

BAR 2023

# FCGAT: Interpretable Malware Classification Method using Function Call Graph and Attention Mechanism

March 3, 2023

Minami Someya<sup>\*†</sup> Yuhei Otsubo<sup>†\*</sup> Akira Otsuka<sup>\*</sup>

\* Institute of Information Security, Japan

† National Police Agency, Japan

# Malware Classification

- Identifies malware family or category

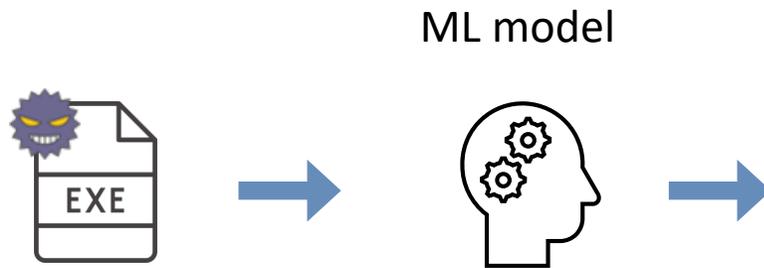
Emotet,  
Trickbot ...

Downloader,  
Ransomware ...

- Usefulness of malware classification:
  - Understands malware behavior
  - Helps with malware analysis
- Limitations of conventional signature-based methods:
  - Cannot keep up with creating pattern files of new malware
- Solution:
  - Use machine learning to classify malware

# Drawbacks in ML Method: Lack of Interpretability

- ML-models are often treated as black boxes

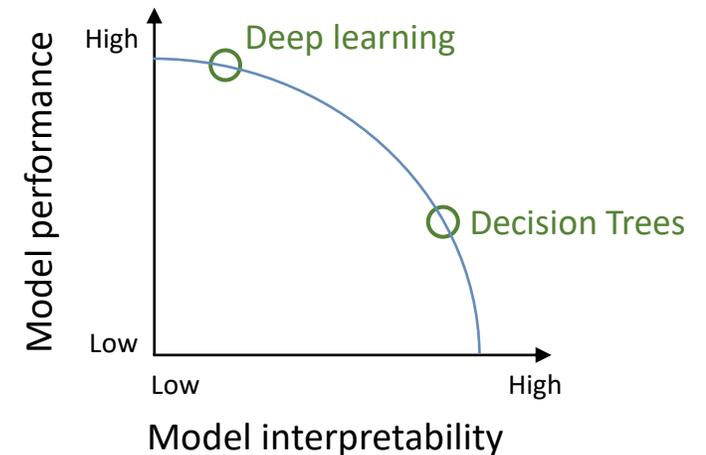


- IcedID
- Trickbot
- njrat

Why Trickbot?  
Where's the malicious code?



- Model interpretability and performance are often in a trade-off relationship



## Research Goals

- Creating classifier that can explain the reasons for malware classification
- Achieving both high classification performance and interpretability

## Solution

- FCGAT: Interpretable Malware Classification Method  
using **F**unction **C**all **G**raph and **A**ttention Mechanism

## Contributions

- Successfully classified malware families with high performance comparable to cutting-edge methods
- Analyzed the explanations and obtained insight into the functions that characterize malware

# Determining Feature Set

Byte	Basic Block	Function
<div data-bbox="405 396 529 482" style="border: 1px solid black; padding: 5px; text-align: center;">5A</div>	<div data-bbox="1016 351 1470 529" style="border: 1px solid black; padding: 5px;"><pre>mov ebx, [ebp+8] jmp short loc_100013</pre></div>	<div data-bbox="1640 339 2265 558" style="border: 1px solid black; padding: 5px;"><pre>funcA push ebp mov ebp, esp ... call OpenMutexA</pre></div>



Which feature set should we use?

We need to consider the interpretability of explanation results.

# Determining Feature Set

Byte	Basic Block	Function
<div style="border: 1px solid black; padding: 5px; width: fit-content; margin: auto;">5A</div>	<div style="border: 1px solid black; padding: 5px; width: fit-content; margin: auto;">                     mov ebx, [ebp+8]                      jmp short loc_100013                 </div>	funcA <div style="border: 1px solid black; padding: 5px; width: fit-content; margin: auto;">                         push ebp                          mov ebp, esp                          ...                          call OpenMutexA                     </div>
<div style="border: 1px solid black; padding: 5px; width: fit-content; margin: auto;">                         3B 52 2B F3 25 68 AD                          33 55 6A 11 9A 4C B4                          2B CA 4B D3 75 98 AD                          45 7A 24 F3 27 68 4D                          3B 5A 2C FB 25 69 CD                          2F 6A 23 F3 5A 2B F3                          34 52 2B F3 25 58 A4                     </div> <div style="background-color: #f0f0f0; padding: 5px; margin-top: 10px;">                         explanation results                          (important feature)                     </div>	<h3 style="text-align: center;">Control Flow Graph (CFG)</h3> <div style="text-align: center;"> <div style="border: 1px solid black; padding: 5px; width: fit-content; margin: auto;">                             mov ebx, [ebp+8]                              jmp short loc_100013                         </div> <div style="display: flex; justify-content: space-around; margin-top: 10px;"> <div style="border: 1px solid black; padding: 5px; width: 40%;">                             cmp byte ptr [ebx], 0                              jnz short loc_100021                         </div> <div style="border: 1px solid black; padding: 5px; width: 40%;">                             pop edi                              pop esi                              ...                         </div> </div> </div>	<h3 style="text-align: center;">Function Call Graph (FCG)</h3> <div style="text-align: center;"> <div style="border: 1px solid black; border-radius: 10px; padding: 5px; width: 100px; margin: auto;">start</div> <div style="display: flex; justify-content: space-around; margin-top: 20px;"> <div style="border: 1px solid black; border-radius: 10px; padding: 5px; width: 100px;">funcA</div> <div style="border: 1px solid black; border-radius: 10px; padding: 5px; width: 100px;">funcB</div> </div> <div style="display: flex; justify-content: space-around; margin-top: 20px;"> <div style="border: 1px solid black; border-radius: 10px; padding: 5px; width: 150px;">OpenMutexA</div> <div style="border: 1px solid black; border-radius: 10px; padding: 5px; width: 150px;">InternetOpenW</div> </div> </div>

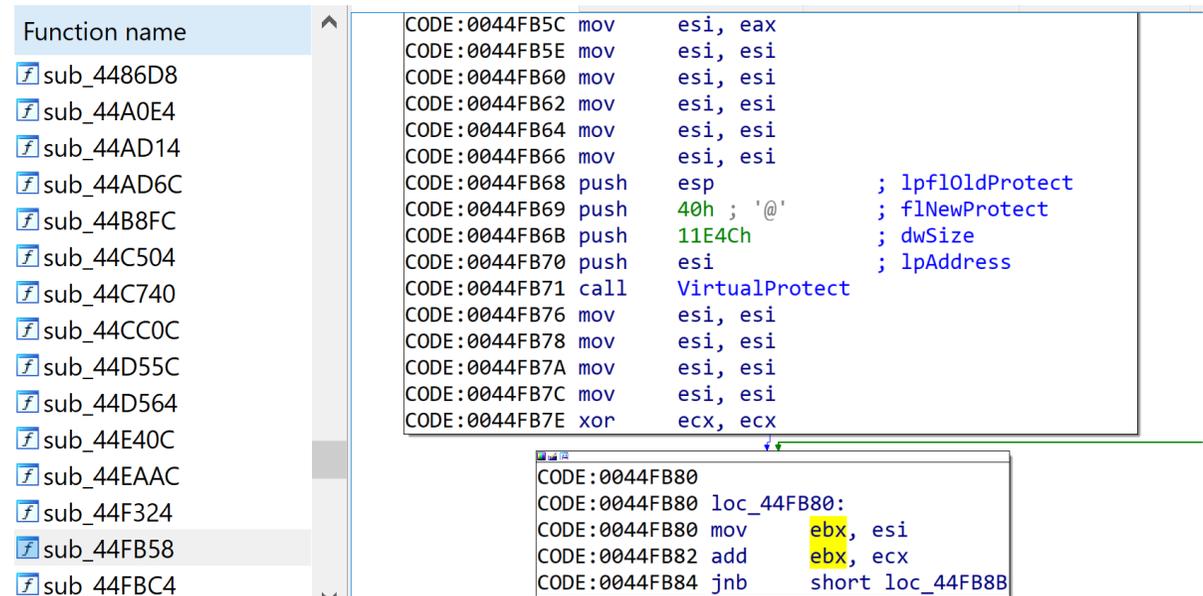
Difficult to interpret  
 The granularity is too small 😞

Easy to interpret 😊

# Why Function-based Feature

- Easier to interpret than byte or basic block
- Often reused in a same malware family
  - Malware is rarely implemented from scratch
- Functions and their relationships are focused on during analysis

→ Function Call Graph (FCG)



The screenshot displays a debugger interface with a list of functions on the left and their assembly code on the right. The function list includes:

- sub\_4486D8
- sub\_44A0E4
- sub\_44AD14
- sub\_44AD6C
- sub\_44B8FC
- sub\_44C504
- sub\_44C740
- sub\_44CC0C
- sub\_44D55C
- sub\_44D564
- sub\_44E40C
- sub\_44EAAC
- sub\_44F324
- sub\_44FB58
- sub\_44FBC4

The assembly code for the selected function (sub\_44FB58) is shown below:

```
CODE:0044FB5C mov     esi, eax
CODE:0044FB5E mov     esi, esi
CODE:0044FB60 mov     esi, esi
CODE:0044FB62 mov     esi, esi
CODE:0044FB64 mov     esi, esi
CODE:0044FB66 mov     esi, esi
CODE:0044FB68 push    esp                ; lpf10ldProtect
CODE:0044FB69 push    40h ; '@'         ; flNewProtect
CODE:0044FB6B push    11E4Ch             ; dwSize
CODE:0044FB70 push    esi                ; lpAddress
CODE:0044FB71 call   VirtualProtect
CODE:0044FB76 mov     esi, esi
CODE:0044FB78 mov     esi, esi
CODE:0044FB7A mov     esi, esi
CODE:0044FB7C mov     esi, esi
CODE:0044FB7E xor     ecx, ecx

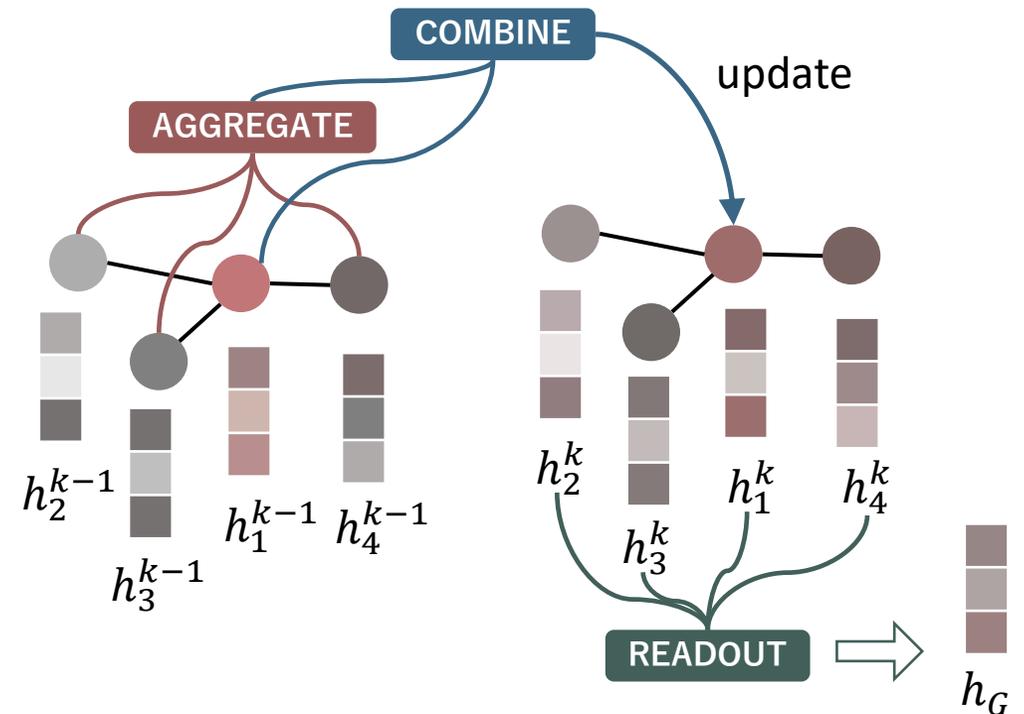
CODE:0044FB80
CODE:0044FB80 loc_44FB80:
CODE:0044FB80 mov     ebx, esi
CODE:0044FB82 add     ebx, ecx
CODE:0044FB84 jnb    short loc_44FB8B
```

# Graph Neural Network (GNN)

- Updating feature to reflect the graph structure

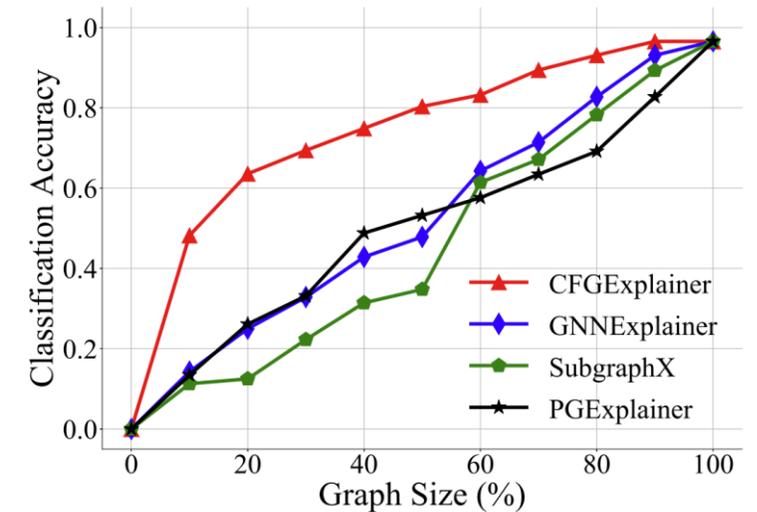
$$\mathbf{h}_v^k \leftarrow \text{COMBINE}(\mathbf{h}_v^{k-1}, \text{AGGREGATE}(\{\mathbf{h}_u^{k-1}, \forall u \in \mathcal{N}(v)\}))$$

- AGGREGATE:**  
Aggregate features of neighboring nodes
- COMBINE:**  
Update the features of the node to be updated using the neighboring nodes
- READOUT:**  
Obtain a representation of the entire graph from the nodes in the graph



# Related Work: CFGExplainer

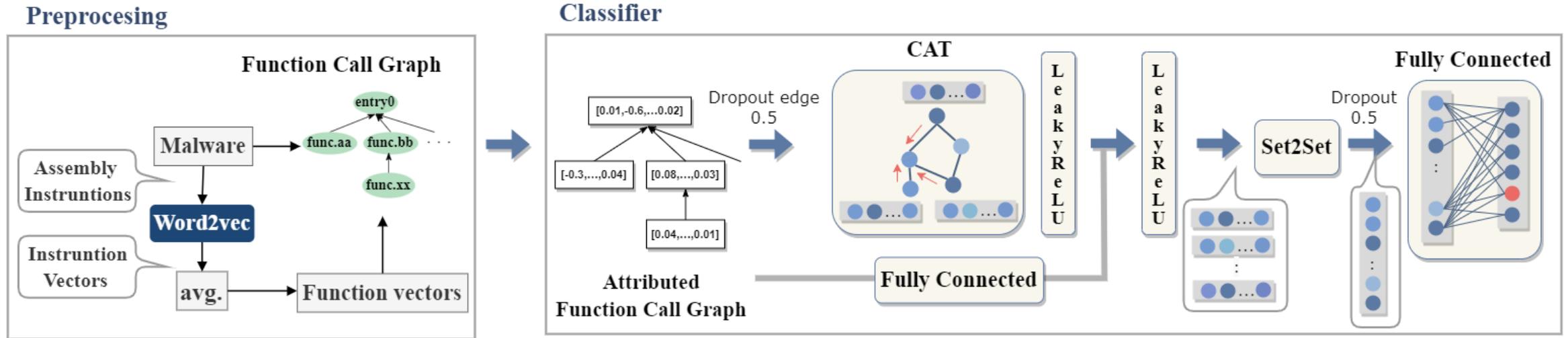
- Explanation method for GNN-based malware classification models
- Uses Control Flow Graph with Basic Blocks as nodes
- Identifies subgraph that contributes most to classification by pruning less important nodes
- Difference from our research:
  - CFGExplainer uses basic block, our FCGAT uses function



(d) Ldpinch

Herath et al. 2022. "CFGExplainer: Explaining Graph Neural Network-Based Malware Classification from Control Flow Graphs." In 2022 52nd Annual IEEE/IFIP International Conference on Dependable Systems and Networks (DSN). IEEE.

# Overview of FCGAT

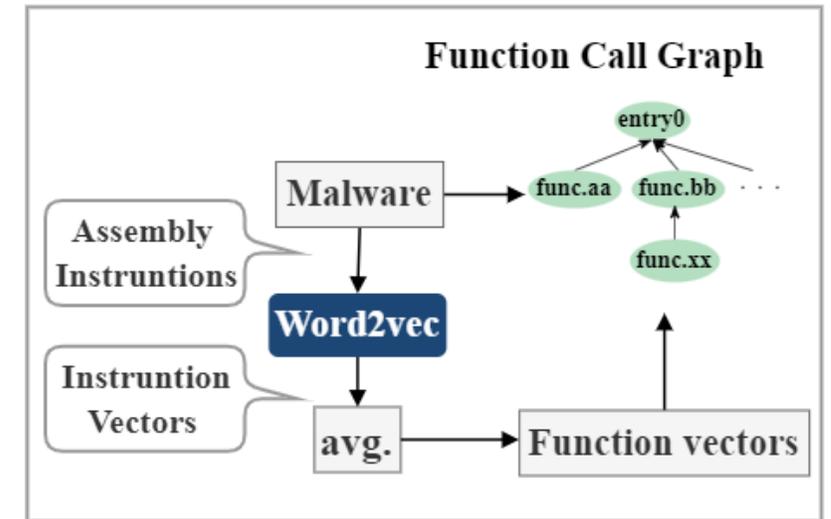


- Preprocessing
  - Creating Function Call Graph
  - Creating feature vector of function (function vector)
- Classifier
  - Malware classification using Graph Neural Network

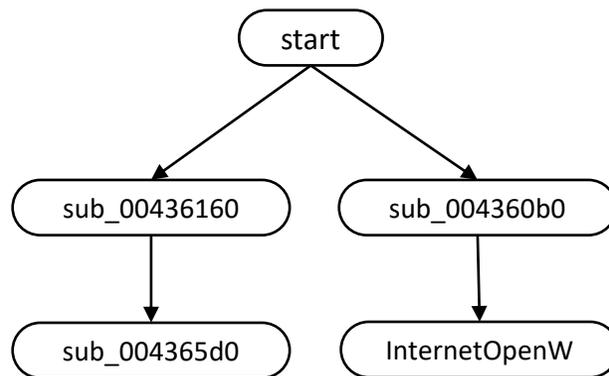
# Preprocessing

- Creating Function Call Graph
  - Using IDA Pro
  - Reversing arrows of FCG
    - The processing of called function is included in that of calling function

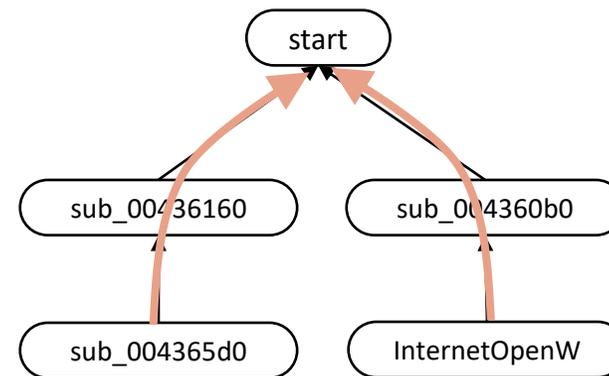
## Preprocessing



General FCG

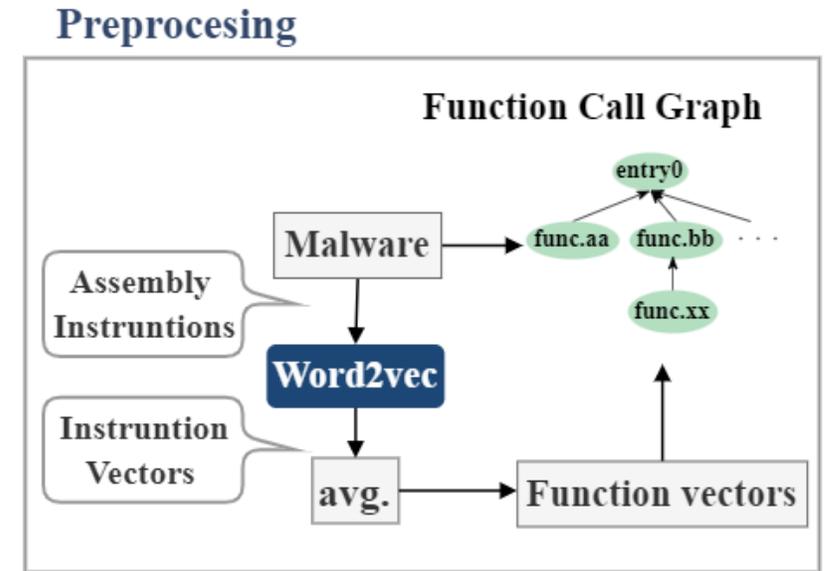


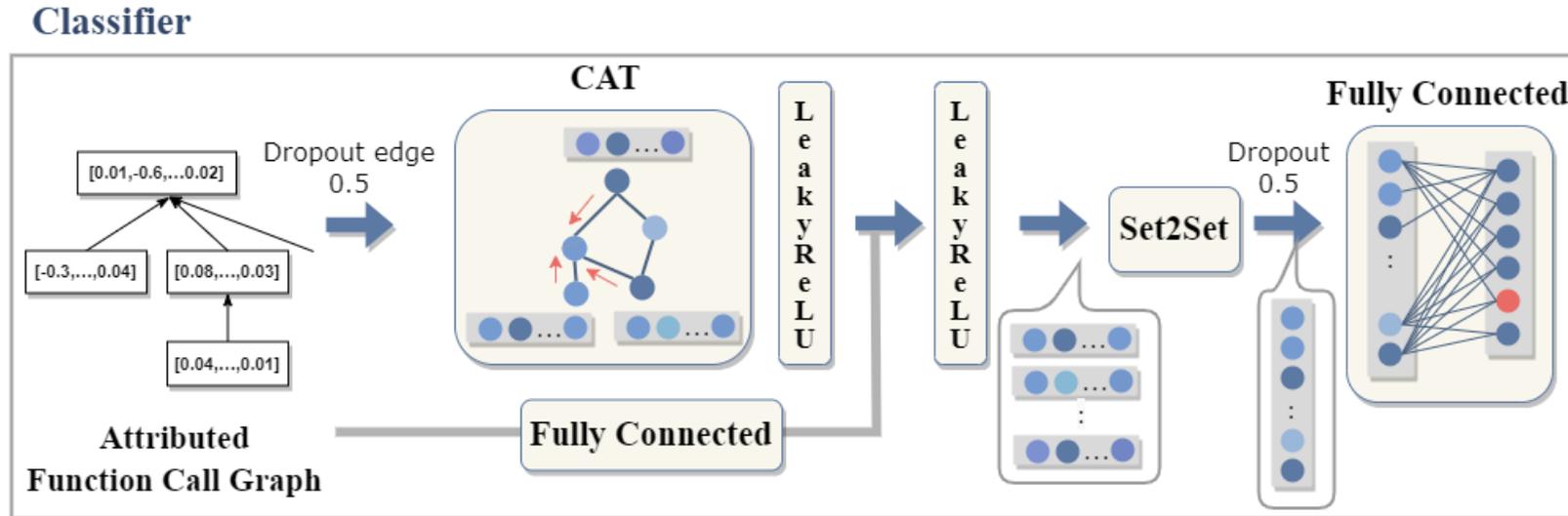
Reverse FCG



# Preprocessing

- Creating Function Call Graph
  - Using IDA Pro
  - **Reversing arrow** of FCG
- Creating function vectors
  - Using **Word2vec** for creating instruction vectors
    - Instruction  $\leftrightarrow$  word, function  $\leftrightarrow$  sentence
  - Averaging instruction vectors in a function to obtain a function vector



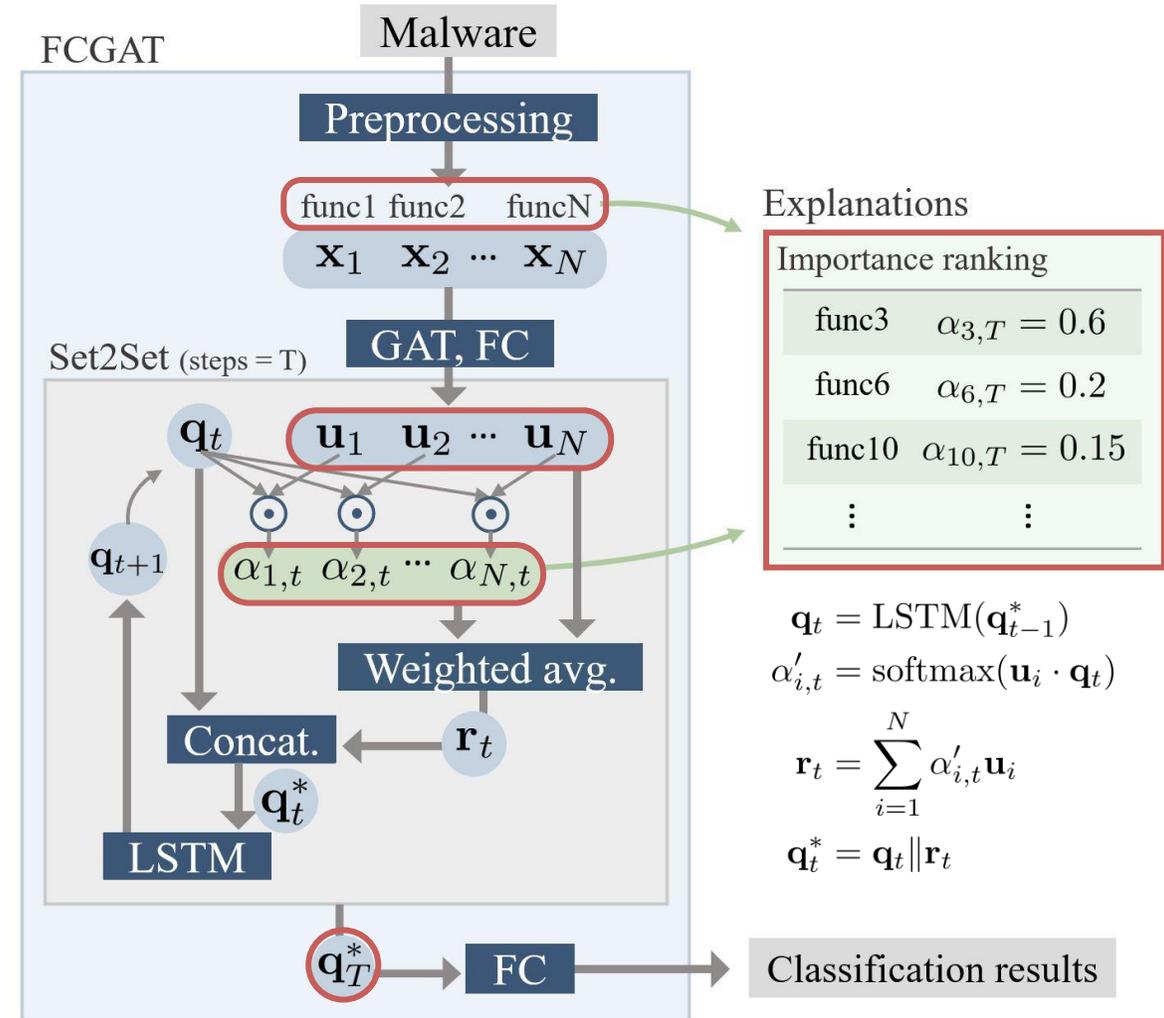


- Graph Attention Network (GAT)
  - Updating function vectors using FCG
- Set2Set (readout process)
  - Next slide
- Fully Connected (final layer)
  - Classifying into a number of classes

$$\mathbf{u}_i = \text{LeakyReLU}(\mathbf{W}_1 \mathbf{x}_{v_i} + \text{LeakyReLU}(\|\sum_{\mathbf{x}_j \in N(v_i)}^{K} \alpha_{ij}^k \mathbf{W}^k \mathbf{x}_j\|))$$
$$\alpha_{ij} = \text{softmax}_j (\text{LeakyReLU}(\mathbf{a}^\top \cdot [\mathbf{W}_2 \mathbf{x}_{v_i} \| \mathbf{W}_2 \mathbf{x}_{v_j}]))$$

# Set2Set (readout process)

- The key to interpretability
- Inputs :  $\mathbf{u}$  Updated function feature  
Outputs:  $\mathbf{q}_T^*$  Feature of the malware
- A larger  $\alpha'_{i,t}$  (attention weight) is assigned to the more important function vector
- Importance ranking of function is provided



- **Classification Performance**
  - Perform malware **family** classification
  - Compare with demonstration results of existing studies by Ma et al.
  - Conduct a replicated experiment of GEMAL (using FCG but not interpretable)
- **Classification Interpretability**
  - Perform malware **category** classification
  - Extract the importance ranking of functions as explanations
  - Confirm the effectiveness of these explanations

Ma, Yixuan et al. 2021. “A Comprehensive Study on Learning-Based PE Malware Family Classification Methods.” In Proceedings of the 29th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering, 1314–25. ESEC/FSE 2021.

Wu, Xiao-Wang et al. 2022. “Embedding Vector Generation Based on Function Call Graph for Effective Malware Detection and Classification.” Neural Computing & Applications.

# Classification Performance

Category	Model	MalwareBazaar dataset				BIG-2015			
		Accuracy	Precision	Recall	F1-Score	Accuracy	Precision	Recall	F1-Score
Image	ResNet-50	96.68	96.91	96.75	96.83	98.42	96.57	95.68	96.08
	VGG-16	96.35	96.58	96.54	96.56	93.94	90.32	81.89	87.27
	Inception-V3	95.83	95.67	95.79	95.73	96.99	93.67	94.46	94.03
	IMCFN	97.38	97.53	97.41	97.47	97.77	95.93	94.81	95.13
Binary	CBOW+MLP	<b>97.81</b>	<b>97.92</b>	<b>98.08</b>	<b>98.00</b>	98.41	97.63	96.67	97.12
	MalConv	95.92	96.04	96.43	96.20	97.02	94.34	92.62	93.33
Disassembly	MAGIC	92.82	88.03	87.36	87.45	98.05	96.75	94.03	95.14
	Word2vec+KNN	95.64	93.34	94.29	93.79	98.07	96.41	96.51	96.45
	MCSC	96.80	94.97	94.51	94.70	97.94	95.97	96.17	96.06
	<b>FCGAT</b>	<b>98.11</b>	<b>98.03</b>	<b>98.27</b>	<b>98.15</b>	<b>99.27</b>	<b>97.93</b>	<b>98.45</b>	<b>98.18</b>
	GEMAL	97.71	97.65	98.00	97.82	<b>99.37</b>	<b>98.26</b>	<b>98.48</b>	<b>98.37</b>

by Ma et al.

by us

FCGAT outperforms all other methods on all metrics

FCGAT is equivalent to the replication experiment of GEMAL

# Classification Interpretability

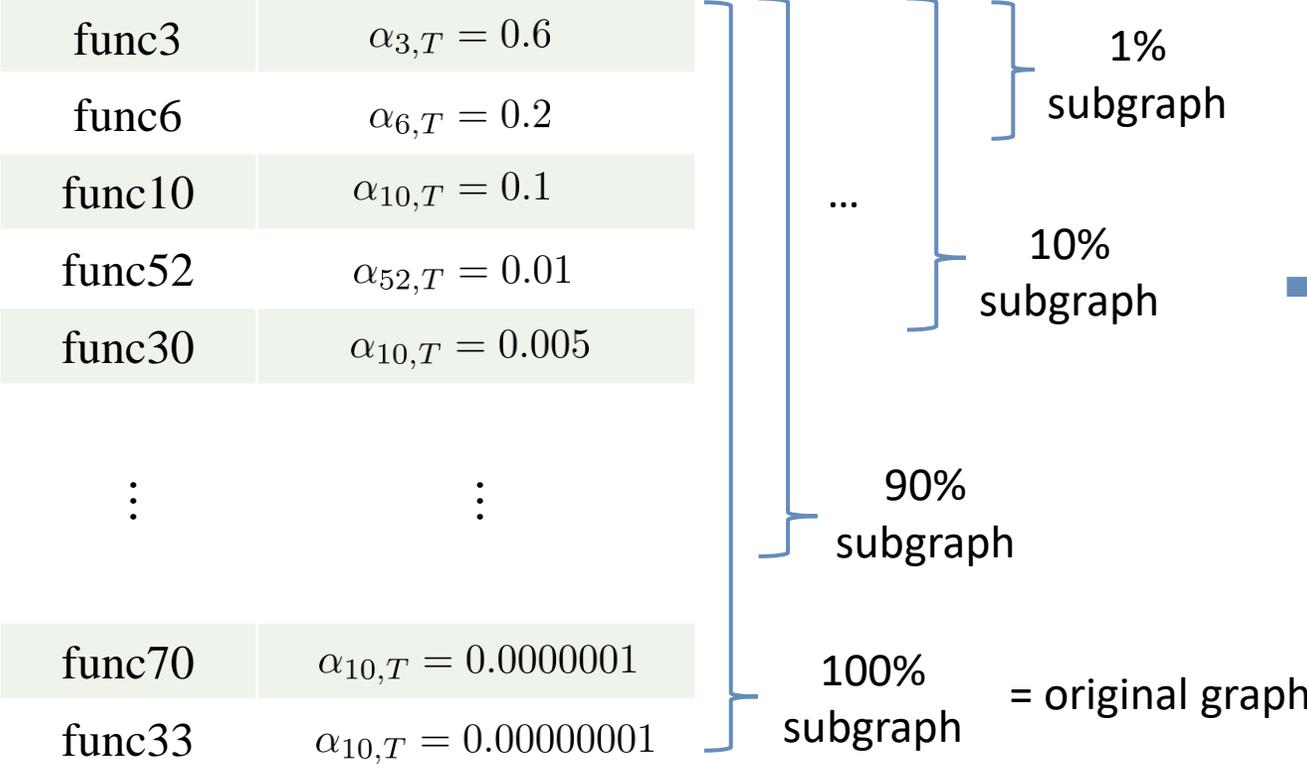
- Malware category classification
  - Malicious behavior is more common in categories
- Dataset: BODMAS-8cat
  - Exclude packed samples detected by peid from BODMAS

Category	Family Counts	Sample Counts
backdoor	31	598
downloader	19	967
dropper	17	397
informationstealer	19	347
ransomware	18	169
trojan	282	15,674
virus	8	93
worm	87	5,124
total	481	23,369

Yang, Limin et al. 2021. “BODMAS: An Open Dataset for Learning Based Temporal Analysis of PE Malware.” In 2021 IEEE Security and Privacy Workshops (SPW), 78–84.

- Measuring the classification performances of subgraphs

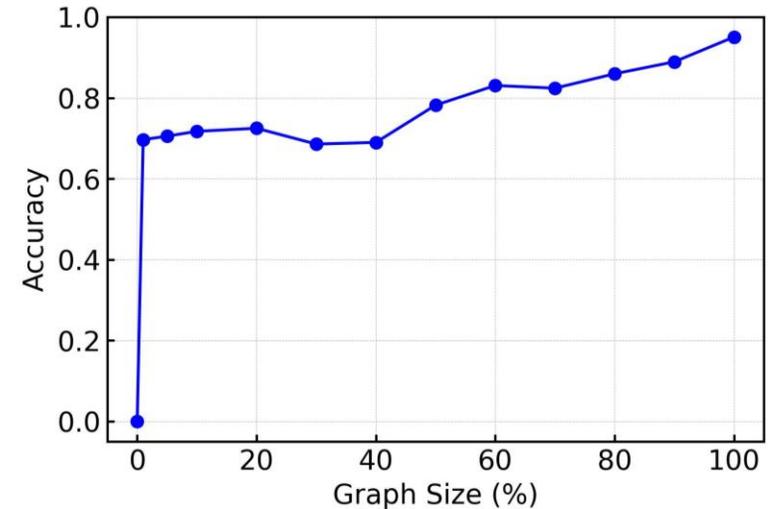
Importance ranking



➔ Evaluate classification performance

# Classification Accuracies of Subgraphs

- Only top 6 functions (average per sample) achieve 69.67% accuracy
- Comparison with existing study
  - CFGExplainer showed 52.39% accuracy in the 10% subgraph
  - FCGAT achieves 71.73% accuracy !
- Using function, malware can be characterized with a small number of nodes



Graph Size (%)	Average Number of Nodes	Accuracy
1	6.2	69.67
5	28.8	70.53
10	56.9	71.73
20	113.4	72.48
30	170.1	68.57
40	226.5	68.99
50	283.0	78.20
60	339.6	83.06
70	396.2	82.40
80	452.7	85.96
90	509.3	88.92
100	565.4	95.06

# Trend Analysis of Malware Categories

Aggregate results for each category of functions with max attention weights

backdoor		downloader		dropper		informationstealer	
CreateFileA	119	InternetOpenW	260	sub_460d8c2	278	GetFocus	202
fclose	105	MessageBoxA	130	GetWindowThreadProcessId	48	mciSendStringA	47
sub_4f847a1	83	mciSendStringA	77	ChecksumMappedFile	17	__dllonexit	11
SetWindowLongA	74	sub_f9da9b9	61	IIDFromString	14	free	10
GetCommandLineA	47	sub_1236153e	48	IsProcessorFeaturePresent	4	CreateFileA	9
sub_a9f1051	29	DispatchMessageA	34	sub_62922e7	3	sub_d0f95e9	8
sub_cf1fee5	18	UpdateWindow	29	__imp_VirtualProtect	2	lstrcatA	7
_fdopen	17	CoUninitialize	27	ShellExecuteA	2	GdipCreateFromHDC	7
ransomware		trojan		virus		worm	
CoRegisterMallocSpy	34	sub_a9f1051	1,316	GetActiveWindow	38	__abnormal_termination	1,314
CoReleaseMarshalData	22	__vbaUI114	1,036	strepv	31	EnterCriticalSection	892
InternetReadFile	19	CloseHandle	709	start_80bc47b	15	sub_f7f9a83	407
ReadFile	12	InternetOpenW	633	?ProcessWndProcException	2	__ZNSt6locale5_ImplC2Ej	255
SetPropA	11	GetSystemDirectoryA	512	nullsub_3	2	__ZNSt6locale5_ImplC1Ej	205
SwitchDesktop	11	sub_acc2cf0	510	CoUninitialize	1	sub_d0f95e9	188
IsProcessorFeaturePresent	10	OleSetMenuDescriptor	472	GetFileVersionInfoSizeW	1	__ZNSt6locale6globalERKS_	182
CreateOleAdviseHolder	7	rtcLowerCaseVar	388	EnableMenuItem	1	SetWindowsHookExA	158

- Some functions reflect category characteristics

# Case Study: Ransomware/ GandCrab

- GandCrab is ransomware that appeared in 2018
- We will see two samples, GandCrab#1 and GandCrab#2

DLL file		EXE file	
GandCrab#1		GandCrab#2	
function	$\alpha'_{i,4}$	function	$\alpha'_{i,4}$
<u>aes_encrypt</u>	0.5984	<u>aes_encrypt</u>	0.0331
<u>aes_decrypt</u>	0.4015	<u>aes_decrypt</u>	0.0084
sub_10007BB0	2.002e-06	SetHandleInformation	0.0080
sub_10003F70	4.166e-07	GetTickCount	0.0080
sub_10004C20	7.423e-10	InitializeCriticalSection	0.0080

- These samples are common to the top two important functions, *aes\_encrypt* and *aes\_decrypt*, which are characteristic of ransomware

- Proposed FCGAT
  - The first study to explain malware classification on a per function basis
- Evaluated classification performance
  - High performance competitive to the latest method
- Confirmed the effectiveness of the explanations
  - Functions that reflect the malware feature
  - A small number of functions characterizing malware

Contact: [mgs227501@iisec.ac.jp](mailto:mgs227501@iisec.ac.jp)