

Natural language processing meets software testing

Michael Ernst

UW CSE

Joint work with Juan Caballero, Alberto Goffi, Alessandra
Gorla, Mauro Pezze, Irfan Ul Haq, and Sai Zhang

March 15, 2016

What is software?

What is software?

- A sequence of instructions that perform some task

What is software?

An engineered object amenable to formal analysis

- A sequence of instructions that perform some task

Formalizations

		$h \in \text{Heap}$	$= \text{Addr} \rightarrow \text{Obj}$
		$\iota \in \text{Addr}$	$= \text{Set of Addresses} \cup \{\text{null}_a\}$
		$o \in \text{Obj}$	$= {}^r\text{Type}, \text{Fields}$
		${}^r\text{T} \in {}^r\text{Type}$	$= \text{OwnerAddr ClassId}\langle{}^r\text{Type}\rangle$
$\text{P} \in \text{Program}$	$::=$	$\overline{\text{Class}}, \text{ClassId}, \text{Expr}$	
$\text{Cls} \in \text{Class}$	$::=$	$\text{class ClassId}\langle\text{TVarId}\rangle$ $\text{extends ClassId}\langle{}^s\text{Type}\rangle$ $\{ \overline{\text{FieldId}} {}^s\text{Type}; \text{Met}$	
		$\text{Fs} \in \text{Fields}$	$= \text{FieldId} \rightarrow \text{Addr}$
		$\iota \in \text{OwnerAddr}$	$= \text{Addr} \cup \{\text{any}_a\}$
		${}^r\Gamma \in {}^r\text{Env}$	$= \overline{\text{TVarId}} {}^r\text{Type}; \overline{\text{ParId}} \text{Addr}$
${}^s\text{T} \in {}^s\text{Type}$	$::=$	${}^s\text{NType} \mid \text{TVarId}$	
${}^s\text{N} \in {}^s\text{NType}$	$::=$	$\text{OM ClassId}\langle{}^s\text{Type}\rangle$	
$u \in \text{OM}$	$::=$	$h, {}^r\Gamma, e_0 \rightsquigarrow h_0, \iota_0$	
$\text{mt} \in \text{Meth}$	$::=$	$\iota_0 \neq \text{null}_a$	
		$h_0, {}^r\Gamma, e_2 \rightsquigarrow h_2, \iota$	
		$h' = h_2[\iota_0.f := \iota]$	
$w \in \text{Purity}$	$::=$	OS-Upd	
$e \in \text{Expr}$	$::=$	$h, {}^r\Gamma, e_0.f = e_2 \rightsquigarrow h'$	
		$\text{Expr.MethId}\langle{}^s\text{Type}\rangle(\text{Expr}) \mid$ $\text{new } {}^s\text{Type} \mid ({}^s\text{Type}) \text{Expr}$	
${}^s\Gamma \in {}^s\text{Env}$	$::=$	$\overline{\text{TVarId}} {}^s\text{NType}; \overline{\text{ParId}} {}^s\text{Type}$	
		GT-Read	
		$\frac{\Gamma \vdash e_0 : N_0 \quad N_0 = _}{\Gamma \vdash e_0.f : N_0 \triangleright f\text{Type}(C_0, f)}$	
		GT-Upd	
		$\frac{\Gamma \vdash e_0 : N_0 \quad N_0 = u_0 \ C_0 \langle _ \rangle \quad T_1 = f\text{Type}(C_0, f) \quad \Gamma \vdash e_2 : N_0 \triangleright T_1 \quad u_0 \neq \text{any} \quad rp(u_0, T_1)}{\Gamma \vdash e_0.f = e_2 : N_0 \triangleright T_1}$	
		OS-Read	
		$\frac{h, {}^r\Gamma, e_0 \rightsquigarrow h', \iota_0 \quad \iota_0 \neq \text{null}_a \quad \iota = h'(\iota_0) \downarrow_2 (f)}{h, {}^r\Gamma, e_0.f \rightsquigarrow h', \iota}$	
		DYN	
		$\frac{\left. \begin{array}{l} h \vdash {}^r\Gamma : {}^s\Gamma \\ h \vdash \iota_1 : \text{dyn}({}^s\text{N}, h, \iota_1) \\ h \vdash \iota_2 : \text{dyn}({}^s\text{T}, \iota_1, h(\iota_1) \downarrow_1) \\ {}^s\text{N} = u_N \ C_N \langle _ \rangle \\ u_N = \text{this}_u \Rightarrow {}^r\Gamma(\text{this}) \\ \text{free}({}^s\text{T}) \subseteq \text{dom}(C_N) \end{array} \right\} \Rightarrow h \vdash \iota_2 : \text{dyn}({}^s\text{N} \triangleright {}^s\text{T}, h, {}^r\Gamma) \quad \begin{array}{l} {}^r\text{T} = \iota' \ _ \langle _ \rangle \quad \iota \vdash {}^r\text{T} \ _ \langle _ \rangle : \iota' \ C \langle \overline{{}^r\text{T}} \rangle \quad \iota \vdash {}^r\text{T} \ _ \langle _ \rangle : \iota' \ C \langle \overline{{}^r\text{T}}_a \rangle \Rightarrow \iota \vdash \overline{{}^r\text{T}} \ _ \langle _ \rangle : \overline{{}^r\text{T}}_a \\ \text{dom}(C) = \overline{X} \quad \text{free}({}^s\text{T}) \subseteq \overline{X} \circ \overline{X'} \end{array}}{\text{dyn}({}^s\text{T}, \iota, {}^r\text{T}, (\overline{X'} \ \overline{{}^r\text{T}}'; _)) = {}^s\text{T}[\iota'/\text{this}, \iota'/\text{peer}, \iota/\text{rep}, \text{any}_a/\text{any}_u, \overline{{}^r\text{T}}/\overline{X}, \overline{{}^r\text{T}}'/\overline{X'}]}$	

What is software?

- A sequence of instructions that perform some task

What is software?

- A sequence of instructions that perform some task
- Test cases
- Version control history
- Issue tracker
- Documentation
- ...

How should it be analyzed?

Analysis of a natural object

- Machine learning over executions
- Version control history analysis
- Bug prediction
- Upgrade safety
- Prioritizing warnings
- Program repair

Natural language in programs

This talk:

1. Variable names:
find undesired variable interactions
2. Error messages and user manuals:
find inadequate diagnostic messages
3. Procedure documentation:
generate test oracles

Undesired variable interactions

```
int totalPrice;  
int itemPrice;  
int shippingDistance;  
totalPrice = itemPrice + shippingDistance;
```

Undesired variable interactions

```
int totalPrice;  
int itemPrice;  
int shippingDistance;  
totalPrice = itemPrice + shippingDistance;
```

- The compiler issues no warning
- A human can tell the abstract types are different

Idea:

- Cluster variables based on usage in program operations
- Cluster variables based on words in variable names

Differences indicate bugs or poor variable names

Clustering based on operations

Abstract type inference [ISSTA 2006]

```
int totalCost(int miles, int price, int tax) {
    int year = 2016;
    if ((miles > 1000) && (year > 2000)) {
        int shippingFee = 10;
        return price + tax + shippingFee;
    } else {
        return price + tax;
    }
}
```

Clustering based on operations

Abstract type inference [ISSTA 2006]

```
int totalCost(int miles, int price, int tax) {  
    int year = 2016;  
    if ((miles > 1000) && (year > 2000)) {  
        int shippingFee = 10;  
        return price + tax + shippingFee;  
    } else {  
        return price + tax;  
    }  
}
```

Clustering based on variable names

Compute variable name similarity

1. Tokenize each variable into dictionary words
 - `in_authskey15` \Rightarrow {"in", "authentications", "key"}
 - Expand abbreviations, best-effort tokenization
2. Compute word similarity
 - For all $w1 \in \text{var1}$ and $w2 \in \text{var2}$, use WordNet or edit distance
3. Combine word similarity into variable name similarity
 - $\text{maxwordsim}(w1) = \text{maximum wordsim}(w1, w2)$ for $w2 \in \text{var2}$
 - $\text{varsim}(\text{var1}) = \text{average maxwordsim}(w1)$ for $w1 \in \text{var1}$

Results

- Found an undesired variable interaction in grep

```
if (depth < delta[tree->label])
    delta[tree->label] = depth;
```
- Loses top 3 bytes of depth
- Not exploitable because of guards elsewhere in program, but not obvious here

Inadequate diagnostic messages

Scenario: user supplies a wrong configuration option
`--port_num=100.0`

Problem: software issues an unhelpful error message

- “... unexpected system failure ...”
- “... unable to establish connection ...”
- Better: “`--port_num` should be an integer”

Goal: detect such problems before shipping the code

Challenges for proactive detection of inadequate diagnostic messages

- How to *trigger a configuration error*?
- How to *determine the inadequacy* of a diagnostic message?

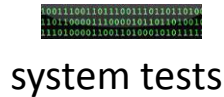
ConfDiagDetector's solutions

- How to *trigger a configuration error*?

- Configuration mutation + run system tests



+



failed tests \approx triggered errors

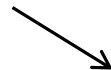
- How to *determine the inadequacy* of a diagnostic message?

- Use a **NLP technique to check its semantic meaning**

Similar semantic meanings?



Diagnostic messages
output by **failed** tests



User manual

Message analysis

- A message is adequate, if it
 - contains the mutated option name or valueOR
 - has a similar semantic meaning with the manual description

Text similarity technique [Mihalcea'06]



A message



Manual description

Has similar semantic meanings, if many words in them have similar meanings

Example:

~~The program goes wrong~~

~~The software fails~~



- *Remove all stop words*
- *For each word in the diagnostic message, tries to find the similar words in the manual*
- *Two sentences are similar, if “many” words are similar between them.*

Results

- Reported 25 missing and 18 inadequate messages in Weka, JMeter, Jetty, Derby
- Validation by 3 programmers:
 - 0% false negative rate
 - 2% false positive rate

Test oracles for exceptional behavior

Exceptional behavior is a significant source of failures, but is under-tested (significantly less coverage)

Goal: create test oracles (= `assert` statements)

Although programmers may not write tests, the programmer does provide other indications: procedure documentation (e.g., Javadoc)

```
/**
 * Checks whether the comparator is now locked
 * against further changes.
 *
 * @throws UnsupportedOperationException if the
 * comparator is locked
 */
protected void checkLocked() {...}
```

Text to code

1. Parse the @throws expression using the Stanford Parser
 - Parse tree, grammatical relations, cross-references
 - Challenges:
 - Often not a well-formed sentence; code snippets as nouns/verbs
 - Referents are implicit, assumes coding knowledge
2. Match each subject to a Java element
 - Pattern matching
 - Semantic similarity
 - Lexical similarity to identifiers, types, documentation
3. Match each predicate to a Java element
4. Create assert statement from expressions and methods

Automatically generated tests

- A test generation tool outputs:
 - Passing tests – useful for regression testing
 - Failing tests – indicates a program bug
- Without a formal specification, tool guesses whether a given behavior is correct
 - False positives: report a failing test that was due to illegal inputs
 - False negatives: fail to report a failing test because it might have been due to illegal inputs
- Results: Reduced false positive test failures in EvoSuite by 1/3 or more

Machine learning + software engineering

- Software is more than source code
- Formal program analysis is useful, but insufficient
- Analyze and generate all software artifacts

A rich space for further exploration