#### ESEC/FSE 2017 PADERBORN, GERMANY, SEPTEMBER 04 - 08









Today

```
var geom2d = function() {
  var t = numeric.sum;
  function r(n, r) {
    this.x = n;
    this.y = r;
  }
  u(r, {
    P: function e(n) {
      return t([ this.x * n.x,
                   this.y * n.y ]);
    }
  });
  function u(n, r) {
    for (var t in r) n[t] = r[t];
    return n;
 }
  return {
   V: r
 };
}();
```

Today

```
var geom2d = function() {
  var t = numeric.sum;
  function r(n, r) {
    this.x = n;
    this.y = r;
  }
  u(r, {
    P: function e(n) {
      return t([ this.x * n.x,
                   this.y * n.y ]);
    }
  });
  function u(n, r) {
    for (var t in r) n[t] = r[t];
    return n;
  }
  return {
    V: r
  };
}();
```

Today

#### Data-driven method + tool

```
var geom2d = function() {
                                      var geom2d = function()
                                        var sum = numeric.sum;
  var t = numeric.sum;
  function r(n, r)
                                        function Vector2d(x, y) {
                                          this.x = x;
    this.x = n;
    this.y = r;
                                          this.y = y;
  }
  u(r, {
                                        mix(Vector2d, ·
                                          P: function dotProduct(vector) {
    P: function e(n) {
      return t( this.x * n.x,
                                            return sum([ this.x * vector.x,
                   this.y * n.y ]);
                                                              this.y * vector.y ]);
  });
                                        });
                                        function mix(dest, src) {
  function u(n, r) {
    for (var t in r) n[t] = r[t];
                                          for (var k in src) dest[k] = src[k];
                                          return dest;
    return n;
  }
  return {
                                        return {
   V: r
                                          V: Vector2d
  };
                                        };
                                      }();
}();
```

Why?

• Programs are (also) written to be read

"Instead of imagining that our main task is to instruct a computer what to do, let us concentrate rather on **explaining to human beings what we want a computer to do**." [Don Knuth]



Why?

- Programs are (also) written to be read
- Well-chosen variable names are critical to source code readability, reusability, maintainability
- Example tasks:
  - reverse engineering binaries
  - reverse engineering obfuscated JavaScript
  - consistent styling in large, distributed teams

Why?

- Programs are (also) written to be read
- Well-chosen variable names are critical to source code readability, reusability, maintainability
- Example tasks:
  - reverse engineering binaries
  - reverse engineering obfuscated JavaScript
  - consistent styling in large, distributed teams

# Why?

Programs

#### JSNice.org

[Predicting program properties from Big Code, ACM POPL 2015]

- Well-chose
   code rea
- Example
  - reverse
  - reverse
  - consiste



Martin Vechev, "Probabilistic Learning From Big Code". Keynote at ISSTA 2016

# Key ingredient



• The "naturalness" of software [Hindle et al, 2011]

#### Natural languages are complex



#### Natural languages are complex

Tiger, Tiger burning bright In the forests of the night What immortal hand or eye, Could frame thy fearful symmetry?



# ...but most utterances are simple & repetitive



#### Can be Rich, Powerful, Expressive







Can be Rich, Powerful, Expressive

...but "in nature" is mostly Simple, Repetitive, Boring

Can be Rich, Powerful, Expressive

..but "in nature" is mostly Simple, Repetitive, Boring Statistical Models

Can be Rich, Powerful, Expressive

..but "in nature" is mostly Simple, Repetitive, Boring Statistical Models



# The "naturalness of software" thesis

Programming Languages are complex...

...but Natural Programs are simple & repetitive.

and this, too, CAN BE EXPLOITED!!

[Hindle et al, 2011]



#### Translate

Turn off instant translation



English Persian Chinese English - detected -	Romanian English Dutch - Translate
I know what you named your × identifiers last summer!	lk weet wat je je id's genaamd afgelopen zomer!
•) //	☆ 「 • ◆ ペ





Minified Source Code









#### What's the relevance of Machine Translation?















$$e_{\text{best}} = \underset{e}{\operatorname{argmax}} p(e \mid f)$$

$$\underset{e}{\text{Eages}} = \underset{e}{\operatorname{argmax}} \frac{p(f \mid e)p(e)}{p(f)}$$

$$= \underset{e}{\operatorname{argmax}} p(f \mid e)p(e) \quad \text{(for a given } f\text{)}$$







Aligned French-English Corpus







#### Translating minified (f) to clear JS (e)



### Translating minified (f) to clear JS (e)

#### GitHub + minifier



EN: 1 know what you named your ídentífiers!

NL: Ik weet wat je je ID's genoemd!

Natural language: non-trivial alignment

- Reordering
- Different length
- Dropped words

EN: 1 know what you named your identifiers! NL: 1k weet wat je je 1D's genoemd! Natural language: non-trivial alignment

- Reordering
- Different length
- Dropped words

EN: I know what you named your identifiers! NL: Ik weet wat je je ID's genoemd!

Natural language: non-trivial alignment

- Reordering
- Different length
- Dropped words

#### function u(n, r) {

function mix(dest, src){

EN: I know what you named your identifiers! NL: Ik weet wat je je ID's genoemd!

Natural language: non-trivial alignment

- Reordering
- Different length
- Dropped words

function u(n, r) { function mix(dest, src){

Minification: straightforward alignment

}



?

function r(n, r) {
 for (var t in r) n[t] = r[t];
 return n;

function r(n, r) {
 for (var t in r) n[t] = r[t];
 return n;
}

```
function r(n, r) {
```

```
for (var t in r) n[t] = r[t];
```

return n;

}



function r(n, r) {

```
for (var t in r) n[t] = r[t];
```

return n;

}





}





mode

#### Evaluation

- Held-out test set: 2,149 files
- Comparison to JSNice [Raychev et al, 2015]
- Metric: % names recovered



#### Evaluation

- Held-out test set: 2,149 files
- Comparison to JSNice [Raychev et al, 2015]
- Metric: % names recovered
- Global vs. local names (globals don't change)

```
var geom2d = function() {
  var t = numeric.sum;
  function r(n, r) {
    this.x = n;
    this.y = r;
  }
```





#### % names recovered (2,149 test files)

![](_page_52_Figure_1.jpeg)

![](_page_53_Figure_0.jpeg)

### Becoming JSNaughty

![](_page_54_Figure_1.jpeg)

### Becoming JSNaughty

![](_page_55_Figure_1.jpeg)

#### % names recovered (2,149 test files)

![](_page_56_Figure_1.jpeg)

#### Examples

![](_page_57_Figure_1.jpeg)

https://github.com/bvasiles/jsNaughty

Clear-text corpus

Aligned clear-text/

minified corpus

- Identifier renaming using SMT, e.g., minified JS, decompiled C
- Generic, mature off-the-shelf technology (Moses)
- Language dependence restricted to tokenization and scope analysis
  - dependency parse in JSNice
- Promising results: ~50% better than JSNice on local names, on average

![](_page_58_Figure_7.jpeg)

training

This material is based upon work supported by the National Science Foundation under Grant No. 1414172

#### Machine translation for code

• Oda et al. (ASE '15): *code to pseudocode* 

```
# Python
if n % 3 == 0:
Pseudo-code:
if n is divisible by 3
```

• Karaivanov et al. (Onward! '14): porting C# to Java

<pre>// C# Console . WriteLine ( "Hello World!" );</pre>	
<pre>// Java System . out . println ( "Hello World!" )</pre>	•

#### Machine translation for code

• Oda et al. (ASE '15): *code to pseudocode* 

```
# Python
if n % 3 == 0:
Pseudo-code:
if n is divisible by 3
```

• Karaivanov et al. (Onward! '14): porting C# to Java

```
// C#
Console . WriteLine ( "Hello World!" );
// Java
System . out . println ( "Hello World!" );
```

 Nguyen et al. (FSE' 13, ASE '15): *porting Java to C#*

```
// Java
public void findResultEdges() {
  for (Iterator it = dirEdgeList.iterator(); it.hasNext();) {
    DirectedEdge de = (DirectedEdge) it.next();...}
}
// C#
public void FindResultEdges() {
  foreach (DirectedEdge de in _dirEdgeList){...}
}
```