

# Stylometry and Machine Learning in 28C3, 29C3, 31C3 ...

- 32C3 – alternative non-stylometric keynote
  - by Rachel
- 32C3: stylometry and machine learning continued
  - by Aylin



# What happened since 31C3?

De-anonymization became easier



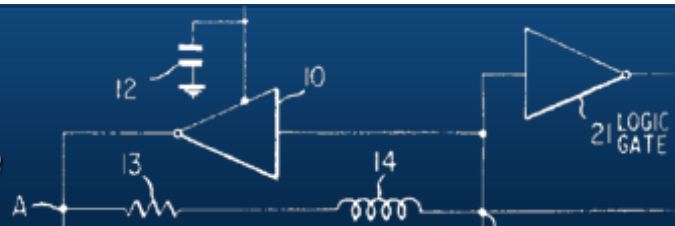
More privacy concerns for programmers



# 32C3: stylometry and machine learning continued

- Paper released today
- Check our blog
  - <https://freedom-to-tinker.com>

**FREEDOM TO TINKER**  
*research and expert commentary on digital technologies in public life*



# Stylistic fingerprints

- Stylometry has been applied to:
  - Fine-art
  - Music
  - Unconventional text
  - Translated text
  - Source code



***“Style expressed in code  
can be  
quantified and characterized.”***

**≈**

***Research ✓***



# Supervised stylometry

- Given a set of documents of known authorship, classify a document of unknown authorship
  - Classifier trained on undisputed text
- Scenario: Alice the Anonymous Blogger vs. Bob the Abusive Employer
  - Alice blogs about abuses in Bob's company
    - Blog posted anonymously (Tor, pseudonym, etc)
  - Bob obtains 5,000 words of each employee's writing
    - Bob uses stylometry to identify Alice as the blogger





**theconnor**: Cisco just offered me a job! Now I have to weigh the utility of a **fatty paycheck** against the **daily commute to San Jose** and **hating the work**.

about 19 hours ago from web · [Reply](#) · [View Tweet](#)

The screenshot shows a Twitter interface with a red background. At the top, the Twitter logo is on the left, and navigation links (Home, Profile, Find People, Settings, Help, Sign out) are on the right. The main content is a reply from user 'timmylevad' to 'theconnor'. The reply text reads: '@theconnor Who is the hiring manager. I'm sure they would love to know that you will hate the work. We here at Cisco are versed in the web.' Below the text, it says 'about 18 hours ago from TweetDeck in reply to theconnor'. The user's profile picture and name 'timmylevad' (Tim Levad) are shown at the bottom left of the tweet. A blue arrow points from the caption below to the user's name.

channel partner advocate for Cisco Alert



# Fingerprints in textual data

- **Stylistic fingerprints**



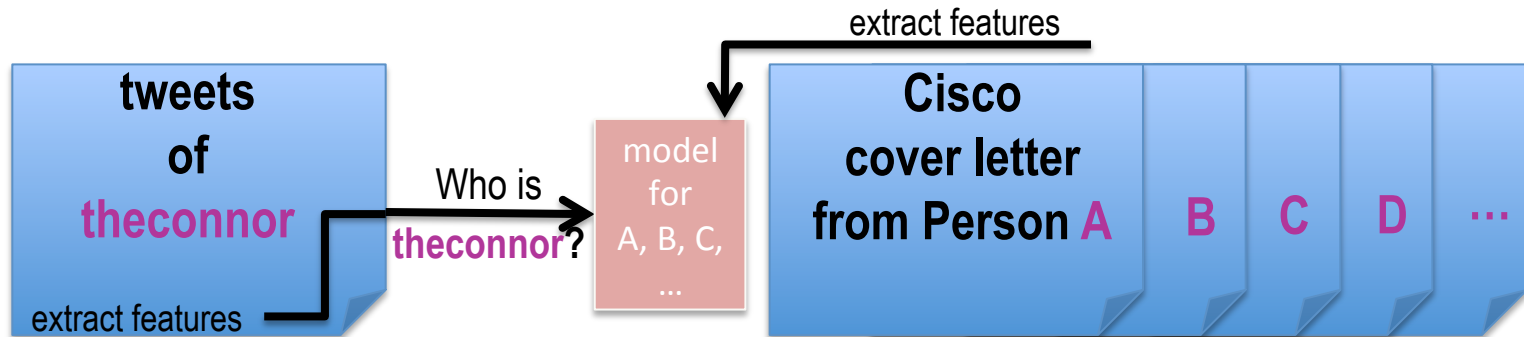
**theconnor**: Cisco just offered me a job! Now I have to weigh the utility of a fatty paycheck against the daily commute to San Jose and hating the work.

about 19 hours ago from web · Reply · View Tweet

Frequency of punctuation

Frequency of function words

- Use stylometric fingerprints to find who “**theconnor**” is.
- Collect the rest of the tweets in the timeline, compare to the cover letters submitted to Cisco and identify **theconnor**.







The image shows a screenshot of a Twitter profile page for the user @theconnor. The profile picture is a black and white illustration of a man in a suit with a cigarette in his mouth. The header is red. Below the profile picture, the statistics are: TWEETS 1,235, FOLLOWING 83, FOLLOWERS 95, and FAVORITES 9. There is a 'Follow' button. The bio area contains the text: '@theconnor's Tweets are protected. Only confirmed followers have access to @theconnor's Tweets and complete profile. Click the "Follow" button to send a follow request.' The name is 'theconnor' with a lock icon, the handle is '@theconnor', the website is 'theconnor.net', and it says 'Joined January 2009'.

theconnor 

@theconnor

theconnor.net

 theconnor.net

 Joined January 2009

@theconnor's Tweets are protected.

Only confirmed followers have access to @theconnor's Tweets and complete profile. Click the "Follow" button to send a follow request.

**theconnor** made her Twitter profile private and deleted all information on her homepage right after the event but it was too late since search engines cache search results which can lead to old information.





Search Twitter



Have an account? [Log in](#)



TWEETS  
**1,235**

FOLLOWING  
**83**

FOLLOWERS  
**95**

FAVORITES  
**9**

**Follow**

**theconnor**

@theconnor

[theconnor.net](#)

[theconnor.net](#)

Joined January 2009

@theconnor's Tweets are protected.

Only confirmed followers have access to @theconnor's Tweets and complete profile. Click the "Follow" button to send a follow request.

theconnor → Connor Riley →



# What about fingerprints in source code?



**zooko**

@zooko

I just heard from an intern at Apple that they disallow her from contributing to open source on her own time. That's illegal, right?

RETWEETS

11

LIKES

6



11:04 AM - 4 Dec 2015



# DE-ANONYMIZING PROGRAMMERS VIA CODE STYLOMETRY



De-anonymizing Programmers via Code Stylometry.

**Aylin Caliskan-Islam, Richard Harang, Andrew Liu, Arvind Narayanan, Clare Voss,  
Fabian Yamaguchi, and Rachel Greenstadt.**

Usenix Security Symposium, 2015

# Source code stylometry

- Everyone learns coding on an individual basis, as a result code in a unique style, which makes de-anonymization possible.
- Software engineering insights
  - programmer style changes while implementing sophisticated functionality
  - differences in coding styles of programmers with different skill sets
- Identify malicious programmers.



# Source code stylometry: Who wrote this code?

- **Scenario 1:**

Alice analyzes a library with malicious source code.

Bob has a source code collection with known authors.

Bob will search his collection to find Alice's adversary.



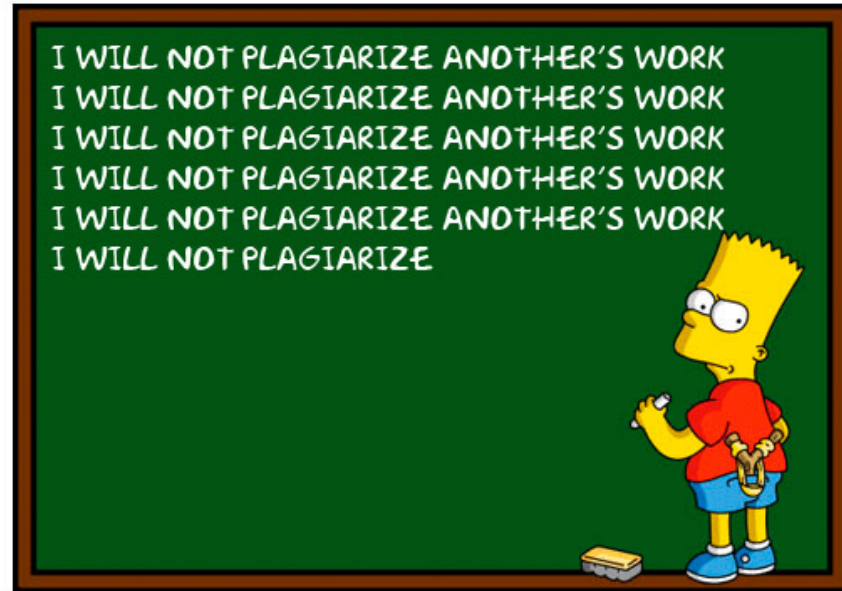
# Source code stylometry: Who wrote this code?

- **Scenario 2:**

Alice got an extension for her programming assignment.

Bob, the professor has everyone else's code.

Bob wants to see if Alice plagiarized.



# Source code stylometry

Iran confirms death sentence for 'porn site' web programmer



No technical difference between security-enhancing and privacy-infringing...





# Comparison to related work

Related Work	Author Size	Instances	Average LOC	Language	Features	Method	Result
MacDonell et al.	7	351	148	C++	lexical & layout	Case-based reasoning	88.0%
Frantzeskou et al.	8	107	145	Java	lexical & layout	Nearest neighbor	100.0%
Elenbogen and Seliya	12	83	100	C++	lexical & layout	C4.5 decision tree	74.7%
Shevertalov et. al.	20	N/A	N/A	Java	lexical & layout	Genetic algorithm	80%
Frantzeskou et al.	30	333	172	Java	lexical & layout	Nearest neighbor	96.9%
Ding and Samadzadeh	46	225	N/A	Java	lexical & layout	Nearest neighbor	75.2%
<b>Ours</b>	<b>35</b>	<b>315</b>	<b>68</b>	<b>C++</b>	<b>lexical &amp; layout &amp; syntactic</b>	<b>Random forest</b>	<b>100.0%</b>
<b>Ours</b>	<b>250</b>	<b>2,250</b>	<b>77</b>	<b>C++</b>			<b>98.0%</b>
<b>Ours</b>	<b>1,600</b>	<b>14,400</b>	<b>70</b>	<b>C++</b>			<b>93.6%</b>



# Comparison to related work

Related Work	Author Size	Instances	Average LOC	Language	Features	Method	Result
Frantzeskou et al.	<u>30</u>	333	172	Java	lexical & layout	Nearest neighbor	<u>96.9%</u>
Ding and Samadzadeh	46	225	N/A	Java	lexical & layout	Nearest neighbor	75.2%
Ours	35	315	68	C++	lexical & layout & syntactic	Random forest	100.0%
Ours	<u>250</u>	2,250	77	C++			<u>98.0%</u>
Ours	1,600	14,400	70	C++			93.6%



# Comparison to related work

Related Work	Author Size	Instances	Average LOC	Language	Features	Method	Result
Frantzeskou et al.	30	333	172	Java	lexical & layout	Nearest neighbor	96.9%
Ding and Samadzadeh	<u>46</u>	225	N/A	Java	lexical & layout	Nearest neighbor	<u>75.2%</u>
Ours	35	315	68	C++	lexical & layout & syntactic	Random forest	100.0%
Ours	250	2,250	77	C++			98.0%
Ours	<u>1,600</u>	14,400	70	C++			<u>93.6%</u>



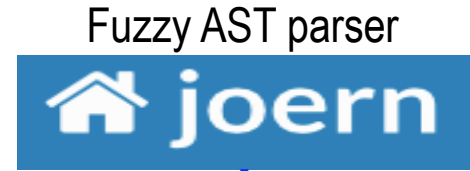
# Source code stylometry

## Machine learning workflow

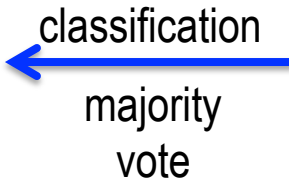
Dataset in CPP ~100,000 users

# code jam

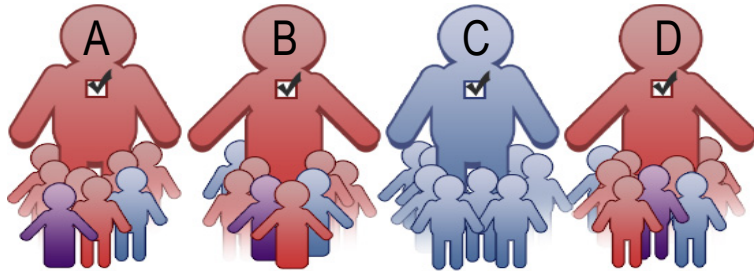
```
System.out.println("hello, world!");
```



Extract features



Random Forest



# code jam

```
System.out.println("hello, world!");
```

- 2008-2014
- Same problems
- Limited time
- Problems get harder
- C++ most common

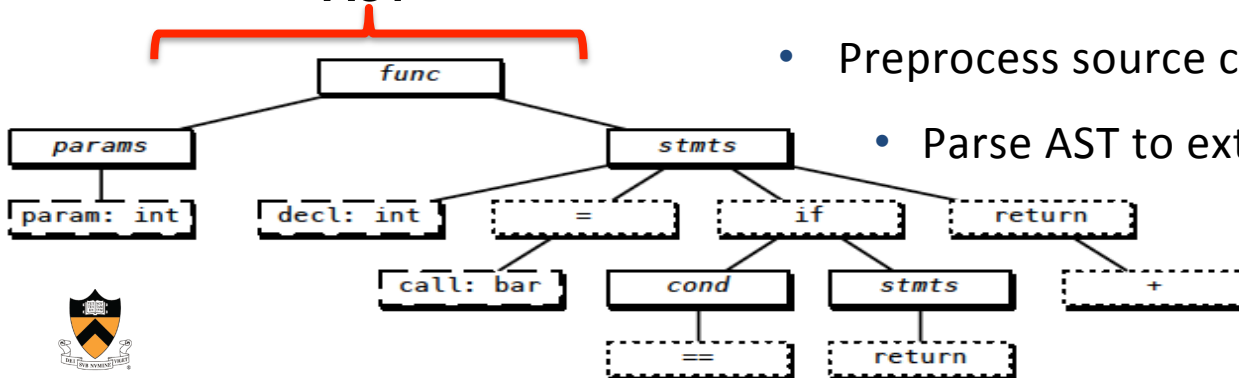


# Source code stylometry

```
int foo(int y)
{
    int n = bar(y);
    if (n == 0)
        return 1;
    return (n + y);
}
```

Source Code

AST



- Stylometry can be used in source code to identify the author of a program.

- Extract layout and lexical features from source code.

- Abstract syntax trees (AST) in code represent the structure of the program.

- Preprocess source code to obtain AST.

- Parse AST to extract coding style features.



# Source code stylometry

## Method

- ① Use random forest as the machine learning classifier,
  - ① avoid over-fitting
  - ② multi-class classifier by nature
- ② K-fold cross validation
- ③ Validate method on a different dataset



# Case 1: Authorship attribution

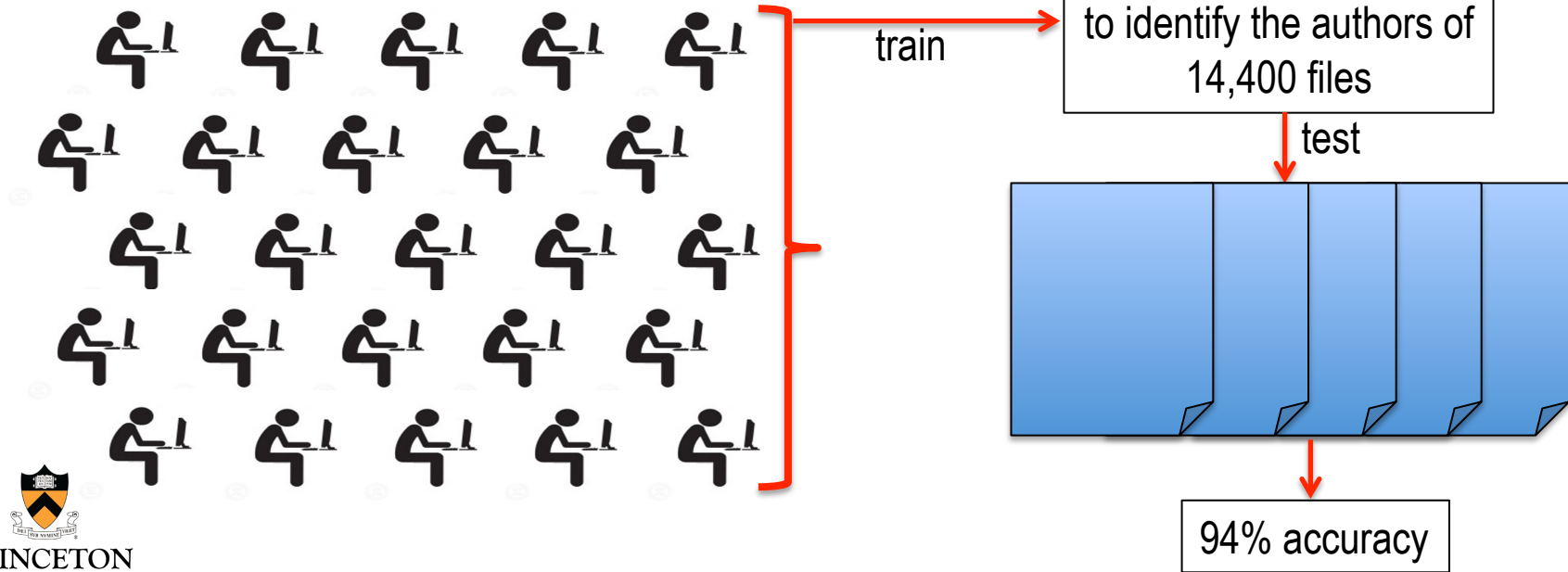
- Who is this anonymous programmer?
- Who is Satoshi?





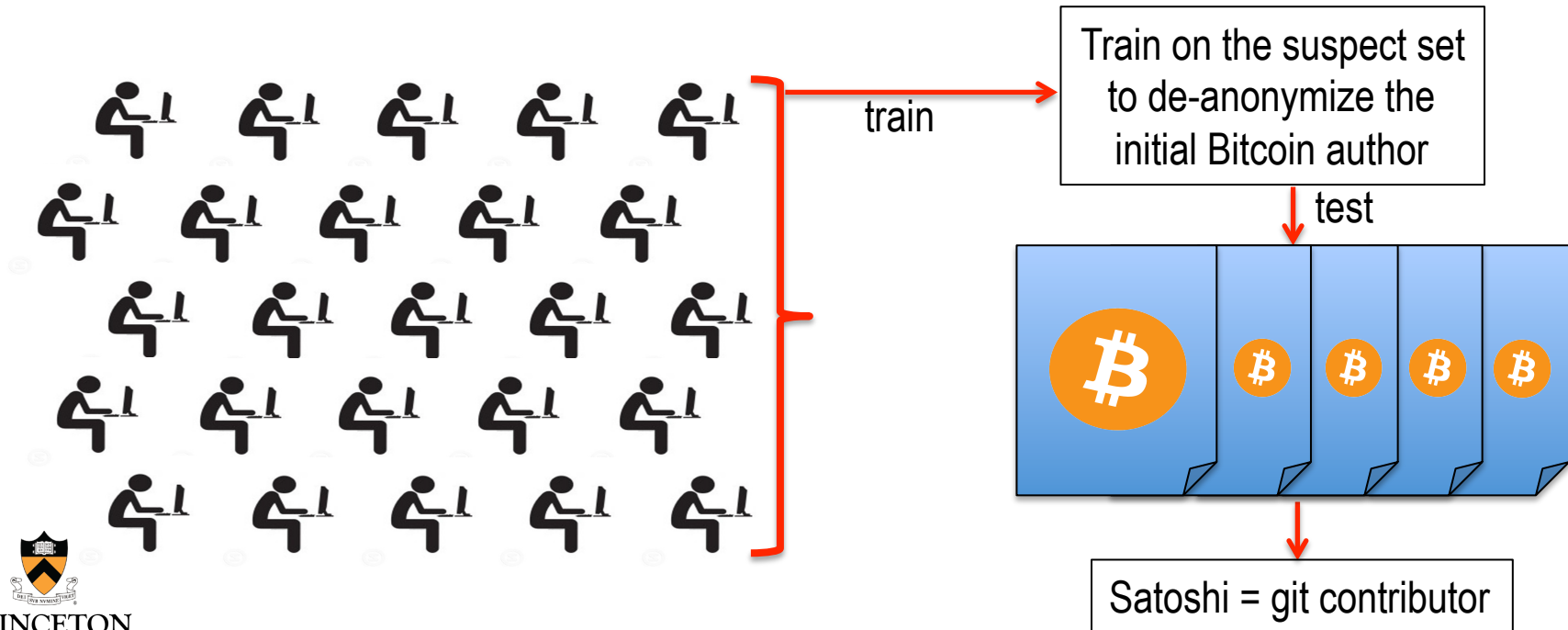
# Case 1: Authorship attribution

- 94% accuracy in identifying 1,600 authors of 14,400 anonymous program files.



# Case 1: Authorship attribution

- If only we had a suspect set for Satoshi...



# Case 2: Obfuscation

- Who is the programmer of this obfuscated source code?
- Code is obfuscated to become unrecognizable.
- Our authorship attribution technique is impervious to common off-the-shelf source code obfuscators.



# Case 2: C++ Obfuscation - STUNNIX

Sample file with C++ code

```
#ifdef __STL_USE_EXCEPTIONS /* this is conditional preprocessing */
extern void __out_of_range (const char *);
#define OUTFRANGE(cond, msg) \
    do { if (cond) __out_of_range (#cond); } while (0)
#else
#include <cassert>
#define OUTFRANGE(cond, msg) assert (!(cond))
#endif

template <class charT, class traits, class Allocator>
basic_string <charT, traits, Allocator>&
basic_string <charT, traits, Allocator>::
replace (size_type pos1, size_type n1,
        const basic_string& str, size_type pos2, size_type n2)
{
    //rather complex body follows
    const size_t len2 = str.length () + 2;

    if (pos1 == 0 && n1 >= length () && pos2 == 0 && n2 >= len2)
```

# Case 2: C++ Obfuscation - STUNNIX

Sample file with C++ code



```
#ifdef __STL_USE_EXCEPTIONS /* this is conditional preprocessing */
extern void __out_of_range (const char *);
#define OUTFRANGE(cond,msg) \
    do { if (cond) __out_of_range (#cond); } while (0)
#else
#include <cassert>
#define OUTFRANGE(cond,msg) assert (!(cond))
#endif

template <class charT, class traits, class Allocator>
basic_string <charT, traits, Allocator>&
basic_string <charT, traits, Allocator>::
replace (size_type pos1, size_type n1,
        const basic_string& str, size_type pos2, size_type n2)
{
    //rather complex body follows
    const size_t len2 = str.length () + 2;

    if (pos1 == 0 && n1 >= length () && pos2 == 0 && n2 >= len2)
```



# Case 2: C++ Obfuscation - STUNNIX

Sample file with C++ code

```

#ifdef z7929401884
extern void za4ldafc42e(const char*);
#define zlc52ffdd48(z22fc207d33, zde05b8b1b0) \
    do { if (z22fc207d33) za4ldafc42e(z22fc207d33); } while ((0x1a1+8313-0x221a))
#else
#include <cassert>
#define zlc52ffdd48(z22fc207d33, zde05b8b1b0) \
    do { if (z22fc207d33) za4ldafc42e(z22fc207d33); } while ((0x1a1+8313-0x221a))
#endif
template<class zd9cfc9cefe, class z9cdf2cd536, Allocator>
::replace(size_type z795f772c7c, size_type z8ad17de27a, size_type za2e5f06cde) {
const size_t z5ldea4lale=str.length()+ (0x12ac+3131-0xlee5); if (z795f772c7c==
(0x455+8190-0x2453)&& zddd43c876a>=length() && z8ad17de27a== (0xc15+4853-0x1f0a)&&
za2e5f06cde>=z5ldea4lale) return operator=(str); zlc52ffdd48(z8ad17de27a>
z5ldea4lale, "\x65\x72\x72\x6f\x72\x20\x69\x6e\x20\x72\x65\x70\x6c\x61\x63\x65");
#ifdef zd943335d79
++ :: z021c346d26. z1534cdbaf9;
#endif

```

Same set of 20 authors  
with 180 program files

Classification  
Accuracy

Original source code

99%

STUNNIX-Obfuscated source code

99%

# Case 2: C Obfuscation - TIGRESS

```
#include<stdio.h>
int main()
{
    int T,test=1;
    double C,F,X,rate,time;
    scanf("%d",&T);
    while(T--)
    {
        scanf("%lf %lf %lf",&C,&F,&X);
        rate=2.0;
        time=0;
        while(X/rate>C/rate+X/(rate+F))
        {
            time+=C/rate;
            rate+=F;
        }
        time+=X/rate;
        printf("Case #d: %lf\n",test++,time);
    }
    return 0;
}
```

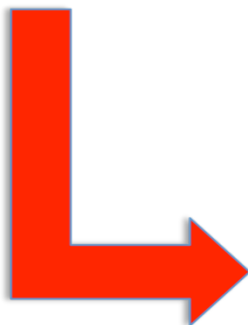


# Case 2: C Obfuscation - TIGRESS

```

#include<stdio.h>
int main()
{
    int T,test=1;
    double C,F,X,rate,time;
    scanf("%d",&T);
    while(T--)
    {
        scanf("%lf %lf %lf",&C,&F,&X);
        rate=2.0;
        time=0;
        while(X/rate>C/rate+X/(rate+F))
        {
            time+=C/rate;
            rate+=F;
        }
        time+=X/rate;
        printf("Case #%d: %lf\n",test++,time);
    }
    return 0;
}

```



```

struct _IO_FILE;
struct timeval {
    long tv_sec ;
    long tv_usec ;
};
enum __1_main_$op {
    __1_main__string
$value_LIT_0$result_REG_1__convert_void_star2void_star
$result_STA_0$left_REG_0__local
$result_STA_0$value_LIT_0__store_void_star
$left_STA_0$right_STA_1__local
$result_STA_0$value_LIT_0__convert_void_star2void_star
$left_STA_0$result_REG_0__local
$result_REG_0$value_LIT_1__convert_void_star2void_star
$result_STA_0$left_REG_0__store_void_star$right_STA_0$left_REG_0 =
46,
    __1_main__local$result_REG_0$value_LIT_1__constant_int
$result_STA_0$value_LIT_0__store_int$right_STA_0$left_REG_0__local
$result_STA_0$value_LIT_0__convert_void_star2void_star
$left_STA_0$result_REG_0__string
$value_LIT_0$result_REG_1__convert_void_star2void_star
$result_STA_0$left_REG_0__store_void_star
$right_STA_0$left_REG_0__local
$result_REG_0$value_LIT_1__convert_void_star2void_star
$result_STA_0$left_REG_0 = 44,
    __1_main__convert_void_star2void_star
$left_STA_0$result_REG_0__load_int
$left_REG_0$result_REG_1__MinusA_int_int2int
$result_REG_0$left_REG_1$right_REG_2__store_int
$left_STA_0$right_REG_0__goto$label_LAB_0 = 161,
    __1_main__local$result_STA_0$value_LIT_0__local
$result_REG_0$value_LIT_1__convert_void_star2void_star
$result_STA_0$left_REG_0__load_double
$left_STA_0$result_REG_0__local
$result_REG_0$value_LIT_1__convert_void_star2void_star
$result_STA_0$left_REG_0__load_double
$left_STA_0$result_STA_0__convert_double2double
$left_STA_0$result_REG_0__local
$result_REG_0$value_LIT_1__convert_void_star2void_star

```





# Case 2: C Obfuscation - TIGRESS

```
#include<stdio.h>
int main()
{
    int T,test=1;
    double C,F,X,rate,time;
    scanf("%d",&T);
    while(T--)
    {
        scanf("%lf %lf %lf",&C,&F,&X);
        rate=2.0;
        time=0;
        while(X/rate>C/rate+X/rate)
        {
            time+=C/rate;
            rate+=F;
        }
        time+=X/rate;
        printf("Case #%d: %lf\n",T,time);
    }
    return 0;
}
```

```
struct _IO_FILE;
struct timeval {
    long tv_sec ;
    long tv_usec ;
};
enum _1_main_$op {
    _1_main_string
$value_LIT_0$result_REG_1__convert_void_star2void_star
$result_STA_0$left_REG_0__local
$result_STA_0$value_LIT_0__store_void_star
```

Same set of 20 authors  
with 180 program files

Classification  
Accuracy

Original C source code

96%

TIGRESS-Obfuscated source code

67%

Random chance of correct de-anonymization

5%



```
$left_STA_0$result_REG_0__load_int
$left_REG_0$result_REG_1__MinusA_int_int2int
$result_REG_0$left_REG_1$right_REG_2__store_int
$left_STA_0$right_REG_0__goto$label_LAB_0 = 161,
    _1_main__local$result_STA_0$value_LIT_0__local
$result_REG_0$value_LIT_1__convert_void_star2void_star
$result_STA_0$left_REG_0__load_double
$left_STA_0$result_REG_0__local
$result_REG_0$value_LIT_1__convert_void_star2void_star
$result_STA_0$left_REG_0__load_double
$left_STA_0$result_STA_0__convert_double2double
$left_STA_0$result_REG_0__local
$result_REG_0$value_LIT_1__convert_void_star2void_star
```



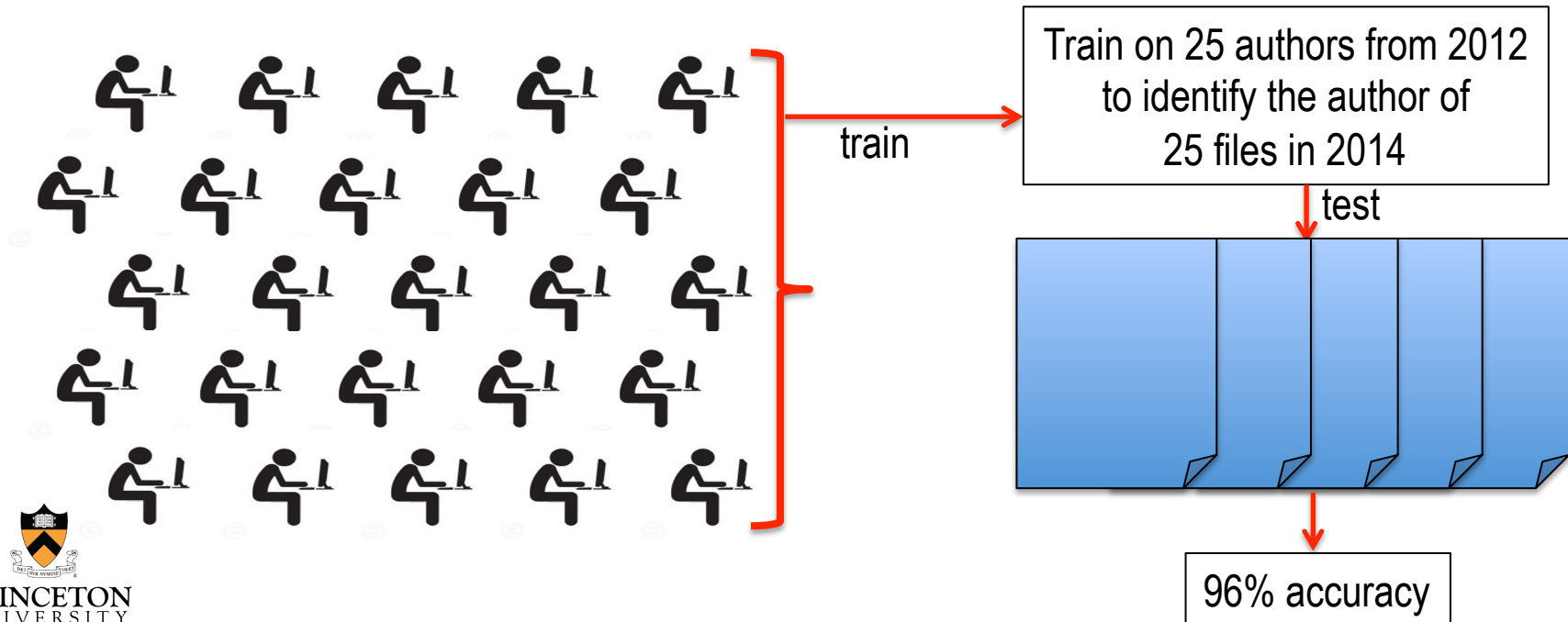
# Case 3: Coding style throughout years

- Is programming style consistent?
- If yes, we can use code from different years for authorship attribution.



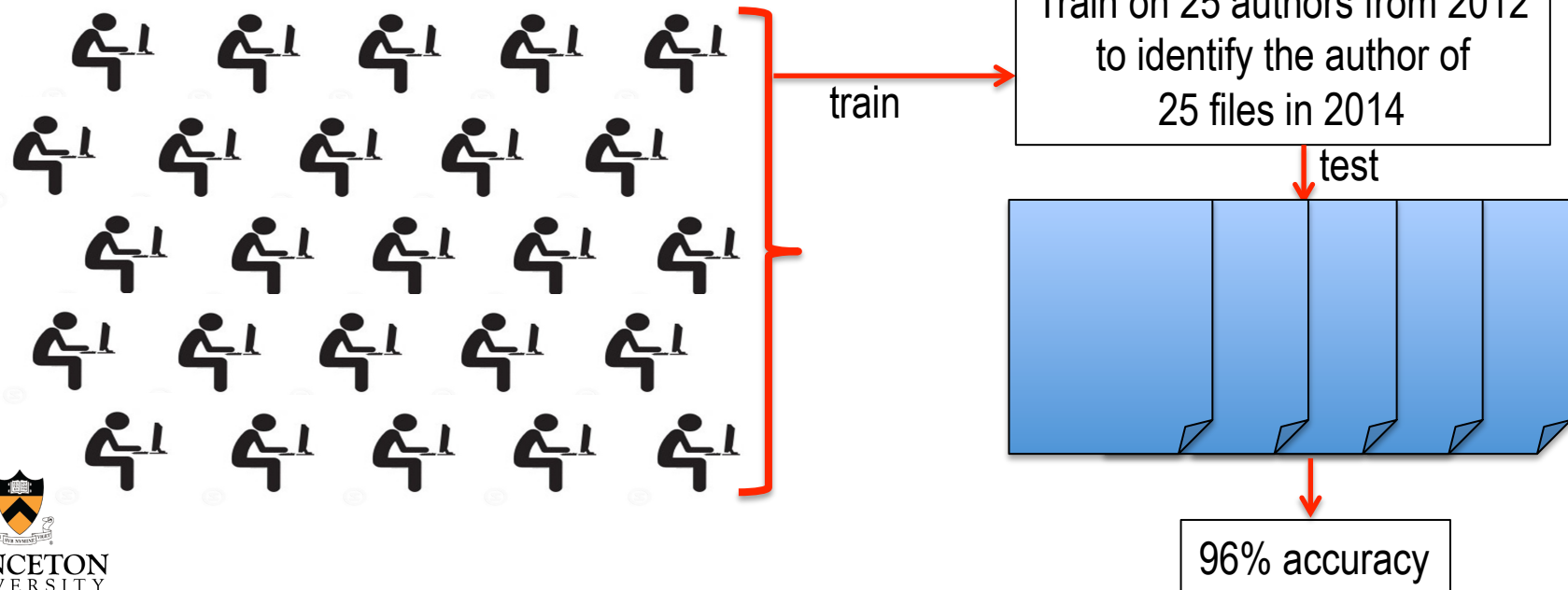
# Case 3: Coding style throughout years

- Coding style is preserved up to some degree throughout years.



# Case 3: Coding style throughout years

- 98% accuracy, train and test in 2014
- 96% accuracy, train on 2012, test on 2014



# Case 4: Generalizing the approach - python

Feature set: Using 'only' the Python equivalents of syntactic features

Application	Programmers	Instances	Result
Python programmer de-anonymization	229	2,061	53.9%
Top-5 relaxed classification	229	2,061	75.7%
Python programmer de-anonymization	23	207	87.9%
Top-5 relaxed classification	23	207	99.5%



# Results

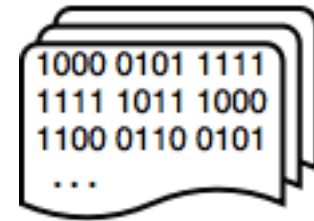
A new principled method with a robust syntactic feature set for de-anonymizing programmers.

Application	Classes	Instances	Accuracy
Stylometric plagiarism detection	250 class	2,250	98.0%
Large scale de-anonymization	1,600 class	14,400	93.6%
Copyright investigation	Two-class	1,080	100.0%
Authorship verification	Two-class/One-class	2,240	91.0%
Open world problem	Multi-class	420	96.0%



# Future work

- Multiple authorship detection
- Multiple author identification
- Anonymizing source code
  - obfuscation is not the answer
- Stylometry in executable binaries
  - Authorship attribution



# What about executable binaries?

Source Code



Compiled code looks cryptic

```
#include <cstdio>
#include <algorithm>
using namespace std;
#define For(i,a,b) for(int i = a; i < b; i++)
#define FOR(i,a,b) for(int i = b-1; i >= a; i--)
double nextDouble() {
    double x;
    scanf("%lf", &x);
    return x;}
int nextInt() {
    int x;
    scanf("%d", &x);
    return x;}
int n;
double a1[1001], a2[1001];
int main() {
    freopen("D-small-attempt0.in", "r", stdin);
    freopen("D-small.out", "w", stdout);
    int tt = nextInt();
    For(t,1,tt+1) {
        int n = nextInt();    . . .
```

```
00100000 00000000 00001000 00000000 00101000 00000000
00000000 00000000 00110100 00000000 00000000 00000000
00000100 00001000 00000000 00000001 00000000 00000000
00000000 00000001 00000000 00000000 00000101 00000000
00000000 00000000 00000100 00000000 00000000 00000000
00000011 00000000 00000000 00000000 00110100 00000001
00000000 00000000 00110100 10000001 00000100 00001000
00000000 00000000 00010011 00000000 00000000 00000000
00000100 00000000 00000000 00000000 00000001 00000000
00000000 00000000 00000001 00000000 00000000 00000000
00000000 00000000 00000000 00000000 00000000 10000000
00000100 00001000 00000000 10000000 00000100 00001000
11001000 00010111 00000000 00000000 11001000 00010111
00000000 00000000 00000101 00000000 00000000 00000000
00000000 00010000 00000000 00000000 00000001 00000000
00000000 00000000 11001000 00010111 00000000 00000000
11001000 10100111 00000100 00001000 11001000 10100111
00000100 00001000 00101100 00000001 00000000 00000000
00000000 00000000 00000000 00010000 00000000 00000000
00000010 00000000 00000000 00000000 11011100 00010111
```

. . .





# Can we identify the author of this binary?

```

00100000 00000000 00001000 00000000 00101000 00000000
00000000 00000000 00110100 00000000 00000000 00000000
00000100 00001000 00000000 00000001 00000000 00000000
00000000 00000001 00000000 00000000 00000101 00000000
00000000 00000000 00000100 00000000 00000000 00000000
00000011 00000000 00000000 00000000 00110100 00000001
00000000 00000000 00110100 10000001 00000100 00001000
00000000 00000000 00010011 00000000 00000000 00000000
00000100 00000000 00000000 00000000 00000001 00000000
00000000 00000000 00000001 00000000 00000000 00000000
00000000 00000000 00000000 00000000 00000000 10000000
00000100 00001000 00000000 10000000 00000100 00001000
11001000 00010111 00000000 00000000 11001000 00010111
00000000 00000000 00000101 00000000 00000000 00000000
00000000 00010000 00000000 00000000 00000001 00000000
00000000 00000000 11001000 00010111 00000000 00000000
11001000 10100111 00000100 00001000 11001000 10100111
00000100 00001000 00101100 00000001 00000000 00000000
00000000 00000000 00000000 00010000 00000000 00000000
00000010 00000000 00000000 00000000 11011100 00010111

```

...



# WHEN CODING STYLE SURVIVES COMPILATION: DE-ANONYMIZING PROGRAMMERS FROM EXECUTABLE BINARIES



PRINCETON  
UNIVERSITY

When Coding Style Survives Compilation: De-anonymizing Programmers from Executable Binaries  
**Aylin Caliskan-Islam, Fabian Yamaguchi, Edwin Dauber, Richard Harang, Konrad Rieck, Rachel Greenstadt, and Arvind Narayanan.**

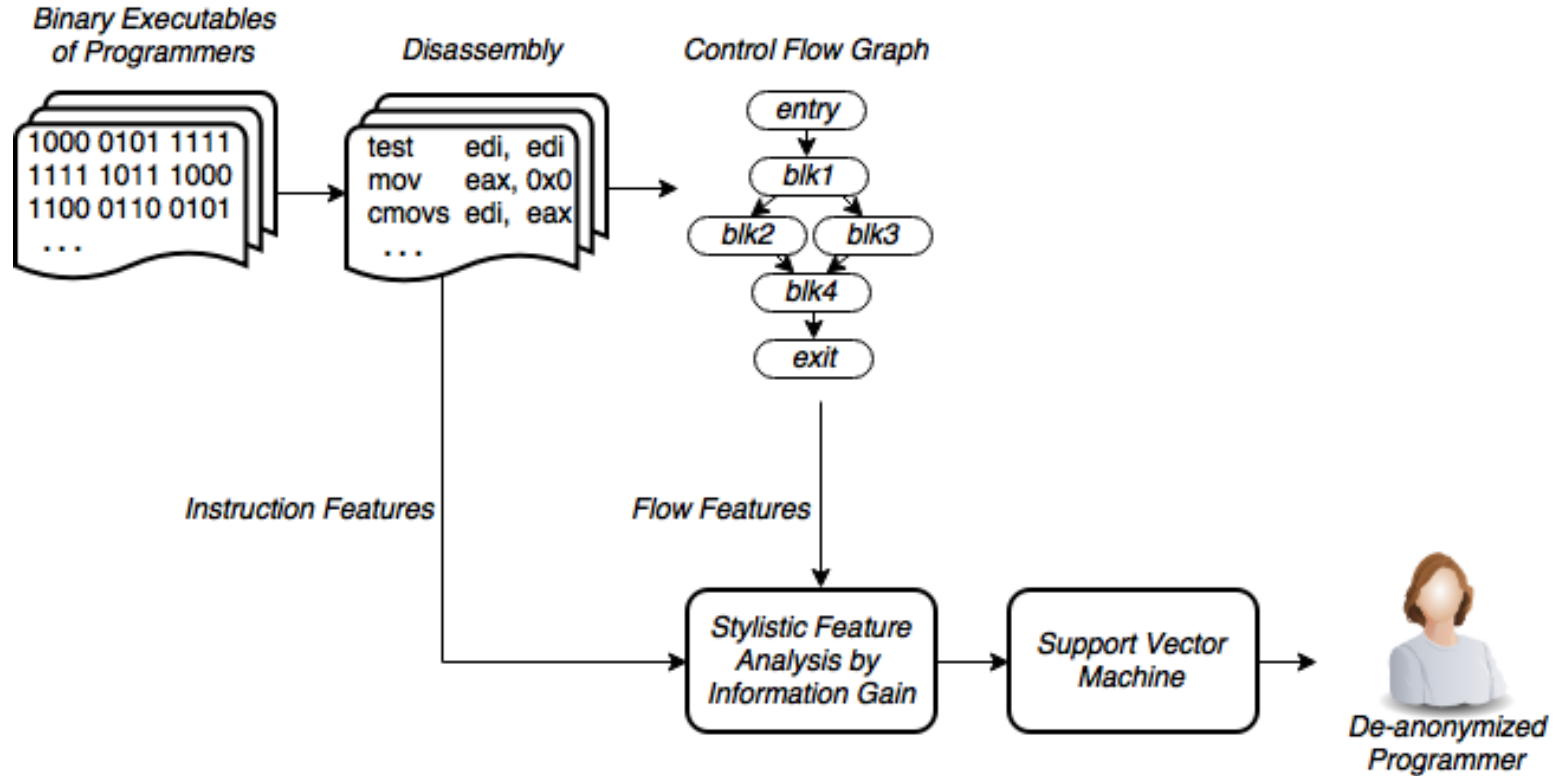
Under Submission, 2015

# Finding the author of an executable binary?

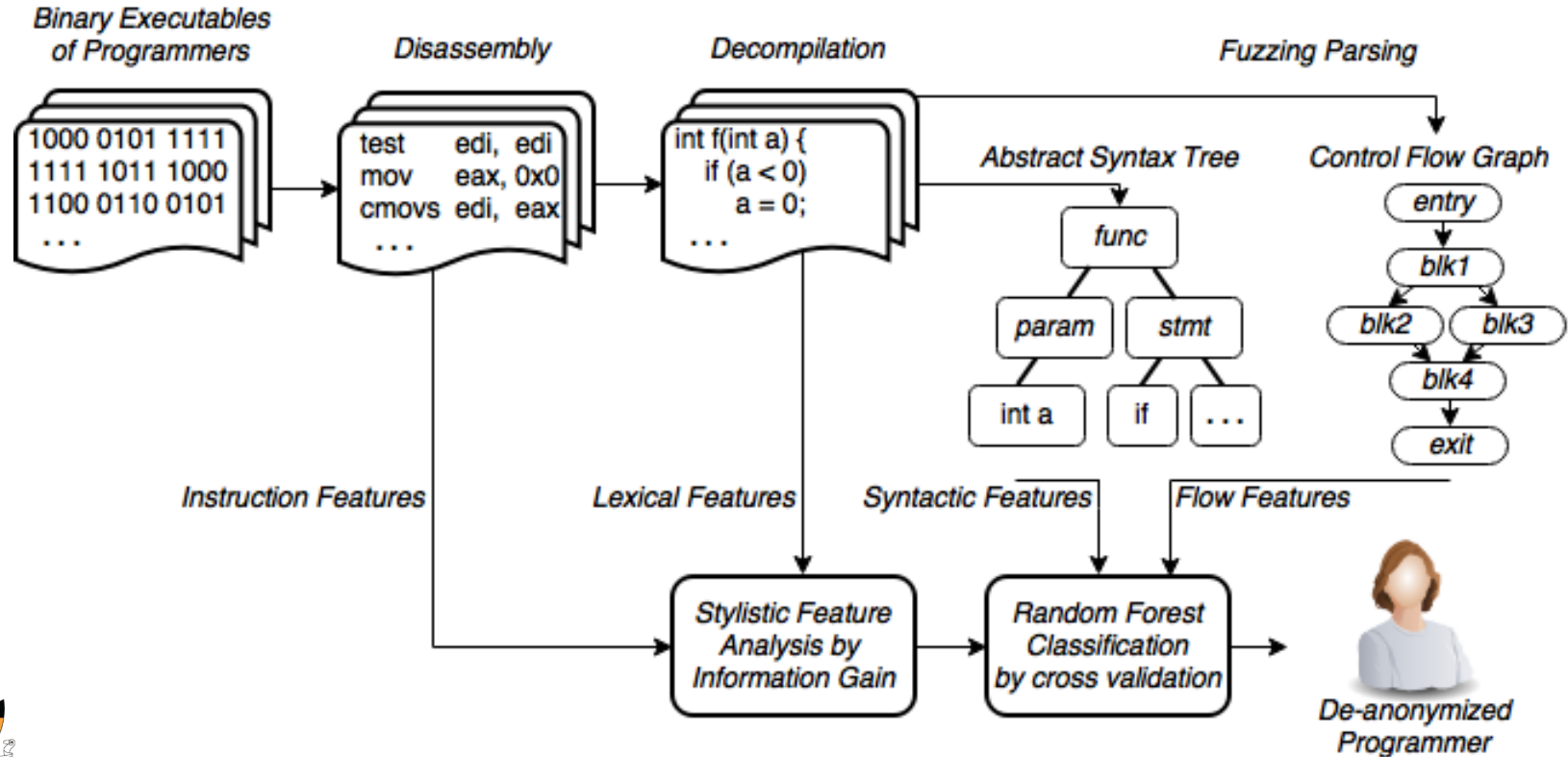
- Coding style in compiled code
- Threat to privacy and anonymity
- Malware classification?



# Related work

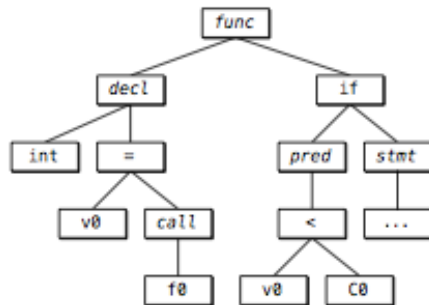


# Our workflow



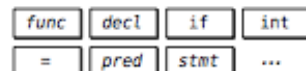
# Features

Abstract syntax tree (AST)

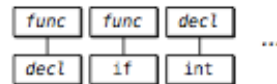


Syntactic features

AST unigrams:

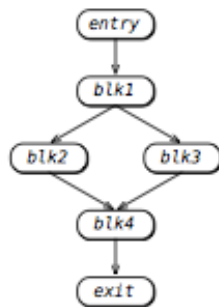


AST bigrams:



AST depth: 5

Control-flow graph (CFG)

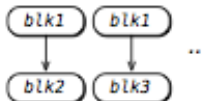


Control-flow features

CFG unigrams:



CFG bigrams:

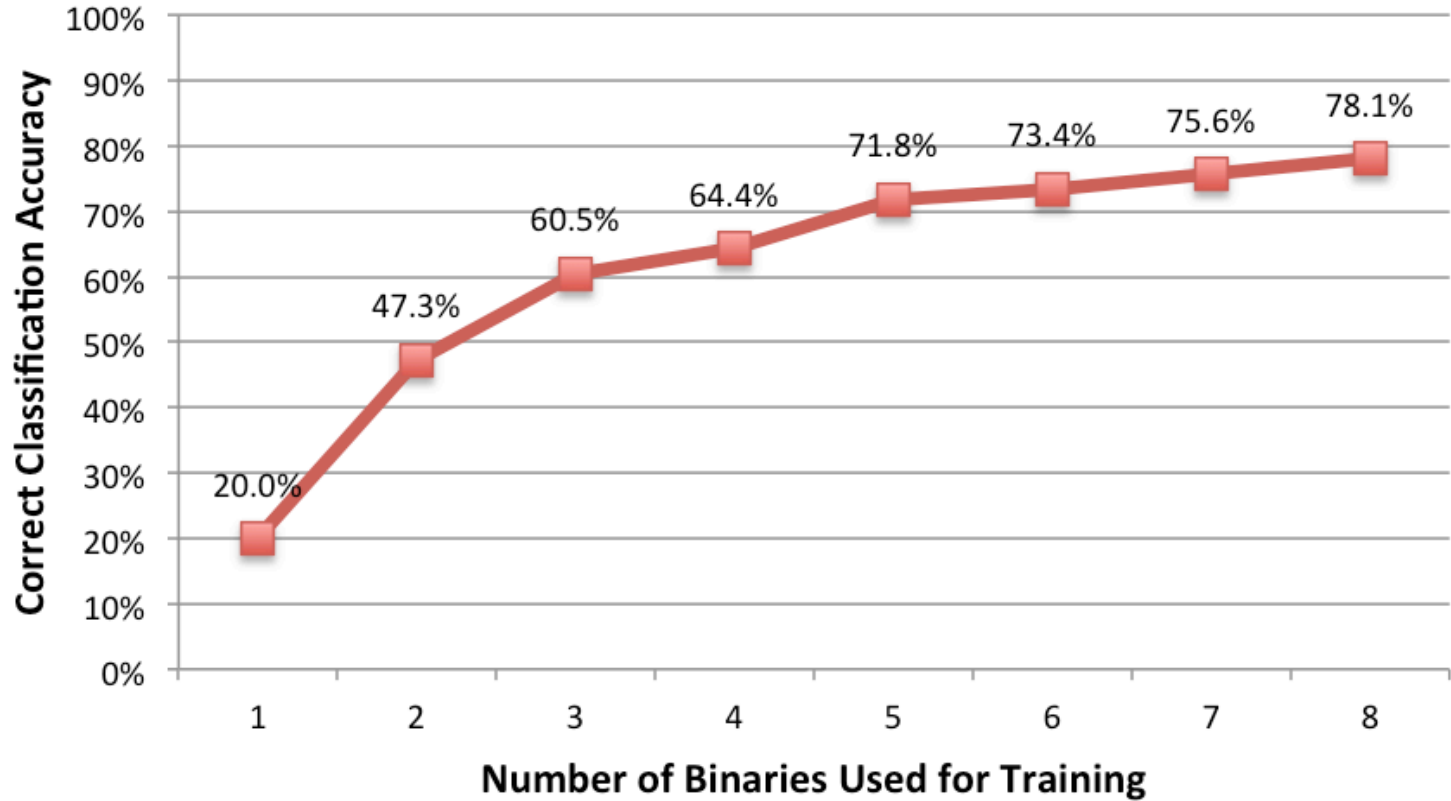


# Feature set from 100 programmers

Feature	Source	Number of Possible Features	Information Gain Features
word unigrams	hex-rays decompiled code	29,686	102
AST node TF*	hex-rays decompiled code	14,663	24
Labeled AST edge TF*	decompiled code	25,941	88
AST node TFIDF**	decompiled code	14,663	8
AST node average depth	decompiled code	14,663	21
C++ keywords	decompiled code	73	5
radare2 disassembly unigrams	radare disassembly	12,629	45
radare2 disassembly bigrams	radare disassembly	33,919	75
ndisasm disassembly unigrams	ndisasm disassembly	532	8
ndisasm disassembly bigrams	ndisasm disassembly	4,570	25
CFG unigrams	Snowman CFG	11,503	5
CFG unigram TFIDF**	Snowman CFG	11,503	10
CFG bigrams	Snowman CFG	38,554	10
Total		201,396	426
<i>TF* = term frequency</i>			
<i>TFIDF** = term frequency inverse document frequency</i>			

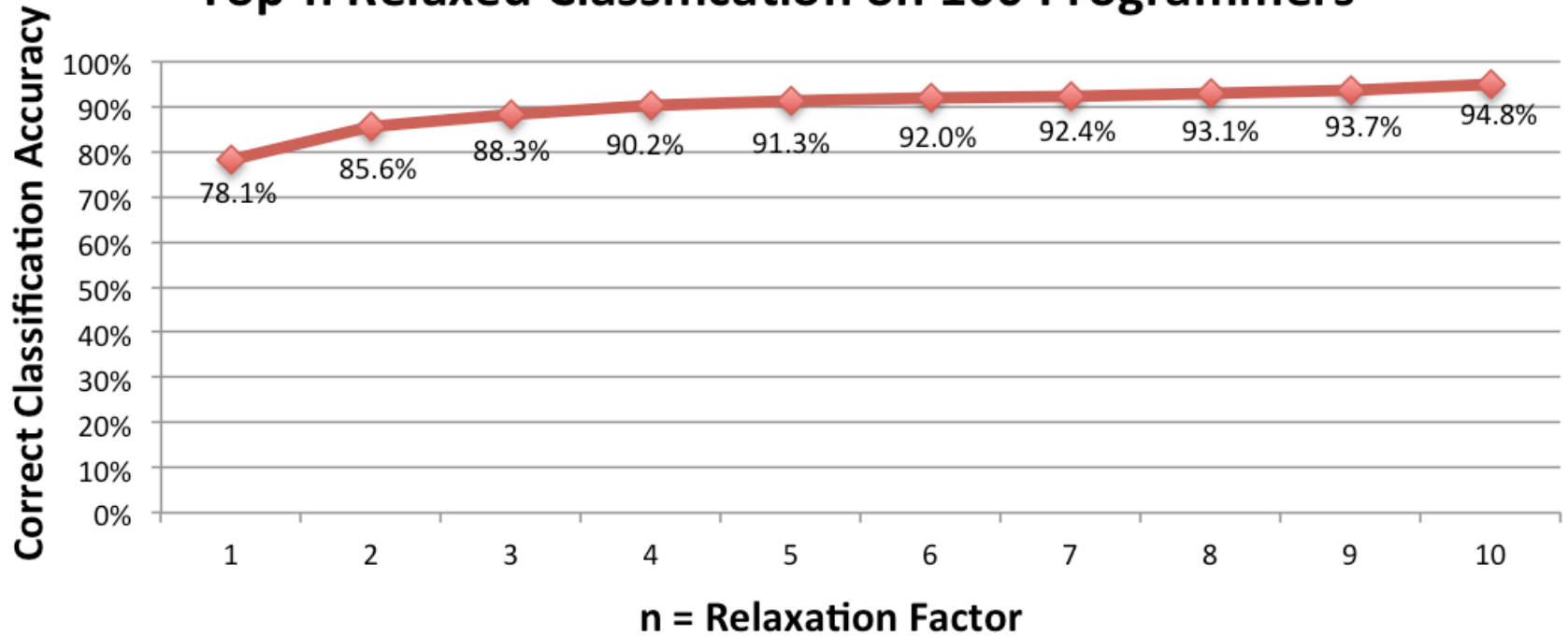


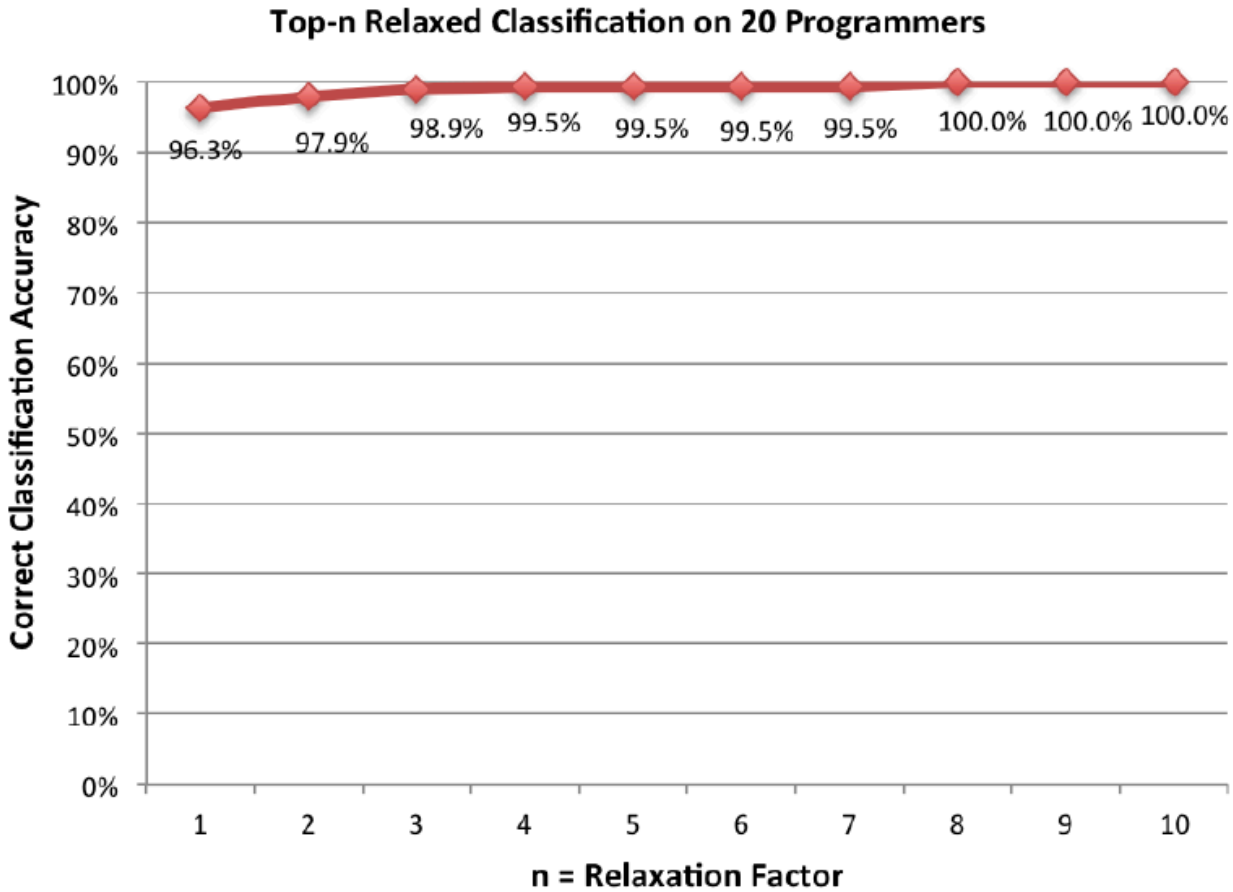
Amount of Training Data Required for Classifying 100 Programmers



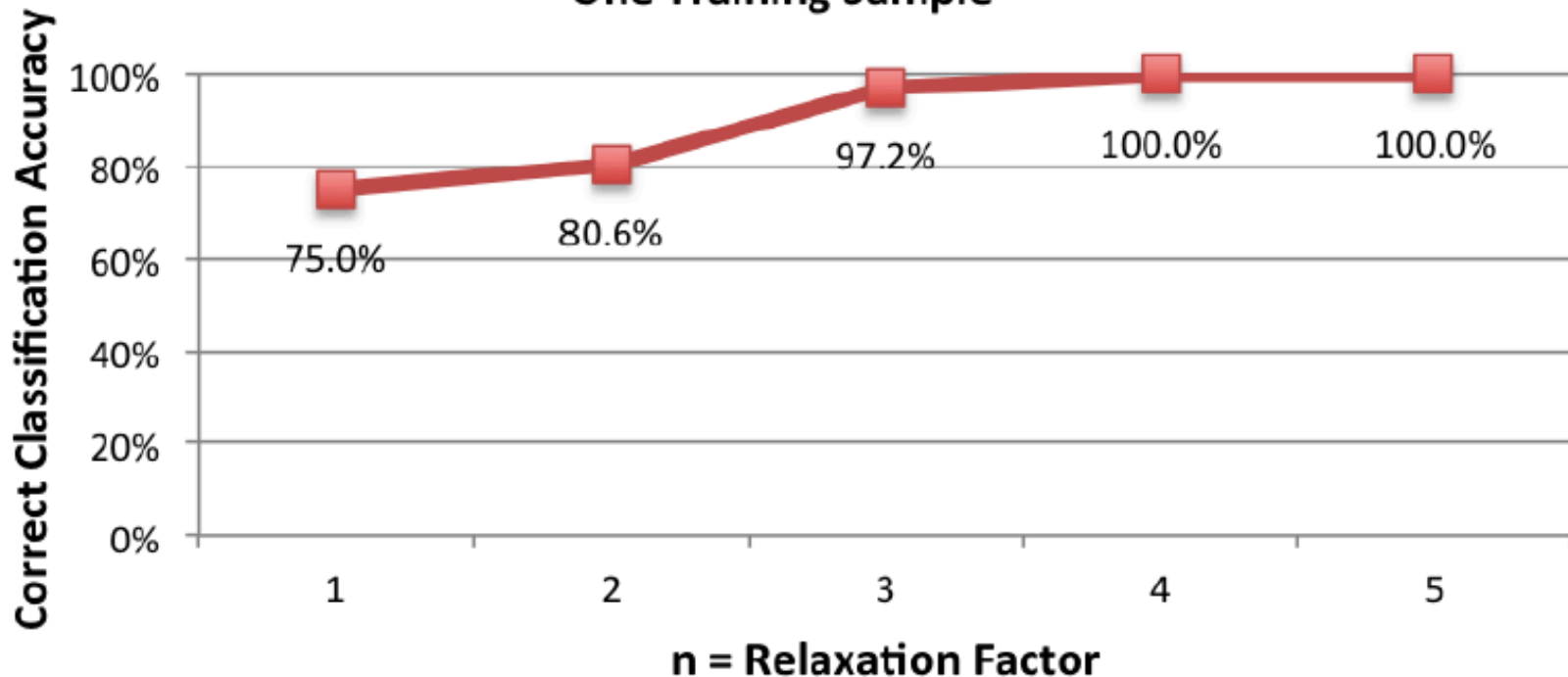


## Top-n Relaxed Classification on 100 Programmers

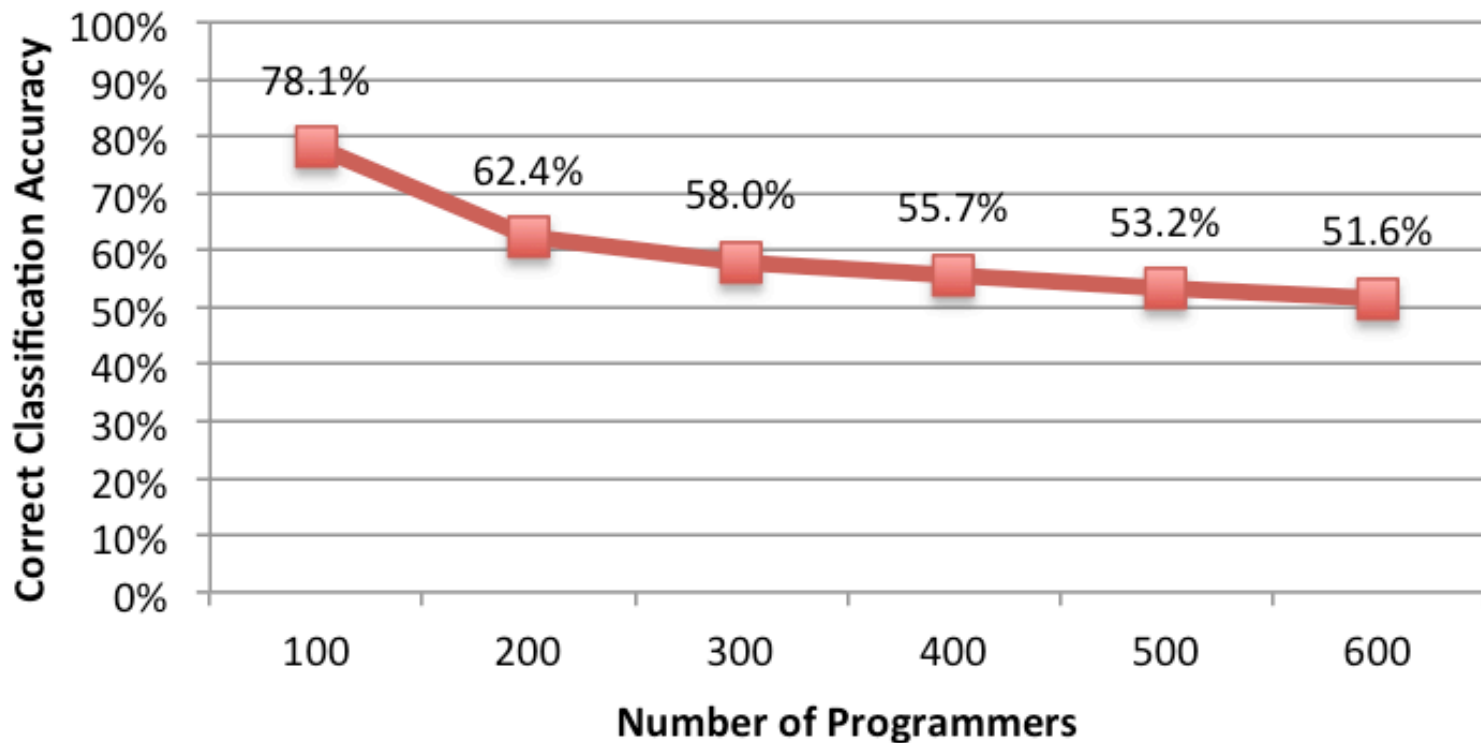




## Top-5 Relaxed Classification of 20 Programmers with One Training Sample



## Large Scale Programmer De-anonymization



<b>Related Work</b>	<b>Number of Programmers</b>	<b>Number of Training Samples</b>	<b>Accuracy</b>
Rosenblum et al.	<b>20</b>	8-16	77%
This work	<b>100</b>	8	78%
Rosenblum et al.	20	<b>8-16</b>	77%
This work	20	<b>2</b>	78%
This work	20	<b>6</b>	95%
This work	20	<b>8</b>	96%
Rosenblum et al.	100	8-16	<b>61%</b>
This work	100	8	<b>78%</b>
Rosenblum et al.	191	8-16	<b>51%</b>
This work	191	8	<b>63%</b>
This work	<b>600</b>	8	<b>52%</b>



# Compiler Optimization

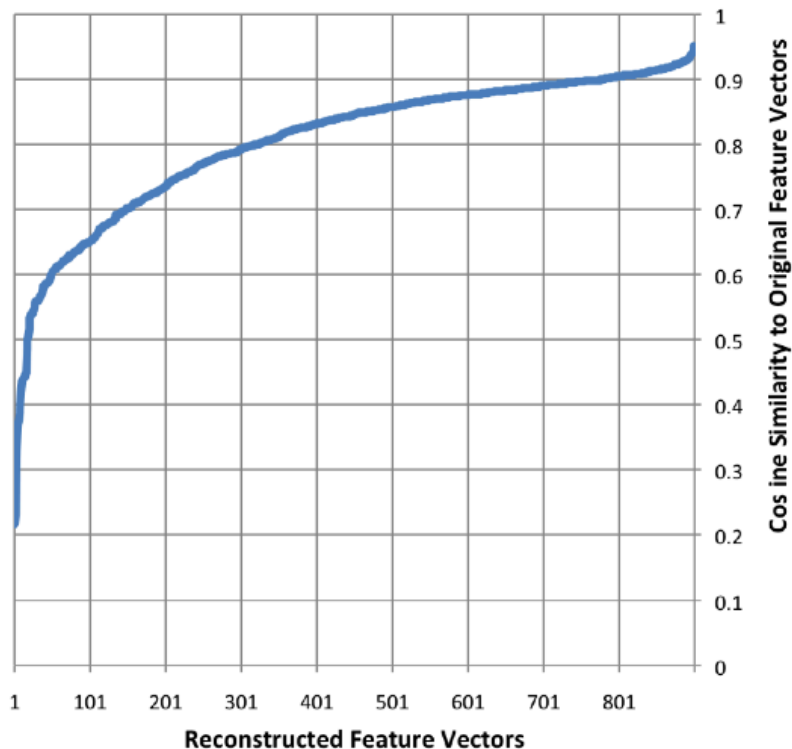
<b>Number of Programmers</b>	<b>Number of Training Samples</b>	<b>Compiler Optimization</b>	<b>Accuracy</b>
100	8	None	78.3%
100	8	Stripped symbols	65.9%
100	8	1	64.2%
100	8	2	61.3%
100	8	3	60.1%

The drop in accuracy is not tragic!



# Reconstructing original features

Cosine Similarity Between  
Original and Reconstructed Feature Vectors

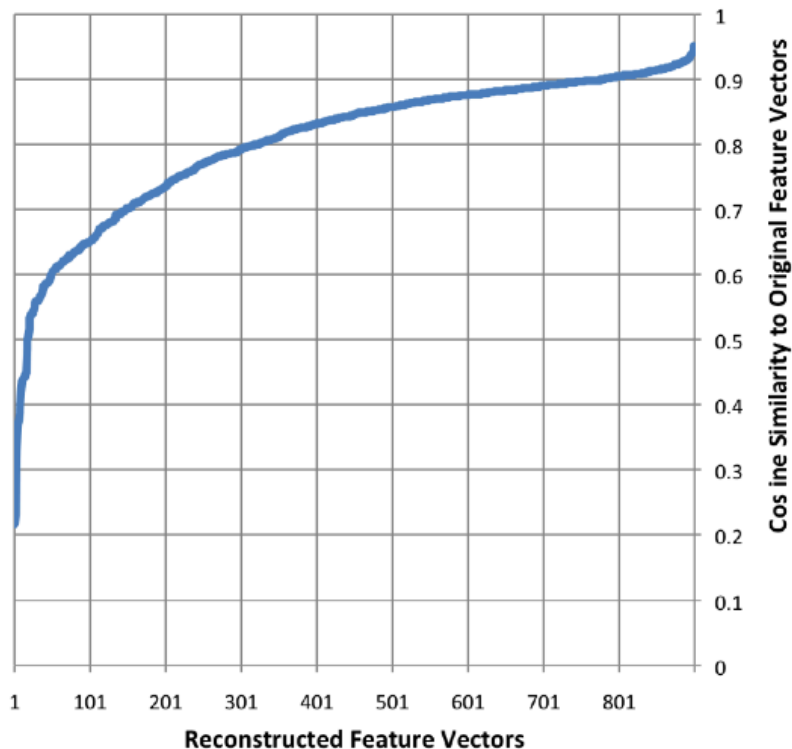


- Original vs predicted features
  - Average cos distance: 0.81
- Original vs decompiled features
  - Average cos distance: 0.35



# Reconstructing original features

Cosine Similarity Between  
Original and Reconstructed Feature Vectors



- Original vs predicted features
    - Average cos distance: 0.81
- ↓
- *This suggests that original features are transformed but not entirely lost in compilation.*





# Insights

More advanced programmers are easier to de-anonymize

<b>A = #authors, F = max #problems completed</b>			
<b>N = #problems included in dataset (<math>N \leq F</math>)</b>			
<b>A = 20</b>			
<b>F = 14</b>	<b>F = 7</b>	<b>F = 12</b>	<b>F = 6</b>
<b>N = 7 easier</b>	<b>N = 7</b>	<b>N = 6 easier</b>	<b>N = 6</b>
<b>Average accuracy after 10 iterations</b>			
<b>88.2%</b>	<b>80.7%<sup>1</sup></b>	<b>86.7%</b>	<b>78.1%<sup>1</sup></b>
<b><sup>1</sup>Drop in accuracy due to programmer skill set.</b>			



# Programmer De-anonymization in the wild

- ✓ Single authored GitHub repositories
- ✓ The repository has at least 500 lines of code

Type	Amount
Authors	49
Repositories	117
Files	782
Repositories / Author	2 – 5
Files / Author	2 – 88

Compile  
repositories

Dataset	Authors	Total Files	Accuracy
GitHub	12	50	62.0%
GCJ	12	50	68.0%



# Future work

- Anonymizing executable binaries
  - optimizations is not the answer
- De-anonymizing collaborative binaries
- Malware family classification



# Available tools

- Programmer de-anonymization
  - <https://github.com/calaylin>
- Jstylo
  - prose authorship attribution framework
- Anonymouth
  - writing anonymization





Dr. Richard Harang



Dr. Konrad Rieck



Dr. Arvind Narayanan

THANKS TO ALL MY CO-AUTHORS



Dr. Clare Voss



Dr. Fabian Yamaguchi



Dr. Rachel Greenstadt

# 934 important features of code stylometry by information gain OUT OF 120,000 FEATURES

Feature Type	Percentage	Count
Word Unigram Frequency	55%	517
AST Node-Bigram Frequency	31%	291
AST Node Average Depth	5%	48
AST Node Frequency	4%	38
AST Node TFIDF	2%	19
C++ Keywords	2%	15
Layout Features	1%	6

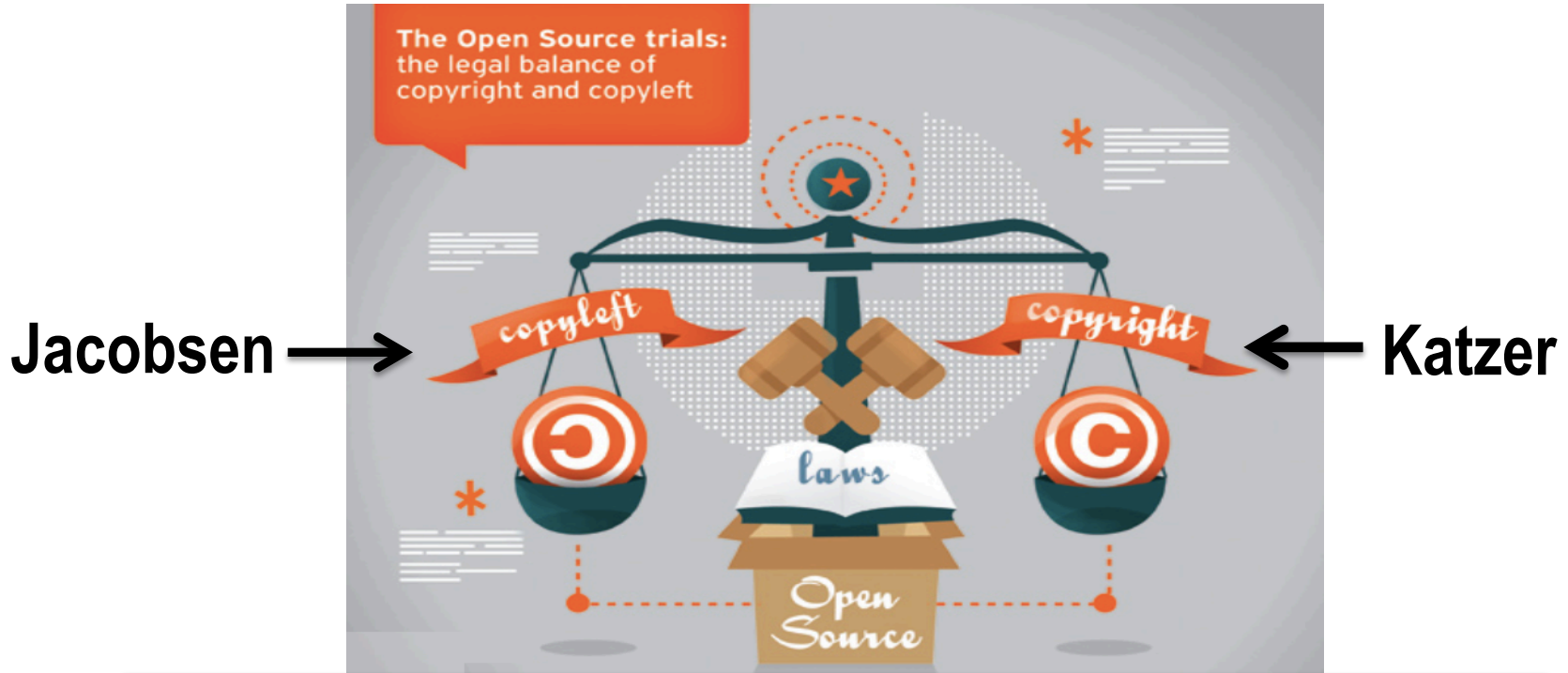


# EXTRA Case: Copyright investigation

- Copyleft programs are free but licensed
- Did this programmer take a copyleft code and distribute it commercially?
  - *Jacobsen vs Katzer (Java Model Railroad Interface)*
- Two-class machine learning classification task
  - Class 1: the copyleft programmer
  - Class 2: the copyright programmer
  - Test: both the copyleft and copyright code



# EXTRA Case: Copyright investigation



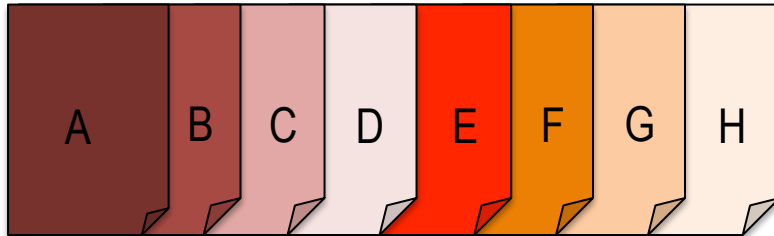
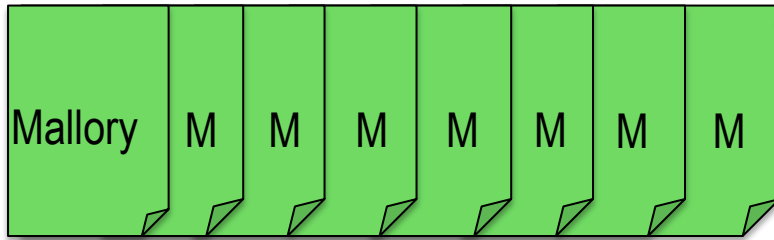
60 pairs of authors each with 9 program files    Classification Accuracy

Two-class task    100%

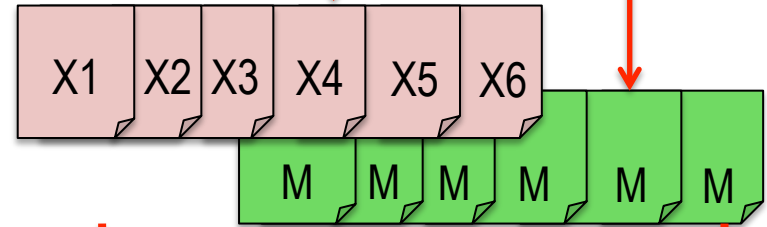


# EXTRA Case: Authorship verification

- Is this source code really written by this programmer?



Train on 8 files from Mallory and one file from authors A, B, C, D, E, F, G, and H.



Test on 6 files that belong to Mallory and 6 files that belong 6 random authors.

91% accuracy in 80 sets of experiments



# Extra Case: Difficult tasks & advanced coders

- Insights about programmers and coding style:
  - Implementing harder functionality makes programming style more unique

Same set of 62 authors	Classification Accuracy
Solving 7 easy problems	98%
Solving 7 more difficult problems	99%



# Extra Case: Difficult tasks & advanced coders

- Insights about programmers and coding style.
  - Better programmers have more distinct coding style

Two sets of 62 authors	Classification Accuracy
Less advanced programmers	97%
More advanced programmers	98%

