

Using Covering Arrays to Construct a Benchmark For Concurrency Testing Tools

==

Using Combinatorial Test Design Techniques for Benchmark Construction

IBM Haifa Research Labs

Based on joint work of
Itai Segall, Eitan Farchi (IBM Haifa Research Labs)
Jeremy Bardbury, Kevin Jalbert and David Kelk
(University of Ontario Institute of Technology)

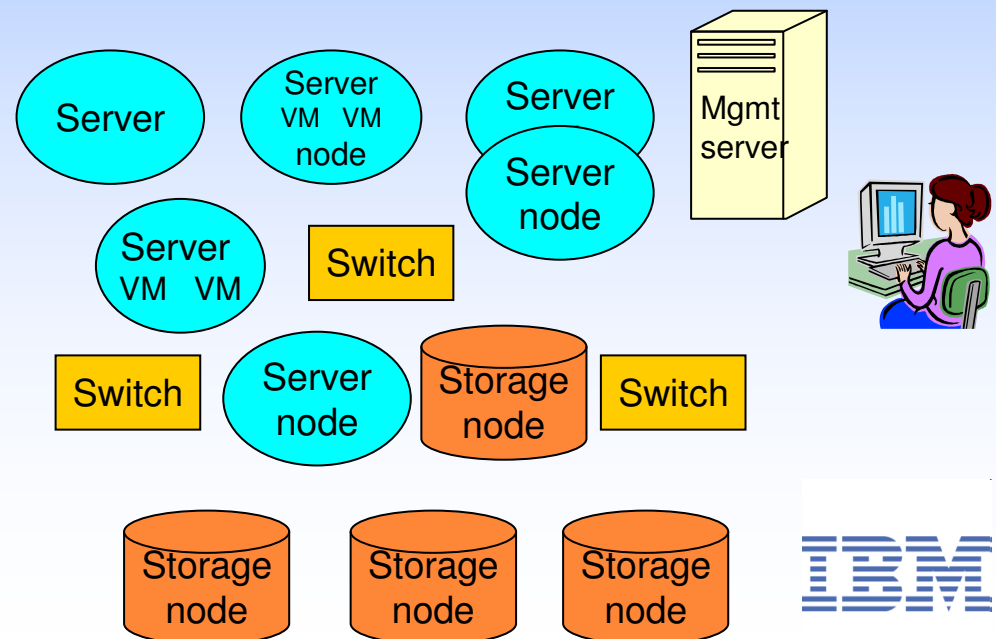


Background – Concurrency Bug Detection Techniques

- Concurrency Testing
 - Coverage of the interleaving space in addition to coverage of the source code, inputs, etc.
- Static Analysis
 - Typically aimed at specific patterns or bugs
- Software Model Checking (exhaustive search)
 - Scalability issues
- No satisfactory benchmark exists
 - Effective
 - Small
 - Systematic

IBM Case Study System

- A distributed application for managing system resources in a networked environment.
 - Already described in IR10.2.
- Consists of:
 - a management server
 - application that communicates with multiple managed clients and with users of the management system.
 - Managed clients
 - physical or virtual resources distributed over a network.
- Managed resources include servers, virtual servers, storage devices, and network devices
- Used by IBM customers for managing IBM hardware and virtual devices, such as servers, Virtual Machines (VMs), switches and storage devices.
- The case study was performed on some new components of a version of this system which is still under development and has not yet been released for customer use.
 - Uses simulated environment



Background – Covering Arrays

- $CA(k,v,t)$ is an array:
 - k columns
 - v possible values in each
 - All combinations of size t appear at least once

Background – Covering Arrays – cont.

- $CA(4,2,2)$:

0	0	0	0
0	1	1	1
1	0	1	1
1	1	0	1
1	1	1	0

Background – Covering Arrays – cont.

- CA(4,2,2):

0	0	0	0
0	1	1	1
1	0	1	1
1	1	0	1
1	1	1	0



Background – Covering Arrays – cont.

- CA(4,2,2):

0	0	0	0
0	1	1	1
1	0	1	1
1	1	0	1
1	1	1	0

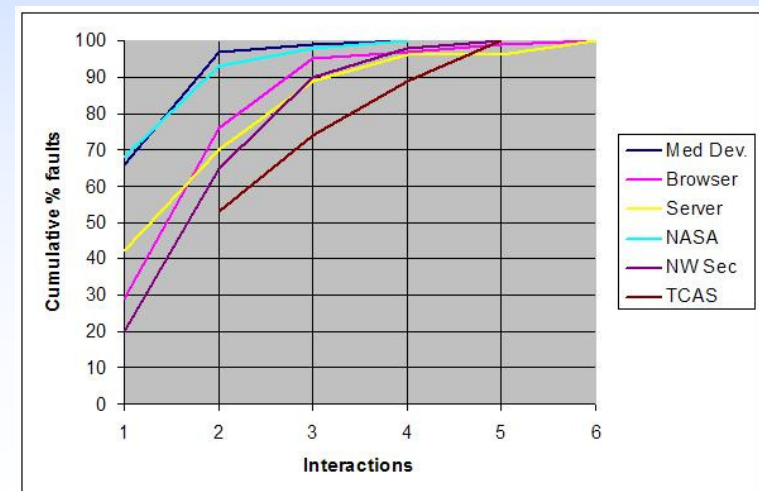


Background – Covering Arrays – cont.

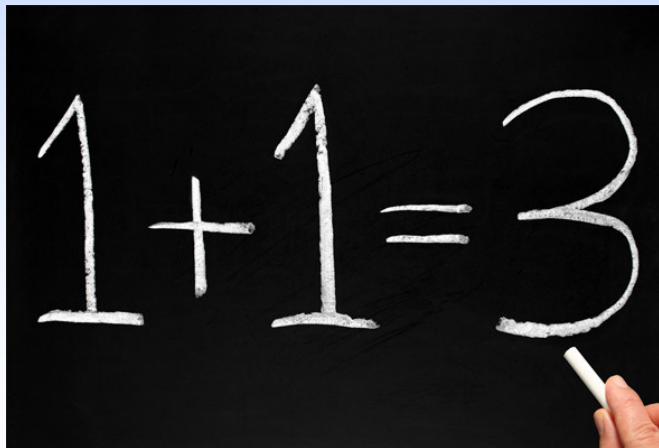
- Extensions to “pure” covering arrays:
 - Restrictions
 - Mixed-values
 - Mixed-strength
 - And many more...

Background – Combinatorial Test Design (CTD)

- CTD = Using covering arrays for testing
- Very well studied for 30 years
- Based on studies that show that software defects are caused by a small combination of attributes



Using Covering Arrays to Construct a Benchmark For Concurrency Testing Tools



Rationale

- Effective benchmark – differentiates between tools
 - For each pair of tools – have at least one task sample on which one outperforms the other
 - Typically, would want also another task sample on which the “other” outperforms the “one”
- (Intuition-only) claim: tool differentiation results from small combinations of factors

Rationale – cont.

- Exhaustive search tools vs. “noise-making” tools
 - Exhaustive tools choke on large programs and/or large # threads and/or large # critical regions
 - Noise-making tools will have difficulty finding certain bug patterns if the bugs are deep and/or rare
- “Noise-making” tools
 - Different heuristics are less/more effective based on combination of # threads and # critical regions

The proposed model

- Program size – # statements
- Program size – # critical regions
- Program size – % statements in critical regions
- Num threads
- Path error density
- Bug depth
- Bug pattern

Each combination here characterizes a (buggy) concurrent program

Pairwise plan for the model

	Program Size - # Statements	Program Size - # Critical Regions	Program Size - % Statements in Critical Regions	# Threads	Path Error Density	Bug Depth	Bug Pattern
1	Small	Small	Large	Small	Medium	High	TwoStageAccess
2	Medium	Small	Medium	Large	VeryLow	Medium	NonAtomicAssumedAtomic
3	Large	Medium	Small	Small	Low	Medium	BlockingCriticalSection
4	Medium	Large	Large	Medium	Low	Low	Interference
5	Small	Medium	Medium	VeryLarge	High	Low	OrphanedThread

... there are 44 such programs in our plan

... but you can't run a table
- or -
finding the needles in the haystack

- Our plan characterizes 44 combinations out of > 14,000 possible in the model
- No chance of finding actual programs that exactly match these



Mutants to the rescue...

- We therefore propose to:
 - Leverage existing benchmarks
 - Apply program mutation to generate new programs



Analyzing existing benchmarks

	Program Size – (sloc, statements)	Program Size – # Critical Regions	Program Size – % Statements in Critical Regions	Bug Pattern
account	Small(139,77)	Medium(7)	Small(11.6883)	NoLock
airlinetickets	Small(61,34)	Small(0)	Small(0)	Interference
allocationvector	Small(163,83)	Small(3)	Large(22.8916)	TwoStageAccess
boundedbuffer	Small(328,192)	Medium(5)	Large(20.3125)	NotifyInsteadOfNotifyAll
bubblesort	Small(236,84)	Small(3)	Medium(11.9048)	NonAtomicAssumedAtomic, OrphanedThread
bubblesort2	Small(98,43)	Small(1)	Medium(6.97674)	Initialization-Sleep

To be continued...



ConMAAn

- 25 mutation operators, in 5 categories:
 - Modify parameters of concurrent methods
 - Modify the occurrence of concurrency method calls (re-moving, replacing, exchanging)
 - Modify concurrency keywords (addition and removal)
 - Switch concurrent objects
 - Modify critical regions (shift, expand, shrink, split)



Zooming out...



	Combinatorial Test Design	Combinatorial Benchmark Design

Zooming out...



	Combinatorial Test Design	Combinatorial Benchmark Design																																										
Rationale	<p>The graph plots 'Considered Tests' (Y-axis, 0-100) against 'Interactions' (X-axis, 1-6). Five data series are shown: Med Dev (blue), Browser (red), Server (green), NASA (yellow), and NW Sec (purple). All series show an increasing trend, with Med Dev and Browser reaching 100% tests by 6 interactions, while others reach it by 5 or 6 interactions.</p> <table border="1"><caption>Approximate data from the graph</caption><thead><tr><th>Interactions</th><th>Med Dev</th><th>Browser</th><th>Server</th><th>NASA</th><th>NW Sec</th></tr></thead><tbody><tr><td>1</td><td>40</td><td>30</td><td>20</td><td>10</td><td>10</td></tr><tr><td>2</td><td>70</td><td>50</td><td>40</td><td>20</td><td>20</td></tr><tr><td>3</td><td>90</td><td>70</td><td>60</td><td>40</td><td>30</td></tr><tr><td>4</td><td>95</td><td>80</td><td>70</td><td>50</td><td>40</td></tr><tr><td>5</td><td>98</td><td>90</td><td>80</td><td>60</td><td>50</td></tr><tr><td>6</td><td>100</td><td>100</td><td>90</td><td>70</td><td>60</td></tr></tbody></table>	Interactions	Med Dev	Browser	Server	NASA	NW Sec	1	40	30	20	10	10	2	70	50	40	20	20	3	90	70	60	40	30	4	95	80	70	50	40	5	98	90	80	60	50	6	100	100	90	70	60	(Intuitively) - differentiation results from small combinations
Interactions	Med Dev	Browser	Server	NASA	NW Sec																																							
1	40	30	20	10	10																																							
2	70	50	40	20	20																																							
3	90	70	60	40	30																																							
4	95	80	70	50	40																																							
5	98	90	80	60	50																																							
6	100	100	90	70	60																																							

Zooming out...



	Combinatorial Test Design	Combinatorial Benchmark Design
Rationale	<p>The graph plots 'Considered Tests' (0-100) against 'Interactions' (1-6). It shows that for 6 interactions, all domains reach 100% of tests. For 5 interactions, Med Dev is at ~95%, Browser at ~90%, Server at ~85%, NASA at ~80%, NW Sec at ~75%, and TCAS at ~70%. For 4 interactions, Med Dev is at ~85%, Browser at ~80%, Server at ~75%, NASA at ~70%, NW Sec at ~65%, and TCAS at ~60%. For 3 interactions, Med Dev is at ~75%, Browser at ~70%, Server at ~65%, NASA at ~60%, NW Sec at ~55%, and TCAS at ~50%. For 2 interactions, Med Dev is at ~65%, Browser at ~60%, Server at ~55%, NASA at ~50%, NW Sec at ~45%, and TCAS at ~40%. For 1 interaction, Med Dev is at ~55%, Browser at ~50%, Server at ~45%, NASA at ~40%, NW Sec at ~35%, and TCAS at ~30%.</p>	(Intuitively) - differentiation results from small combinations
Models	Describe points of variability in tests – potential causes of defects	Describe points of variability in task samples – potential causes of variability in tool effectiveness



Zooming out...



	Combinatorial Test Design	Combinatorial Benchmark Design
Rationale	<p>The graph plots 'Considered Tests' (0-100) against 'Interactions' (1-6). It shows that for 6 interactions, Med Dev and Browser consider nearly all tests, while other domains like TCAS consider significantly fewer tests.</p>	(Intuitively) - differentiation results from small combinations
Models	Describe points of variability in tests – potential causes of defects	Describe points of variability in task samples – potential causes of variability in tool effectiveness
Analyzing existing artifacts	Tests commonly need to be abstracted to the level of the model	Samples in existing benchmarks need to be analyzed



Zooming out...



	Combinatorial Test Design	Combinatorial Benchmark Design
Rationale	<p>The graph plots 'Considered Tests' (0-100) against 'Interactions' (1-6). The curves show that as interactions increase, the number of tests considered approaches 100. Med Dev reaches 100% tests at 2 interactions, while other domains like Browser and Server reach 100% at 3 interactions.</p>	(Intuitively) - differentiation results from small combinations
Models	Describe points of variability in tests – potential causes of defects	Describe points of variability in task samples – potential causes of variability in tool effectiveness
Analyzing existing artifacts	Tests commonly need to be abstracted to the level of the model	Samples in existing benchmarks need to be analyzed
Generation		Mutations

Thanks!

Questions ?