



Statistical De-obfuscation for Android

Petar Tsankov, ETH Zurich



Benjamin
Bichsel



Veselin
Raychev



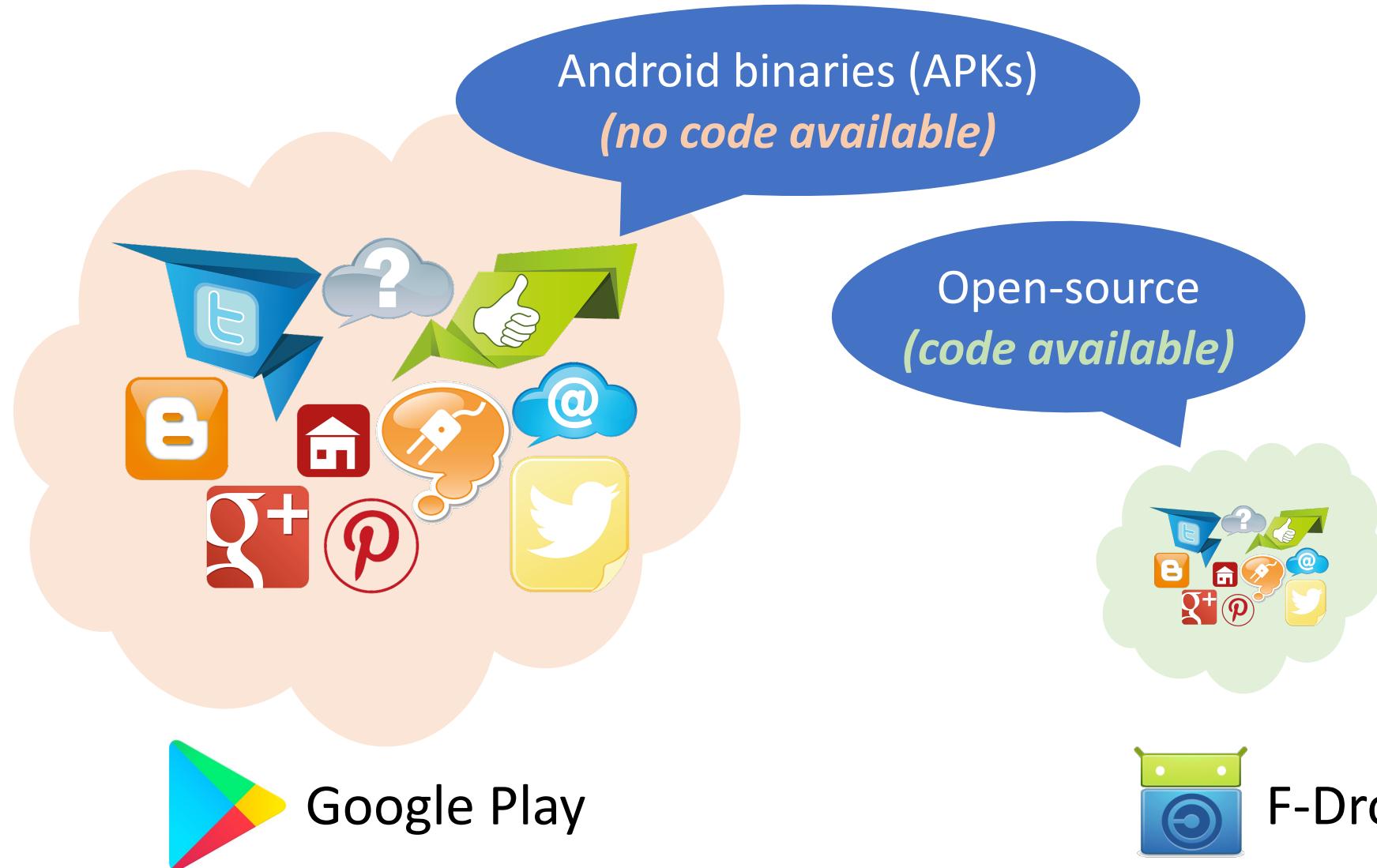
Petar
Tsankov



Martin
Vechev

DeGuard
Team

Why De-obfuscate Android Applications?



Why De-obfuscate Android Applications?



Google Play

2.6M
APKs

Which ones use
vulnerable libraries?

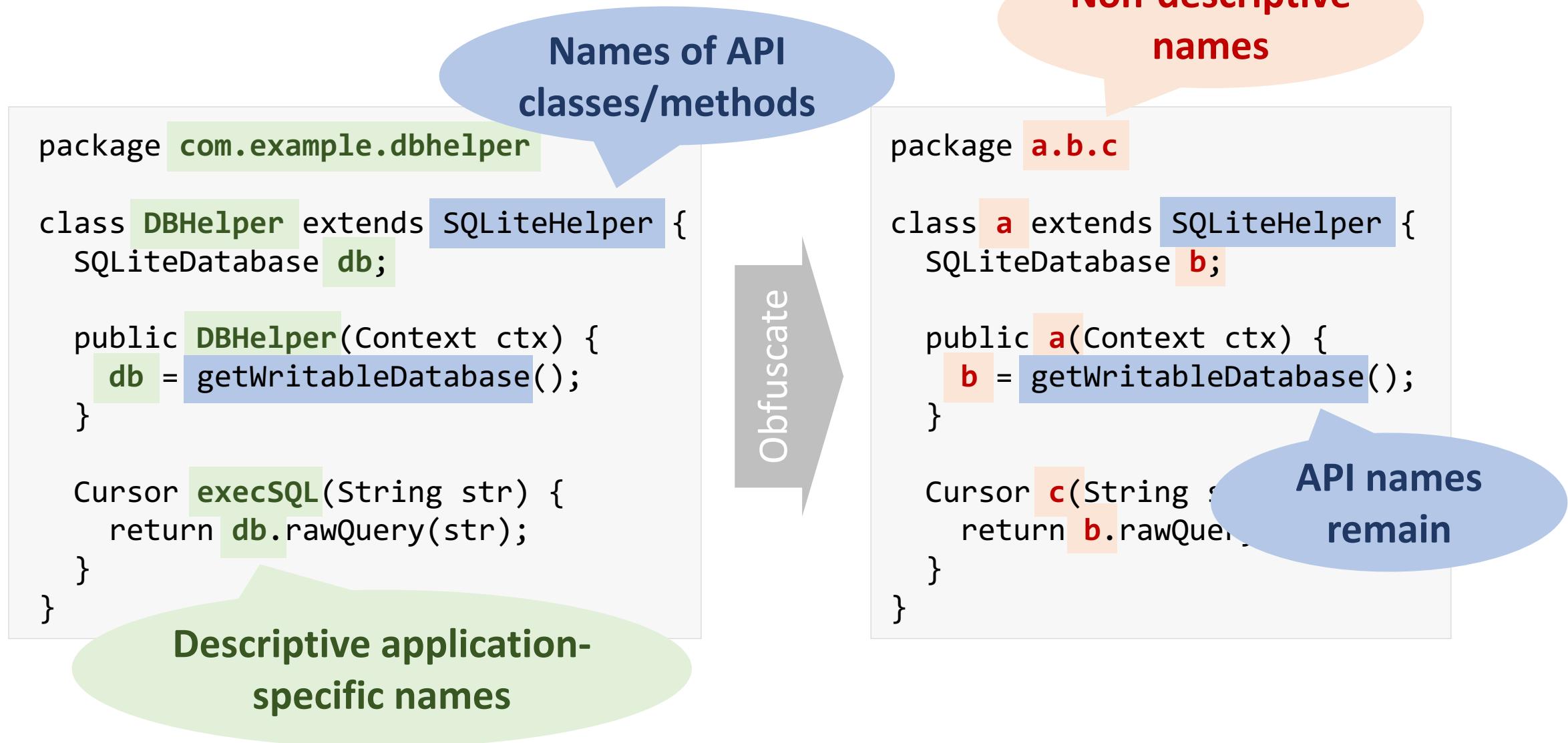


5K
APKs



F-Droid

Layout Obfuscation in Android



Layout Obfuscation in Android

```
package com.example;

class DBHelper extends SQLiteOpenHelper {
    public DBHelper(Context context) {
        super(context, "mydatabase", null, 1);
    }

    Cursor execSQL(String sql) {
        return db.rawQuery(sql, null);
    }
}
```

Names of API

Non-descriptive names



Security Challenges



Code inspection



Third-party library detection

... many others

Descriptive application-specific names

```
teHelper {
    tx) {
        database();
    }
}
```

API names remain

Layout Obfuscation in Android

```
package com.example.dbhelper  
  
class DBHelper extends SQLiteHelper {  
    SQLiteDatabase db;  
  
    public DBHelper(Context context) {  
        db = getWritableDatabase();  
    }  
  
    Cursor execSQL(String str) {  
        return db.rawQuery(str);  
    }  
}
```

Names of API
classes/methods

Descriptive application-
specific names

```
package a.b.c  
  
class a extends SQLiteHelper {  
    SQLiteDatabase b;  
  
    Context ctx) {  
        b = getWritableDatabase();  
    }  
  
    Cursor c(String str) {  
        return b.rawQuery(str);  
    }  
}
```

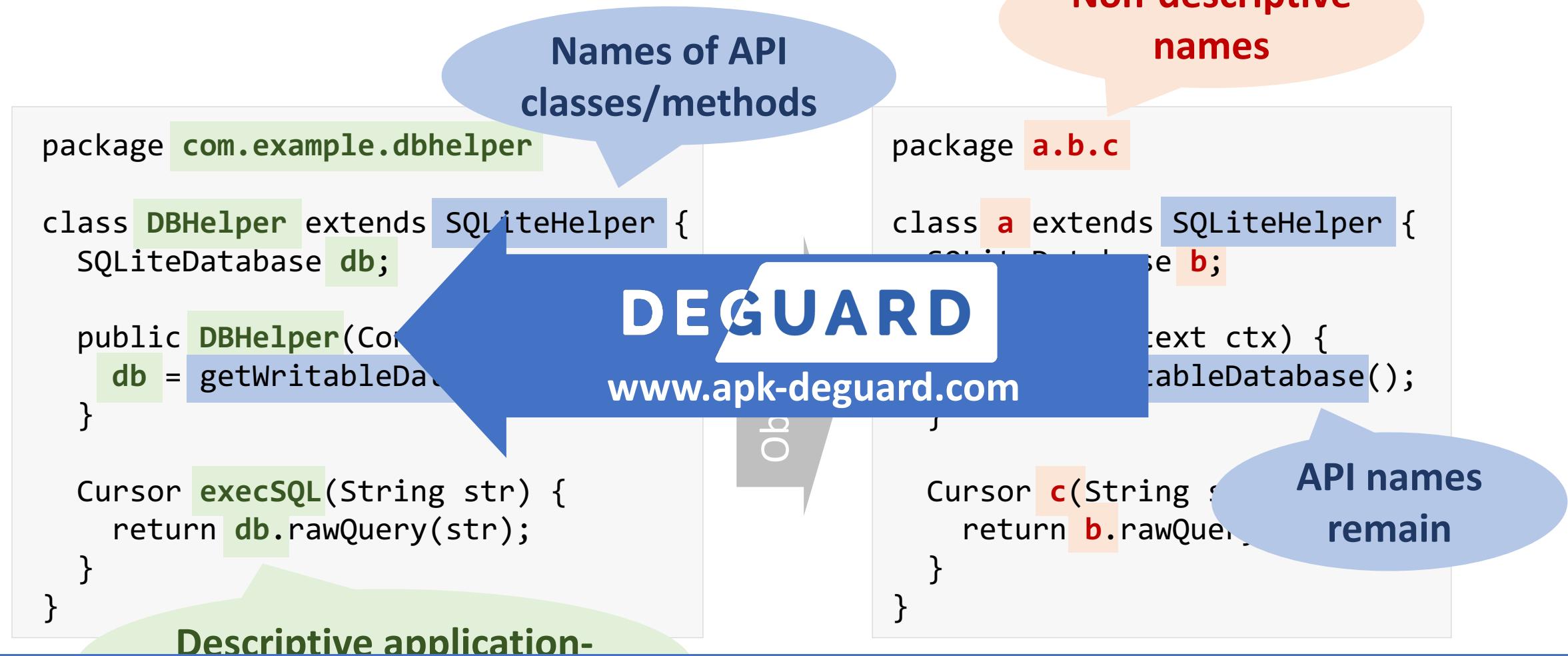
Non-descriptive
names

Can we reverse
layout obfuscation



API names
remain

Layout Obfuscation in Android



Yes, with roughly 80% accuracy!

Demo

Released in October 2016, so far: > 100GB distinct APKs de-obfuscated

Reddit posts/comments

evantarka WillowTree • 3 points 2 days ago

Nice! This should help with debugging issues in play services or other libs that are obfuscated just to make my life harder a bit easier.

[permalink](#) [embed](#) [pocket](#)

oleeEncantado 2 points 4 days ago

Works quite well, I've tried on some small games.

[permalink](#) [embed](#) [pocket](#)

Tycon712 • 3 points 2 days ago

Can someone tell me what the point of using Proguard is if there are tools out there like this?

[permalink](#) [embed](#) [pocket](#)

theheartbreakpug • 6 points 2 days ago

As far as I know, this is brand new. I asked the creator of ProGuard a week ago how hard it is to unobfuscate code after it's run through proguard. He said it strips all the names out of the code so it's essentially impossible. I'm super impressed by what they've done here.

•
•
•



Tweets

DJ-AR\$-IN @dharshin • Oct 17

Deobfuscate #proguard 'ed APKs: apk-deguard.com. The paper on the inner workings: srl.inf.ethz.ch/papers/deguard... #Android #MobileSecurity

[link](#) [retweet](#) [comment](#) [like](#) [more](#)



Brian Carpenter @geeknik • 9h

Android Deobfuscation with Machine Learning reverses the effects of ProGuard [PDF] srl.inf.ethz.ch/papers/deguard...

[link](#) [retweet](#) 3 [comment](#) [like](#) 4 [more](#)



rvivek @vivek_310 • Oct 17

"Android Deobfuscation with Machine Learning reverses the effects of ProGuard" #security #feedly



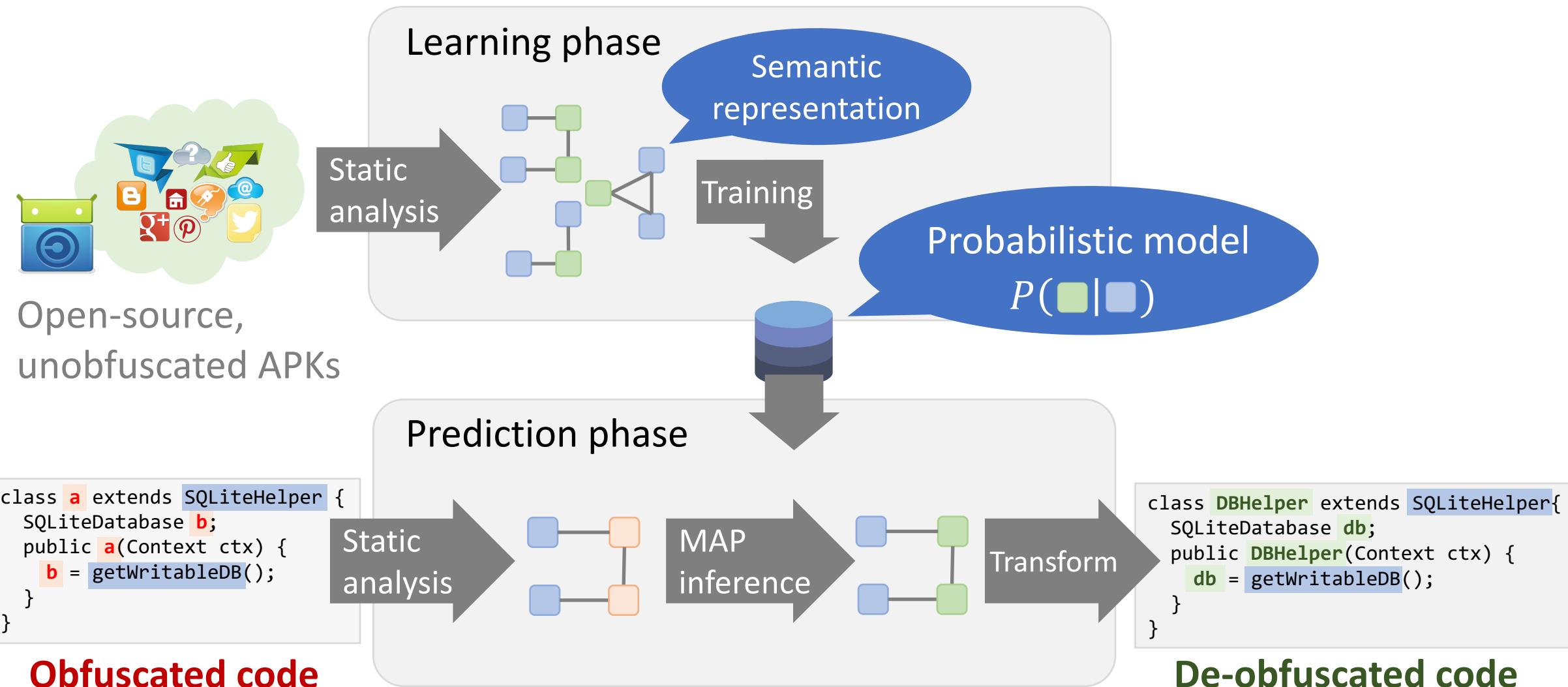
Android Deobfuscation with Machine... • /r/Reverse...

8 points and 1 comments so far on reddit

•
•
•

How Does DeGuard Work?

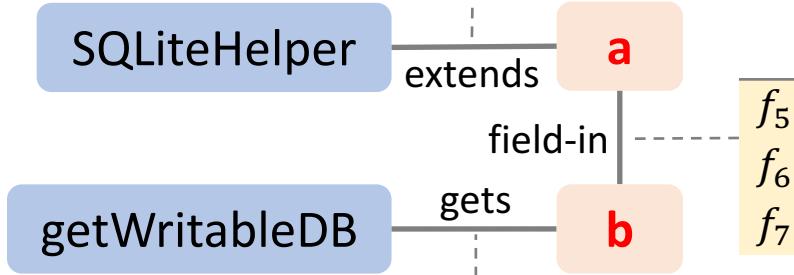
DeGuard: System Overview



Probabilistic Graphical Models

Probabilistic Graphical Models

	name1	name2	weight
f_1	SQLiteHelper	DBUtils	0.3
f_2	SQLiteHelper	DBHelper	0.2



	name1	name2	weight
f_3	getWritableDatabase	db	0.7
f_4	getWritableDatabase	instance	0.4

SQLiteHelper, getWritableDatabase Known variables

a, b Unknown variables

f_1, f_2, \dots, f_7 Feature functions

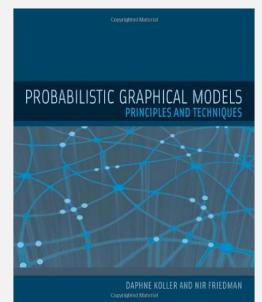
	name1	name2	weight
f_5	DBUtils	instance	0.5
f_6	DBHelper	db	0.4
f_7

```
class a extends SQLiteHelper {  
    SQLiteDatabase b;  
    public a(Context ctx) {  
        b = getWritableDatabase();  
    }  
}
```

Graph + features define a **probabilistic graphical model**

$$P(a, b | \text{SQLiteHelper}, \text{getWritableDatabase})$$

$$= \frac{1}{Z} \exp(0.3 \cdot f_1(\text{SQLiteHelper}, a) + 0.2 \cdot f_2(\text{SQLiteHelper}, a) + \dots)$$



Probabilistic Graphical Models

	name1	name2	weight
f_1	SQLiteHelper	DBUtils	0.3
f_2	SQLiteHelper	DBHelper	0.2

SQLiteHelper

a

extends

getWritableDatabase

b

field-in

gets

	name1	name2	
f_3	getWritableDatabase	db	
f_4	getWritableDatabase	instan	

SQLiteHelper, getWritableDatabase Known variables

a, b

Unknown variables

f_1, f_2, \dots, f_7 Feature functions

Next

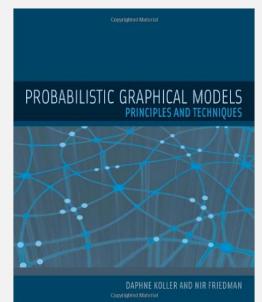
How are the features and their weights learned?

```
class a extends SQLiteHelper {  
    SQLiteDatabase b;  
    public a(Context ctx) {  
        b = getWritableDatabase();  
    }  
}
```

a probabilistic graphical model

$$P(a, b | \text{SQLiteHelper}, \text{getWritableDatabase})$$

$$= \frac{1}{Z} \exp(0.3 \cdot f_1(\text{SQLiteHelper}, a) + 0.2 \cdot f_2(\text{SQLiteHelper}, a) + \dots)$$



Learning

Learning



Unobfuscated
APKs



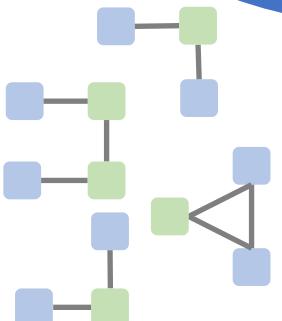
Feature
templates

28 templates

Static
analysis

	name1	name2
f_1	SQLiteHelper	DBUtils
f_2	SQLiteHelper	DBHelper
f_3	getWritableDatabase	db
f_4	getWritableDatabase	instance
f_5	DBUtils	instance
f_6	DBHelper	db
f_7

Features (with
candidate names)



Actual graphs have
> 1,000 nodes

Dependency graphs

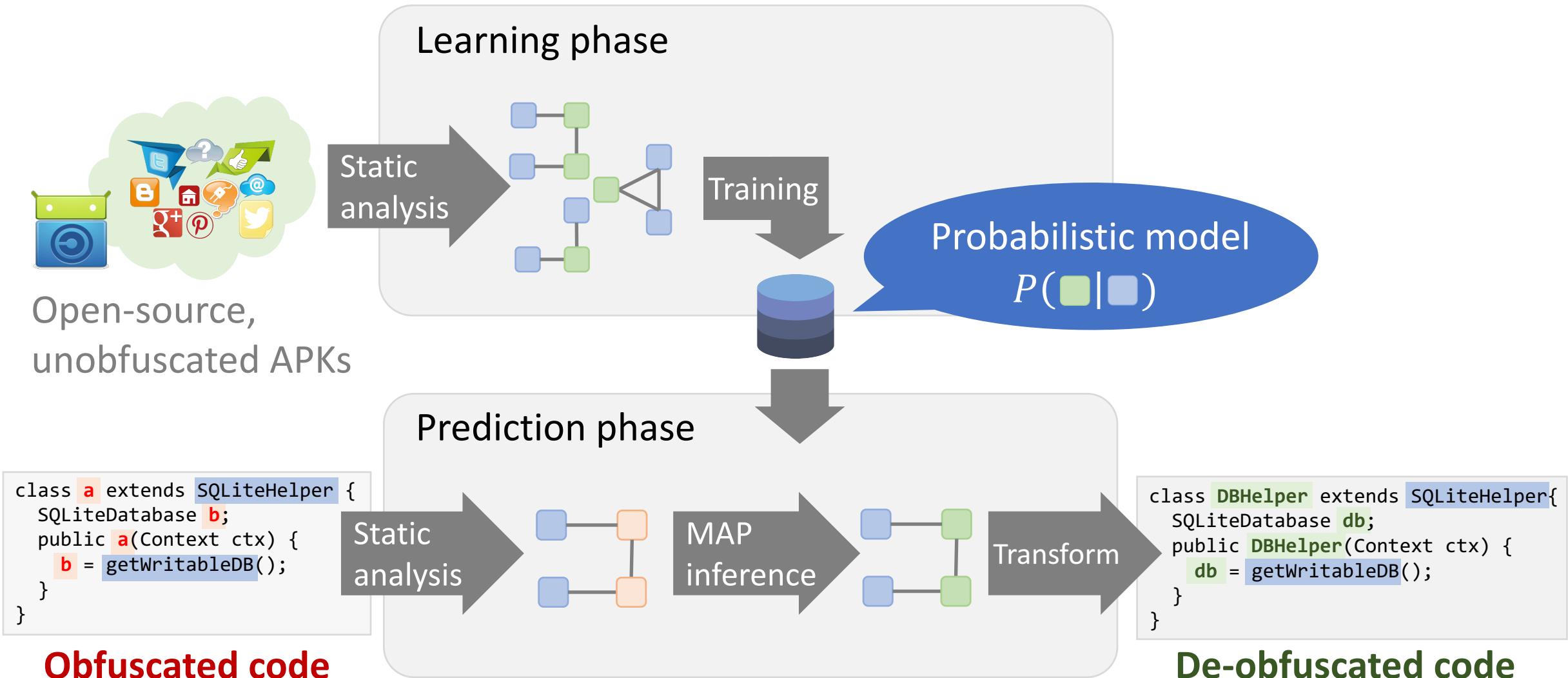
Train
model

	name1	name2	weight
f_1	SQLiteHelper	DBUtils	0.3
f_2	SQLiteHelper	DBHelper	0.2
f_3	getWritableDatabase	db	0.7
f_4	getWritableDatabase	instance	0.4
f_5	DBUtils	instance	0.5
f_6	DBHelper	db	0.4
f_7

> 100,000

Compute **weights** that
maximize $P(\vec{O} = \vec{o}_i | \vec{K} = \vec{k}_i)$ for
all training samples (\vec{o}_i, \vec{k}_i)

DeGuard: System Overview

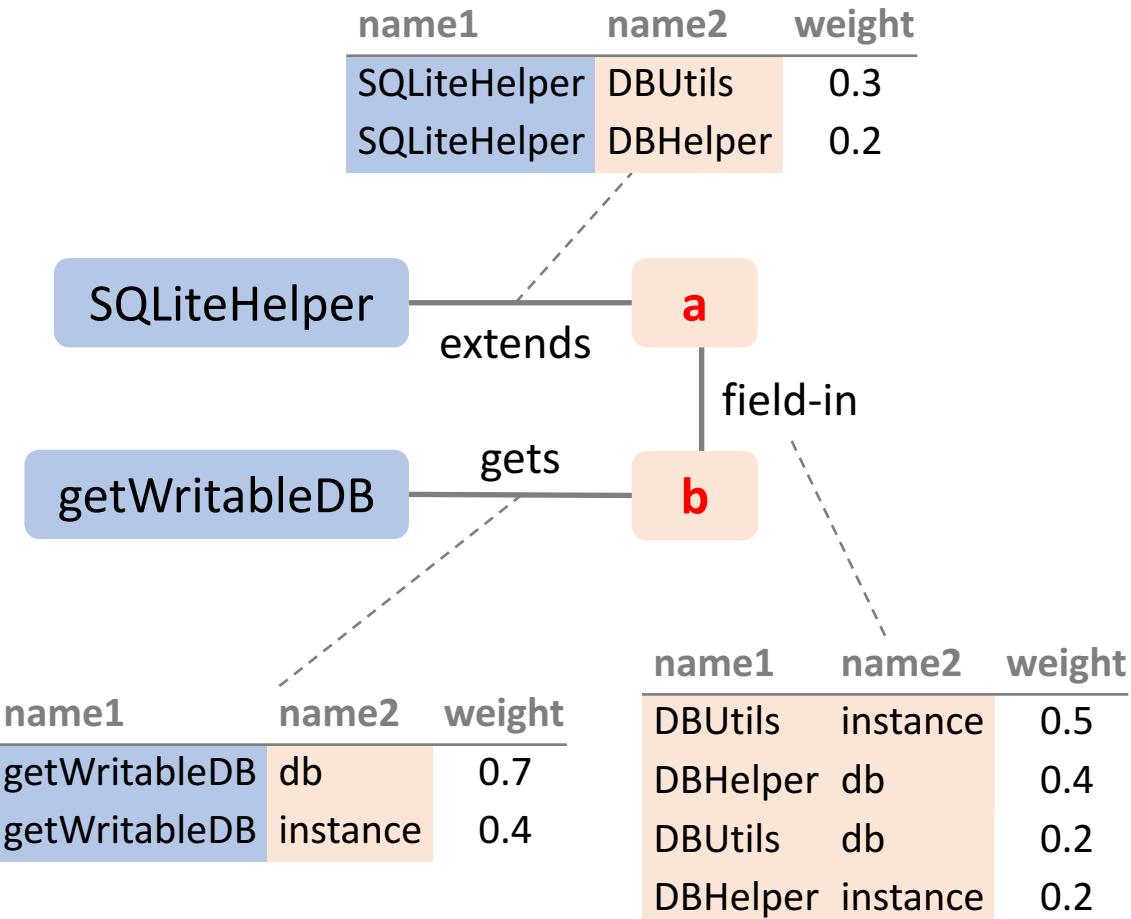


Prediction Phase

```
class a extends SQLiteHelper {  
    SQLiteDatabase b;  
    public a(Context ctx) {  
        b = getWritableDatabase();  
    }  
}
```

Obfuscated Code

Static analysis



Prediction Phase

```
class a extends SQLiteOpenHelper  
    SQLiteDatabase b;  
    public a(Context ctx  
        b = getWritableDatabase()  
    }  
}
```

Obfuscation

MAP Inference

$$\vec{o} = \underset{\vec{o}' \in \Omega}{\operatorname{argmax}} P(\vec{o} = \vec{o}' | \vec{k} = \vec{k})$$

Candidate assignment \vec{o} $P(\vec{o} | \vec{k})^*$

a = DBUtils b = instance 1.2

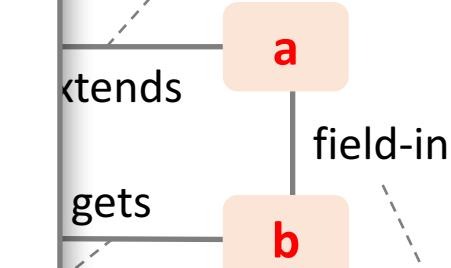
a = DBHelper b = db 1.3

a = DBUtils b = db 0.8

a = DBHelper b = instance 1.2

*Non-normalized

1	name2	weight
Helper	DBUtils	0.3
Helper	DBHelper	0.2



name1	name2	weight
DBUtils	instance	0.5
DBHelper	db	0.4
DBUtils	db	0.2
DBHelper	instance	0.2

Prediction Phase

```
class a extends SQLiteOpenHelper  
    SQLiteDatabase b;  
    public a(Context ctx  
        b = getWritableDatabase()  
    }  
}
```

Obfuscation

MAP Inference

$$\vec{o} = \underset{\vec{o}' \in \Omega}{\operatorname{argmax}} P(\vec{o} = \vec{o}' | \vec{k} = \vec{k})$$

Candidate assignment \vec{o} $P(\vec{o} | \vec{k})^*$

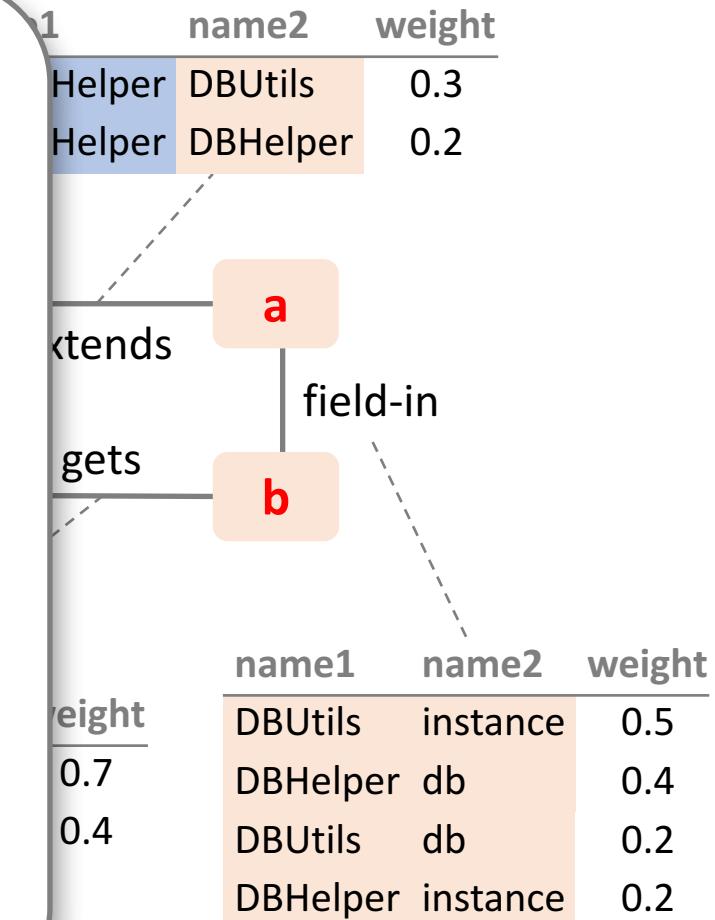
a = DBUtils b = instance 1.2

a = DBHelper b = db 1.3

a = DBUtils b = db 0.8

a = DBHelper b = instance 1.2

*Non-normalized



Prediction Phase

```
class a extends SQLiteHelper {  
    SQLiteDatabase b;  
    public a(Context ctx) {  
        b = getWritableDatabase();  
    }  
}
```

Obfuscated Code

```
class DBHelper extends SQLiteHelper {  
    SQLiteDatabase db;  
    public DBHelper(Context ctx) {  
        db = getWritableDatabase();  
    }  
}
```

Deobfuscated Code

Static analysis

Semantically
the same?

Transform

name1	name2	weight
SQLiteHelper	DBUtils	0.3
SQLiteHelper	DBHelper	0.2

SQLiteHelper

DBHelper

ritableDB

db

extends

field-in

gets

name1	name2	weight
getWritableDatabase	db	0.7
getWritableDatabase	instance	0.4

name1	name2	weight
DBUtils	instance	0.5
DBHelper	db	0.4
DBUtils	db	0.2
DBHelper	instance	0.2

Semantics-Preserving De-obfuscation Constraints



Freely renaming fields/variables/methods
may **change** the application's **semantics**

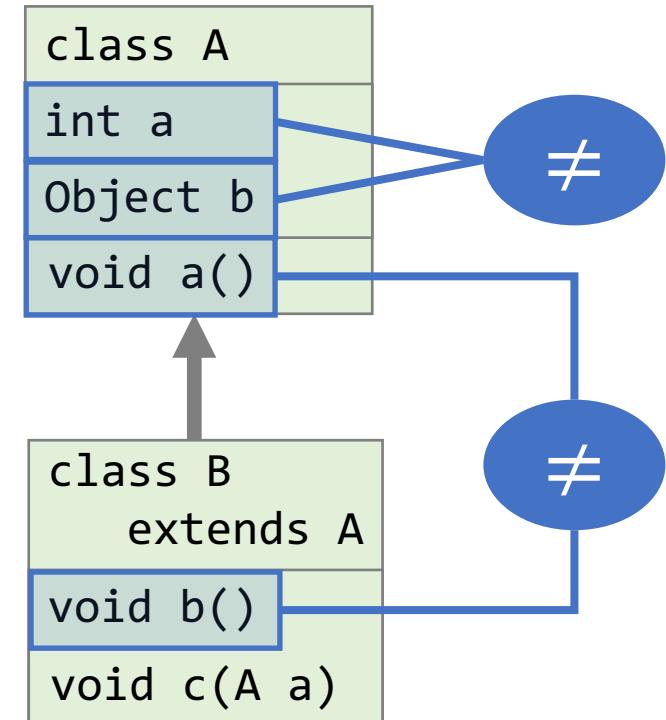
</> Syntactic constraints

e.g. “*Fields within a class must have distinct names*”

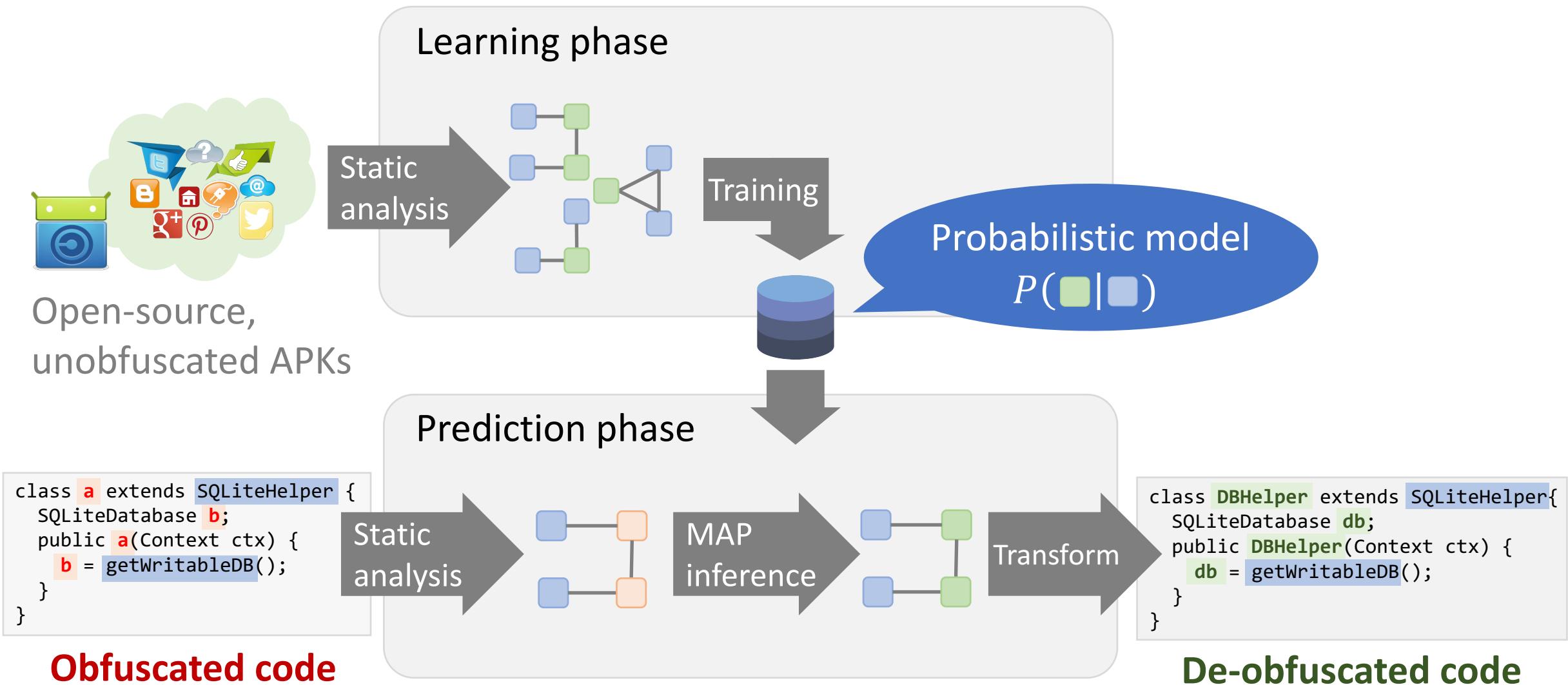


Semantic constraints

e.g. “*Method overloads must be preserved*”



DeGuard: System Overview



DeGuard Implementation

DeGuard Implementation

Static Analysis



- Static analysis framework for Java and Android

Learning and MAP Inference



- Scalable open-source framework for structured prediction
- Open-source:
<http://nice2predict.org>
- Training data: 2K open-source, unobfuscated Android applications



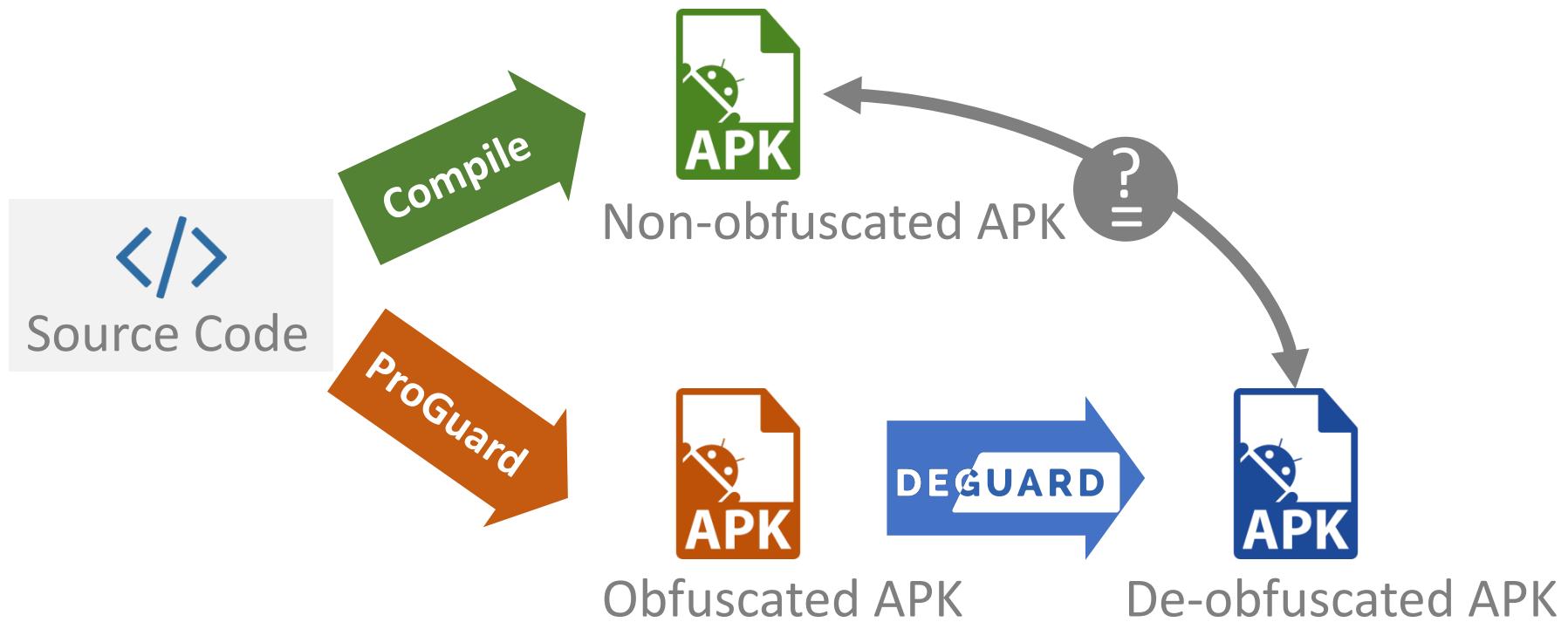
A screenshot of a web browser displaying the DeGuard website at www.apk-deguard.com. The page has a dark blue header with the DeGuard logo and social media links. Below the header, there's a section titled "Statistical Deobfuscation for Android" with a brief description. At the bottom, there's a form with "Select APK File" and "Upload" buttons, along with links for "Try on Sample APK" and "or see a demo in action".

www.apk-deguard.com

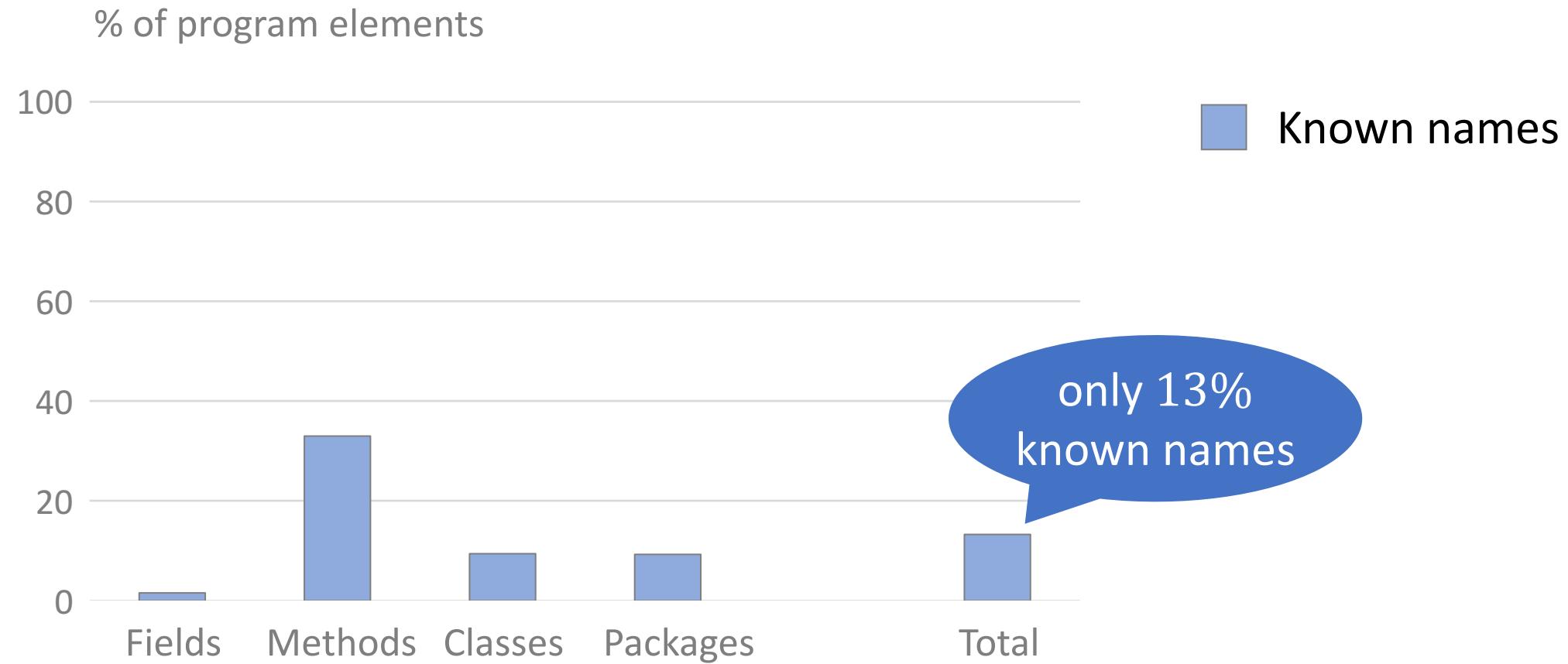
Evaluation

1. Can DeGuard reverse ProGuard?
2. Can DeGuard detect third-party libraries?
3. Is DeGuard useful for malware inspection?

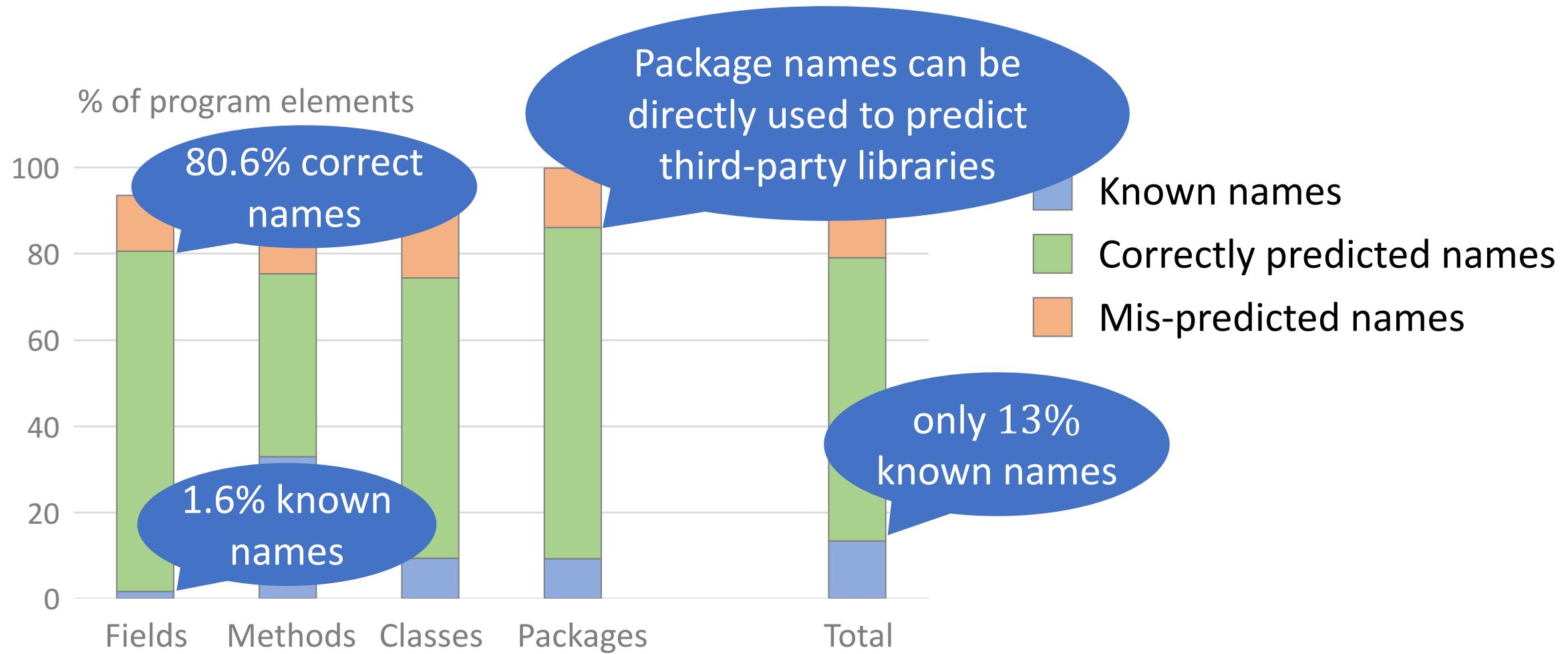
ProGuard Experiment



After Obfuscation

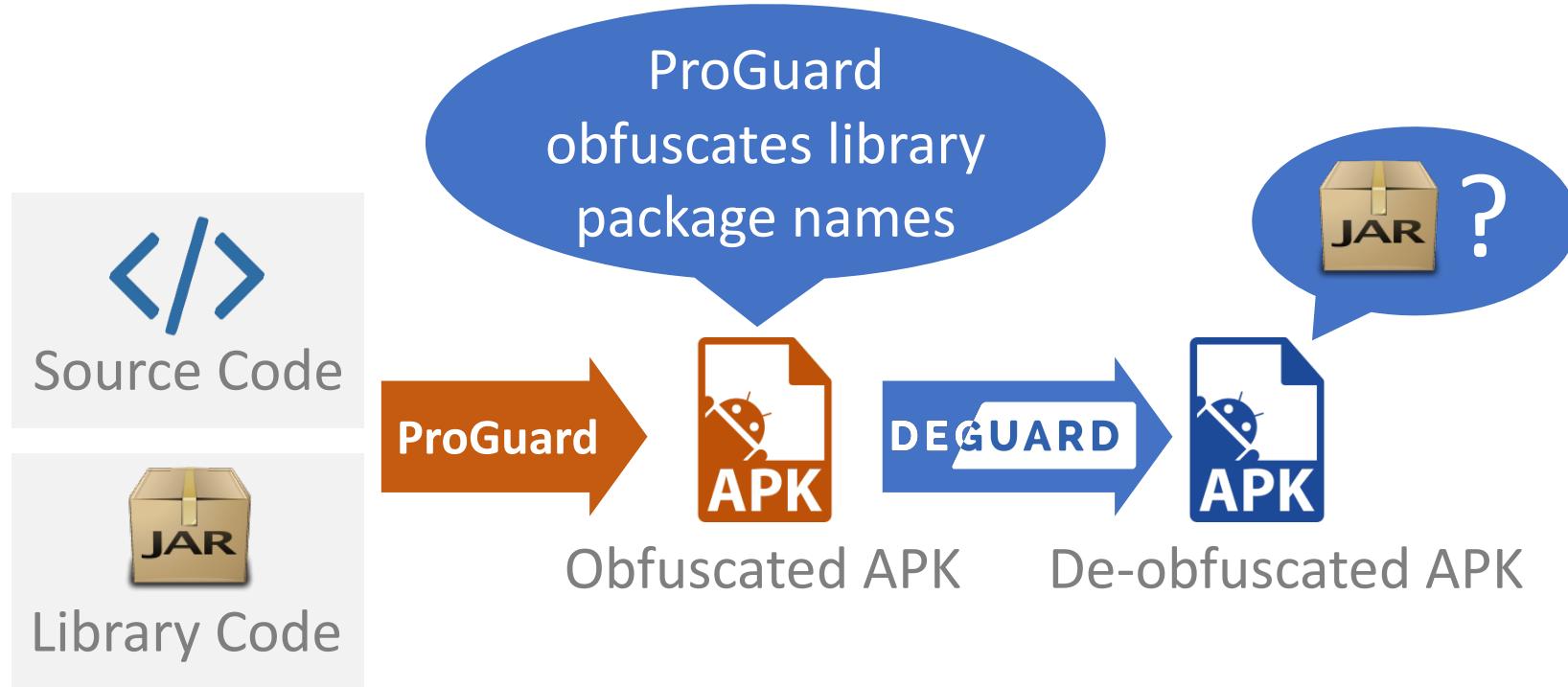


Can DeGuard Reverse ProGuard?



80% of the names are identical to the original ones

Can DeGuard Detect Third-Party Libraries?



Precision: 93.1%

Recall: 91%

Is DeGuard Useful for Malware Inspection?



We de-obfuscated all samples from
the Android Malware Genome Project

```
class d {  
    String a = System.getProperty(..)  
    char[] b;  
    byte [] c;  
    byte[] a(String) {...}  
}
```

Malware Sample



```
class Base64 {  
    String NL = S...  
    char[] ENC;  
    byte [] DEC;  
    byte[] decode(String) {...}  
}
```

Base64
Decoder

De-obfuscated Malware Sample



Reveals string
decoders



Reveals classes that handle
sensitive data (e.g. Location)

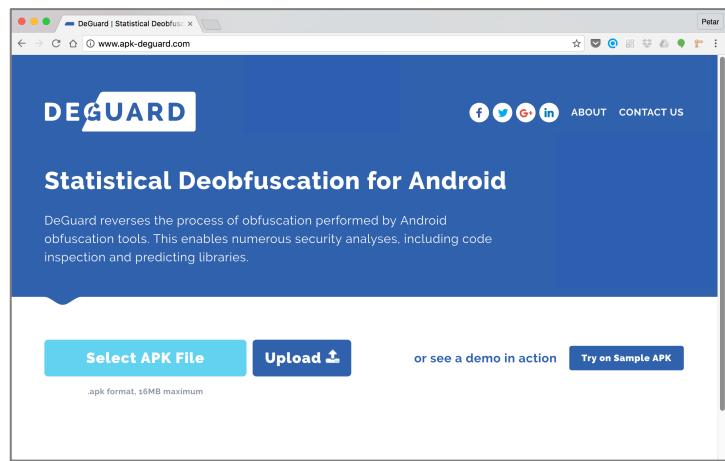
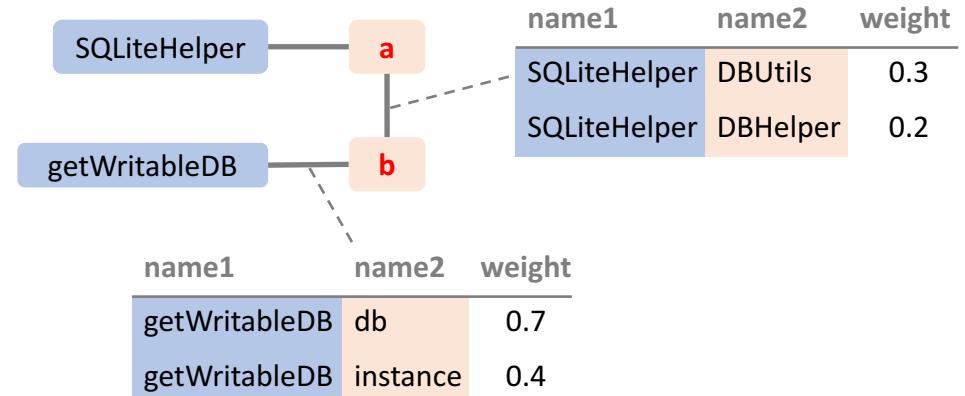


Hard to handle heavily-obfuscated
code (e.g. reflection)

Summary

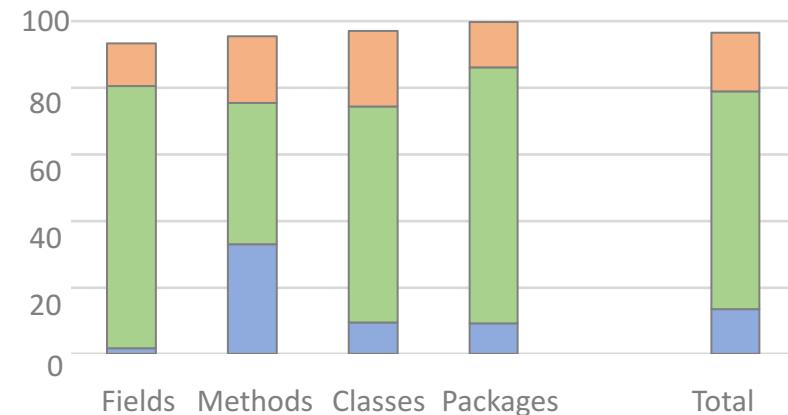
The screenshot shows two panels of Java code. The left panel is the original code with some identifiers highlighted in green (e.g., package names, class names, method names). The right panel is the deobfuscated code where these identifiers have been replaced by placeholder names in orange (e.g., 'a.b.c', 'a', 'b', 'c'). A large blue arrow points from the right panel back to the left, containing the word 'DEGUARD' in white.

```
package com.example.dbhelper  
class DBHelper extends SQLiteHelper {  
    SQLiteDatabase db;  
    public DBHelper(Context ctx) {  
        super(ctx, "mydb", null, 1);  
        db = getWritableDatabase();  
    }  
    Cursor execSQL(String str) {  
        return db.rawQuery(str);  
    }  
}  
  
package a.b.c  
class a extends SQLiteHelper {  
    SQLiteDatabase b;  
    public a(Context ctx) {  
        super(ctx, "mydb", null, 1);  
        b = getWritableDatabase();  
    }  
    Cursor c(String str) {  
        return b.rawQuery(str);  
    }  
}
```



Try online: www.apk-deguard.com

Probabilistic Models



High Prediction Accuracy

For more info: <http://plml.ethz.ch> / <http://srl.inf.ethz.ch>