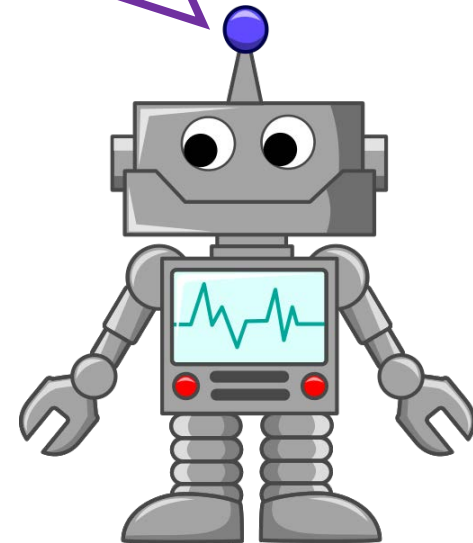


Natural language is a
programming language

Michael D. Ernst

UW CSE

Joint work with Arianna Blasi, Juan Caballero,
Sergio Delgado Castellanos, Alberto Goffi,
Alessandra Gorla, Xi Victoria Lin, Deric Pang,
Mauro Pezzè, Irfan Ul Haq, Kevin Vu,
Chenglong Wang, Luke Zettlemoyer, and Sai
Zhang



Questions about software

- How many of you have used software?
- How many of you have written software?

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

What is software?

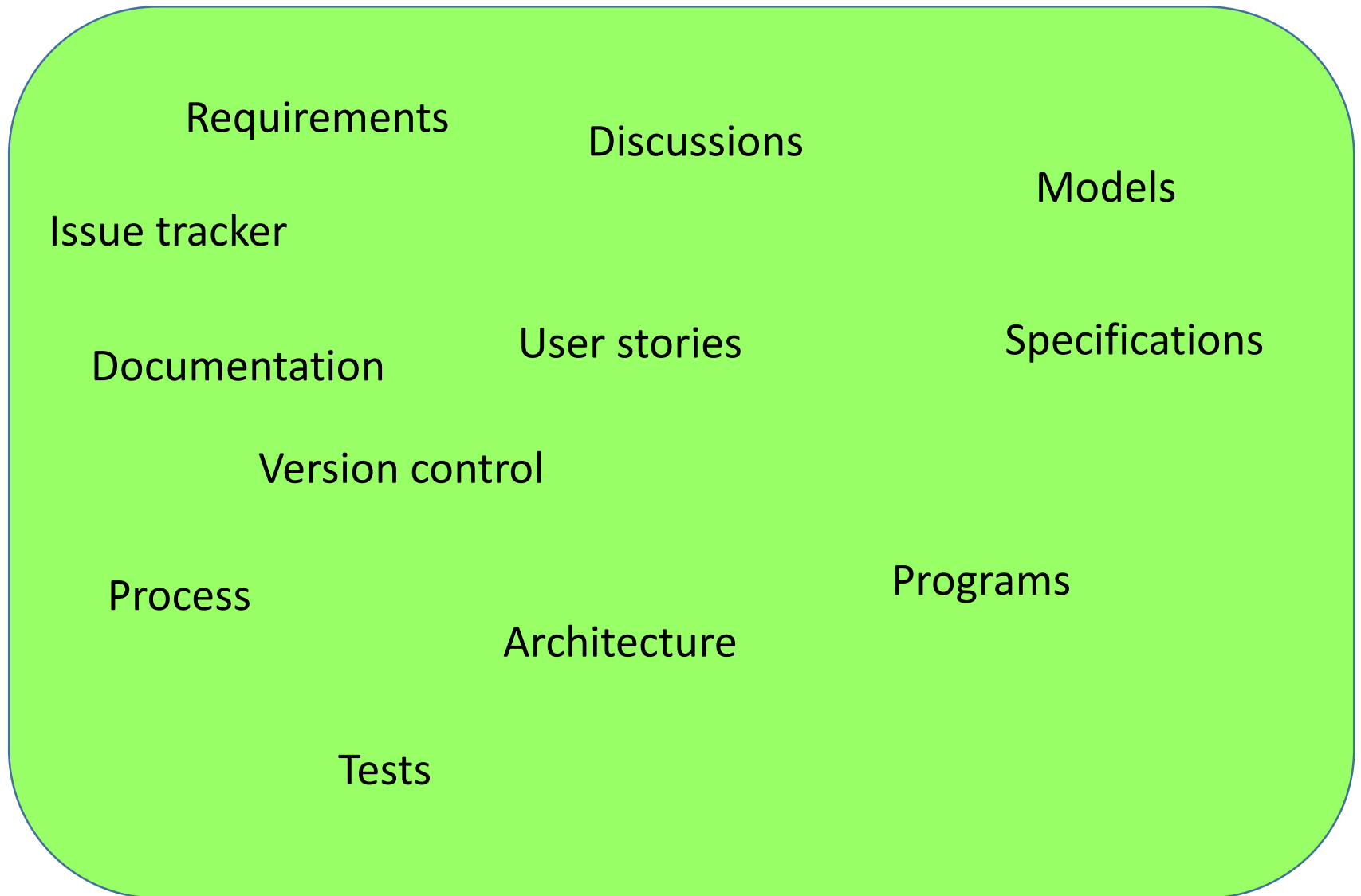
- A sequence of instructions that perform some task

What is software?

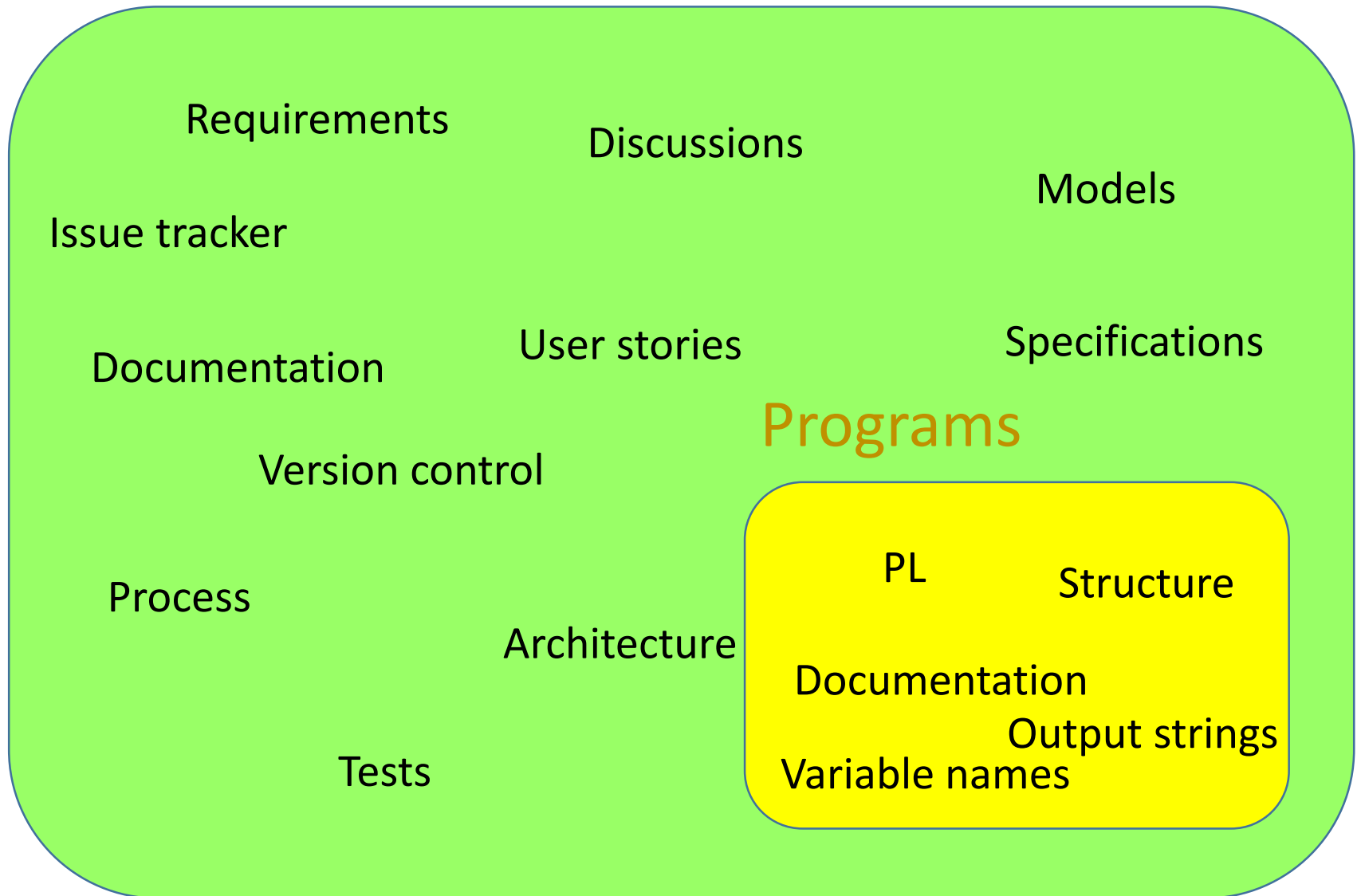
- A sequence of instructions that perform some task
- Test cases
- Version control history
- Issue tracker
- Documentation
- ...

How should it be analyzed?

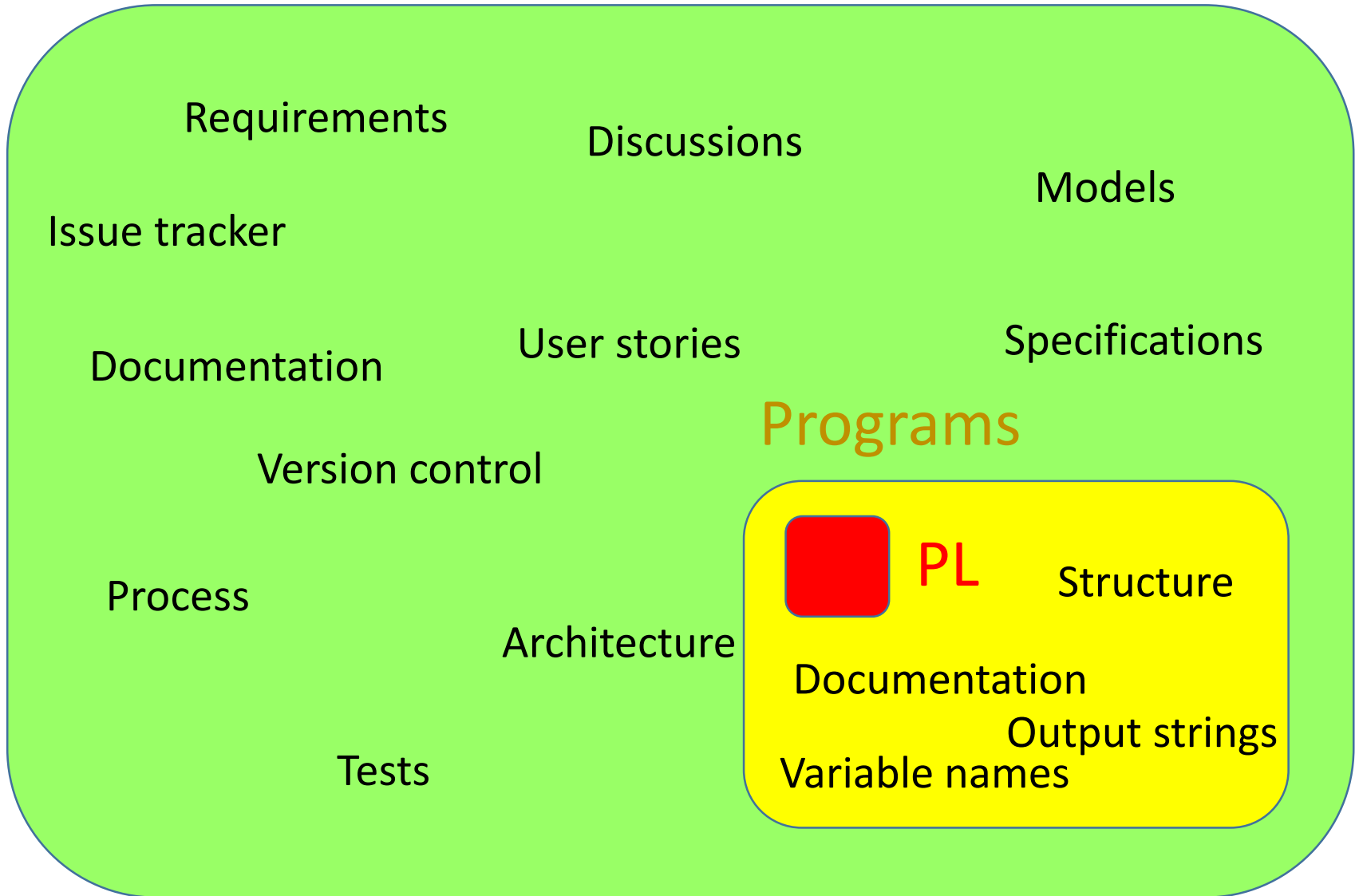
Programming



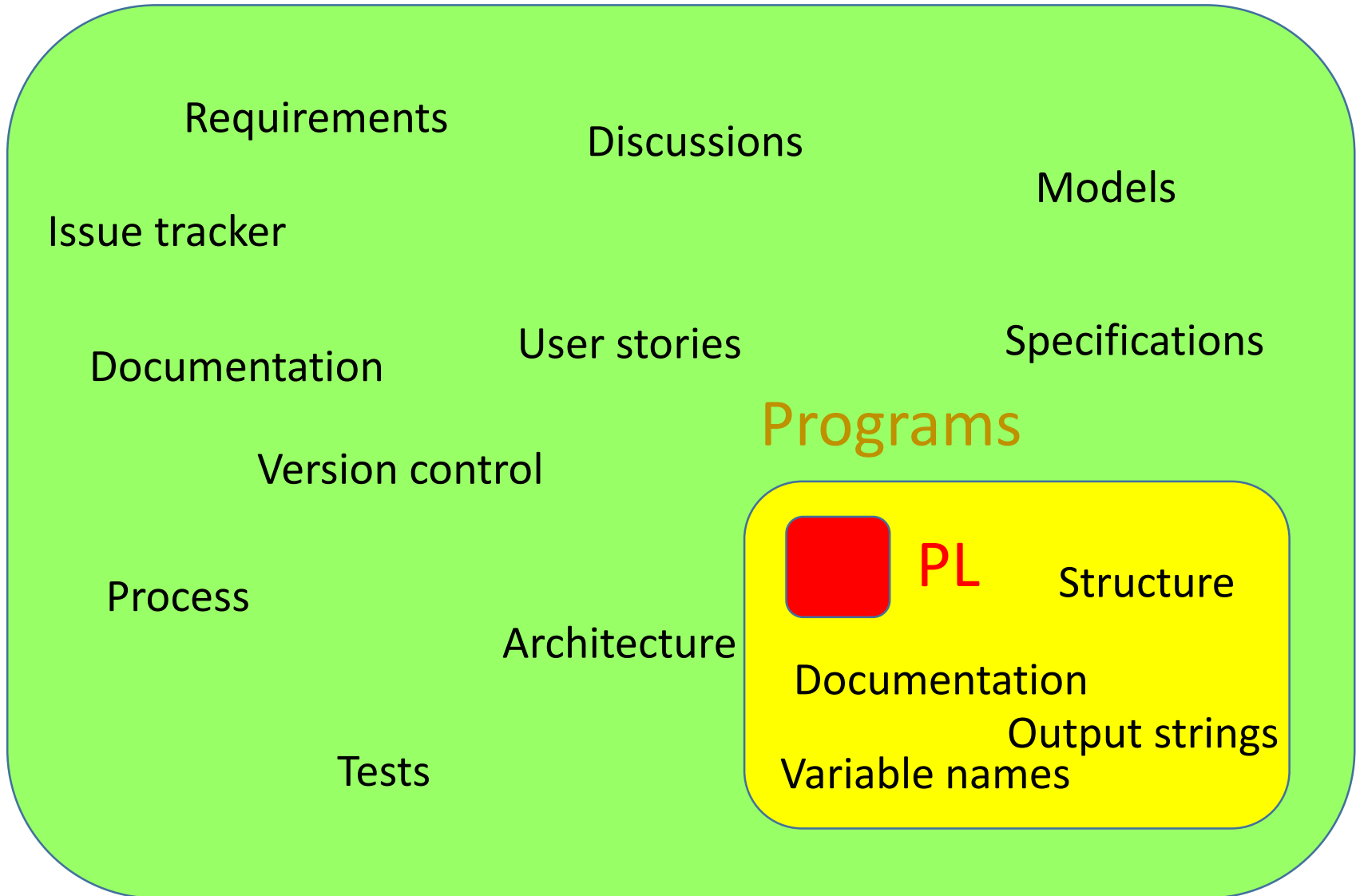
Programming



Programming



Programming



Analysis of a natural object

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



Specifications are needed; Tests are available but ignored

- **Specs are needed.** Many papers start: “Given a program and its specification...”
- **Tests are ignored.** Formal verification process:
 - Write the program
 - Test the program
 - Verify the program, *ignoring* testing artifacts

Observation: Programmers embed semantic info in tests

Goal: translate tests into specifications

Approach: machine learning over executions

Dynamic detection of likely invariants



<https://plse.cs.washington.edu/daikon/>
[ICSE 1999]

- Observe values that the program computes
- Generalize over them via machine learning
- Result: invariants (as in **asserts** or specifications)
 - **$x > \text{abs}(y)$**
 - **$x = 16*y + 4*z + 3$**
 - array **a** contains no duplicates
 - for each node **n**, **$n = n.\text{child}.\text{parent}$**
 - graph **g** is acyclic
- Unsound, incomplete, and **useful**

Programming

Requirements

Discussions

Models

Issue tracker

```
Formalizations h C Heap == Addr -> Obj
               i C Addr == Set of Addresses (null+)
               a C Obj == Type, Fields
               t C Type == OwnerAddr ClassId*Type*
               p C Field == FieldId -> Addr
               d C OwnerAddr == Addr (any)
               s C *VarId *Type; FieldId Addr
               h, T, eo == M, so
               OS-Heap == OS-Heap
               OS-Heap == OS-Heap
               OS-Heap == OS-Heap
```

Documentation

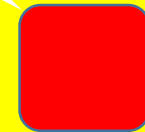
Specifications

Programs

Version control

Process

Architecture



PL

Structure

Documentation

Output strings

Variable names

Tests

Applying NLP to software engineering

Problems

inadequate

diagnostics

incorrect

operations

missing

tests

unimplemented

functionality

Analyze
existing
code

Generate
new
code

NL sources

error

messages

variable

names

code

comments

user

questions

NLP techniques

document

similarity

word

semantics

parse

trees

translation

Applying NLP to software engineering

Problems

NL sources

NLP techniques

inadequate
diagnostics

error
messages

document
similarity

incorrect
operations

variable
names

word
semantics

[ISSTA 2015]

missing
tests

code
comments

parse
trees

unimplemented
functionality

user
questions

translation

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”

Hard for end users to diagnose

Goal: detect such problems *before* shipping the code

- Better message: “`--port_num` should be an integer”

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



configuration

+



system tests



failed tests \approx triggered errors
(We know the root cause.)

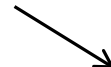
- **How to *determine the inadequacy* of a diagnostic message?**

- Use a NLP technique to check its semantic meaning

Similar semantic meanings?



Diagnostic messages
output by failed tests



User manual
(Assumption: a manual,
webpage, or man page exists.)

When is a message adequate?

- Contains the mutated **option name** or value [Keller'08, Yin'11]

Mutated option:

```
--percentage-split
```

Diagnostic message:

```
"the value of percentage-split should be > 0"
```

- Similar **semantic meaning** as the manual description

Mutated option:

```
--fnum
```

Diagnostic message:

```
"Number of folds must be greater than 1"
```

User manual description of `--fnum`:

```
"Sets number of folds for cross-validation"
```

Classical document similarity: TF-IDF + cosine similarity

1. Convert document into a real-valued vector
 2. Document similarity = vector cosine similarity
- Vector length = dictionary size, values = term frequency (TF)
 - Example: [2_{classical}, 8_{document}, 3_{problem}, 3_{values}, ...]
 - Problem: frequent words swamp important words
 - Solution: values = TF x IDF (inverse document frequency)
 - $IDF = \log(\text{total documents} / \text{documents with the term})$

Problem: does not work well on very short documents

Text similarity technique [Mihalcea'06]



A message



Manual description

The documents have similar semantic meanings if many words in them have similar meanings

Example:

~~The~~ program goes wrong

~~The~~ software fails



1. Remove all stop words.
2. For each word in the diagnostic message, try to find similar words in the manual.
3. 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
 - Tool says message is adequate, humans say it is inadequate
 - 2% false positive rate
 - Tool says message is inadequate, humans say it is adequate
 - Previous best: 16%

Related work

Configuration error diagnosis techniques

- Dynamic tainting [Attariyan'08], static tainting [Rabkin'11], Chronus [Whitaker'04]

Troubleshooting an exhibited error rather than detecting inadequate diagnostic messages

Software diagnosability improvement techniques

- PeerPressure [Wang'04], RangeFixer [Xiong'12], ConfErr [Keller'08] and Spex-INJ [Yin'11], EnCore [Zhang'14]

Requires source code, usage history, or OS-level support

Applying NLP to software engineering

Problems

NL sources

NLP techniques

inadequate
diagnostics

error
messages

document
similarity

incorrect
operations

variable
names

word
semantics

missing
tests

code
comments

parse
trees

unimplemented
functionality

user
questions

translation

[WODA 2015]

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 words in variable names
 - Cluster variables based on usage in program operations
- Differences indicate bugs or poor variable names

Undesired interactions

`distance`

`itemPrice`

`tax_rate`

`miles`

`shippingFee`

`percent_complete`

Undesired interactions

`distance` ↔ `itemPrice`

`itemPrice + distance`

`tax_rate`

`miles`

`shippingFee`

`percent_complete`

Undesired interactions

`int`

`distance ↔ itemPrice`

`miles`

`shippingFee`

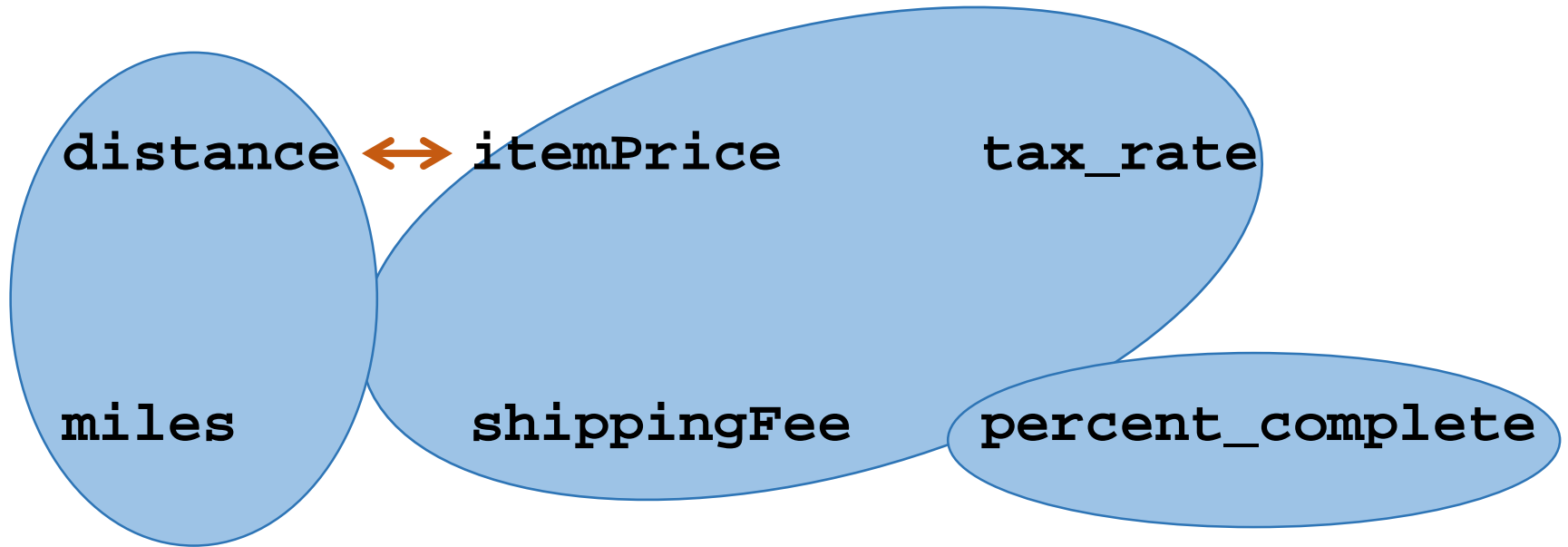
`float`

`tax_rate`

`percent_complete`

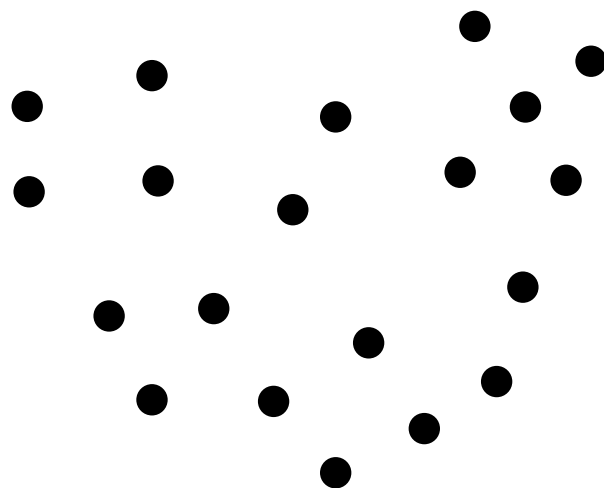
Program types don't help

Undesired interactions



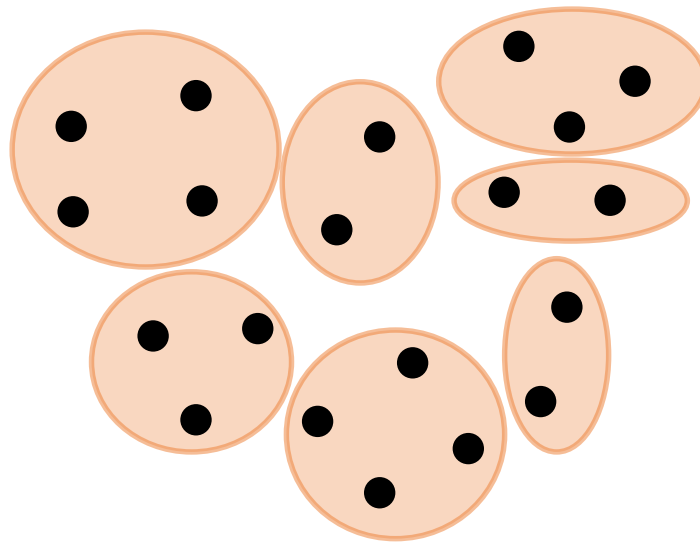
Language indicates the problem

Variables



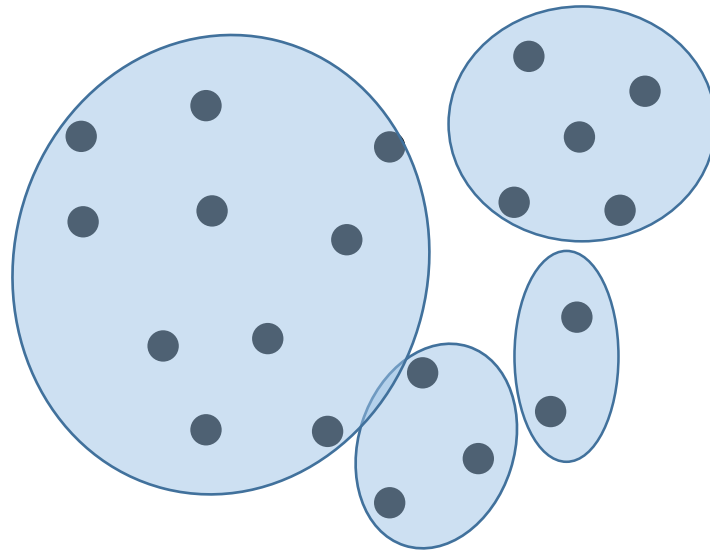
Variable clustering

Cluster based on
interactions:
operations



Variable clustering

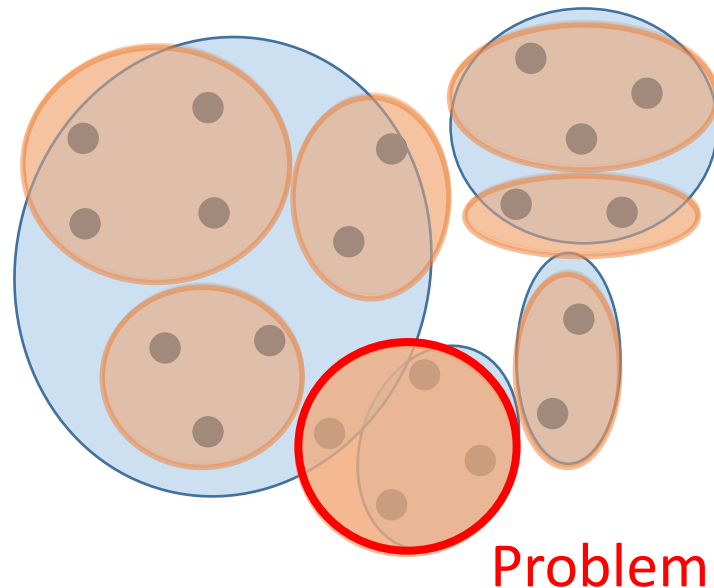
Cluster based on
language:
variable names



Variable clustering

Cluster based on
interactions:
operations

Cluster based on
language:
variable names



Actual algorithm:

1. Cluster based on **operations**
2. Sub-cluster based on **names**
3. Rank an **operation cluster** as suspicious if it contains well-defined **name sub-clusters**

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 for var_1 and var_2

1. **Tokenize** each variable into dictionary words
 - `in_authskey15` \Rightarrow {"in", "authentications", "key"}
 - Expand abbreviations, best-effort tokenization
2. Compute **word similarity**
 - For all $w_1 \in \text{var}_1$ and $w_2 \in \text{var}_2$, use WordNet (or edit distance)
3. Combine word similarity into **variable name similarity**
 - $\text{maxwordsim}(w_1, \text{var}_2) = \max_{w_2 \in \text{var}_2} \text{wordsim}(w_1, w_2)$
 - $\text{varsim}(\text{var}_1, \text{var}_2) = \text{average}_{w_1 \in \text{var}_1} \text{maxwordsim}(w_1, \text{var}_2)$

Results

- Ran on grep and Exim mail server
- Top-ranked mismatch indicates 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

Related work

- Reusing identifier names is error-prone [Lawrie 2007, Deissenboeck 2010, Arnaoudova 2010]
- Identifier naming conventions [Simonyi]
- Units of measure [Ada, F#, etc.]
- Tokenization of variable names [Lawrie 2010, Guerrouj 2012]

Applying NLP to software engineering

Problems

inadequate
diagnostics

incorrect
operations

missing
tests

unimplemented
functionality

NL sources

error
messages

variable
names

code
comments

user
questions

NLP techniques

document
similarity

word
semantics

parse
trees

translation

[ISSTA 2016]

Test oracles (`assert` statements)

A test consists of

- an **input** (for a unit test, a sequence of calls)
- an **oracle** (an `assert` statement)

Programmer-written tests

- often trivial oracles, or too few tests

Automatic generation of tests:

- inputs are easy to generate
- oracles remain an open challenge



Goal: create test oracles
from what programmers already write

Automatic test generation

- Code under test:

```
public class FilterIterator implements Iterator {  
    public FilterIterator(Iterator i, Predicate p) {...}  
    public Object next() {...}  
    ...  
}
```

`/** @throws NullPointerException if either
 * the iterator or predicate are null */`

- Automatically generated test:

```
public void test() {  
    FilterIterator i = new FilterIterator(null, null);  
    i.next();  
}
```

 **Throws NullPointerException!**

Did the tool discover a bug?

It could be:

1. Expected behavior
2. Illegal input
3. Implementation bug

Automatically generated tests

- A test generation tool outputs:
 - Failing tests – indicates a program bug
 - Passing tests – useful for regression testing
- Without a specification, the 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

Programmers write code comments

Javadoc is standard procedure documentation

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

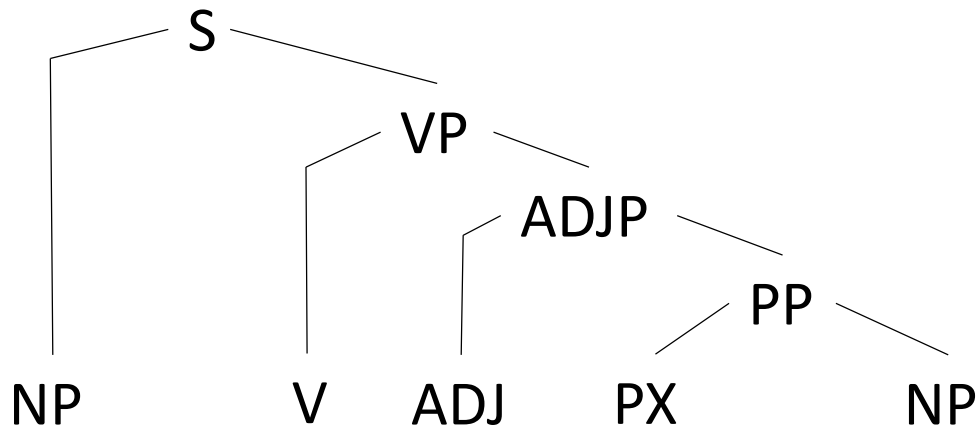
Javadoc comment and assertion

```
class MyClass {  
    ArrayList allFoundSoFar = ...;  
  
    boolean canConvert(Object arg) { ... }  
  
    /** @throws IllegalArgumentException if the  
     *   element is not in the list and is not  
     *   convertible. */  
    void myMethod(Object element) { ... }  
}
```

Condition for exception: `myMethod` should throw iff ...

```
( !allFoundSoFar.contains(element)  
  && !canConvert(element) )
```

Nouns = objects, verbs = operations



The element is greater than the current maximum.



`elt` `compareTo()` > 0 `currentMax`

`elt.compareTo(currentMax) > 0`

Text to code: Toradocu algorithm

1. Parse `@param`, `@return`, and `@throws` expressions 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
 - Lexical similarity to identifiers, types, documentation
3. Match each predicate to a Java element
4. Create assert statement from expressions and methods

Results

Accuracy on 857 Javadoc tags:

- 97% precision
- 72% recall

Can tune parameters to favor either metric

Pre-processing and pattern-matching are important

Discovered specification errors

Improving test generation tools:

- Reduced false positive test failures in EvoSuite by $\geq 1/3$
- Also improved Randoop, but by less

Related work

Heuristics

- JCrasher, Crash'n'Check [Csallner'04, Csallner'05]
- Randoop [Pacheco'07]

Specifications

- ASTOOT [Doong'94]
- Models, contracts, ...

Properties

- Cross-checking oracles [Carzaniga'14]
- Metamorphic testing [Chen'13]
- Symmetric testing [Gotlieb'03]

Natural language documentation

- iComment, aComment, @tComment [Tan'07, Tan'11, Tan'12]

Applying NLP to software engineering

Problems

inadequate
diagnostics

incorrect
operations

missing
tests

unimplemented
functionality

NL sources

error
messages

variable
names

code
comments

user
questions

NLP techniques

document
similarity

word
semantics

parse
trees

translation

Machine translation

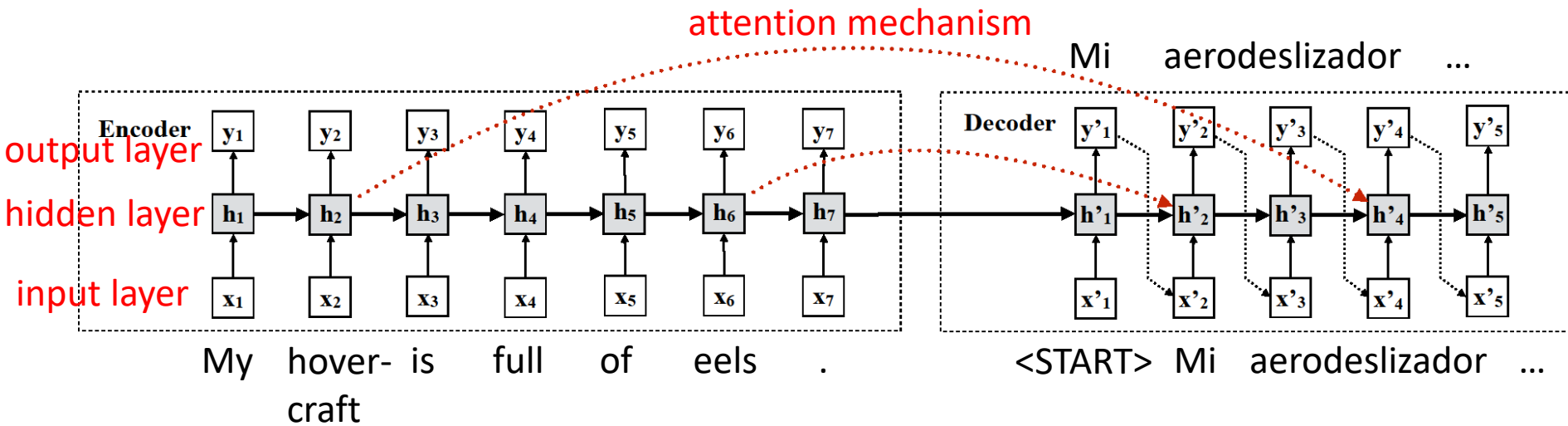
English: “My hovercraft is full of eels.”

Spanish: “Mi aerodeslizador está lleno de anguilas.”

English: “Don’t worry.”

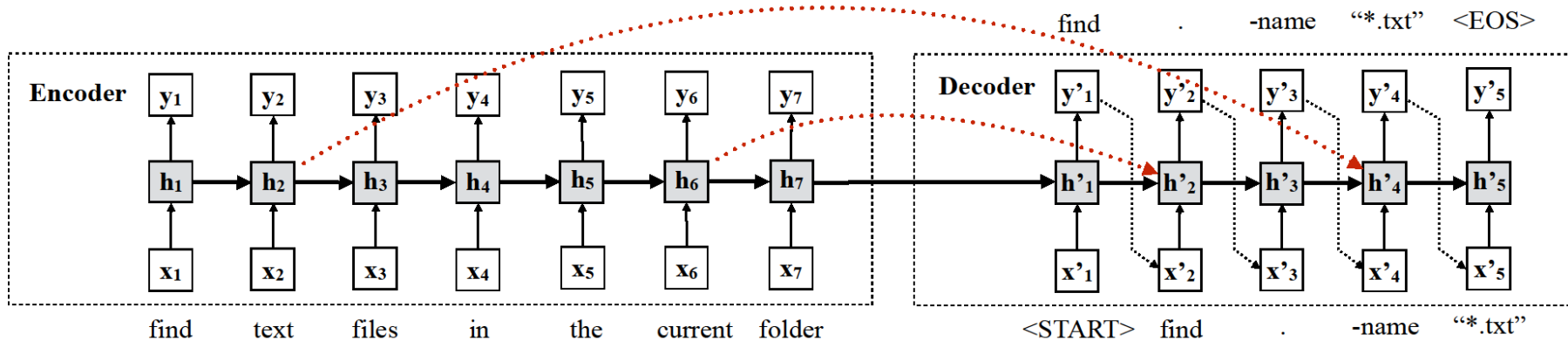
Spanish: “No te preocupes.”

Sequence-to-sequence recurrent neural network translators



Input, hidden, and output functions are inferred from training data using probability maximization.

Tellina: text to commands



- Training data: ~5000 $\langle \text{text, command} \rangle$ pairs
 - Collected manually from webpages, plus cleaning
- 17 file system utilities, > 200 flags, 9 types of constants
 - Compound commands: `()`, `&&`, `||`
 - Nesting: `|`, `$()`, `<()`
 - Strings are opaque; no command interpreters (**awk**, **sed**)
 - No bash compound statements (**for**)

Results

Accuracy for Tellina's first output:

- Structure of command (without constants): 69%
- Full command (with constants): 30%

User experiment:

- Tellina makes users 22% more efficient
 - Even though it rarely gives a perfect command
- Qualitative feedback
 - Most participants wanted to continue using Tellina (5.8/7 Likert scale)
 - Partially-correct answers were helpful, not too hard to correct
 - Output bash commands are sometimes non-syntactic or subtly wrong
 - Needs explanation of meaning of output bash commands

Related work

Neural machine translation

- Sequence-to-sequence learning with neural nets [Sutskever 2014]
- Attention mechanism [Luong 2015]

Semantic parsing

- Translating natural language to a formal representation [Zettlemoyer 2007, Pasupat 2016]

Translating natural language to DSLs

- If-this-then-that recipes [Quirk 2015]
- Regular expressions [Locascio 2016]
- Text editing, flight queries [Desai 2016]

Other software engineering projects

- Analyzing programs before they are written
- Gamification (crowd-sourcing) of verification
- Evaluating and improving fault localization
- Pluggable type-checking for error prevention

- ... many more: systems, synthesis, verification, etc.



UW is hiring! Faculty, postdocs, grad students

Applying NLP to software engineering

Problems

inadequate
diagnostics

incorrect
operations

missing
tests

unimplemented
functionality

NL sources

error
messages

variable
names

code
comments

user
questions

NLP techniques

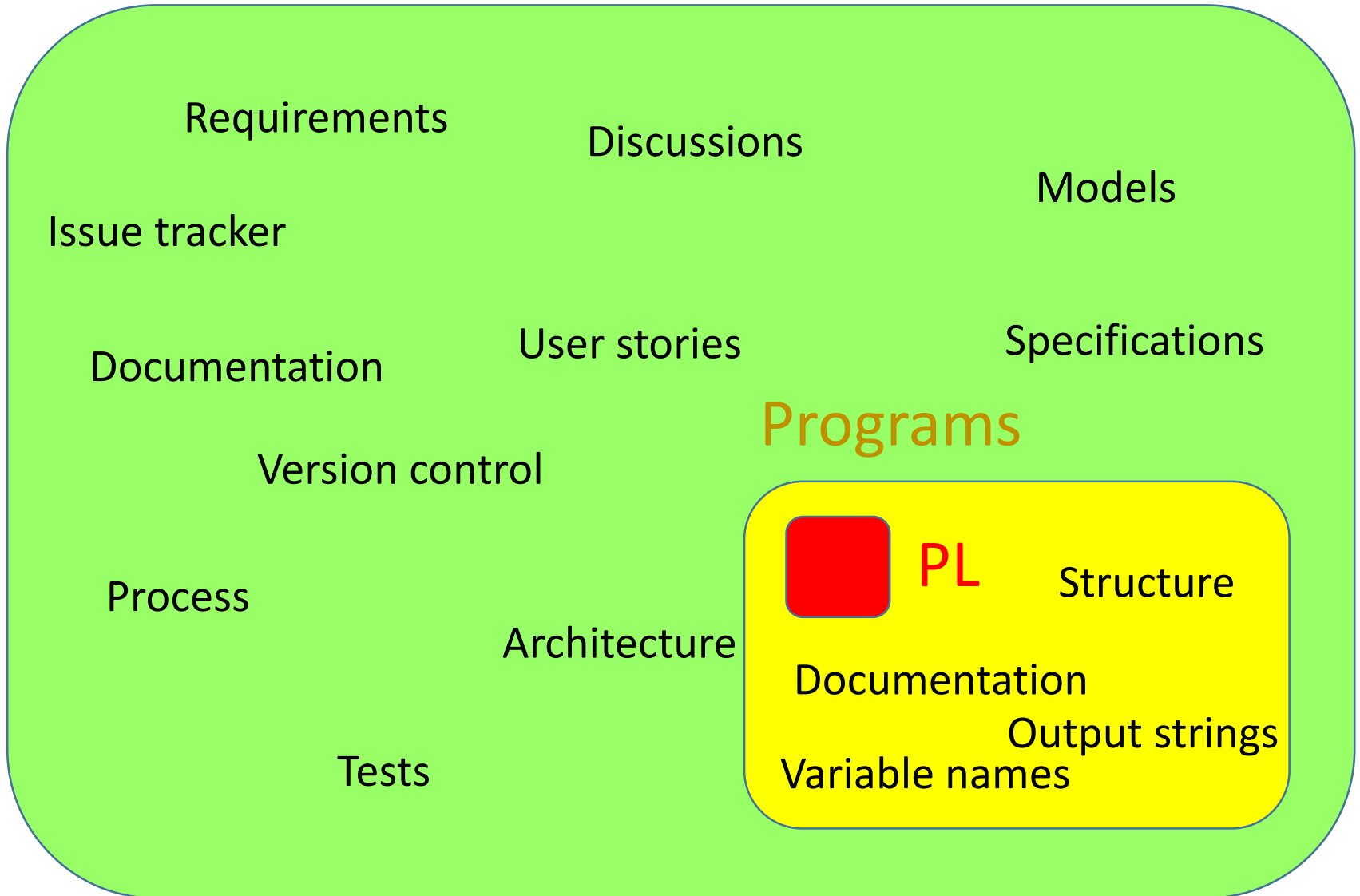
document
similarity

word
semantics

parse
trees

translation

Programming



Analyzing text

- iComment [Tan 2007]: pattern matching for null
- N-gram models: code completion [Hindle 2011], predict variable names, whitespace [Allemanis 2014]
- Mining variable names [Pollock et al.]
- Code → comments [Sridhara 2010]
- DARPA Big Mechanism (read cancer papers)
- JSNice [Raychev 2015]: learn rules for identifiers and types

Analyzing other artifacts by machine learning over the program

- Tests (dynamic invariant detection)
- Mining software repositories
- Defect prediction
- Code completion
- Clone detection
- ... many, many 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

Programming

