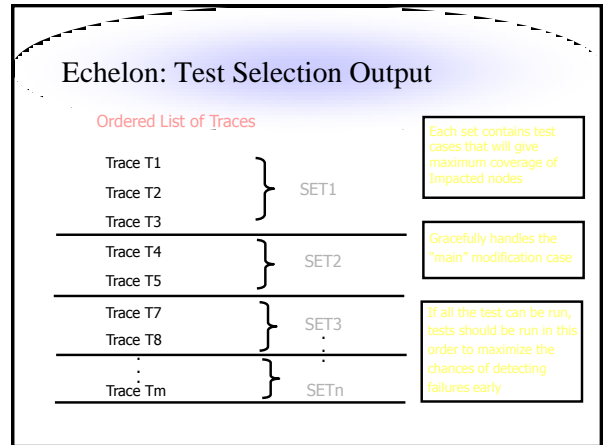
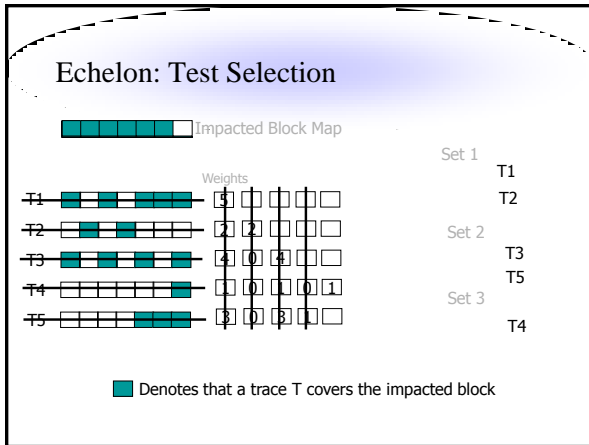
Coverage Impact Analysis

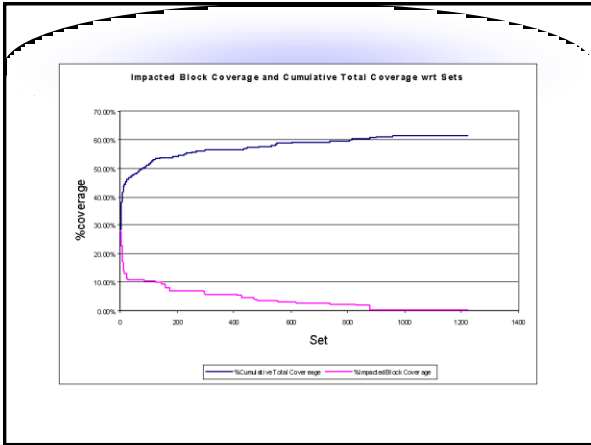
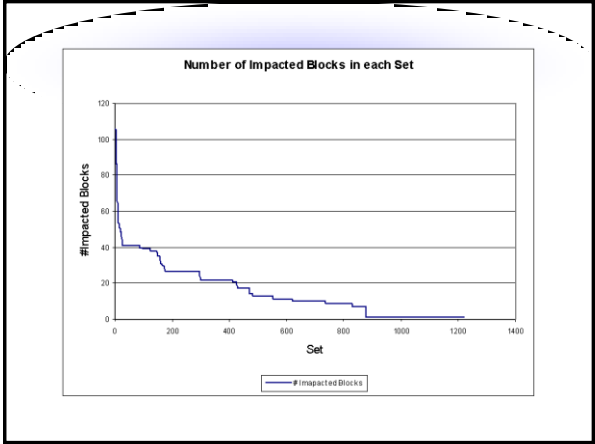
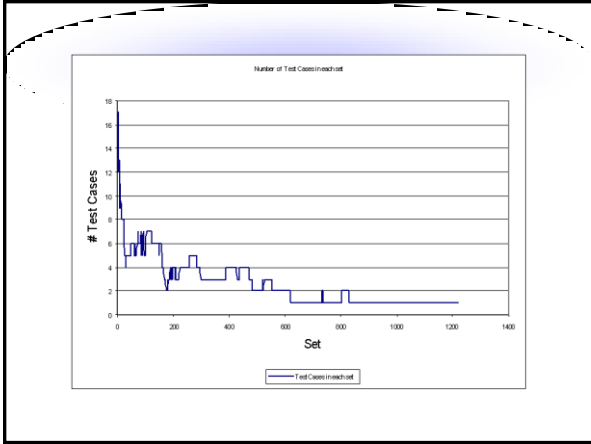
- Limitations - New node may not be executed
 - If there is a path from successor to predecessor
 - If there are changes in control path due to data changes

Echelon : Test Case Prioritization

- Detects minimal sets of test cases that are likely to cover the impacted blocks (old changed and new blocks)
 - Input is traces (test cases) and a set of impacted blocks
 - Uses a greedy iterative algorithm for test selection





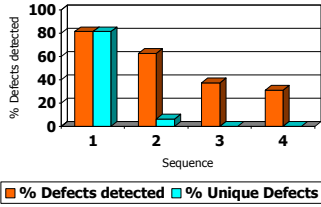


Effectiveness of Echelon

- Important Measure of effectiveness is early defect detection
- Measured % of defects vs. % of unique defects in each sequence
- Unique defects are defects not detected by the previous sequence

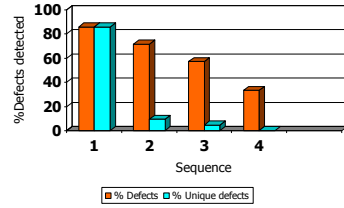
Effectiveness of Echelon

Defects detected in each sequence

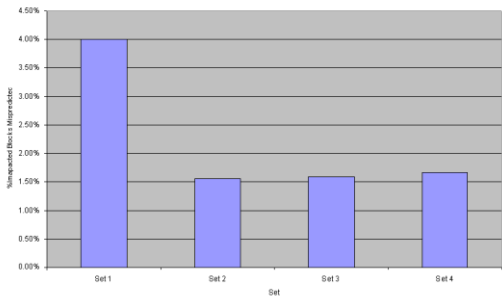


Effectiveness of Echelon

Defects detected in each sequence

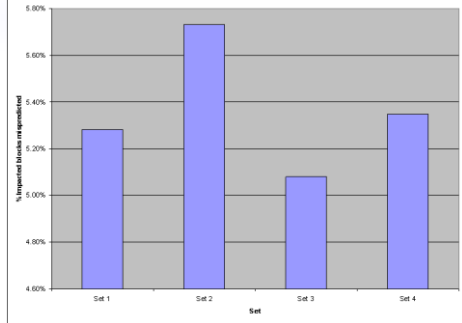


Mispredict due to Limitations

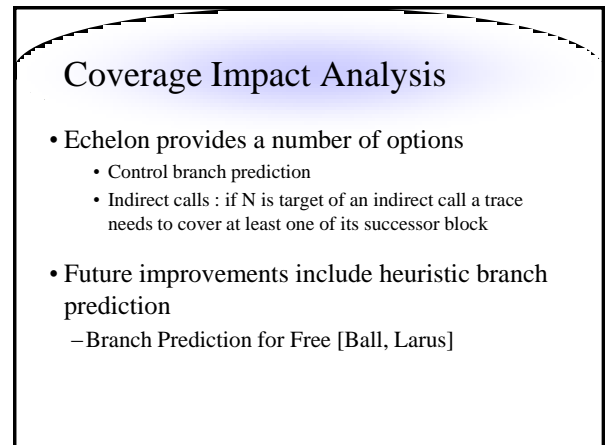
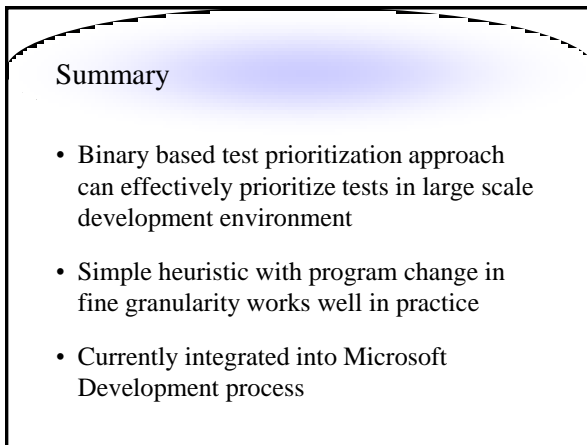
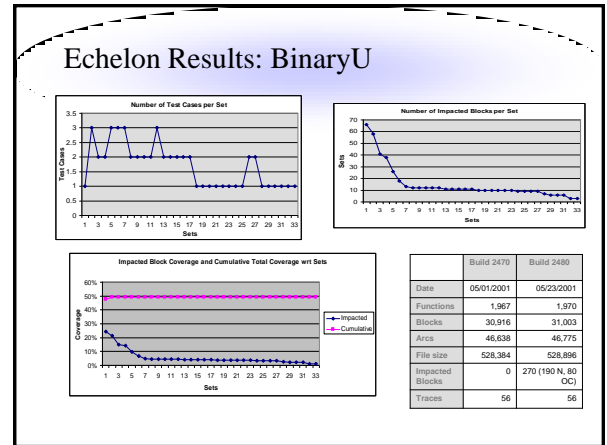
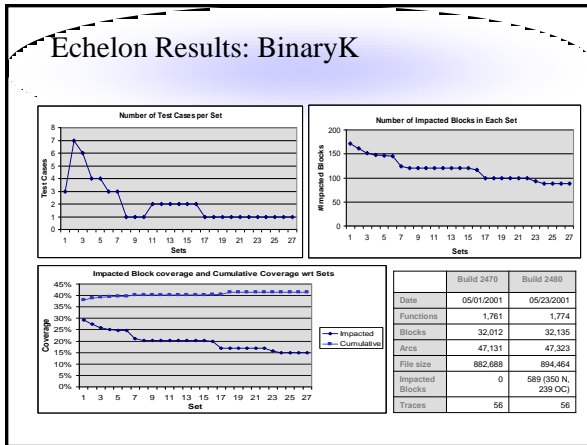


Blocks predicted hit that were not hit

Mispredict due to conservative approach



Blocks predicted not hit that were actually hit
(Blocks were target of indirect calls are being predicted as not hit)



Echelon: Test Selection

- Options

- Calculations of weights can be extended, e.g. traces with great historical fault detection can be given additional weights
- Include time each test takes into calculation
- Print changed (modified or new) source code that may not be covered by any trace
- Print all source code lines that may not be covered by any trace