Automated Analysis of the Results of Performance Regression Tests

SPEC DevOps Performance WG Friday, 05/29/2015

Zhen Ming (Jack) Jiang

Software Construction, AnaLysis and Evaluation (SCALE) Lab York University, Canada

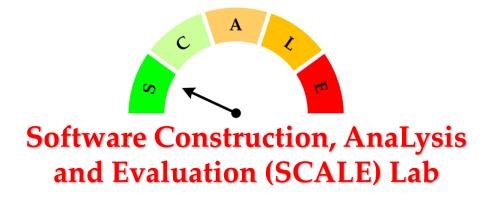




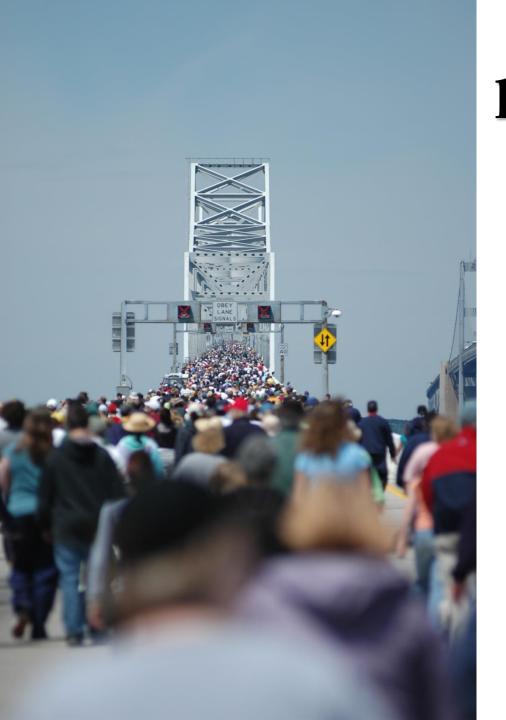


Zhen Ming (Jack) Jiang

- Assistant Professor, Department of Electrical Engineering and Computer Science, Lassonde School of Engineering, York University. Toronto, ON, Canada.
- Head of <u>Software Construction</u>, <u>AnaLytics and Evaluation (SCALE) Lab</u>
- Co-founder and co-organizer for the International Workshop on Large-Scale Testing (LT)
- Research Interests:
 - Software Performance Engineering
 - Mining Software Engineering Data







Most field problems for large scale systems are rarely functional instead they are load-related



















Software Testing



Unit Testing



(Functional) **Regression Testing**

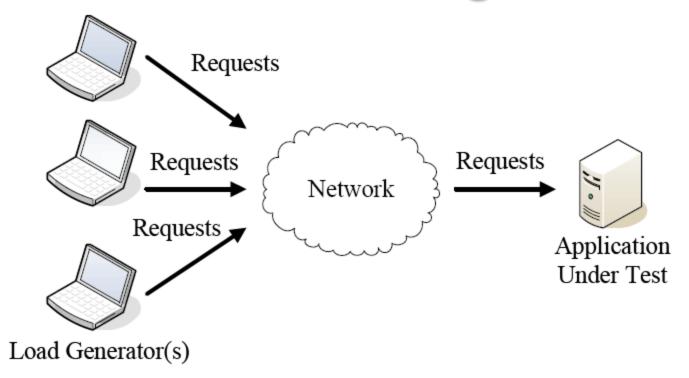
System Integration Testing



Non-Functional Testing

Performance Regression Testing

Performance Regression Testing

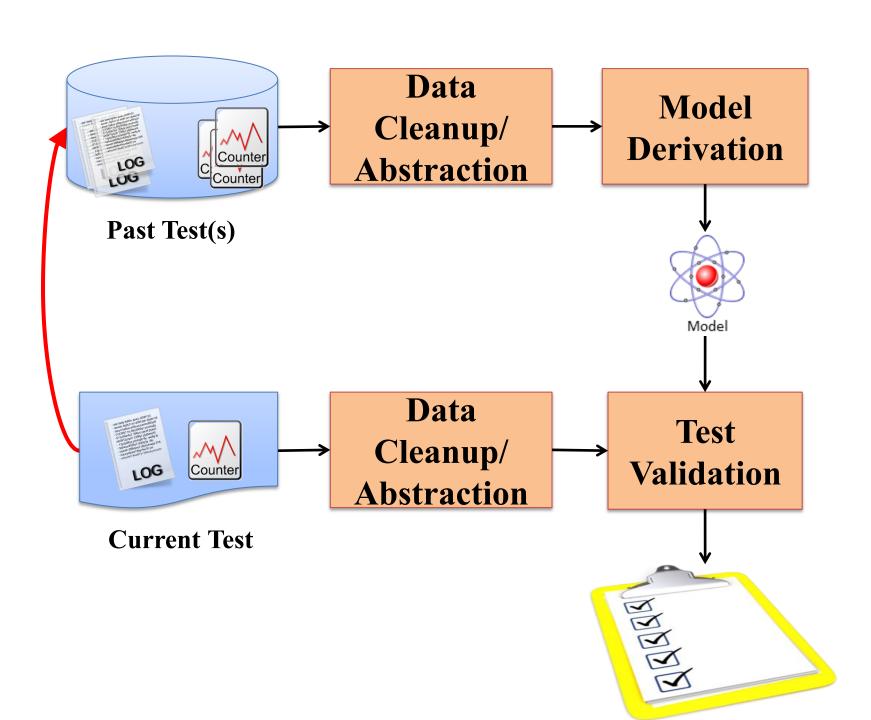


Mimics multiple users repeatedly performing the same tasks

Take hours or even days

Is the system ready for release?





Studied Systems













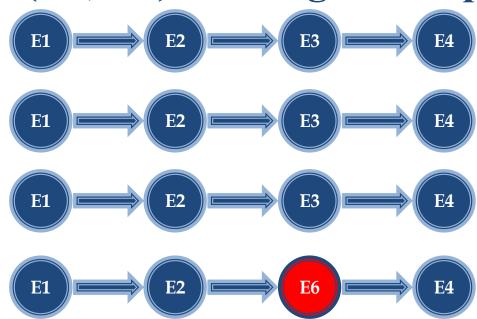
Execution Logs

#	Log Lines
1	time=1, thread=1, session=1, receiving new user registration request
2	time=1, thread=1, session=1, inserting user information to the database
3	time=1, thread=2, session=2, user=Jack, browse catalog=novels
4	time=1, thread=2, session=2, user=Jack, sending search queries to the database
5	time=3, thread=1, session=1, user=Tom, registration completed, sending confirmation email to the user
6	time=3, thread=2, session=2, database connection error: session timeout
7	time=4, thread=1, session=1, fail to send the confirmation email, number of retry = 1
8	time=6, thread=2, session=2, user=Jack, successfully retrieved data from the database
9	time=7, thread=2, system health check
10	time=8, thread=1, session=1, registration email sent successfully to user=Tom
11	time=9, thread=2, session=3, user=Tom, browse catalog=travel
12	time=10, thread=2, session=3, user=Tom, sending search queries to the database
13	time=10, thread=3, session=4, user=Jim, updating user profile
14	time=11, thread=3, session=4, user=Jim, database error: deadlock

Automated Functional Analysis



- (acquire_lock, release_lock)
- (open_inbox, close_inbox)
- If we see (E2, E6) this might be a problem



Dell DVD Store

#	Z-Stat	Kinds	Min	Max	Total		Event		
						Session	IID=19420. Entering purchase for simple quantity queries		
Fre	q		Sa	mple			Details (Sort by Freq)		
87,	528 (9	98%)		2logs. 2logs.			E ₁₃ > SessionID=19420, Entering purchase for simple quantity queries E ₁₄ > SessionID=19420, Initial purchase, update cart		
1,4	36 (<	[1%]				E ₁₃ > SessionID=16242, Entering purchase for simple quantity queries E ₁₃ > SessionID=16242, Entering purchase for simple quantity queries			
358	3 (<1%	(o)		2logs. 2logs.		020	E ₁₃ > SessionID=13496 Entering nurchase for simple quantity queries re commit		
3 (4)	34.73	2	317	16,273	ğ	99.99	9% reduction in viewed log lines		
E22	20.65	2	1	3,857					

Automated Performance Analysis

- Each load test has periods of peak/light load
 - Fluctuations of response time







- Same workload applied across different tests
 - Similar response time distributions

Performance Analysis Report - Visual Comparison

Event Sequences	Deviation	Throughput
<u>E1→E2→E3→E4</u>	0.8	5,000 5,002
<u>E1→E3→E5→E7→E4</u>	<u>0.5</u>	4,500 4,498
E1 .F2 .FE .F7 .F9 .F4	0.2	4,000

Throughput	Deviation	Event Sequences
	Response Time Distribution Previous Test Current Test Throughput Deviation Sequences Throughput Sequences	
	Throughput Deviation Sequences	
5,000	5,000 5,002 <u>0.8</u> <u>1 0.8</u> <u>1 E1→E2→E3→E4</u>	
5,002	4,500 4,498 0.5 E1→E3→E5→E7→E4 1,000 Evolution	<u>E1→E2→</u> <u>E3→E4</u>
	4,000 4,010 Previous Test Current Test Current Test	
	0 1000 2000 3000	
	Start time (in sec)	
4,500 4,498	<u>0.5</u>	<u>E1→E3→E5→</u> <u>E7→E4</u>
4,000 4,010	<u>0.3</u>	$\begin{array}{c} E1 \rightarrow E3 \rightarrow E5 \rightarrow \\ E7 \rightarrow E2 \rightarrow E4 \end{array}$

	Throughput	Deviation	Event Sequences
	5,000 5,002	0.8	<u>E1→E2→E3→E4</u>
L	4,500 4,498	<u>0.5</u>	<u>E1→E3→E5→E7→E4</u>
	4 000		

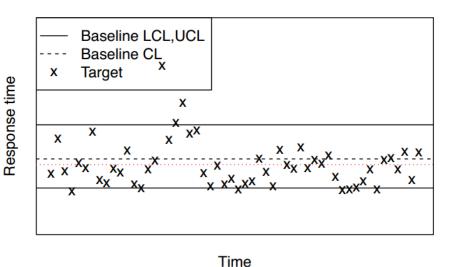
Performance Analysis Report - Step-wise Performance Diagnosis

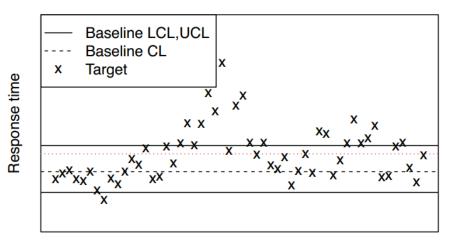
Throughput	Deviation	Event Sequences				
		<u>E1→E2→E3→E4</u>				
		Average	Deviation	Event Sequences Stepwise		
				E1: sessionId=2, time=1, user=Tom, login successfully		
5,000 5,002	0.8	1.0 1.1	0.1	E2: sessionId=2, time=1, user=Tom, browse catalog, catalog=bird		
		0.5 1.8	0.9	E3: sessionId=2, time=2, user=Tom, add item, item=parrot, quantity=1		
		0.5 0.5	0.1	E4: sessionId=2, time=3, user=Tom, user checkout successfully		
4,500						
4,498						
4,000		discovered a MySQL performance bug				
4,010	We					

Performance Counters

	Α	В	С	D	Е
1	Time	Disk Reads/sec	Disk Writes/sec	Page Faults/sec	Memory
2	2/29/08 16:58	0.049986394	0.000723659	0.003876542	3534848
3	2/29/08 17:01	0	0	0	3534848
4	2/29/08 17:04	0.060612225	0.027551011	0.016530607	3534848
5	2/29/08 17:07	0	0	0	3534848
6	2/29/08 17:10	0	0	0	3534848
7	2/29/08 17:13	0.060733302	0.027606046	0.016563628	3534848
8	2/29/08 17:16	0	0	0	3534848
9	2/29/08 17:19	0.060727442	0.027603383	0.01656203	3534848
10	2/29/08 17:22	0	0	0	3534848
11	2/29/08 17:25	0	0	0	3534848
12	2/29/08 17:28	0	0	0	3534848
13	2/29/08 17:31	0	0	0	3534848
14	2/29/08 17:34	0.121368621	0.055167555	0.038617289	3534848
15	2/29/08 17:37	0	0	0	3534848
16	2/29/08 17:40	0	0	0	3534848
17	2/29/08 17:43	0	0	0	3534848
18	2/29/08 17:46	0	0	0	3534848
19	2/29/08 17:49	0	0	0	3534848
20	2/29/08 17:52	0	0	0	3534848
21	2/29/08 17:55	0.121392912	0.055178596	0.033107158	3534848
22	2/29/08 17:58	0.060592703	0.027542138	0.02203371	3534848

Deriving Performance Ranges Using Control Charts





Normal Operations

Suspicious Test

Time

- Derive control charts from the past good tests
- Flag new tests as anomalous if there are many violations in the control charts

Deriving Frequent Itemset from Past Test(s)

Time	DB read/sec	Throughput	Request Queue Size
10:00	Medium	Medium	Low
10:03	Medium	Medium	Low
10:06	Low	Medium	Medium
10:09	Medium	Medium	Low
10:12	Medium	Medium	Low
10:15	Medium	Medium	Low

Deriving Frequent Itemset from Past Test(s)

Time	DB read/sec	Throughput	Request Queue Size	
10:00	Medium	Medium	Low	
10:03	Medium	Medium	Low	
10:06	Low	Medium	Medium	
10:09	Medium	Medium	Low	
10:12	Medium	Medium	Low	
10:15	Medium	Medium	Low	

DB read/sec
Medium

Throughput **Medium**

Request Queue size <u>Low</u>

Deriving Frequent Itemset from Past Test(s)

DB read/sec
Medium

Throughput **Medium**

Request Queue size <u>Low</u>

DB read/sec
Medium



Throughput **Medium**



Request Queue size **Low**

DB read/sec Medium



Request Queue Size <u>Low</u>



Throughput **Medium**

Request Queue Size <u>Low</u>



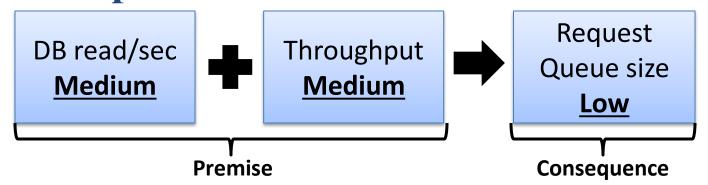
Throughput **Medium**



DB read/sec
Medium

Automated Derivation of Performance Rules

• We have derived from the past test(s) a set of performance rules:



Flags tests where the rule does not hold



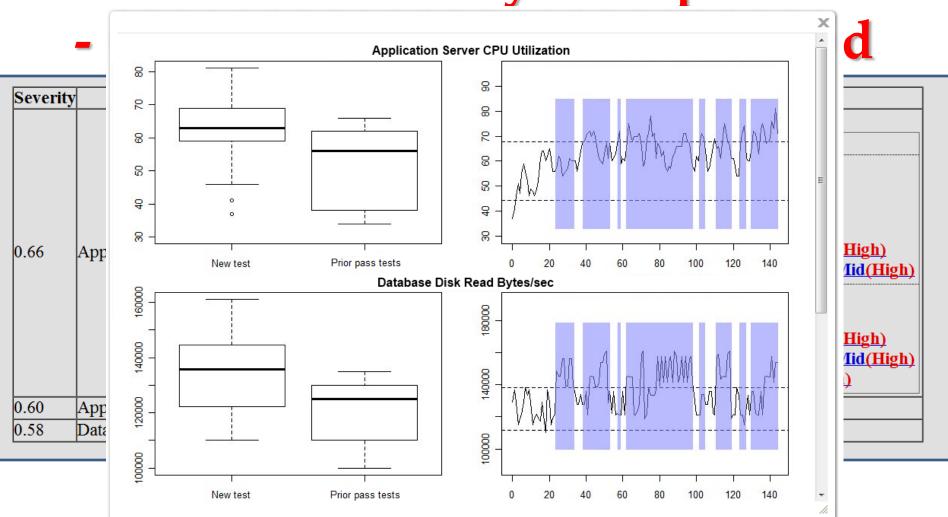
Severity	Performance Regression	Symptoms
0.66	Application Server CPU Utilization	Show Rules
0.60	Application Server Memory Utilization	Show Rules
0.58	Database Disk Read Bytes/sec	Show Rules

Counter Analysis Report

Severity	Performance Regression			Symptoms
		Hide R Graph	<mark>ules</mark> Conf. Change	Expected Correlation
0.66	Application Server CPU Utilization	<u>.</u>	0.79	Database Logical Disk Reads/sec=Mid Database Memory Page Writes/sec=Mid Database CPU Utilization=Mid Database Memory Page Reads/sec=Mid Application Server CPU Utilization=Mid(High) Application Server Memory Utilization=Mid(High)
		≥	0.65	Database Logical Disk Reads/sec=Mid Database Memory Page Reads/sec=Mid Application Server CPU Utilization=Mid(High) Application Server Memory Utilization=Mid(High) Database Disk Read Bytes/sec=Mid(High)
0.60	Application Server Memory Utilization	Show R	<u>Rules</u>	
0.58	Database Disk Read Bytes/sec	Show R	<u>Rules</u>	

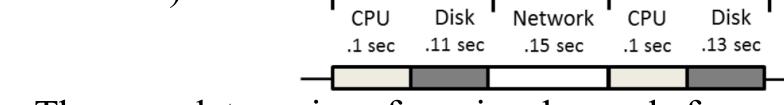
Severity		Symptoms
	Application Server CPU Utilization	Show Rules
0.60	Application Server Memory Utilization	Show Rules
0.58	Database Disk Read Bytes/sec	Show Rules

Counter Analysis Report

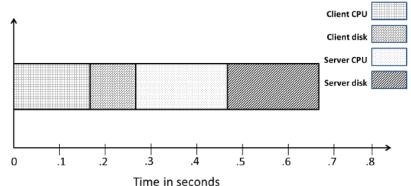


Deriving Transaction Profiles Using Queuing Theory

• Every transaction visits various resources (e.g., CPU and disk)



• The complete series of service demands for one transaction is known as the Transaction Profile (TP)



TP for transaction type A
TP should be stable across different tests.

Otherwise, there is an anomaly

[Ghaith et al., CSMR 2013, STVR 2015]

Correlating Logs and Counters

#	Log Lines
1	time=1, thread=1, session=1, receiving new user registration request
2	time=1, thread=1, session=1, inserting user information to the database
3	time=1, thread=2, session=2, user=Jack, browse catalog=novels
4	time=1, thread=2, session=2, user=Jack, sending search queries to the database
5	time=3, thread=1, session=1, user=Tom, registration completed, sending confirmation email to the user
6	time=3, thread=2, session=2, database connection error: session timeout
7	time=4, thread=1, session=1, fail to send the confirmation email, number of retry = 1
8	time=6, thread=2, session=2, user=Jack, successfully retrieved data from the database
9	time=7, thread=2, system health check
10	time=8, thread=1, session=1, registration email sent successfully to user=Tom
11	time=9, thread=2, session=3, user=Tom, browse catalog=travel
12	time=10, thread=2, session=3, user=Tom, sending search queries to the database
13	time=10, thread=3, session=4, user=Jim, updating user profile
14	time=11, thread=3, session=4, user=Jim, database error: deadlock

	А	В	С	D	Е
1	Time	Disk Reads/sec	Disk Writes/sec	Page Faults/sec	Memory
2	2/29/08 16:58	0.049986394	0.000723659	0.003876542	3534848
3	2/29/08 17:01	0	0	0	3534848
4	2/29/08 17:04	0.060612225	0.027551011	0.016530607	3534848
5	2/29/08 17:07	0	0	0	3534848
6	2/29/08 17:10	0	0	0	3534848
7	2/29/08 17:13	0.060733302	0.027606046	0.016563628	3534848
8	2/29/08 17:16	0	0	0	3534848
9	2/29/08 17:19	0.060727442	0.027603383	0.01656203	3534848
10	2/29/08 17:22	0	0	0	3534848
11	2/29/08 17:25	0	0	0	3534848
12	2/29/08 17:28	0	0	0	3534848
13	2/29/08 17:31	0	0	0	3534848
14	2/29/08 17:34	0.121368621	0.055167555	0.038617289	3534848
15	2/29/08 17:37	0	0	0	3534848
16	2/29/08 17:40	0	0	0	3534848
17	2/29/08 17:43	0	0	0	3534848
18	2/29/08 17:46	0	0	0	3534848
19	2/29/08 17:49	0	0	0	3534848
20	2/29/08 17:52	0	0	0	3534848
21	2/29/08 17:55	0.121392912	0.055178596	0.033107158	3534848
22	2/29/08 17:58	0.060592703	0.027542138	0.02203371	3534848

Combining Logs and Counters

00:01, USER starts a conversation with USER

00:08, 15MB

00:00, 5MB

00:01, USER says MSG to USER

00:16, 15MB

00:02, USER says MSG to USER

00:24, 5MB

00:11, USER says MSG to USER

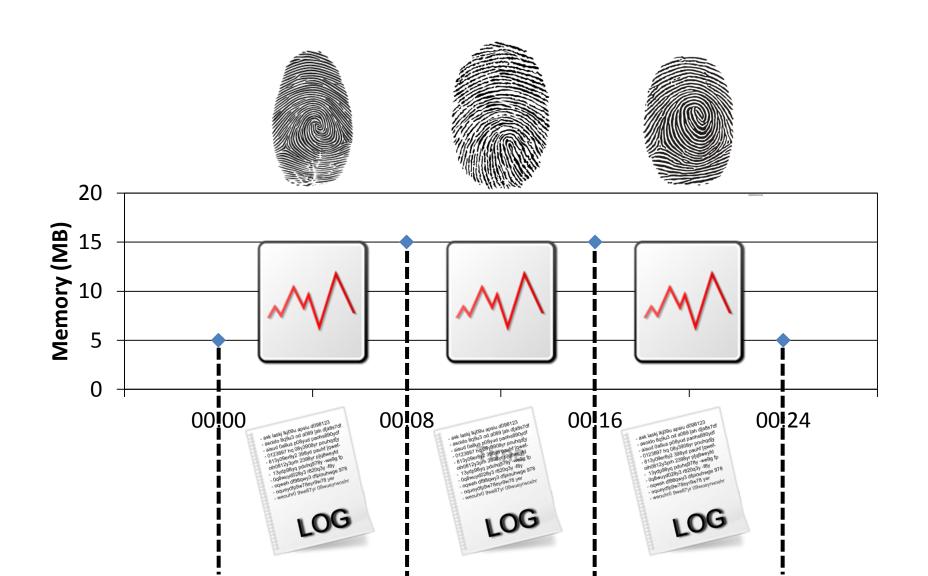
00:12, USER says MSG to USER

00:18, USER ends a conversation with USER





Diagnosing Memory-Related Problems



Creating A Research Community



LT: International Workshop on Large-Scale Testing





<u>Large-Scale Testing</u> includes all different objectives and strategies of testing large-scale software systems using *load*. Examples of large-scale testing include live upgrade testing, load testing, high availability testing, operational profile testing, performance testing, reliability testing, stability testing and stress testing.

Workshop: http://lt2016.eecs.yorku.ca/

Contact: zmjiang@cse.yorku.ca





