# A security analyst's guide to selecting training data

Luke Dramko



### DIRE

int v1  $\longrightarrow$  int size



### DIRE int v1 $\longrightarrow$ int size

### DIRTY \_QWORD \* → struct point \*





### DIRTY \_QWORD \* → struct point \*





- DIRTY \_QWORD \* → struct point \*







- DIRTY \_QWORD \* → struct point \*



Idioms (coming soon!) \*(\_QWORD \*) (a1 + 4120) ---- ctrl->used\_size



# We studied how the **characteristics of training data** influence these machine learning models.



# We use DIRE in our case study.

```
void *file_mmap(int v1, int v2) {
    void * v3;
   v3 = mmap(0, v2, 1, 2, v1, 0);
   if (v3 == (void *) -1) {
        perror("mmap")
        exit(1)
    return v3;
}
             v1
             v2
             v3
```



# DIRE is a machine-learning based tool which automatically suggests variable names in decompiled code.

```
void *file_mmap(int v1, int v2) {
    void * v3;
    v3 = mmap(0, v2, 1, 2, v1, 0);
    if (v3 == (void *) -1) {
        perror("mmap")
        exit(1)
    return v3;
             v1
             v2
             v3
```



# DIRE is trained on large amounts of code from GitHub.





# Which **C** programs do we select?



















We looked at ways to select and augment data to improve performance, especially on rare and underrepresented variable names.

Domains

Helping DIRE perform better on domain-specific identifiers.



http\_response
socket +acc
...



lexeme ast

+acc

Dataset Quality

Carefully selecting training data to improve generalizability.









We looked at ways to select and augment data to improve performance, especially on rare and underrepresented variable names.

Dataset Quality

Carefully selecting training data to improve generalizability.





# Subproblem - Given a binary, how can we identify its domain?

🖵 gpac/(	gpac						
<> Code	<ul> <li>Issues 12</li> </ul>	រៀ Pull requests 2	Actions	Projects	🕮 Wiki	(!) Security	
	ಳ	양 master → 양 16 branches ♡ 11 tags					
<b>jeanlf</b> fixed potential poc compute bug in vvc				~	/ e4		
		.github		try "Issue and Pu	Ill Request te	emplates"	
		applications		misc help and al	ias improvem	ents	
		build		port qjs-libc to n	nsvc		
		debian		Install app icon i	n XDG hicolo	r icon theme	
		extra_lib/include		enable http2 on	windows		

Hand-selected 7 domains from among the GitHub tags in our dataset

# We derived domain labels from GitHub tags.

	Q Notifications	Selected Doma
Go to file Code -	About	kernel
9352 yesterday 🕚 14,855 commits	Modular Multimedia framework for packaging, streaming and playing your favorite content.	game
6 years ago yesterday 2 months ago	gpac.io streaming mp4 hl graphics vr tiling mpag-dash second	shell
2 months ago 5 months ago	mov mp4box prores mpe -ts gpac atsc3 cenc	network
		compiler
e most comma	SU	embedded
		graphics













•













### Strategy: Remove parts of the code





Strategy: Remove parts of the code Create Classifier





Strategy: Remove parts of the code Create Classifier Evaluate





Strategy: Remove parts of the code Create Classifier Evaluate Repeat





Strategy: Remove parts of the code Create Classifier Evaluate Repeat Observe which parts of code are most predictive of software domain.



# In aggregate, string literals are highly indicative of software domain.

- The approach that led to the best performance was to extract the string literals run the classifier only on the string literals.
- Test Accuracy: 83.54%



Dimensionality-reduced visualization of embeddings generated (originally 64 dimensions)



We now know that string literals inform domain.

We can use this information to make DIRE predict domain-specific identifiers more accurately

# We modified the DIRE model to use string literals directly instead of placeholders.

### **Orignal DIRE Model**

result = gf\_log("String");

### **Domain-Aware DIRE Model**

result = gf\_log("[Compositor] Failed to allocate svg gradient stack\n");



We borrowed categorical labeling techniques from machine translation to help inform our model of the domain of the data.

### **Machine Translation**

<english> I went to the store.

<japanese>私は店に行った。

Labeled data is rare — Train mostly on data with no known domain (label <unknown>)

### **DIRE's Variable Renamings**

<graphics> \_\_int64 \_\_fastcall compositor\_init\_s...

<compiler> unsigned \_\_int64 \* \_\_fastcall endian...



# Each modification applied to DIRE's data serves a different purpose.



Provides general, *global* information about software domain to each function in the binary (even those without string literals).

Provides more specific information about a particular function's role and/or the role of certain variables in that function.



# With these modifications, DIRE performs better across all domains.

### Function Body not in Train Accuracy

Domain	Control (Original)	With Modifications	Accuracy Increase
unknown	35.89%	43.85%	7.9
kernel	39.99%	46.07%	6.0
game	22.40%	32.51%	10.1
shell	19.63%	24.49%	4.8
compiler	22.97%	27.42%	4.4
network	20.13%	25.60%	5.4
embedded	46.19%	54.33%	<b>8.</b> I
graphics	18.58%	26.67%	8.0
Mean (Domains)	27.13%	33.87%	6.74



# We improved DIRE to predict domain-specific variable names more accurately.



ML models can struggle on rare or underrepresented concepts like domain specific variable names.

<graphics> \_\_int64 \_\_fastcall compositor\_init\_s...

<compiler> unsigned \_\_int64 \*\_\_fastcall endian...

We modified DIRE's data to help it perform better on The helps DIRE consistently outperform the original domain-specific variable names.



In aggregate, string literals are highly indicative of a given binary's software domain.





# Bonus Slides - Methodology

# We split each dataset into train and test sets.



Domain-specific Data (kernel, network, etc.)



# We train a control model and a model with interventions applied.





# We test each model against each test set.

for each  $\{O\}$  in  $\{O, V\}$ 







# Bonus Slides - Training Set Quality

How do we select "good" training data for DIRE and machine learning models of code in general?



### Variable names are labels for concepts present in code

```
listNode *
listSearchKey(list *list, void *key) {
    listIter iter;
    listNode *node;
    listRewind(list, &iter);
    while((node = listNext(&iter)) != NULL) {
        if (list->match) {
            if (list->match(node->value, key)) {
                return node;
        } else {
            if (key == node->value) {
                return node;
        }
    return NULL;
}
```

### Entities: list, key, iter, etc.

```
ngx_buf_t *
ngx_create_temp_buf(ngx_pool_t *pool, size_t size) {
    ngx_buf_t *b;
    b = ngx_calloc_buf(pool);
    if (b == NULL) {
        return NULL;
    }
    b->start = ngx_palloc(pool, size);
    if (b->start == NULL) {
        return NULL;
    b->pos = b->start;
    b->last = b->start;
    b->end = b->last + size;
    b->temporary = 1;
    return b;
}
```

### Entities: pool, size, start, etc.





Names are labels for the entities present in code



### More diverse labels

Names are labels for the entities present in code



Names are labels for the entities present in code

More diverse labels — More types of code are covered



**Names** are **labels** for the entities present in code

More diverse labels — More types of code are covered







fd fd fd len fd len fd len

### fd len str name

Name

### fd len fd len str name fd len str name Name



fd fd fd len fd len

Entropy 1.49 bits

### fd len str name

Name

### fd len fd len str name fd len str name Name



fd fd fd len fd len fd len

Entropy 1.49 bits

## fd len str name

Name

### fd len fd len str name fd len str name Name

Entropy





fd fd fd len fd len fd len

Entropy 1.49 bits

# fd len str name

Name

Entropy



### fd len fd len str name fd len str name Name

Entropy

1.97<sub>bits</sub>

# To test our hypothesis, we tried training a model on code that has a high number of stars \_\_\_\_\_ on GitHub

- We suspected high-star repositories may be more diverse
- High-star repositories are much less likely to be duplicates

Star 130k

## • We suspected high-star repositories likely have better names





Repositories



















![](_page_58_Picture_3.jpeg)

![](_page_59_Figure_1.jpeg)

![](_page_59_Picture_2.jpeg)

![](_page_59_Picture_5.jpeg)

![](_page_60_Figure_1.jpeg)

![](_page_60_Picture_2.jpeg)

Treatment group

Control group

Low-Star Dataset

![](_page_60_Picture_7.jpeg)

![](_page_61_Figure_1.jpeg)

After separating out testing and validation sets, we trained a DIRE model on high-star code and low-star code

![](_page_61_Picture_3.jpeg)

# We tested each model against each test set.

Experiment I

Experiment 2

Experiment 3

Experiment 4

![](_page_62_Picture_5.jpeg)

![](_page_62_Figure_6.jpeg)

![](_page_63_Figure_0.jpeg)

![](_page_64_Figure_0.jpeg)

![](_page_64_Picture_2.jpeg)