Advance in Software Engineering Research

# Service Analytics: Concept and Applications

楼建光
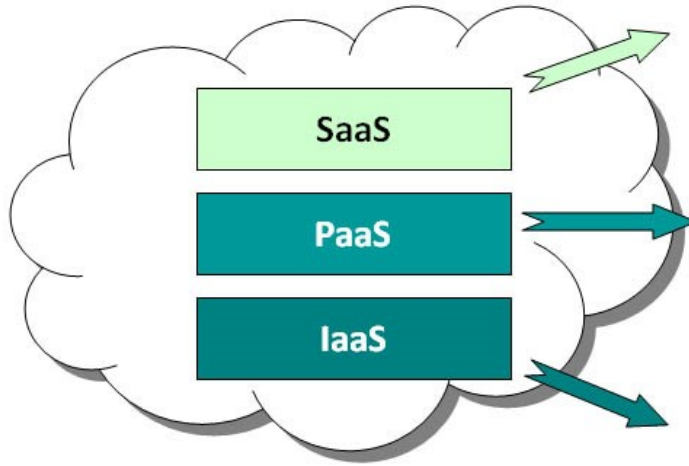Software Analytics Group, Microsoft Research
Dec 10, 2014

Microsoft
**Research**

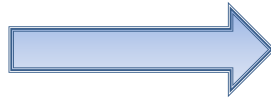*Microsoft*

# Cloud Era

| Who Uses It | What Services are available | Why use it? |
|---|---|---|
| Business Users | EMail, Office Automation, CRM, Website Testing, Wiki, Blog, Virtual Desktop ... | To complete business tasks |
| Developers and Deployers | Service and application test, development, integration and deployment | Create or deploy applications and services for users |
| System Managers | Virtual machines, operating systems, message queues, networks, storage, CPU, memory, backup services | Create platforms for service and application test, development, integration and deployment |

SaaS

PaaS

IaaS

# Software is changing...



On-premise
License
Small Scale

Online Services
Subscription
Large Scale

Microsoft
Research @ 20 Years

# How software is built & operated is changing

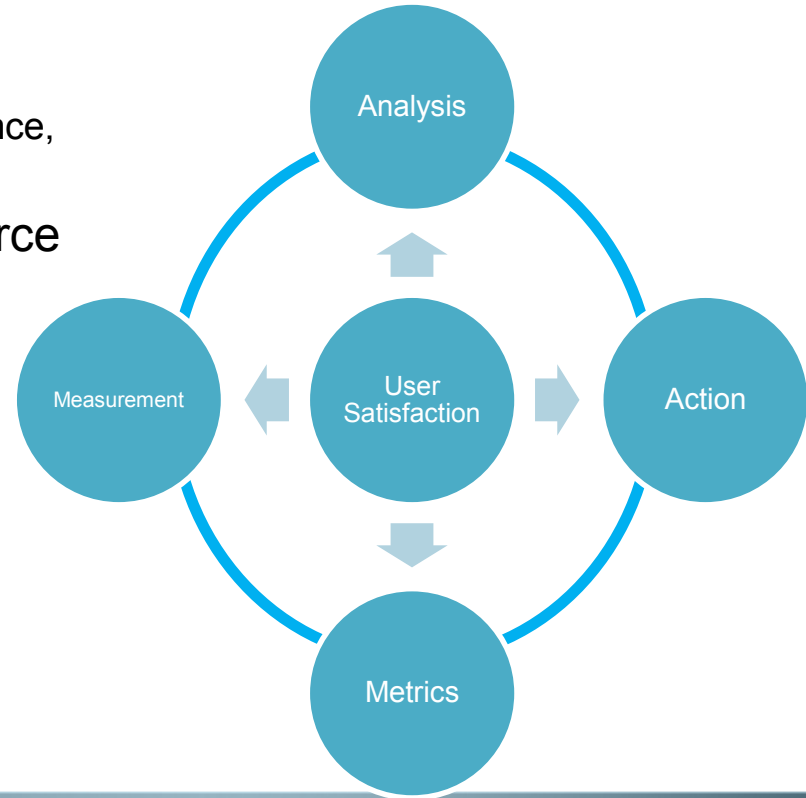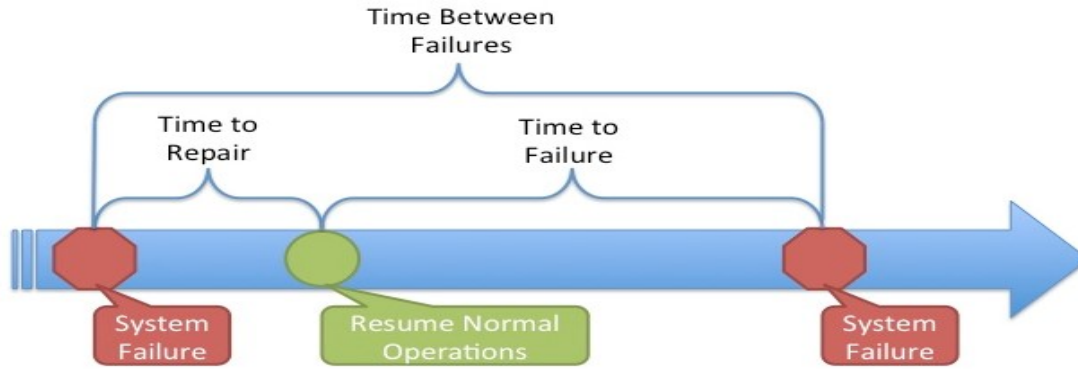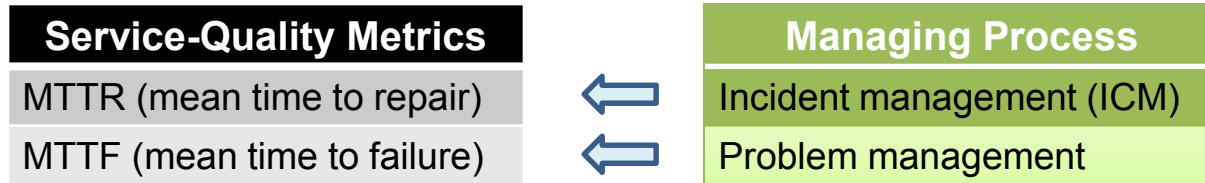| | |
|---|---|
| Code centric | User centric |
| In-lab testing | Debugging in the large |
| Experience & gut-feeling | Data-driven decision making |
| Centralized development | Distributed development |
| Long product cycle | Continuous release |
| … | … |

# User-Centric Service

- Aspects of user satisfaction
  - Usability, reliability, availability, performance, security, privacy, power consumption, …
- User satisfaction as a key driving force for success
  - Prioritization guideline
  - Optimization target
  - Design goal
- Data-driven user satisfaction
  - Metrics
  - Measurement
  - Analysis
  - Action
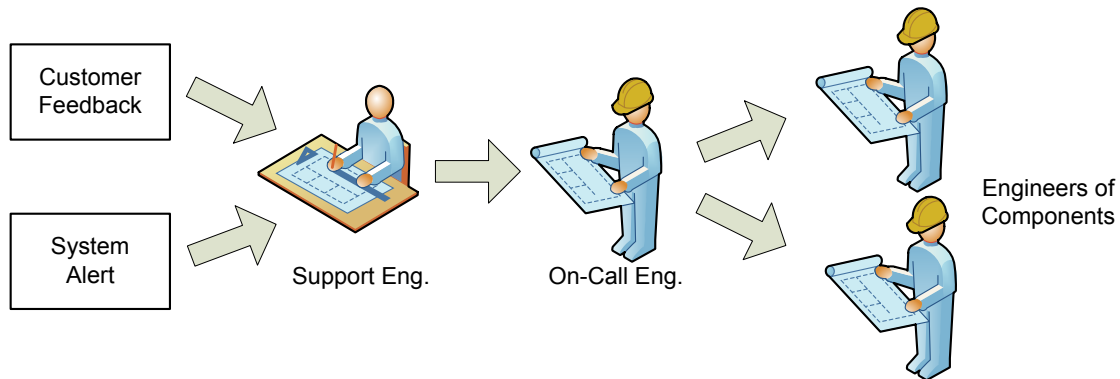
# Reality – 故障不可避免!?

- Google
  - Gmail 在2012/4， 2012/6/07， 2013/8/17, 2014/1/24 等多次发生故障，影响超3300万人
  - 最近一次Google搜索故障，2014/8/27，搜啥都是车祸图片
- Microsoft
  - 2014/11/18 Windows Azure故障
- 微信
  - 2011/12/14, 2013/4/10, 2013/7/22, 2013/8/20, 2014/10/20

# Service Quality Management

| Service-Quality Metrics | | Managing Process |
|---|---|---|
| MTTR (mean time to repair) | ⇐ | Incident management (ICM) |
| MTTF (mean time to failure) | ⇐ | Problem management |



Time Between Failures

Time to Repair — Time to Failure

System Failure — Resume Normal Operations — System Failure

Microsoft Research @ 20 Years

# Incident Management



Customer Feedback → Support Eng. → On-Call Eng. → Engineers of Components

System Alert → Support Eng.

# Incident Management: An Example



缺点：
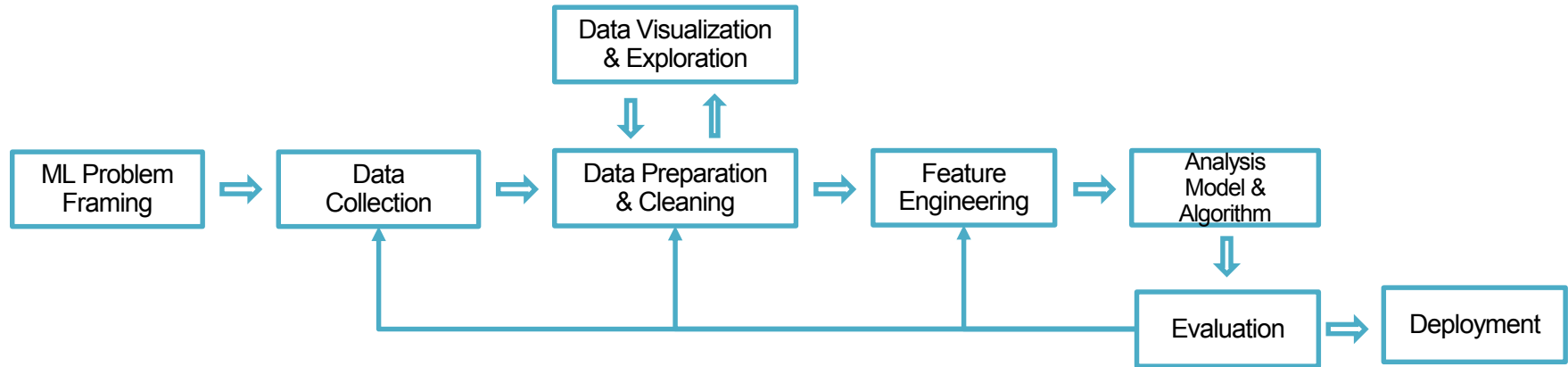1. 没有自动化
2. 故障只等着用户汇报

Microsoft Research @ 20 Years

# What is the Key?

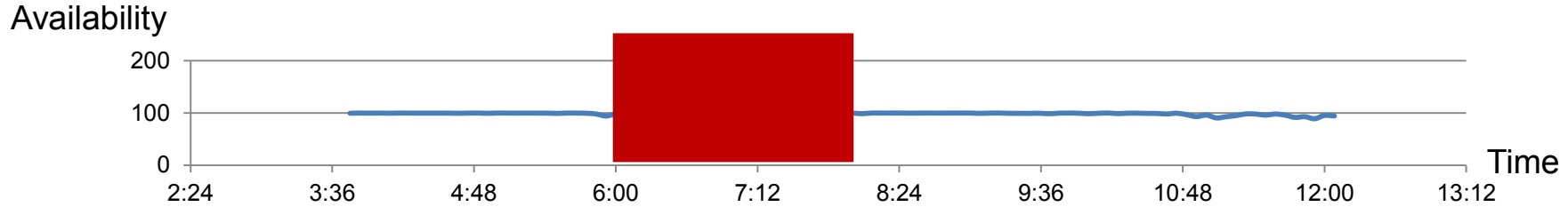Service engineering processes are moving to data-driven

# Formulation: Service Analytics

Service analytics is to enable service *practitioners* to perform *data exploration and analysis* in order to quickly conduct service management tasks.
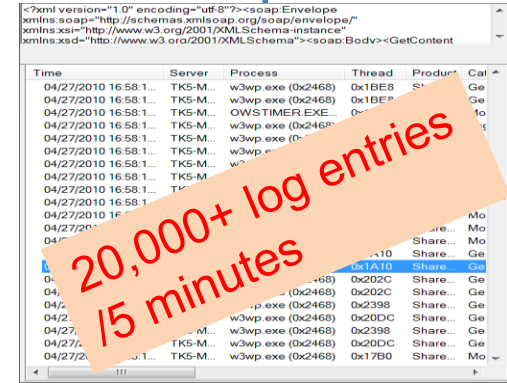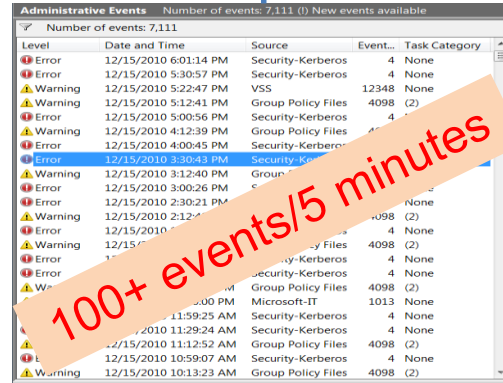
# Logs Generated by An Online Service

Availability



500+ metrics/5 minutes

100+ events/5 minutes

20,000+ log entries /5 minutes

System Resource Measurements

System Events

Transaction Processing event logs

Microsoft Research @ 20 Years

# Important Scenarios

**Problem Detection**

Detect potential issues based on system logs, events, counters, usage data, and customer support records

**Problem Localization & Diagnosis**

Identify the problem site for a service live site issue, or provide information to help pinpoint the potential causes

**Problem Categorization & Prioritization**

Categorize issues and failures to help understand the trend and prioritize management tasks

# Example 1. mining invariants for service problem detection

# Background

- Logs are the major source for telemetry and diagnosis
- Manually inspecting logs is not feasible
  - Large scale of system
  - High complexity of system
- Traditional rule/keyword based log analysis tools:
  - Heavily depend on the knowledge of operators
  - Difficult to keep rules updated when components are frequently revised or upgraded

# Linear Program Invariant

- A predicate always holds the same value under different normal executions.
  - For example:



$$count(A) = count(B) = count(E)$$

$$count(B) = count(C) + count(D)$$

# Invariant and Execution Path

$$count(A) = count(B) = count(E)$$

$$count(B) = count(C) + count(D)$$



Sequential Execution

Execution Branch

Linear invariants reflect the properties of execution path.

# Invariant Violation and Anomaly(1)

- A violation of invariant often indicates a system problem.



$$count(Enter) \neq count(Leave)$$

**Problem
on Critical Section Operations**

# Invariant Violation and Anomaly(2)

- Violated invariants often give diagnosis cues.

$count(A) > count(B)$

$count(B) > count(C) + count(D)$



Sequential Execution

Execution Branch

# Formulation of Invariant

- A linear invariant can be presented as a linear equation:

$$a_0 + a_1 x_1 + a_2 x_2 + \cdots + a_m x_m = 0$$

where $x_i$ is the message count of message i.

- Given a set of logs, we have

$$\boldsymbol{X}\theta = \begin{bmatrix} 1 & x_{11} & x_{12} & \cdots & x_{1m} \\ 1 & x_{21} & x_{22} & \ddots & x_{2m} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & x_{n1} & x_{n2} & \cdots & x_{nm} \end{bmatrix} \theta = 0$$

where

$$\theta = [a_0, a_1, a_2, \cdots, a_m]^T$$

# What Is A Meaningful Invariant?
## -- Sparse Non-zero Coefficients

$$c(B) = c(C) + c(D)$$

$$c(A) = c(B)$$

are more meaningful than

$$c(A) + 3c(B) - 2c(E) - 2c(C) - 2c(D) = 0$$

Any vector in the Null Space of **X** is an invariant;
Only sparse invariants are interested.

# What Is A Meaningful Invariant?
# --  Integer Coefficients

Elementary work flow structures can be interpreted by integer invariants.



Sequential            Branch                        Join

Integer invariants are easy to be understood by human operators.

# Problem Statement

- Due to noise pollution, mining invariants is to find integer sparse solutions of regression.

$$X\theta = \begin{bmatrix} 1 & x_{11} & x_{12} & ... & x_{1m} \\ 1 & x_{21} & x_{22} & \ddots & x_{2m} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & x_{n1} & x_{n2} & \cdots & x_{nm} \end{bmatrix} \theta = 0 \implies arg\ min\|X\theta\|_0$$
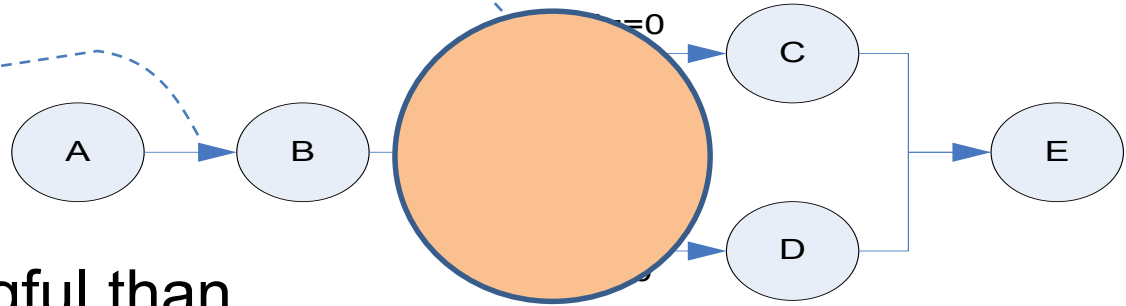
  - Challenges:
    - A typical integer sparse regulation problem (NP-Hard)
    - Traditional method is to relax 0-norm to 1-norm. However, it cannot guarantee to find all invariants.

# Learning Invariant Overview



Parsing → Message Count Vector → Invariant → Anomalies

**Parsing**

```
type=1, PV(1)=172.23.67.0:4635, PV(2)=00001
type=1, PV(1)=172.23.67.1:4635, PV(2)=00002
type=1, PV(1)=172.23.67.2:4635, PV(2)=00003
type=2, PV(1)=00001, PV(2)=0, PV(3)=57717
type=3, PV(1)=00001, PV(2)=0, PV(3)=70795
type=4, PV(1)=\tmp\dfs\name\current\edits,
PV(2)=1049092, PV(3)=2057, PV(4)=0
……

type=1, PV(1)=172.23.67.1:4635, PV(2)=00002
type=1, PV(1)=172.23.67.2:4635, PV(2)=00003
type=2, PV(1)=00001, PV(2)=0, PV(3)=57717
type=3, PV(1)=00001, PV(2)=0, PV(3)=70795
type=4, PV(1)=\tmp\dfs\name\current\edits,
PV(2)=1049092, PV(3)=2057, PV(4)=0
……
```

**Message Count Vector**

```
00001:  [ 1,1,1,2,5,3]
00002:  [1,3,3,3,6,3]
00003:  [1,2,2,2,4,2]
00004:  [ 1,1,1,2,5,3]
00006:  [1,3,3,3,6,3]
00007:  [1,2,2,2,4,2]
00007:  [1,2,2,2,4,2]
00007:  [1,2,2,2,4,2]
--------------------------
… …
```

**Invariant**

$$\theta_1 = [-1,1,0,0,0,0,0],$$
$$\theta_2 = [0,0,1,-1,0,0,0],$$
$$\theta_3 = [0,0,0,0,-1,1,-1]$$
……

**Anomalies**

```
00091:  [ 1,1,0,0,0,0]
violates the
invariant θ₂
00732: [1,3,3,3,4,2]
violates the
invariant θ₃
… …
--------------------------
… …
```

Four Steps:
Auto log parsing, Message Grouping and Counting, Search Invariants, and Anomaly Detection

# Example 2. Healing Online Service Systems via Mining Historical Issue Repositories

# Motivation



Issue Detection → Diagnosis → Healing Action

Issue Repository

Incident Management Process

When a new issue occurred, how to leverage past diagnosis efforts, to identify proper healing action for the new issue?

Microsoft
Research @ 20 Years

# A Simple Example of An Issue

- Symptoms
  - Describing the particular sign and phenomena of the issue

- Solution
  - Recording diagnostic steps and resolution

| ISSUE | Symptoms |
|---|---|
| | • Title:         Browse Homepage failed |
| | • Time:         2012/06/25 13:04:33 |
| | • Datacenter:      XXX |
| | • Type:         Availability |
| | • Traces:  Transaction logs |
| | Solution |
| | • Diagnosis:      SQL connection timeout, SQL-001 blue screen |
| | • Healing action:   Reboot SQL-001. |

Simplified example of an issue

# Characteristics of Logs
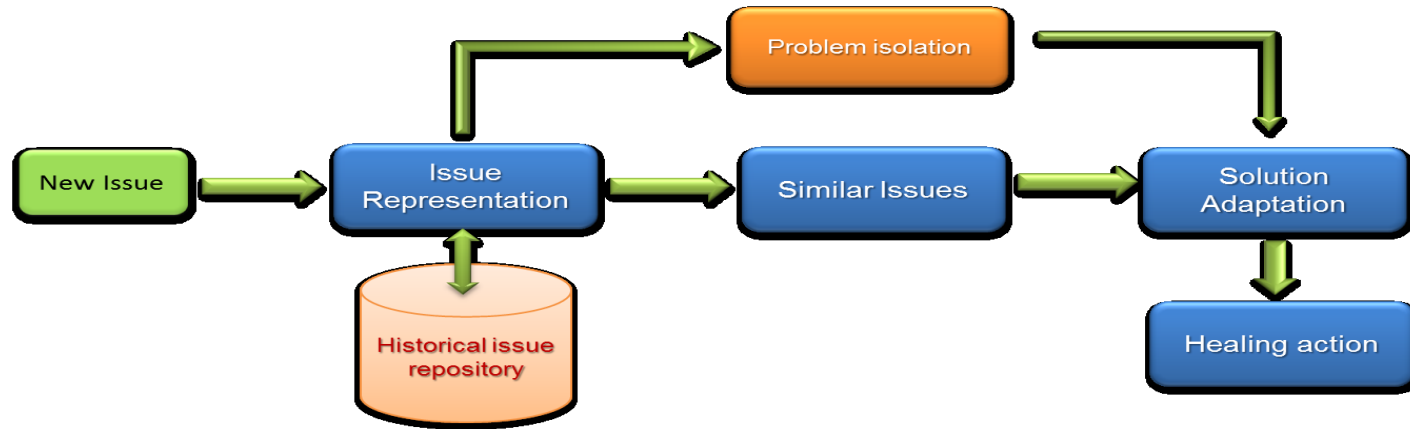
- Highly redundant events

  E.g., 6 events x1 ~ x6 indicate the authentication failure
  - Bias issue representation

- Many irrelevant events to failure

  E.g., event d indicates "SQL usage detection"
  - BUT Relevant to issues, e.g., appearing in only SQL-related issues
  - Downgrading discrimination of issue representation
    e.g., one type of SQL issue needs to reboot SQL; another type needs to patch SQL

| Time | Event | TX ID | Message |
|---|---|---|---|
| 1 | a | A | A entering |
| 1 | b | A | created cookie |
| 2 | c | A | Site = * |
| 3 | d | A | Detected SQL usage |
| 6 | y1 | A | SQL-Exception |
| 6 | Z | A | A leaving |
| 1 | a | B | B entering |
| 3 | x1 | B | B is not sign |
| 3 | x2 | B | building authentication |
| 4 | x3 | B | create sign |
| 4 | x4 | B | create cookie |
| 4 | x5 | B | B does not valid |
| 4 | x6 | B | redirecting B |
| 5 | z | B | B leaving |

**Illustration of transaction logs**
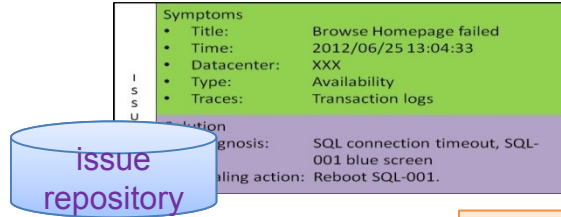
Microsoft Research @ 20 Years

# Our Approach

- Issue-signature extraction
  - Address the challenges posed by logs
- Similarity-metric definition
  - Cosine similarity based on Generalized Vector Space Model (GVSM)
- Healing-action adaptation
  - Structured healing action + fault localization

Microsoft
Research @ 20 Years

# Signature Extraction

## An Issue

Symptoms
- Title: Browse Homepage failed
- Time: 2012/06/25 13:04:33
- Datacenter: XXX
- Type: Availability
- Traces: Transaction logs

Solution
- Diagnosis: SQL connection timeout, SQL-001 blue screen
- Healing action: Reboot SQL-001.

Issue repository

## Issue Signature

| Index | Term (Event Set) | DMI |
|-------|------------------|------|
| 1 | X G T | 0.34 |
| 2 | N S O Y | 0.21 |
| 3 | B C | 0.07 |

Parsing log messages

## Log Sequences

| Time | Event | TX ID | Message |
|------|-------|-------|---------|
| 1 | a | A | A entering |
| 1 | b | A | created cookie |
| 2 | c | A | Site = * |
| 3 | d | A | Detected SQL usage |
| 6 | y1 | A | SQL-Exception |
| 6 | Z | A | A leaving |
| 1 | a | B | B entering |

| Time | Event | TX I |
|------|-------|------|
| 1 | a | A |
| 1 | b | A |
| 2 | c | A |
| 3 | d | A |
| 6 | y1 | A |
| 6 | z | A |

| Time | Event |
|------|-------|
| 1 | a |
| 1 | b |
| 2 | c |
| 3 | d |

## Concept Lattice

{1,2,
{W,X,

{1,2,3,4}
{W,X,G,O,Y,S}

{1,3,5}
W,X,G,O,Y,A}

{W,X,

{W,X,G,

{W,X,G,O,

Formal concept analysis
- Reduce redundancy
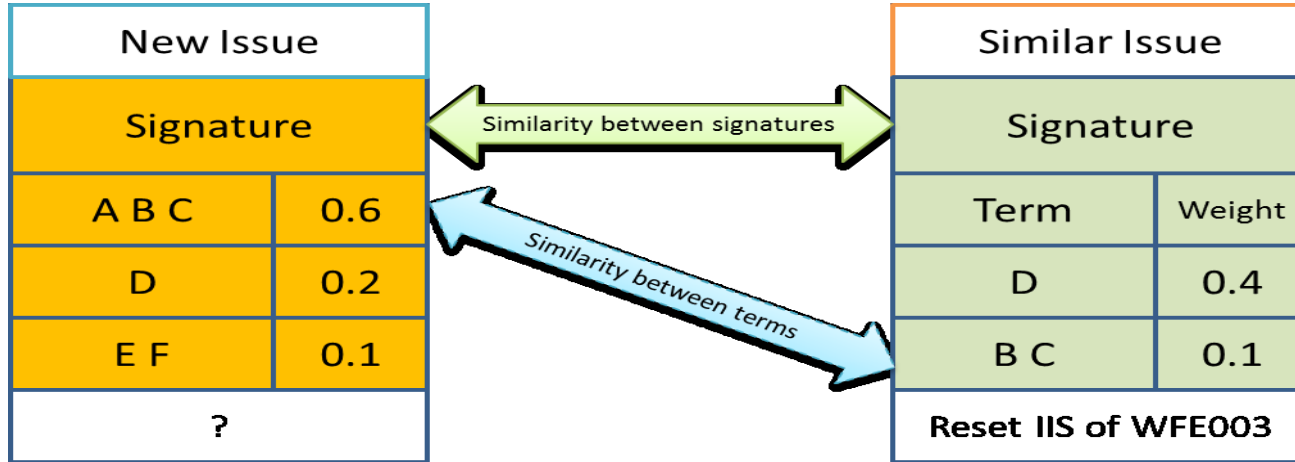- Extract trunk/branch relationship of execution paths

Contrast analysis
- Reduce weak-discrimination
- Measure correlation with Delta Mutual Information (DMI)

ASE 2012, Essen

# Issue Comparison

| New Issue | |
|---|---|
| **Signature** | |
| A B C | 0.6 |
| D | 0.2 |
| E F | 0.1 |
| **?** | |

*Similarity between signatures*

*Similarity between terms*

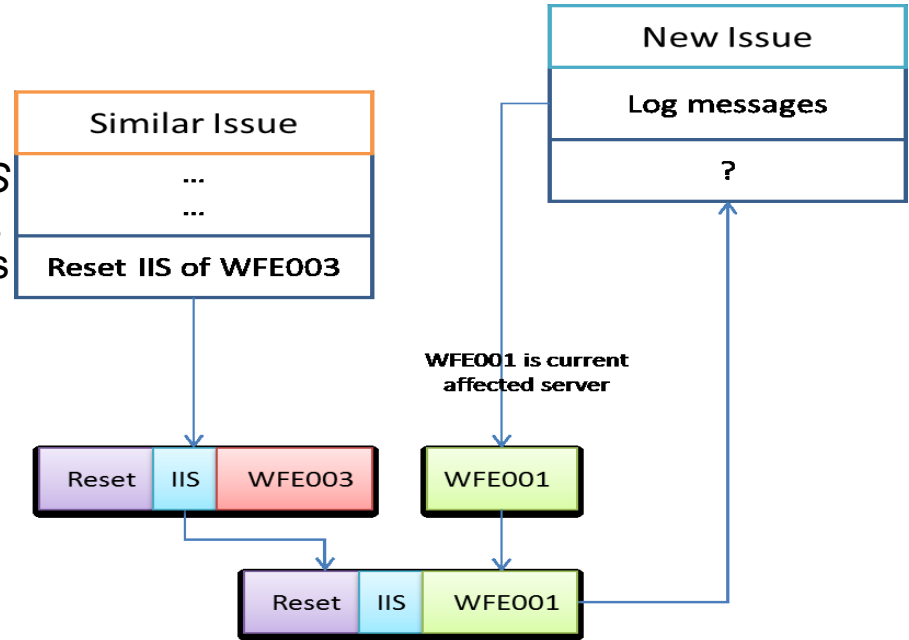| Similar Issue | |
|---|---|
| **Signature** | |
| Term | Weight |
| D | 0.4 |
| B C | 0.1 |
| **Reset IIS of WFE003** | |

Similarity definition: (Generalized Vector Space Model)
- Similarity between terms
- Similarity between signatures
  - Combine term similarity
  - Encode importance of term using DMI as weights

# Healing-Action Adaptation

- Triple structure
  - $< verb, target, locaton >$
- *Verb* & *Target*
  - E.g., *"recycle + AppPool", "Reset + IIS*
  - Extracted from retrieved similar issues by analyzing their solution descriptions
- Location
  - Specific machine/server name, e.g., SQL23524-001
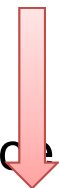  - Obtained by applying fault-localization techniques

Microsoft **Research** @ 20 Years

# Evaluation

---Internal production service: ServiceX

- 332 issues collected in time period: 11/01/2011~02/18/2012

- 146 issues with documented healing actions and recorded logs

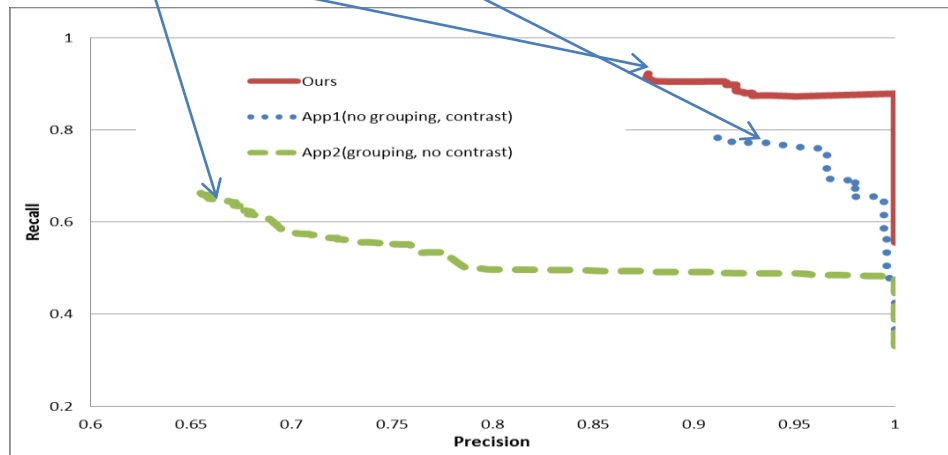69 issues on service upgrade

**77** issues on service interruption or degradation
  – used in evaluation

- Effect of our techniques on overall effectiveness
  – Approach1: Ignore highly-correlated phenomenon (Mutual information + VSM)
  – Approach2: Ignore weakly-discriminative phenomenon (FCA + TF-IDF + VSM)
  – Our approach: FCA + contrast analysis



**Overall ROC curves**

Research @ 20 Years

# Summary

- Mission of Service Analytics
  - Utilize data-driven approach to help create highly performing, user friendly, and efficiently built & operated online services

- Service Analytics is naturally tied with state of engineering practice of service

- Empowering future software practitioners with data analytics mindset & skills

# Advertisement

- We are recruiting!
  - Software analytics researchers (Full-time employee, visiting researchers)
  - Software analytics interns

# Q & A

http://research.microsoft.com/groups/sa/