

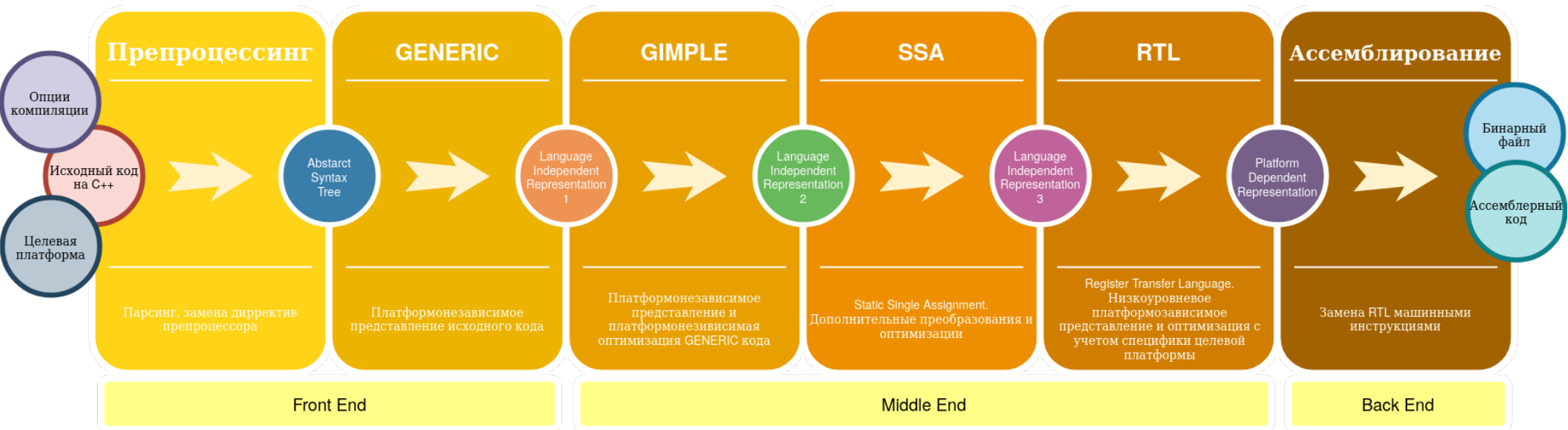
Classification of functions by ML methods based on their binary representation

24.12.2022

The essence of the task

- There is a source code of the program in a compiled programming language (PL)
- There is a binary file received after compilation
- After disassembling, there is the code of the Assembler program and the code of the functions used in the program
- **Problem:** let there be an unknown function in an unknown binary. After disassembling, the Assembler code was obtained. **What does this function do?**

Compilation stages of GCC



Limitations

Ограничения

1

Язык программирования

C++



2

Тип бинарных файлов

Linux ELF



3

Платформа

Intel x86



4

Библиотеки

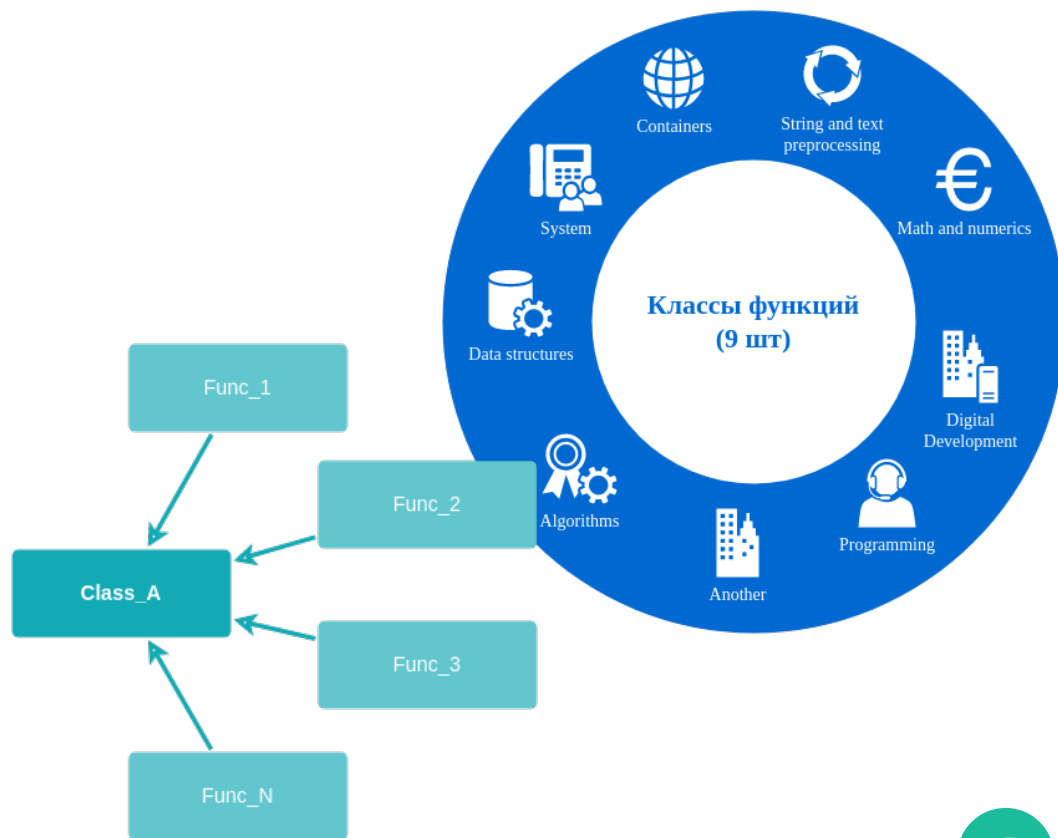
Boost, STL, Asio, Silicon



5

Компиляторы

GCC, Clang



Different compilers

The screenshot displays the Compiler Explorer interface in Mozilla Firefox. The browser address bar shows <https://godbolt.org>. The Compiler Explorer logo is in the top left, and the title bar reads "Compiler Explorer - Mozilla Firefox".

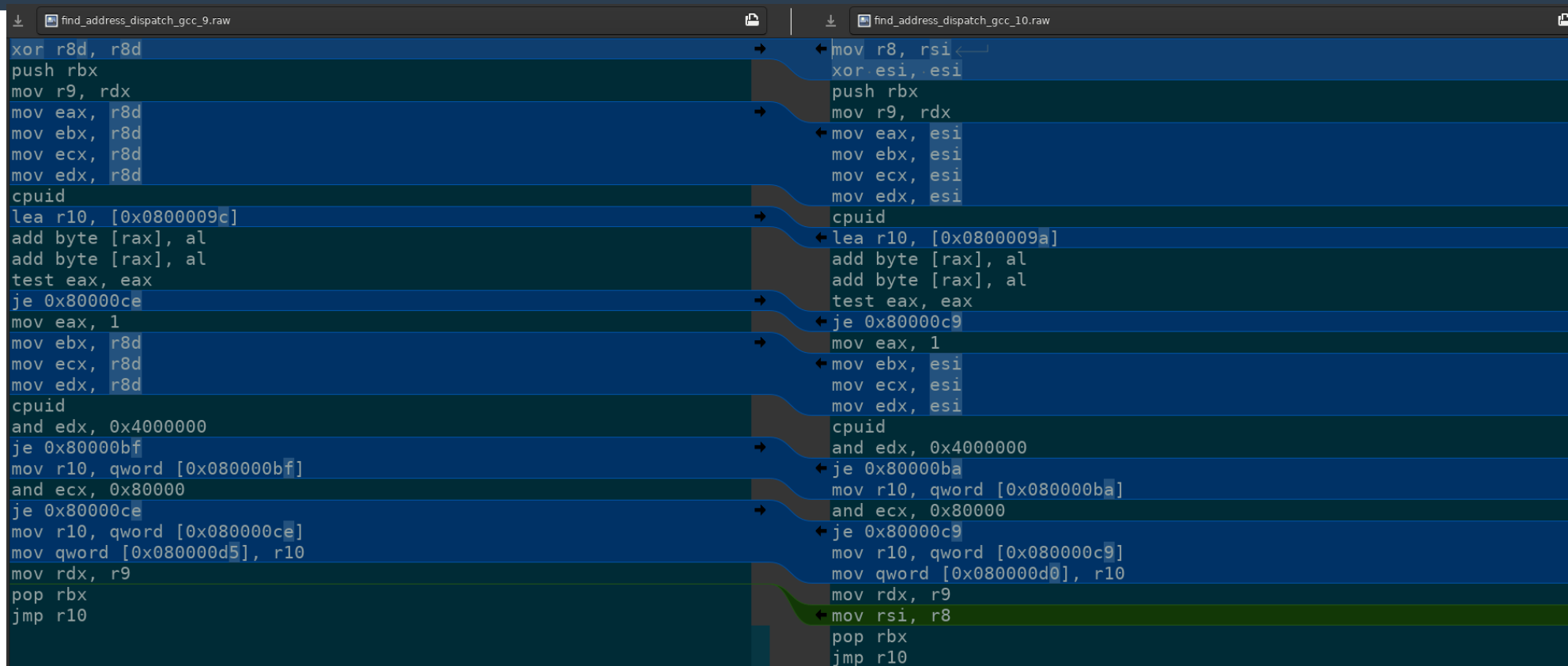
The interface is divided into three main panes:

- Source Code Pane (Left):** Displays the C++ source code for `poco/ASCIIEncoding.cpp`. The code includes a 7-bit ASCII text encoding class with methods `convert`, `queryConvert`, and `sequenceLength`.
- Assembly Pane (Middle):** Shows the assembly output for the `convert` method using `x86-64 gcc 10.2`. The assembly includes instructions like `push rbp`, `mov rbp, rsp`, `mov QWORD PTR [rbp-8], rdi`, and `ret`.
- Assembly Pane (Right):** Shows the assembly output for the `convert` method using `x86-64 clang 10.0.1`. The assembly includes instructions like `push rbp`, `mov rbp, rsp`, `mov qword ptr [rbp-16], rdi`, and `ret`.

At the bottom of the interface, there is a status bar with the following information:

- Output (0/0) x86-64 gcc 10.2 i - cached (8408B)
- Output (0/0) x86-64 clang 10.0.1 i - 442ms (25433B)

Different versions of the same compiler



gcc_9.raw	gcc_10.raw
xor r8d, r8d	mov r8, rsi
push rbx	xor esi, esi
mov r9, rdx	push rbx
mov eax, r8d	mov r9, rdx
mov ebx, r8d	mov eax, esi
mov ecx, r8d	mov ebx, esi
mov edx, r8d	mov ecx, esi
cuid	mov edx, esi
lea r10, [0x0800009c]	cuid
add byte [rax], al	lea r10, [0x0800009a]
add byte [rax], al	add byte [rax], al
test eax, eax	add byte [rax], al
je 0x80000ce	test eax, eax
mov eax, 1	je 0x80000c9
mov ebx, r8d	mov eax, 1
mov ecx, r8d	mov ebx, esi
mov edx, r8d	mov ecx, esi
cuid	mov edx, esi
and edx, 0x4000000	cuid
je 0x80000bf	and edx, 0x4000000
mov r10, qword [0x080000bf]	je 0x80000ba
and ecx, 0x80000	mov r10, qword [0x080000ba]
je 0x80000ce	and ecx, 0x80000
mov r10, qword [0x080000ce]	je 0x80000c9
mov qword [0x080000d5], r10	mov r10, qword [0x080000c9]
mov rdx, r9	mov qword [0x080000d0], r10
pop rbx	mov rdx, r9
jmp r10	mov rsi, r8
	pop rbx
	jmp r10

Getting raw data

- Compilation of the Boost library by various versions of gcc compilers (10.2.1, 9.3.0) with various compilation options (O0, O1, O2, O3, Os, Og, Ofast) using the b2 build system
- Disassembling the received files using Radare2 (r2)
- Extraction of function names, number of instructions, Assembly code, block sequence graph in gml format

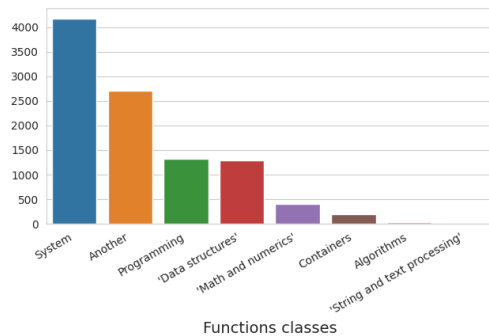
142438 lines

10213 functions

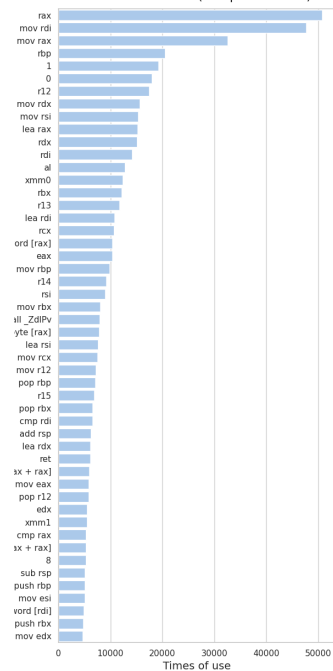
9 classes

EDA

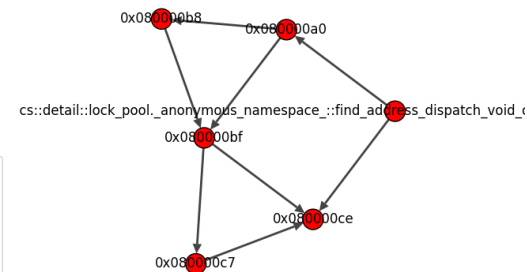
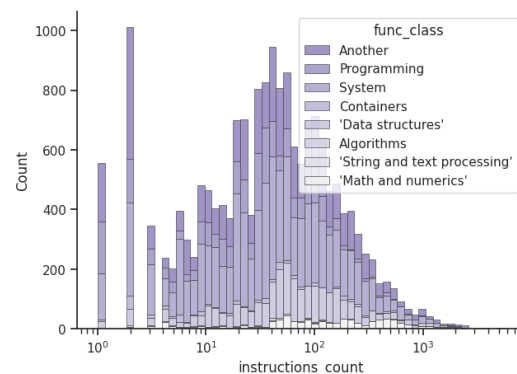
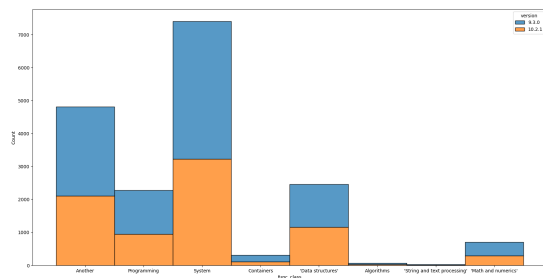
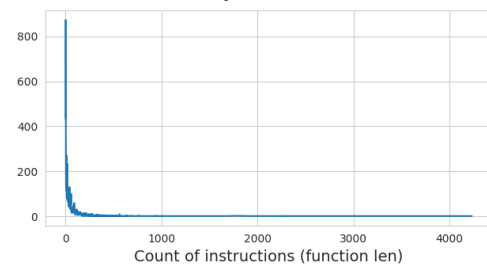
Count functions in classes distribution



Top-50 most used instructions and commands (not precleaned)



Funcs by len distribution



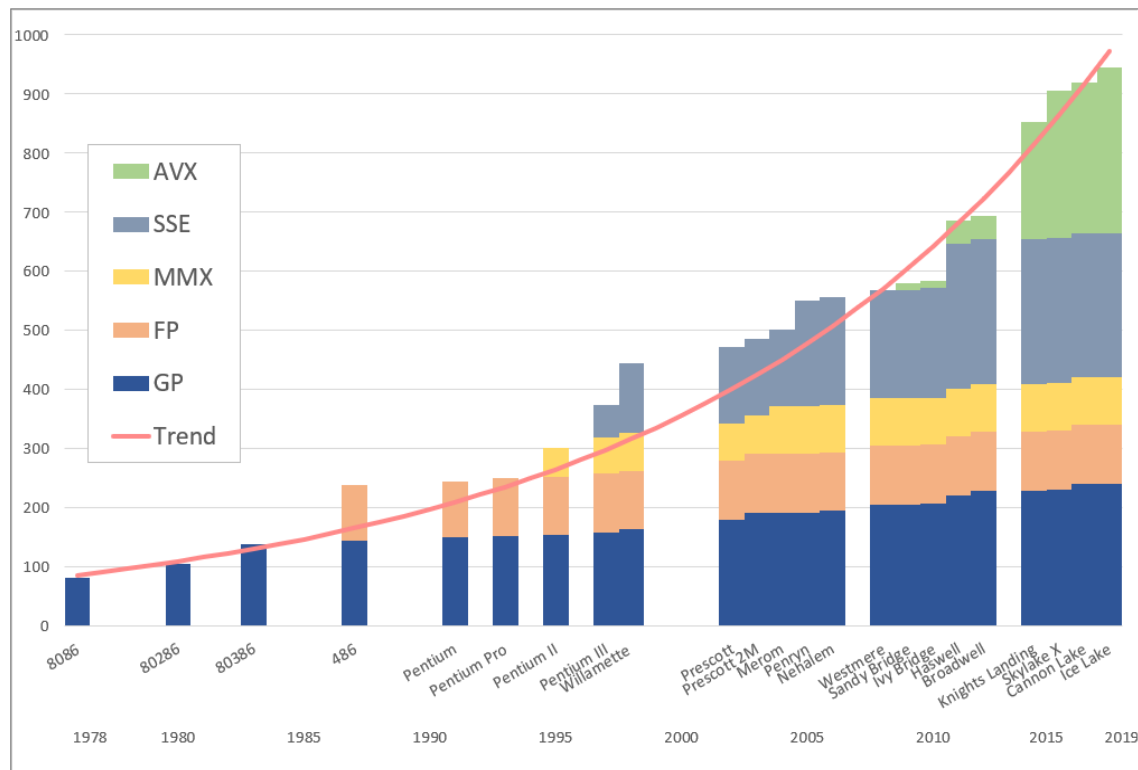
EDA



230000

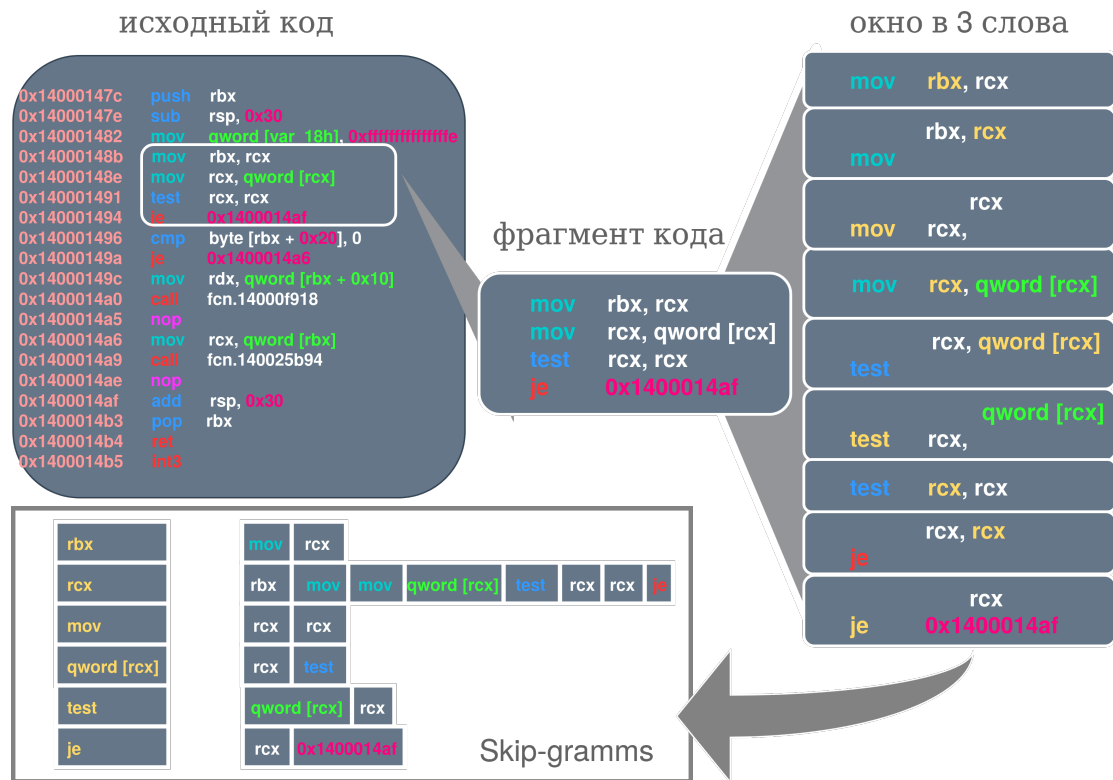
6000

~18000



<https://habr.com/ru/post/503486/>

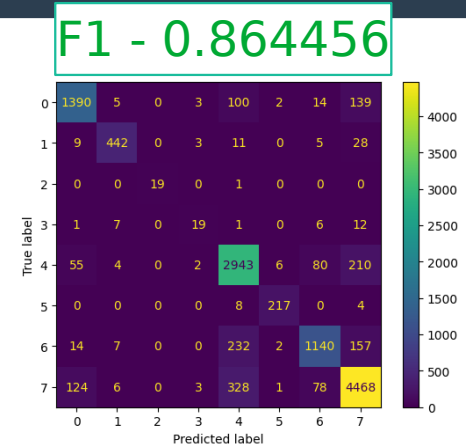
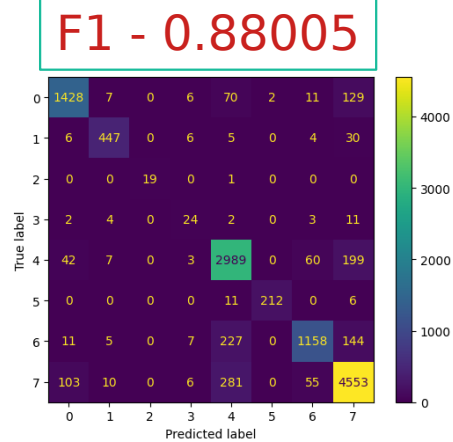
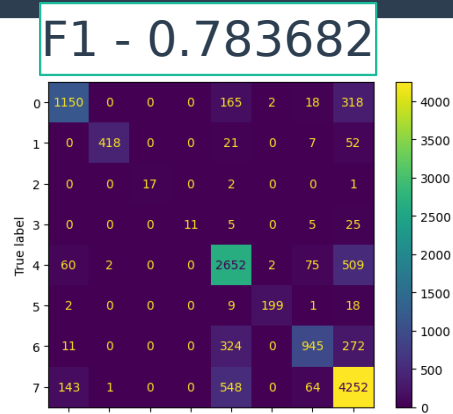
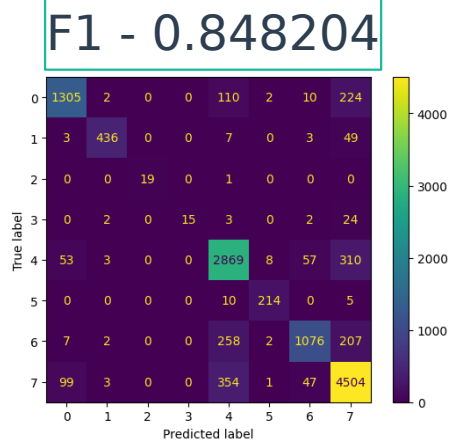
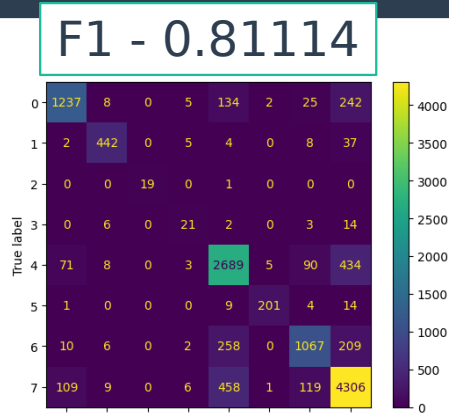
CountVectorizer and TfidfVectorizer



LogisticRegression with CountVectorizer

y='Data structures' top features		y='Math and numerics' top features		y='String and text processing' top features		y='Algorithms' top features		y='Another' top features		y='Containers' top features		y='Programming' top features		y='System' top features	
Weight ²	Feature	Weight ²	Feature	Weight ²	Feature	Weight ²	Feature	Weight ²	Feature	Weight ²	Feature	Weight ²	Feature	Weight ²	Feature
+1.352	and rdi	+1.555	st	+0.640	lock	+0.880	lea	+1.845	<BIAS>	+1.528	qword addr	+1.600	ret	+2.378	<BIAS>
+1.350	and rdi addr	+1.453	loc	+0.579	jne addr	+0.822	rax addr	+1.526	rax addr push	+1.435	sym	+1.265	<BIAS>	+2.352	rax call qword
+1.247	rbp rax call	+1.232	xmm1	+0.553	r8d	+0.749	esi	+1.378	rsi qword rbx	+1.356	boostfcn	+1.134	addr mov	+1.504	byte rax je
+1.214	rax and	+1.076	fcn call fcn	+0.538	ecx	+0.633	lea rsi	+1.311	rdi qword addr	+1.091	cmp qword	+1.067	dword	+1.225	fcn nop ret
+1.212	fcn lea rax	+1.073	rsp addr	+0.535	add	+0.597	rdi qword rbp	+1.213	ret	+1.056	add rdi	+1.057	rax qword rsi	+1.216	fcn lea rsi
+1.063	mov edx dword	+1.023	fld	+0.524	jne	+0.593	rax addr ret	+1.196	addr jmp fcn	+1.028	jmp boostfcn	+1.057	close	+1.166	addr rax mov
+1.030	methodmethod	+0.969	ja	+0.520	jne addr mov	+0.551	r14 mov	+1.189	rsp call	+1.028	jmp boostfcn	+1.049	rdi addr lea	+1.150	addr rax lea
+0.983	movsx	+0.944	fcn call	... 1790 more positive ...		+0.544	r14		methodmethod		nop	+1.045	jne addr add	+1.107	add rdi addr
+0.935	call text	+0.921	comisd	... 63724 more negative ...		+0.542	al	+1.143	rax rax lea	+1.016	call sym	+1.027	fcn test rax	+1.102	rdi addr rcx
... 19786 more positive ...		+0.878	sub rsp addr	-0.572	addr jmp	+0.518	qword rbp	+1.124	je addr lock	+1.004	boostfcn nop	+1.001	addr add rsp	... 28143 more positive ...	
... 45728 more negative ...		+0.845	xword	-0.582	push	... 1541 more positive 16569 more positive ...		+1.004	boostfcn nop word	+0.980	word	... 37371 more negative ...	
-0.912	je addr lock	... 10419 more positive ...		-0.584	eax ret	... 63973 more negative 48945 more negative ...				+0.972	jmp fcn	-1.152	al je addr
-0.982	lea rsi addr	... 55095 more negative ...		-0.664	jmp fcn	-0.555	qword	-1.235	varfcn lea rsi	+0.990	dword addr	... 12273 more positive ...		-1.165	rdi rbp mov
-0.989	rsi qword rbx	-0.848	movq	-0.827	eax	-0.709	rsp	-1.248	rdi rax call	+0.916	rsi rdi	... 53241 more negative ...		-1.166	lea rdi addr
-1.038	rdx addr add	-0.849	qword rax	-0.811	fcn	-0.811	ret	-1.252	rax rax call	+0.894	qword addr ret	-1.036	je addr lea	-1.345	dword addr
-1.043	fcn lea rdi	-0.884	cmp	-1.548	ret	-0.910	jmp	-1.264	jmp addr nop	... 6809 more positive ...		-1.042	or	-1.411	rdi addr lea
-1.330	rdx qword varfcn	-0.925	movdqa	-3.077	<BIAS>	-2.221	<BIAS>	-1.366	rax call qword	... 58705 more negative ...		-1.077	addr add byte	-1.838	rax call fcn
								-1.501	byte rdi addr	-0.995	edi	-1.159	rdi addr test		
										-1.583	byte				

LogisticRegression vs RandomForestClassifier with CountVectorizer vs TfidfVectorizer



LogisticRegression	CountVectorizer	1-2
LogisticRegression	TfidfVectorizer	1-2
LogisticRegression	CountVectorizer	1-3
LogisticRegression	TfidfVectorizer	1-3
RandomForestClassifier	CountVectorizer	1-3

Feature extraction

- Working with a graph (python graph)
- Working separately with operands and operators
- Working with addresses
- Prologue and epilogue of functions
- Stack
- Optimization options
- More compilers and their versions
- More libraries

```
print(*g.get_adjacency())
print(g.vcount())
print(g.incident(0))
print(g.degree(0))
print(g.diameter())
print(g.girth())
print(g.radius())
print(g.average_path_length())
print(g.transitivity_avglocal_undirected())
print(g.laplacian())
```

Where it can be used, in what form to implement

- Traffic filter plugin in IDS Suricata
- Plugin for Radare2 for the initial analysis of unknown functions

End