

# o-glassesX: Compiler Provenance Recovery with Attention Mechanism from a Short Code Fragment

Yuhei Otsubo<sup>\*†</sup>, Akira Otsuka<sup>†</sup>, Mamoru Mimura<sup>‡†</sup>,  
Takeshi Sakaki<sup>§</sup>, and Hiroshi Ukegawa<sup>\*</sup>

<sup>\*</sup>National Police Agency, Tokyo, Japan

<sup>†</sup>Institute of information Security, Kanagawa, Japan

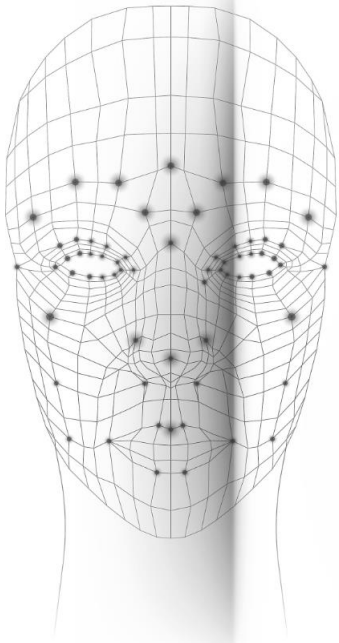
<sup>‡</sup>National Defense Academy, Kanagawa, Japan

<sup>§</sup>The University of Tokyo, Tokyo, Japan

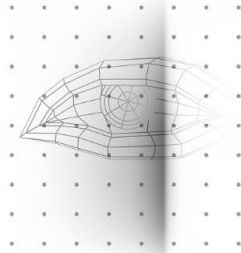
# Introduction

# Forensic Scientists

## FACE RECOGNITION SYSTEM



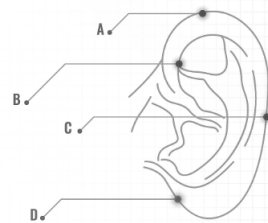
## EYES SCANNER



## EAR IDENTIFICATION SYSTEM

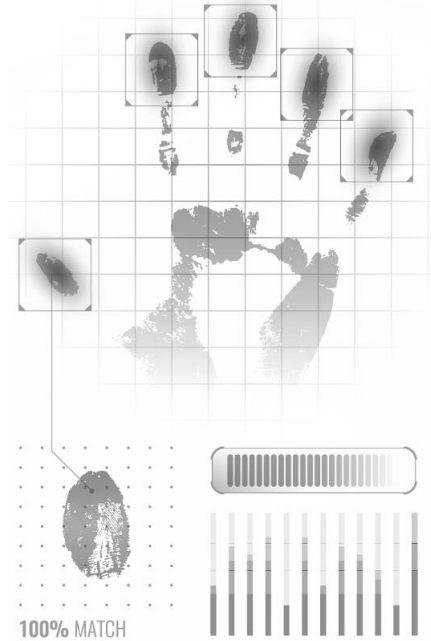
SERIAL: 124386576796870887

Learn more about all our... comprehensive... and other...  
Detailed description of ear identification system, including...  
Detailed description of ear identification system, including...



## HAND RECOGNITION SYSTEM

SERIAL: 124386576796870845



## FINGERPRINT RECOGNITION SYSTEM

SERIAL: 124386576796870867

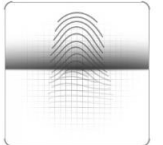
### CURL-SHAPED



OBJECT: A-2107/2  
MODE: SEARCH  
STATUS: COMPLETE



### WAVY-SHAPED



OBJECT: A-2107/3  
MODE: SCANNING  
STATUS: IN PROGRESS...

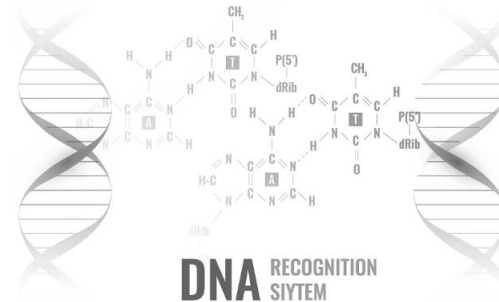
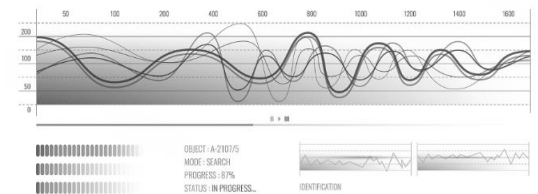
### ARCH-SHAPED LOOK



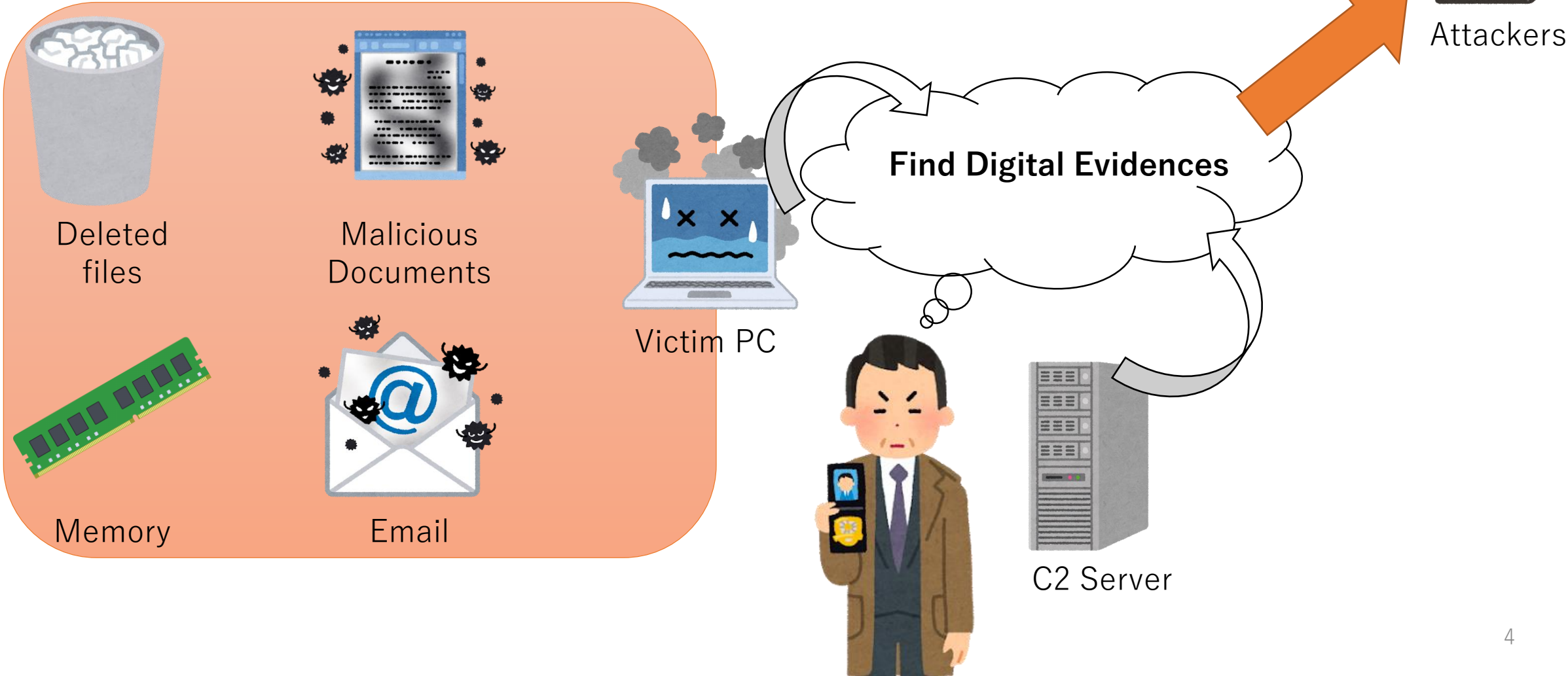
OBJECT: A-2107/1  
MODE: SCANNING  
STATUS: IN PROGRESS...

## VOICE IDENTIFICATION SYSTEM

SERIAL: 124386576796870867



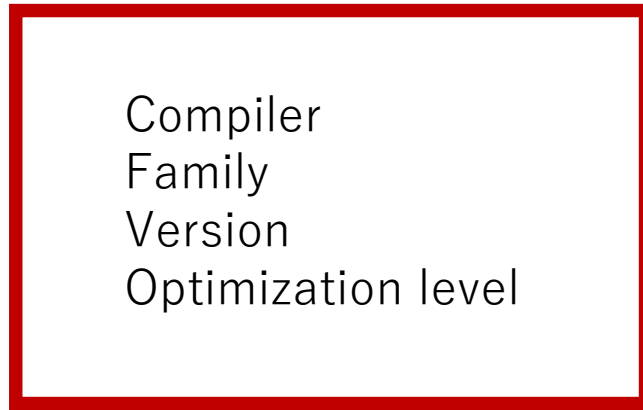
# Computer Forensics



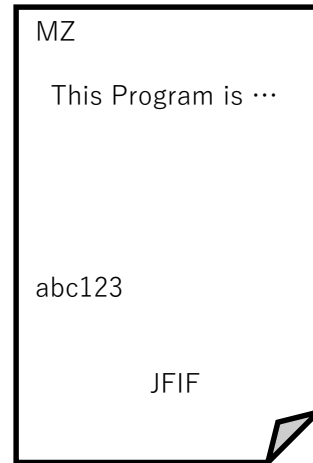
# Author Identification



## Compiler Provenance



## Malicious EXE



Specific Strings

Secret key

Resources

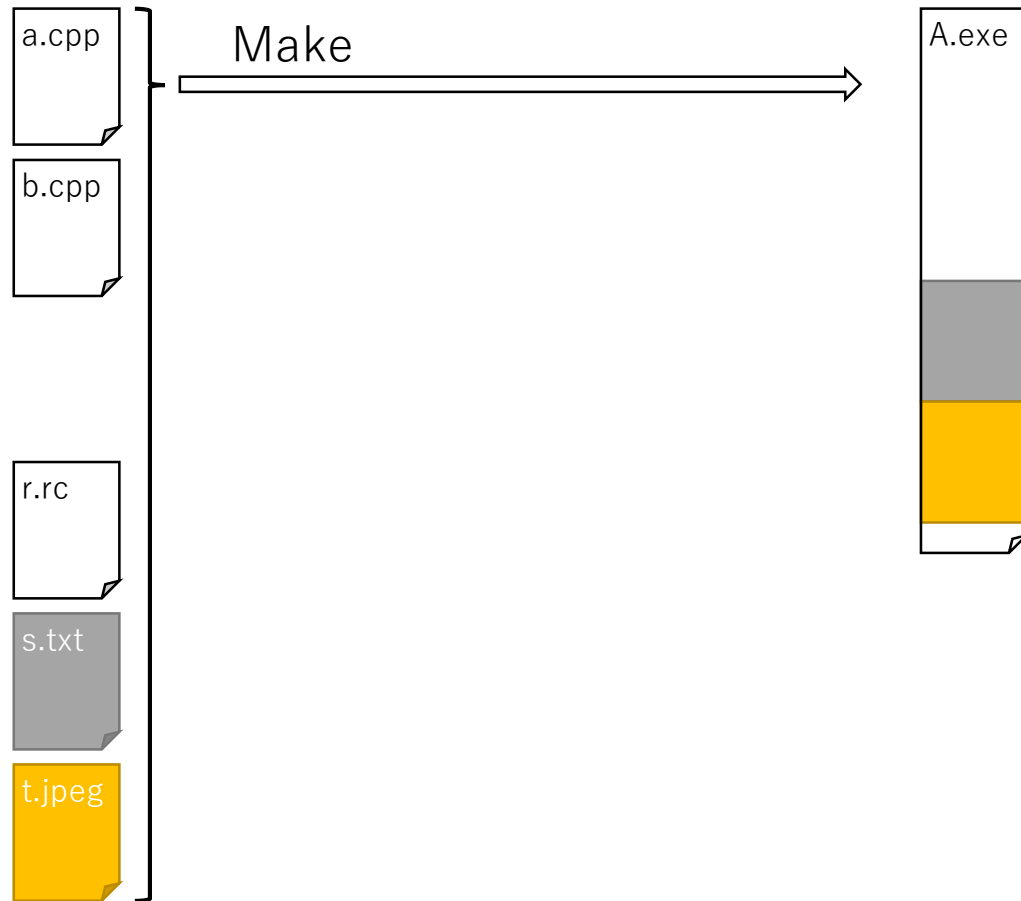
Language

Static Link Libraries

Compile time

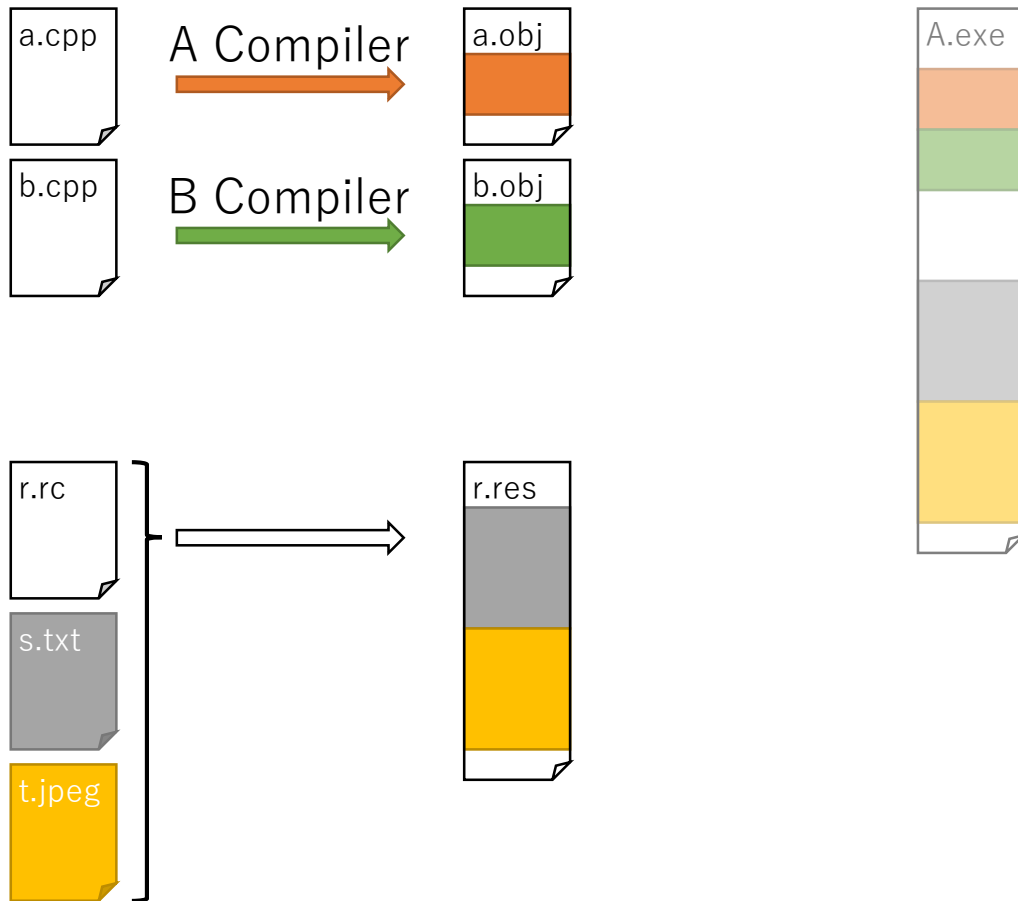
API name

# Multiple Compiler Binary



What is the truth label of **A.exe**?

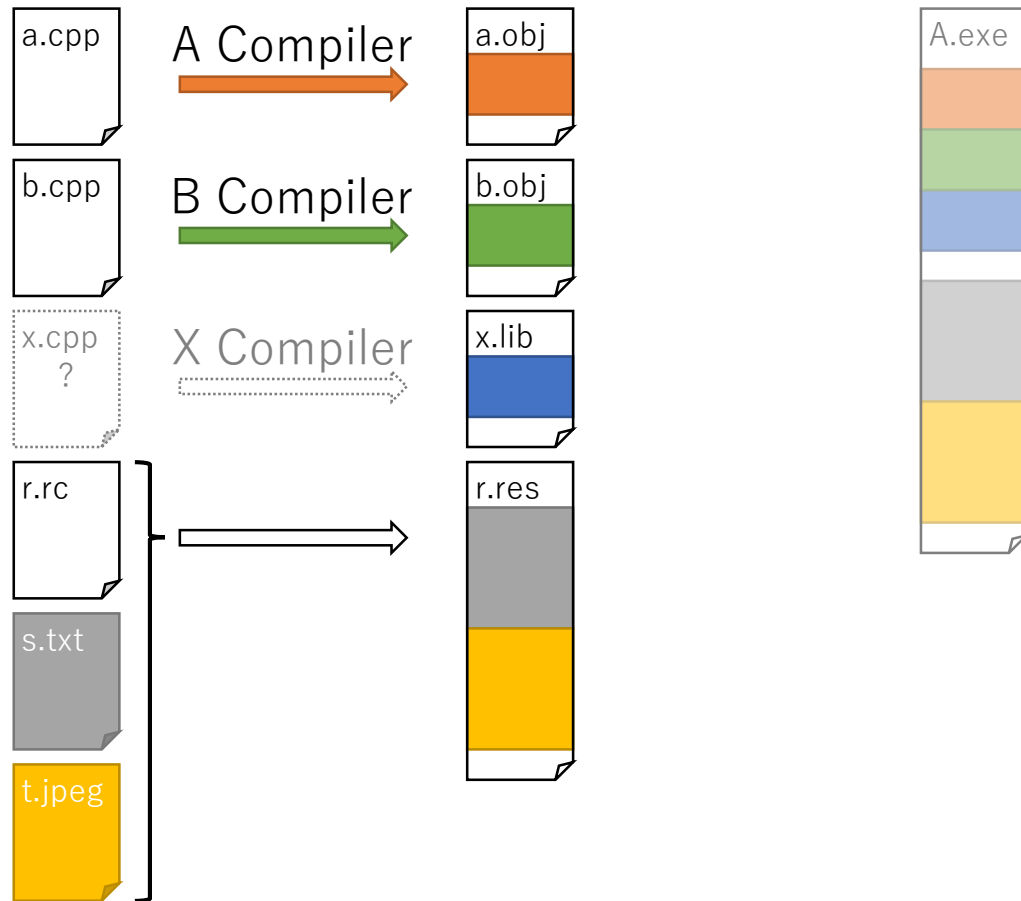
# Multiple Compiler Binary



What is the truth label of **A.exe**?

- **A Compiler?**
- **B Compiler?**

# Multiple Compiler Binary

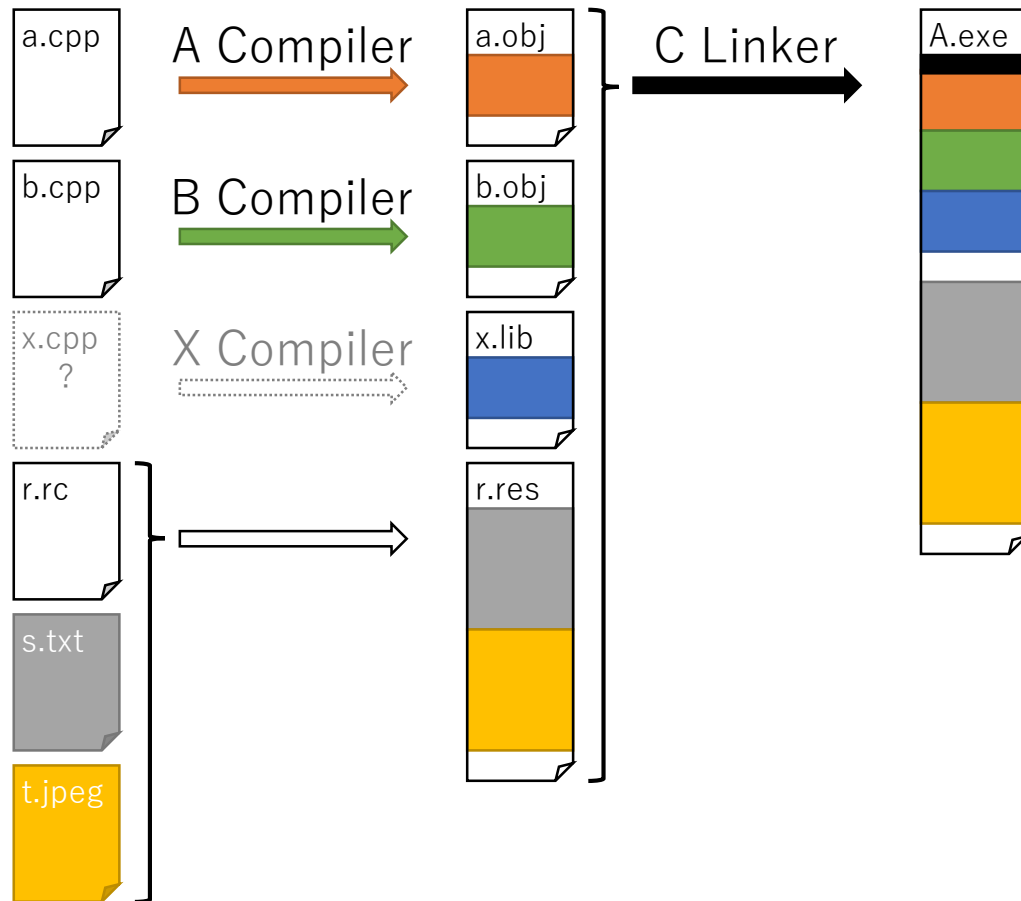


What is the truth label of **A.exe**?

- **A Compiler?**
- **B Compiler?**
- **X Compiler?**



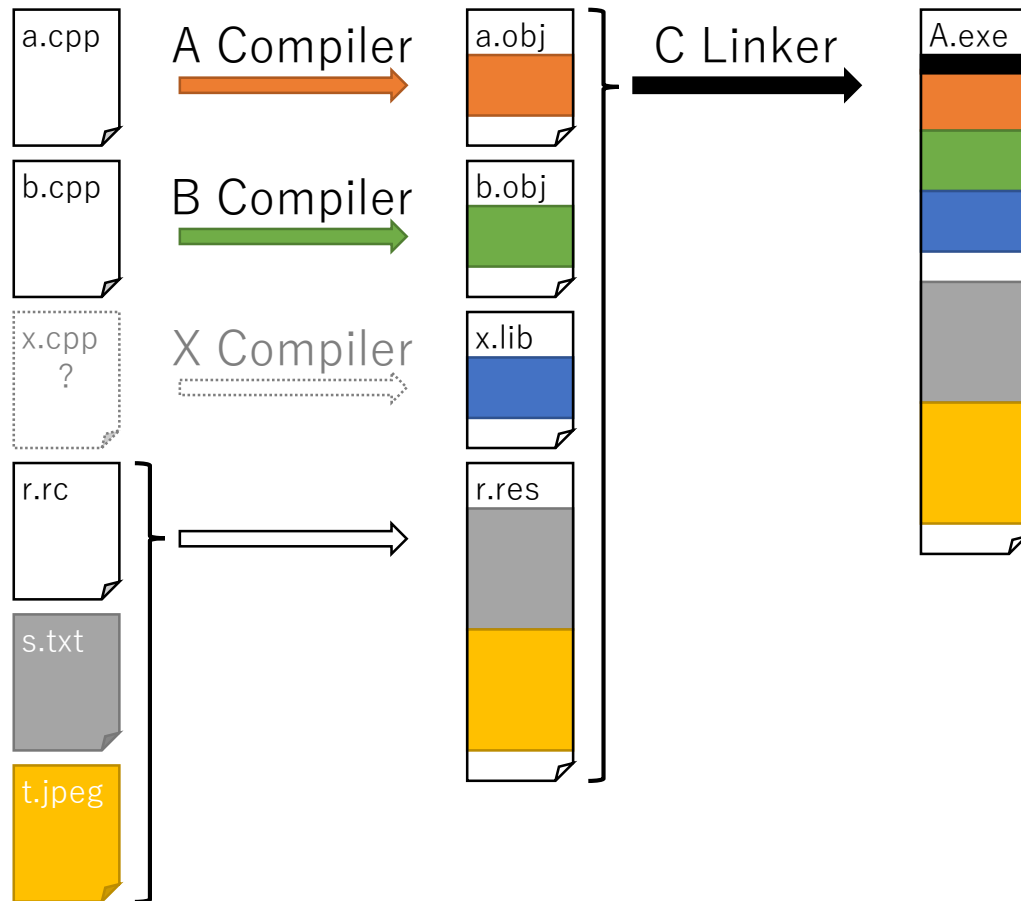
# Multiple Compiler Binary



What is the truth label of **A.exe**?

- **A Compiler?**
- **B Compiler?**
- **X Compiler?**
- **C Linker?**

# Multiple Compiler Binary



What is the truth label of **A.exe**?

- **A Compiler?**
- **B Compiler?**
- **X Compiler?**
- **C Linker?**

What is the truth label of **a.obj**?

- **A Compiler**

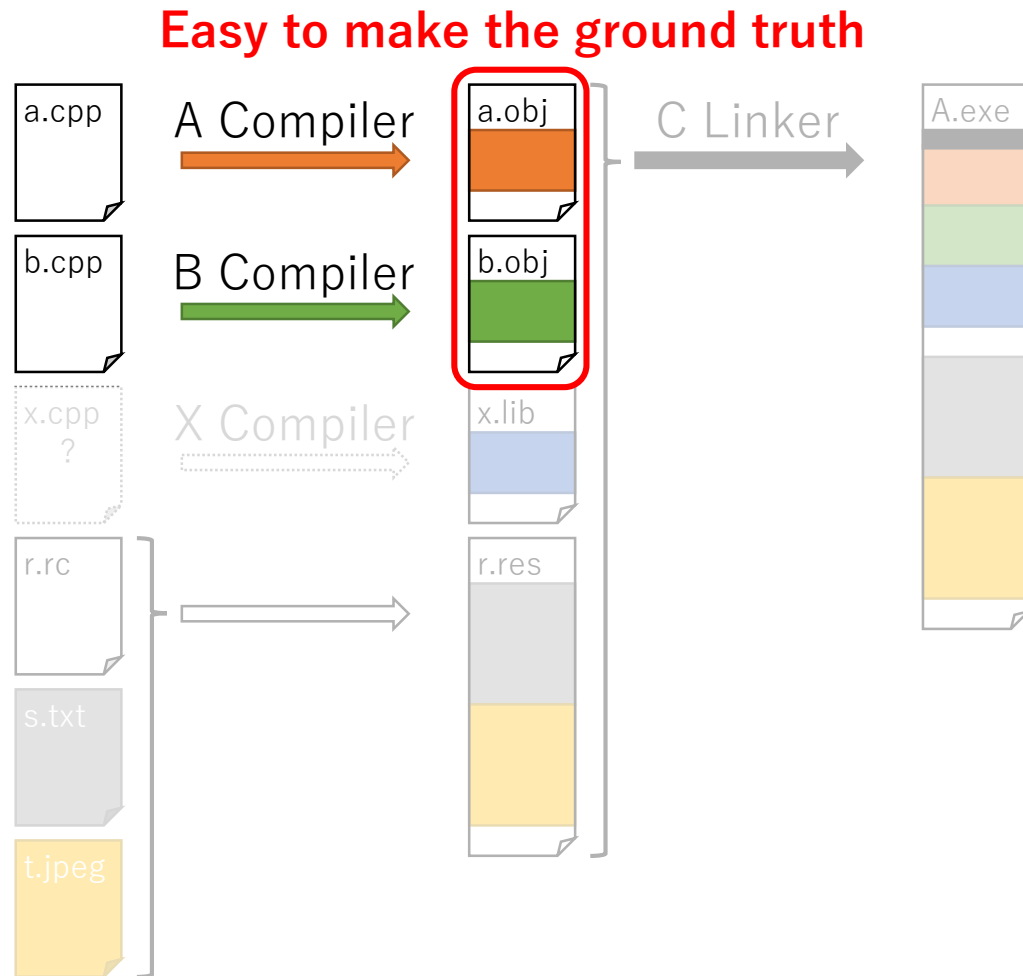
What is the truth label of **b.obj**?

- **B Compiler**

What is the truth label of **x.lib**?

- **Hmm... I think VC, because MS provide it!**

# Multiple Compiler Binary



What is the truth label of A.exe?

- A Compiler?
- B Compiler?
- X Compiler?
- C Linker?

What is the truth label of **a.obj**?

- **A Compiler**

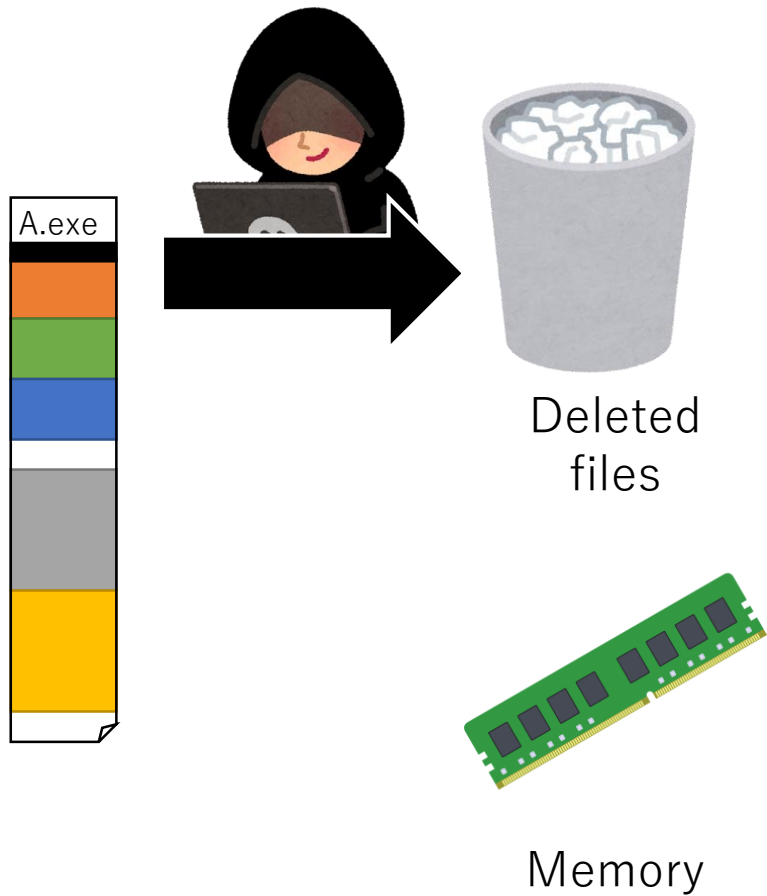
What is the truth label of **b.obj**?

- **B Compiler**

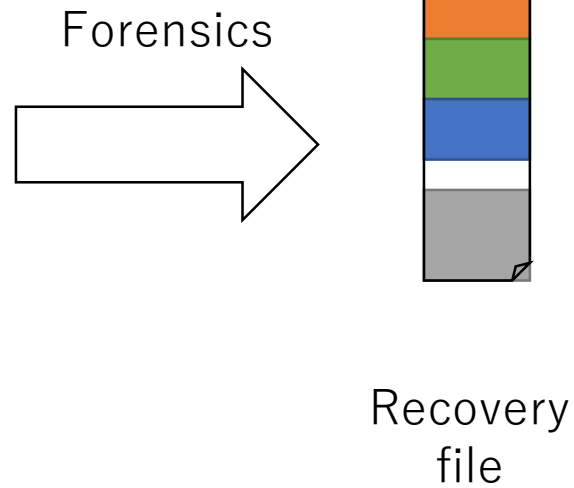
What is the truth label of x.lib?

- Hmm... I think VC, because MS provide it!

# Fragmented Files



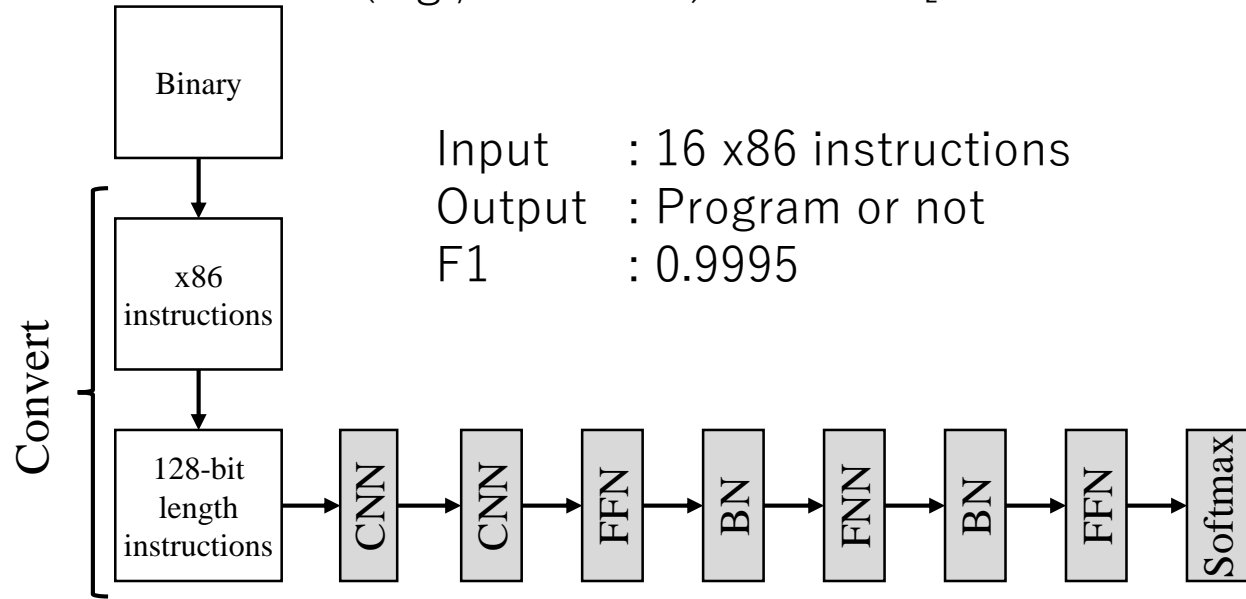
**Collect as much of the attacker's trace as possible even from fragmented files.**



# Preliminaries

# o-glasses

x86 code (e.g., shellcode) detector [arXiv.1806.05328]

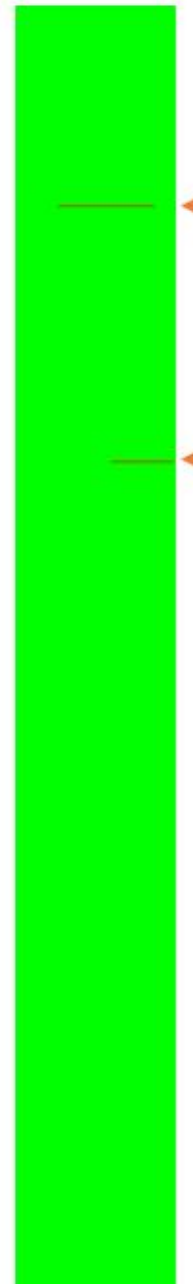


Grayscale

o-glasses  
(1d-CNN)

Correct  
data

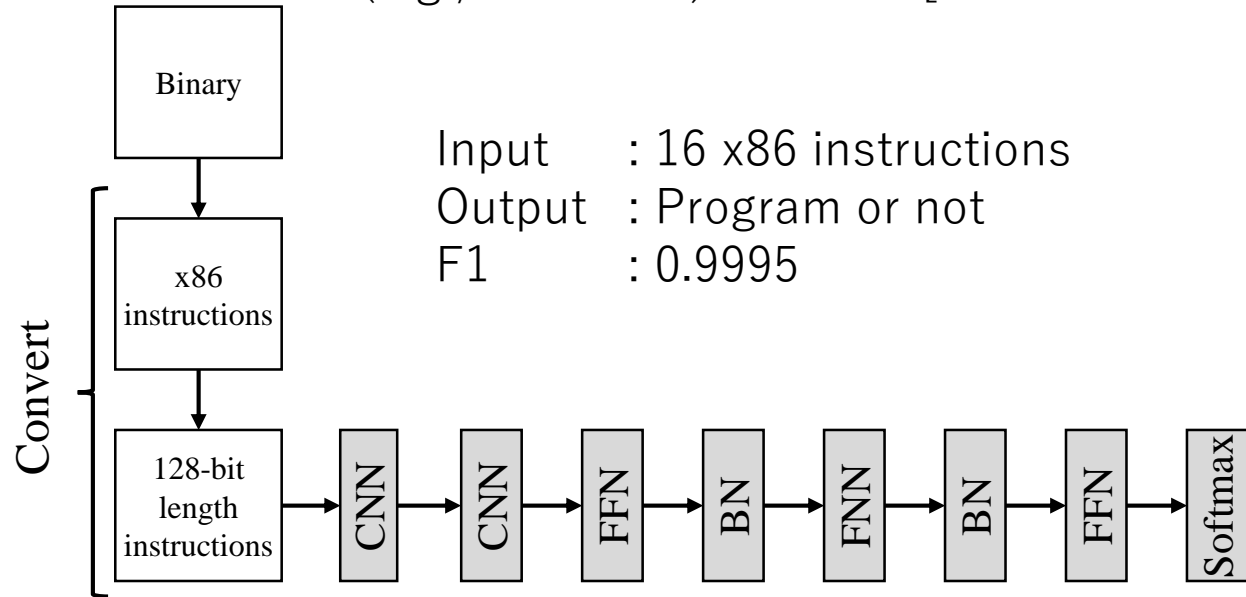
Executable file (compressed)



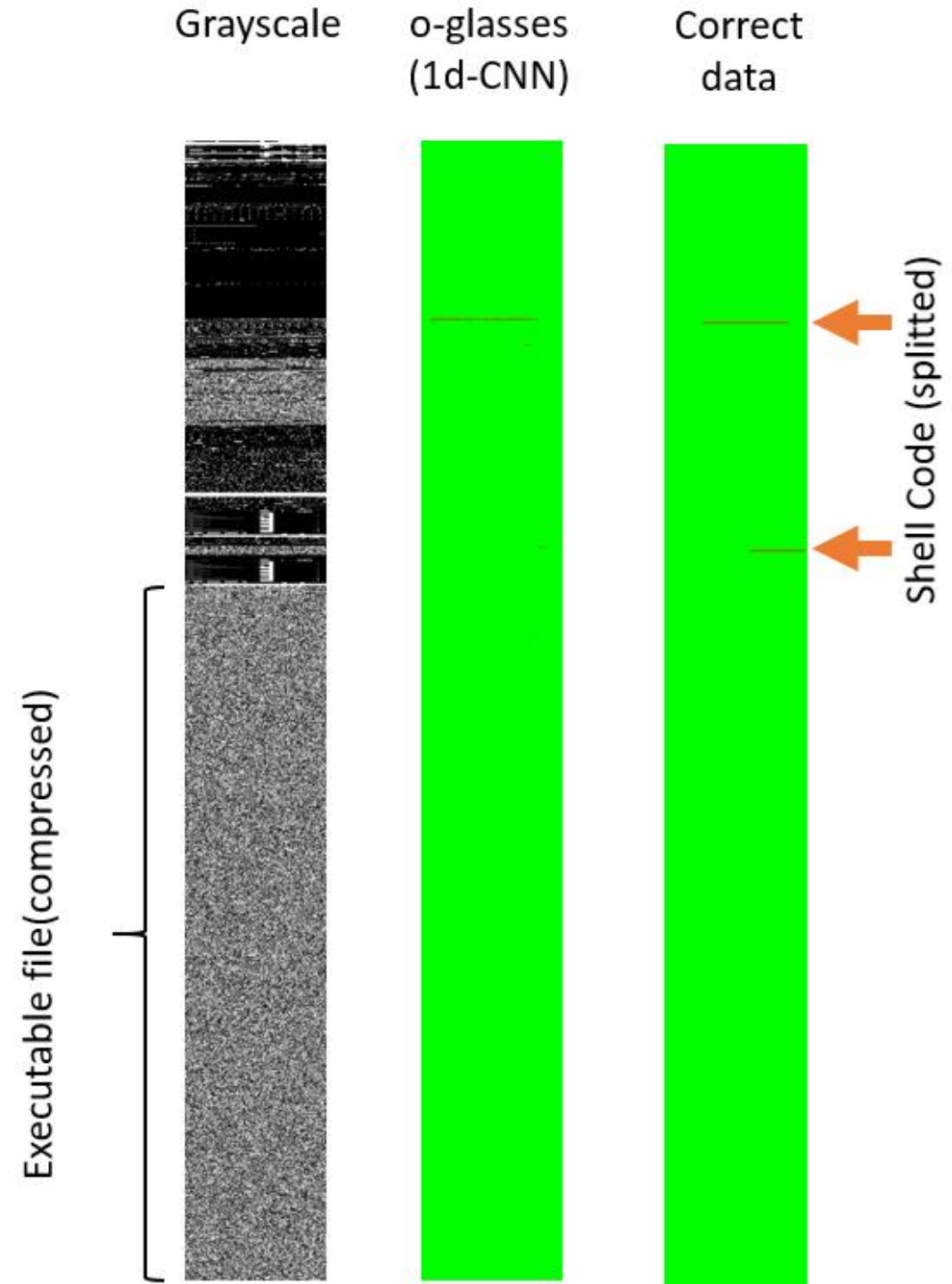
Shell Code (splitted)

# o-glasses

x86 code (e.g., shellcode) detector [arXiv.1806.05328]



- Applying to compiler identification
- Black Box Problem



# Attention Is All You Need

[Łukasz Kaiser et al., NIPS, 2017]

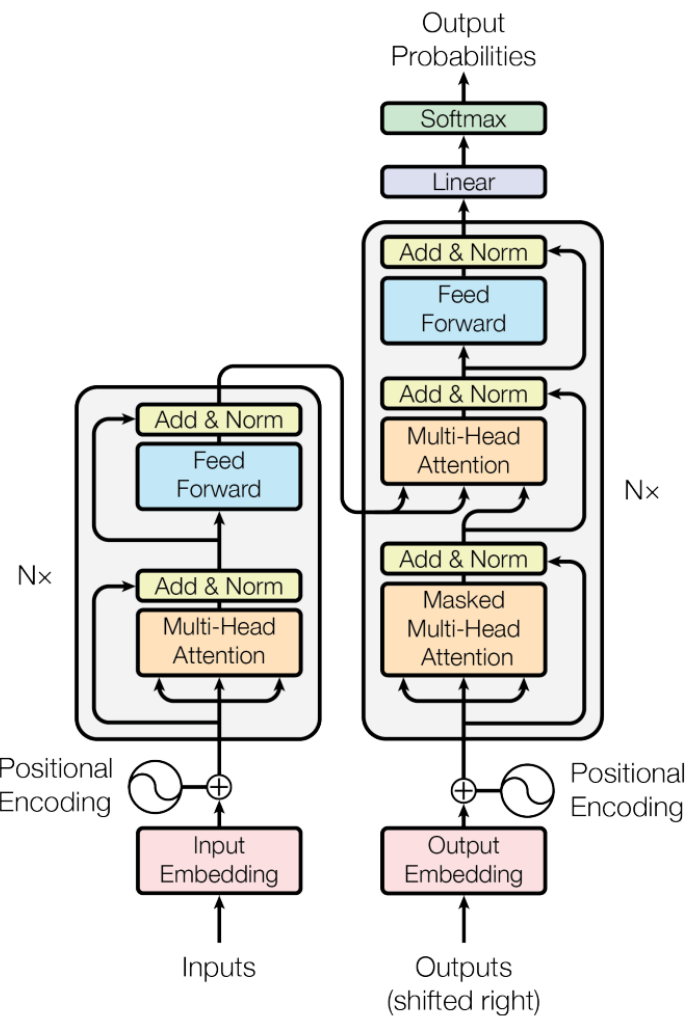
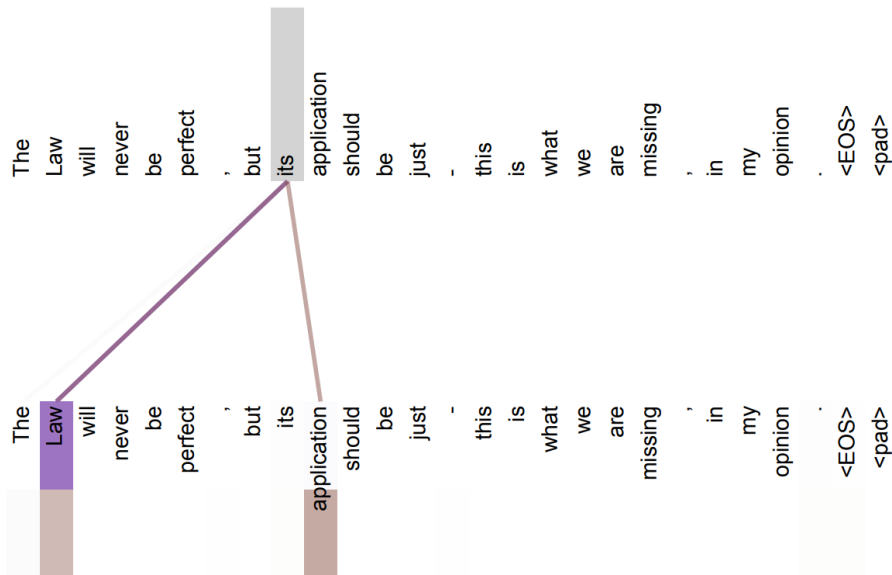


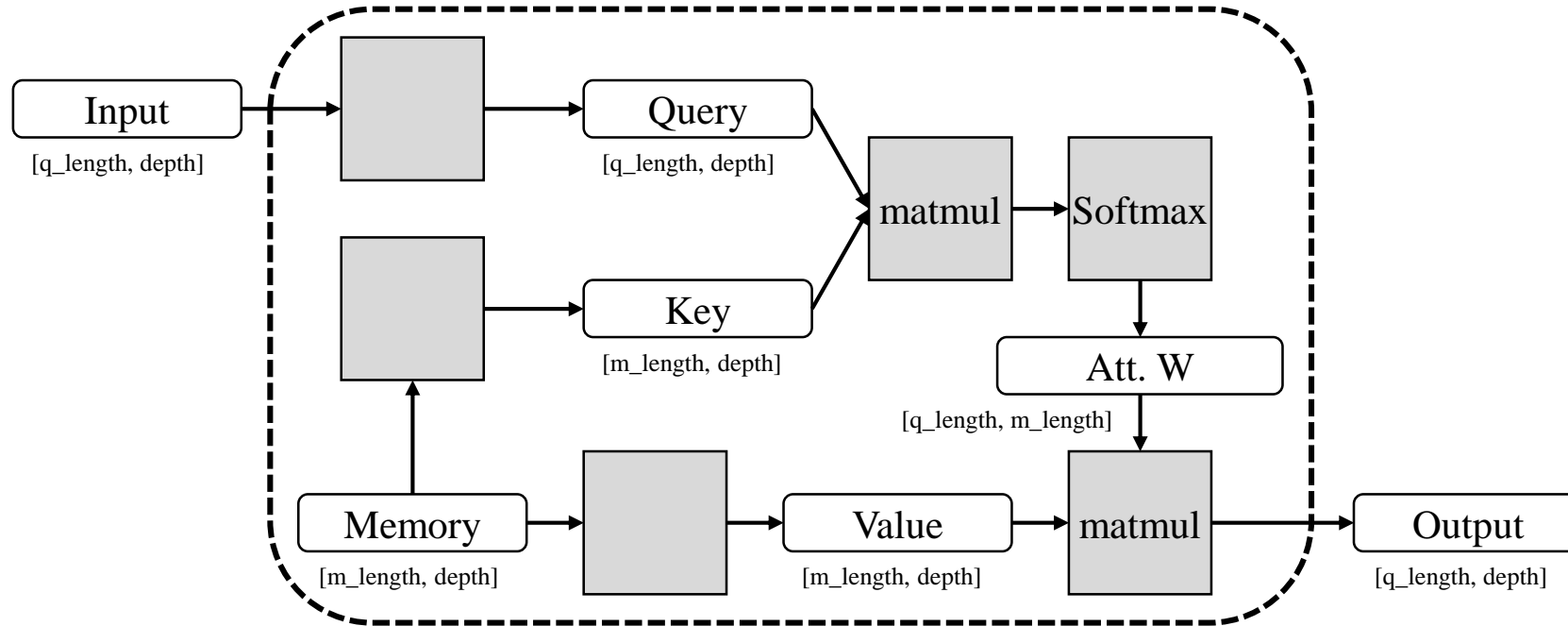
Table 2: The Transformer achieves better BLEU scores than previous state-of-the-art models on the English-to-German and English-to-French newstest2014 tests at a fraction of the training cost.

Model	BLEU		Training Cost (FLOPs)	
	EN-DE	EN-FR	EN-DE	EN-FR
ByteNet [15]	23.75			
Deep-Att + PosUnk [32]		39.2		$1.0 \cdot 10^{20}$
GNMT + RL [31]	24.6	39.92	$2.3 \cdot 10^{19}$	$1.4 \cdot 10^{20}$
ConvS2S [8]	25.16	40.46	$9.6 \cdot 10^{18}$	$1.5 \cdot 10^{20}$
MoE [26]	26.03	40.56	$2.0 \cdot 10^{19}$	$1.2 \cdot 10^{20}$
Deep-Att + PosUnk Ensemble [32]		40.4		$8.0 \cdot 10^{20}$
GNMT + RL Ensemble [31]	26.30	41.16	$1.8 \cdot 10^{20}$	$1.1 \cdot 10^{21}$
ConvS2S Ensemble [8]	26.36	<b>41.29</b>	$7.7 \cdot 10^{19}$	$1.2 \cdot 10^{21}$
Transformer (base model)	27.3	38.1	<b><math>3.3 \cdot 10^{18}</math></b>	
Transformer (big)	<b>28.4</b>	<b>41.0</b>	$2.3 \cdot 10^{19}$	

Figure 1: The Transformer - model architecture.



# Basic of Attention



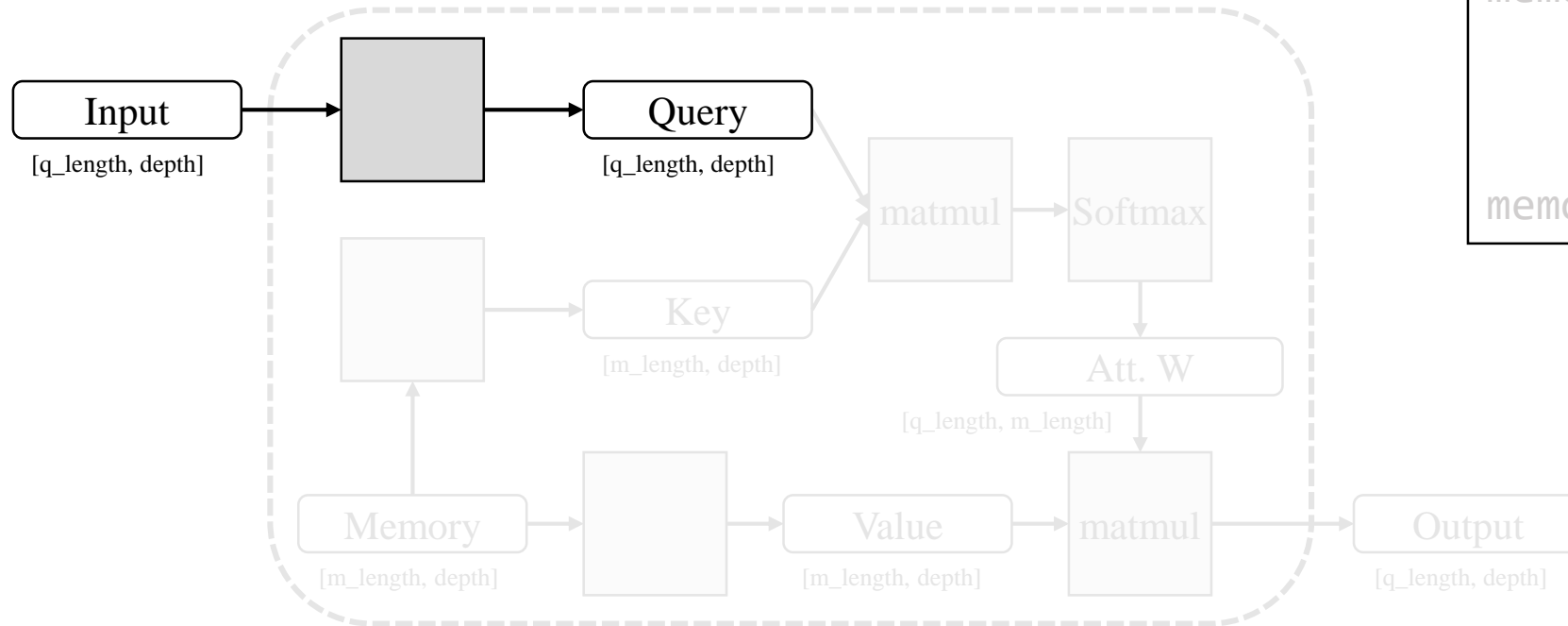
# Basic of Attention

Input

[q\_length, depth]

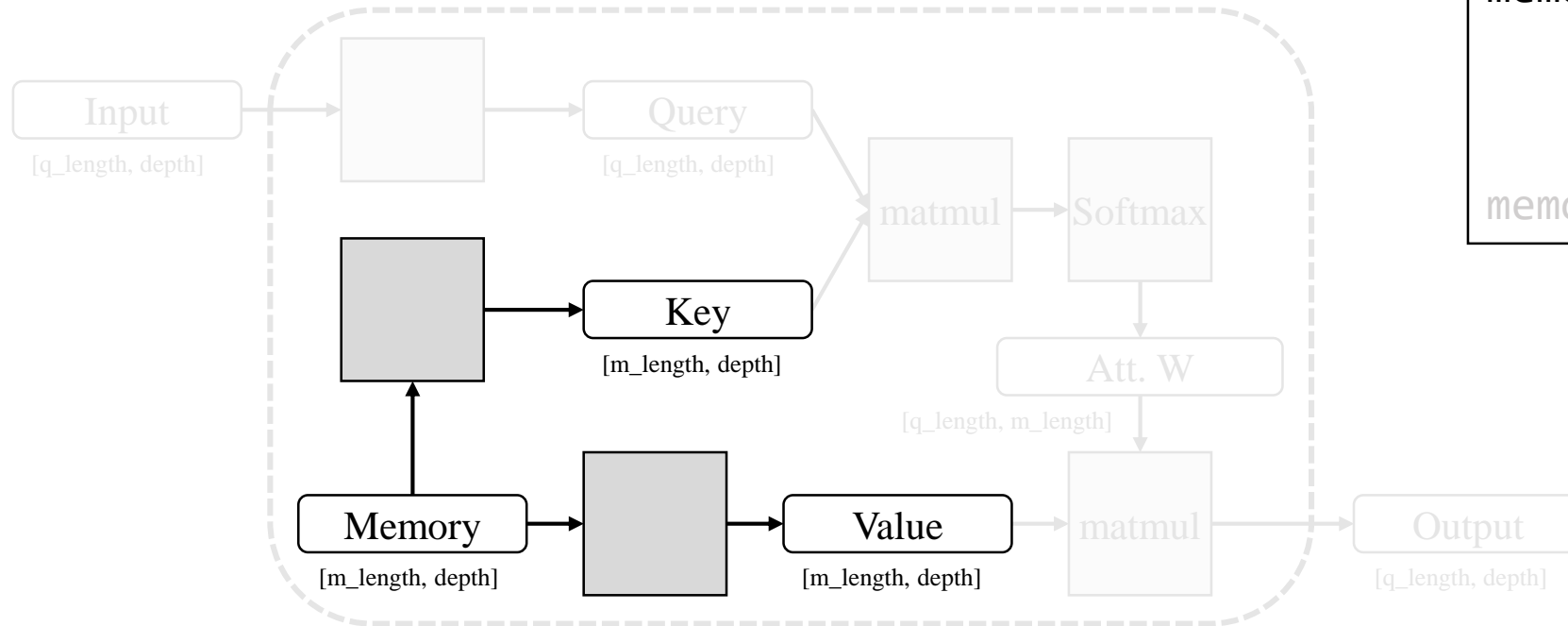
```
query = 'key2'  
memory = {'key1': 'value2',  
          'key2': 'value2',  
          'key3': 'value3',  
          'key4': 'value4'}  
memory[query] = 'value2'
```

# Basic of Attention



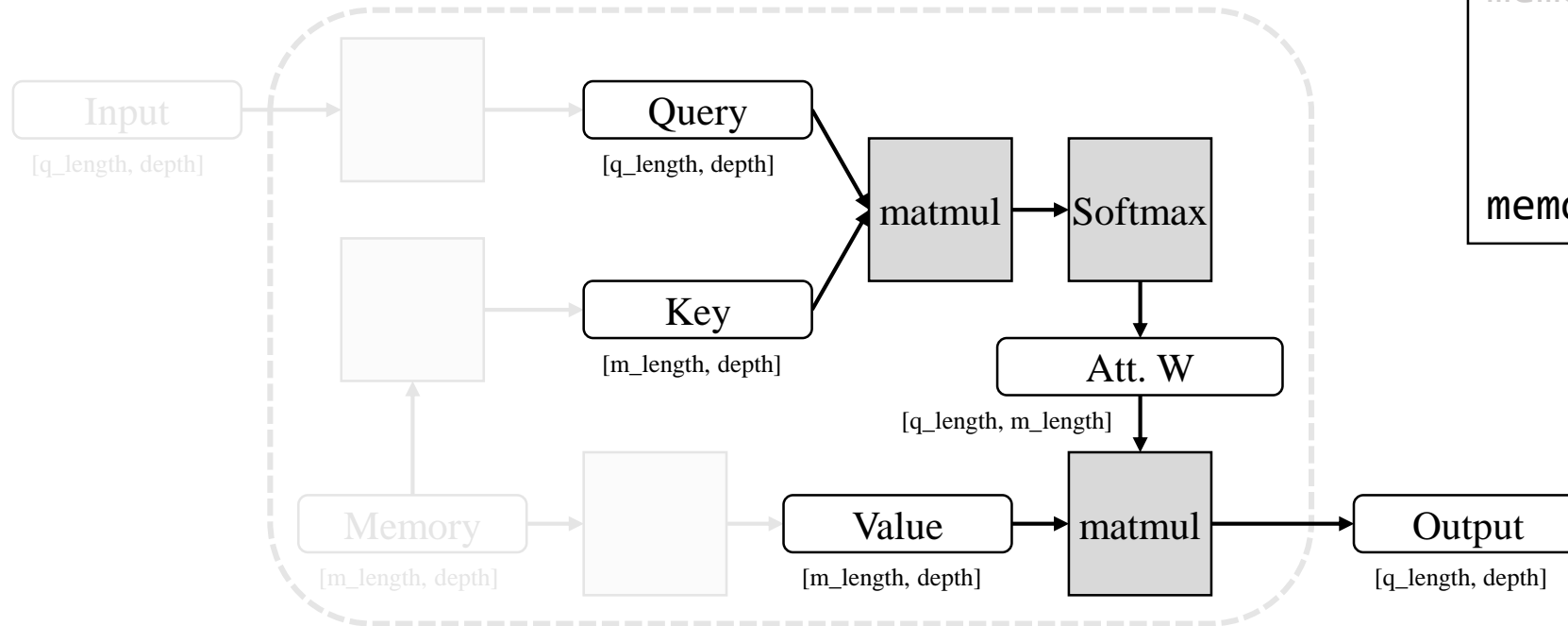
```
query = 'key2'  
memory = {'key1': 'value2',  
          'key2': 'value2',  
          'key3': 'value3',  
          'key4': 'value4'}  
memory[query] = 'value2'
```

# Basic of Attention



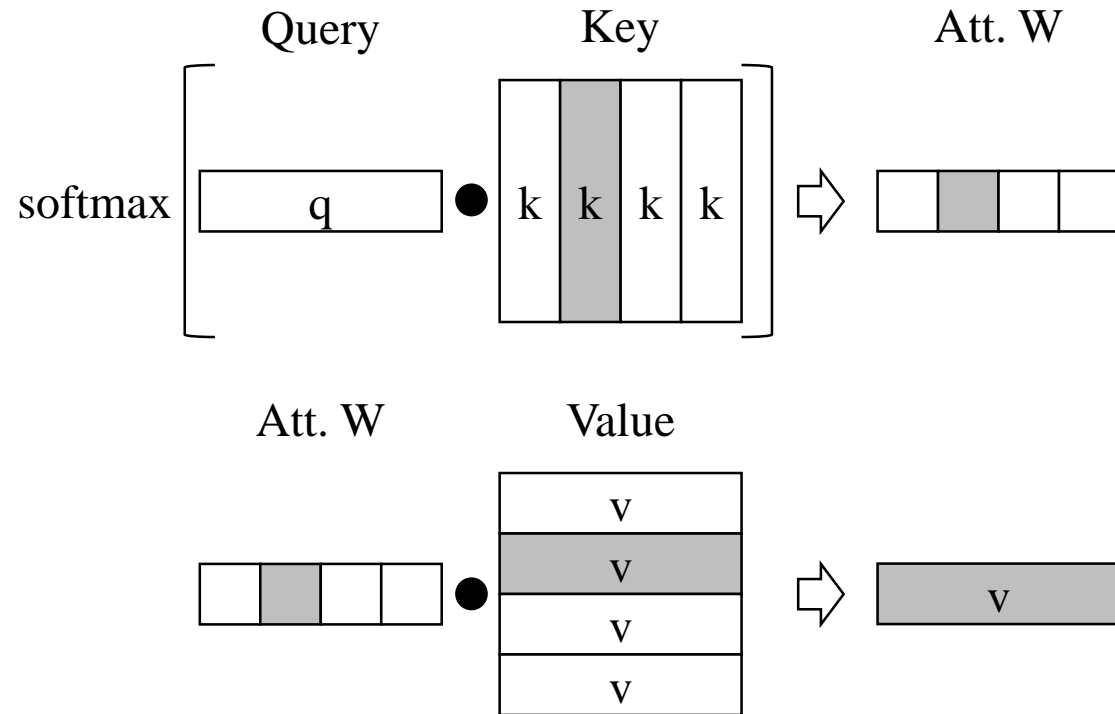
```
query = 'key2'  
memory = {'key1': 'value2',  
          'key2': 'value2',  
          'key3': 'value3',  
          'key4': 'value4'}  
memory[query] = 'value2'
```

# Basic of Attention

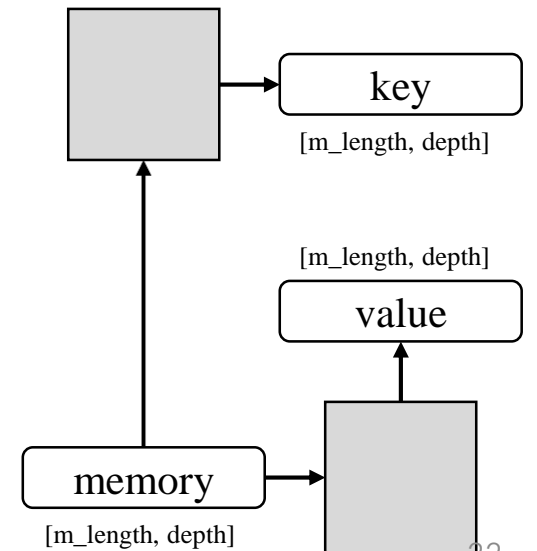


```
query = 'key2'  
memory = {'key1': 'value2',  
          'key2': 'value2',  
          'key3': 'value3',  
          'key4': 'value4'}  
memory[query] = 'value2'
```

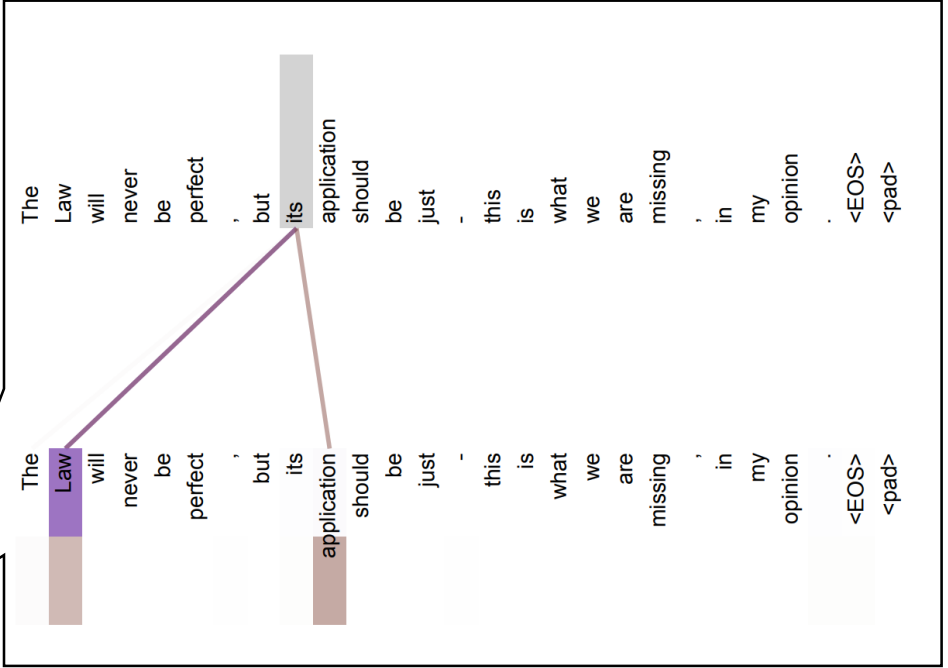
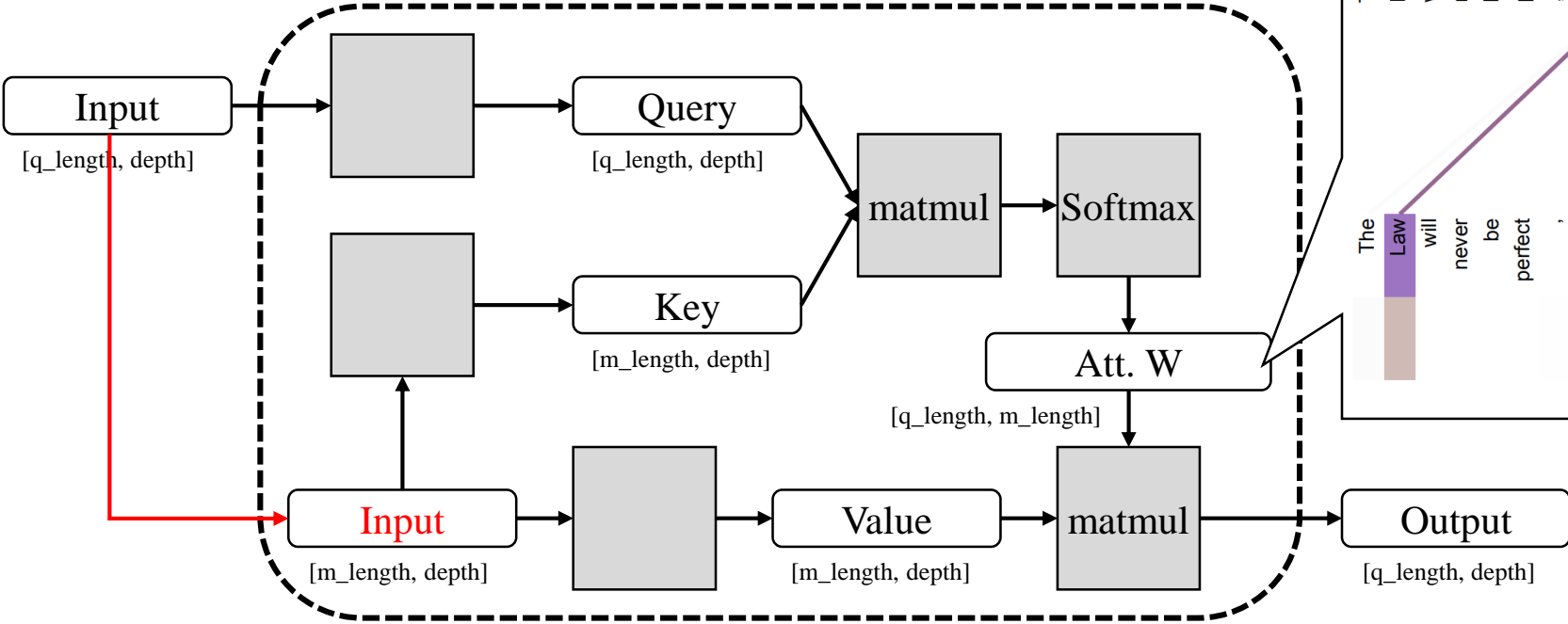
# Dot-Product Attention vs. Dictionary Object



```
query = 'key2'  
memory = {'key1': 'value2',  
          'key2': 'value2',  
          'key3': 'value3',  
          'key4': 'value4'}  
memory[query] = 'value2'
```



# Self-Attention



# Positional Encoding (PE)

PE (Positional Encoding) adds information about the word position to the input word vectors for learning the context of words.

$$Y(X) = X + \alpha PE$$

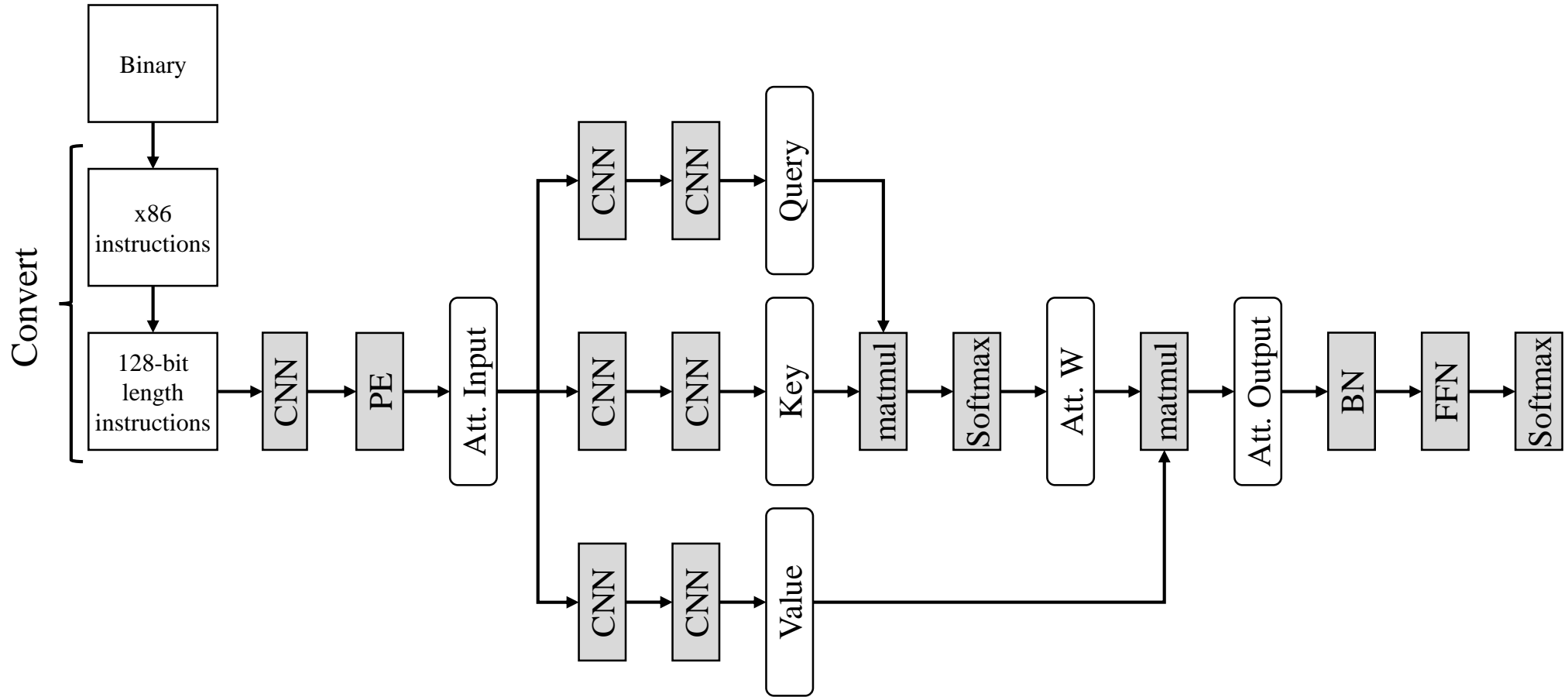
$$PE_{(pos,2i)} = \sin(pos/10000^{2i/d_{model}})$$

$$PE_{(pos,2i+1)} = \cos(pos/10000^{2i/d_{model}})$$

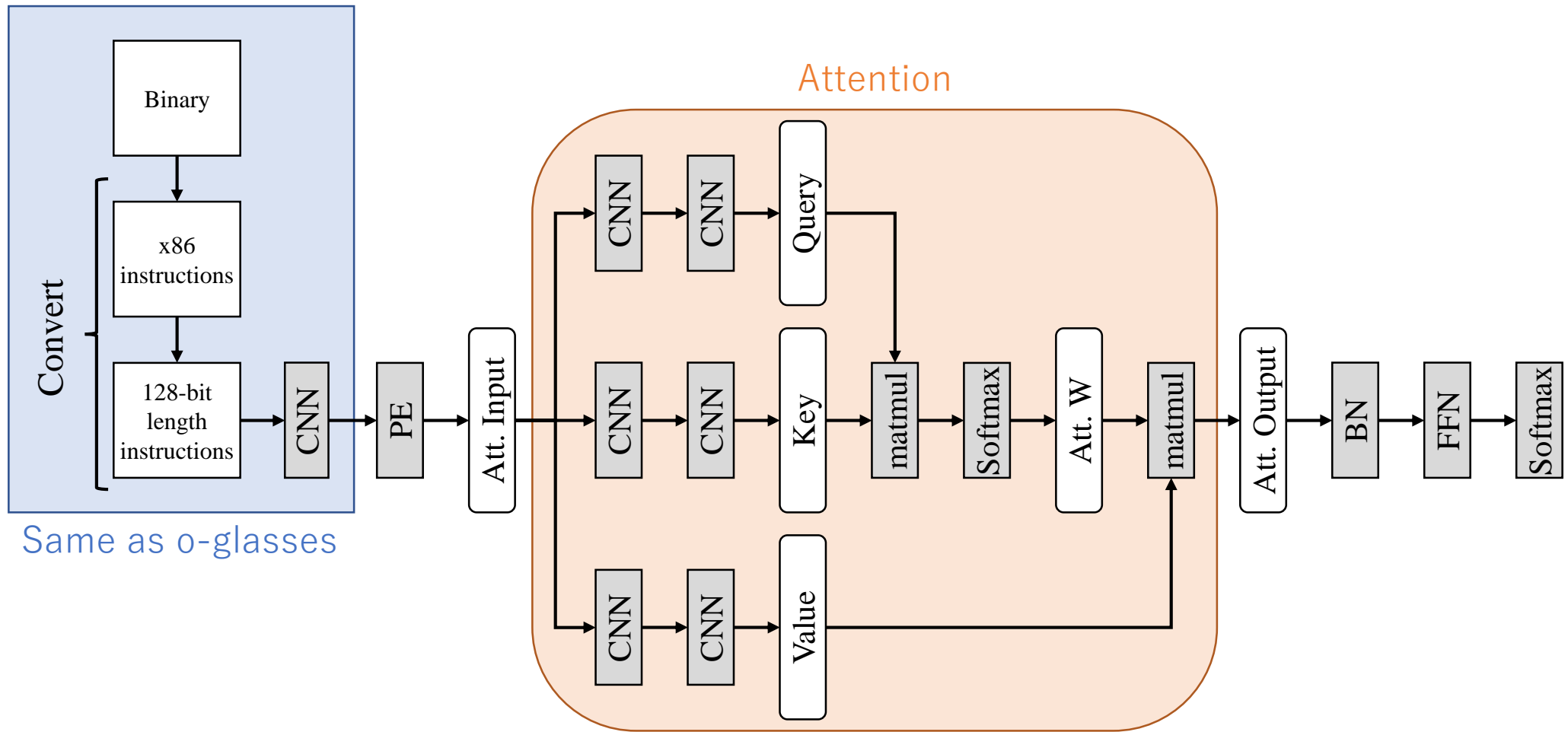


# Proposed Method

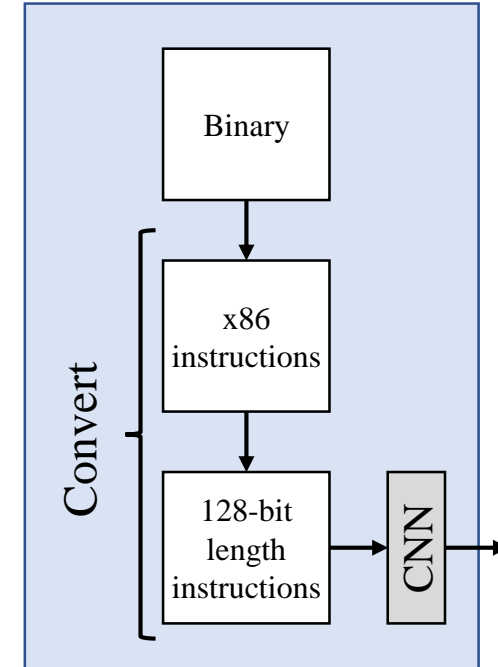
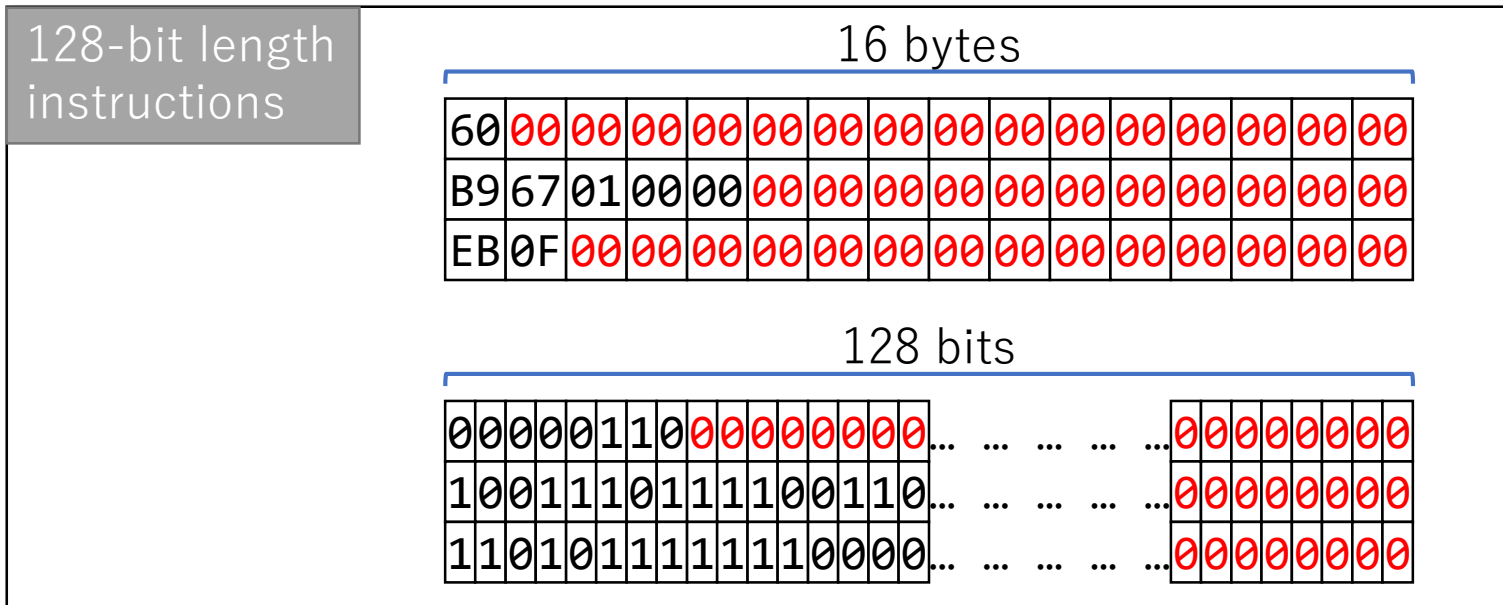
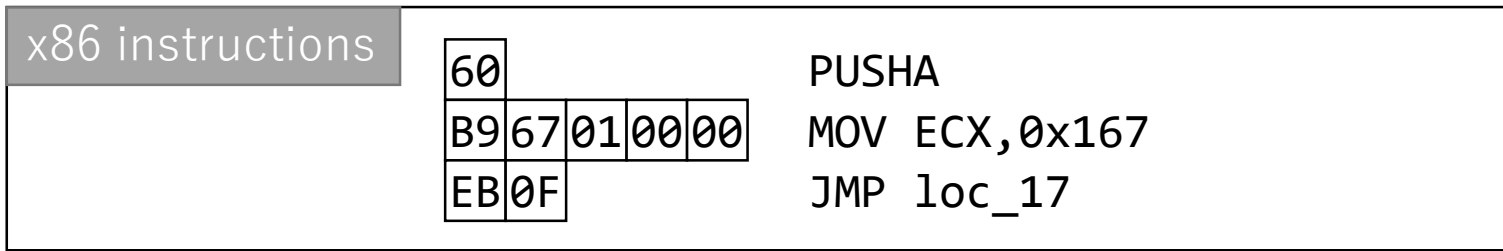
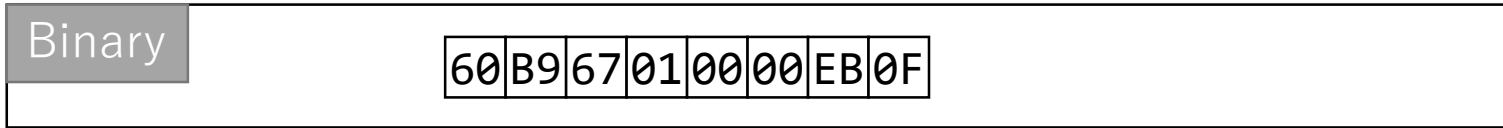
# o-glassesX



# o-glassesX



# Preprocessing details

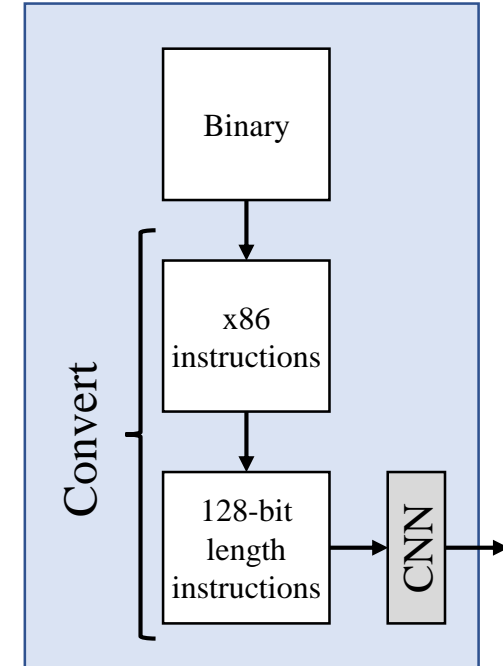
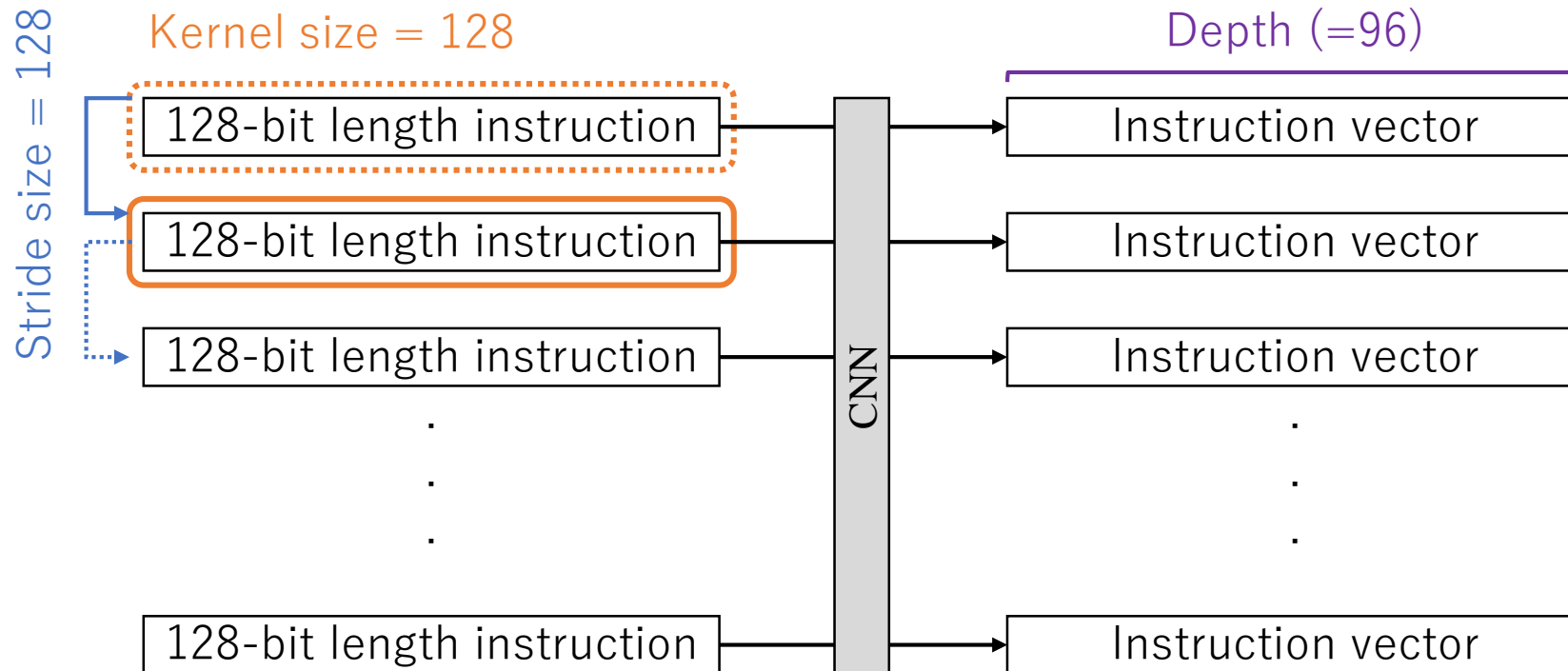


Same as o-glasses

# The 1<sup>st</sup> CNN Layer

Each unit in CNN has specially local connections to the input units, called a Kernel. Every kernel shares the weight parameters with the others in the same layer.

Each kernel covers a single instruction by adjusting the hyperparameters.



Same as o-glasses

# Evaluation

# Training Dataset

Collecting source code files from GitHub  
Compiling various compilers and options

**Total** : 19 labels  
**Compiler** : 4 families  
 Visual C++, GCC, Clang and Intel C++ Compiler  
**Opt. level** : 2 types  
 maximum or not  
**CPU Arc.** : 2 types  
 x86 or x86-64

TABLE II. DATASET COMPILATION SETTINGS.

Compiler family	Version	Architecture		Optimization Level	
		x86	x86-64	Low	High
VC	2003	✓		-Od	-Ox
	2017	✓	✓	-Od	-Ox
GCC	6.3.0	✓	✓	-O0	-O3
Clang	5.0.2	✓	✓	-O0	-O3
ICC	19.0.0.117	✓	✓	-O0	-O3

		Label	#Binaries	#Code
Program	VC2017	x86 VC17,32,none(Od)	1,170	369,605
		x86 VC17,32,max(Ox)	1,147	255,143
		x86-64 VC17,64,none(Od)	1,456	540,568
		x86-64 VC17,64,max(Ox)	1,242	542,020
Program	VC2003	x86 VC03,32,none(Od)	1,350	292,277
		x86 VC03,32,max(Ox)	1,306	270,743
		x86-64 -	-	-
		x86-64 -	-	-
Program	GCC	x86 GCC,32,none(O0)	2,111	227,004
		x86 GCC,32,max(O3)	1,844	239,821
		x86-64 GCC,64,none(O0)	1,582	283,276
		x86-64 GCC,64,max(O3)	1,580	287,775
Program	Clang	x86 Clang,32,none(O0)	1,205	101,024
		x86 Clang,32,max(O3)	1,196	86,521
		x86-64 Clang,64,none(O0)	1,892	332,278
		x86-64 Clang,64,max(O3)	1,883	246,500
Program	ICC	x86 ICC,32,none(Od)	1,761	1,494,677
		x86 ICC,32,max(Ox)	1,724	1,161,499
		x86-64 ICC,64,none(Od)	1,796	1,419,705
		x86-64 ICC,64,max(Ox)	1,728	1,046,958
Others			101	912,855
Total			28,074	10,110,249

# 4-fold Cross Validation (Input: 64 instructions)

Train \ Predict	VC						GCC				Clang				ICC				Others	R	P	F1	
	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19				
VC03,32,none(Od)	1	.9604	.0355	.0026	.0002	.0002	.0002	.0001	.0000	.0000	.0000	.0001	.0000	.0000	.0000	.0001	.0002	.0000	.0000	.0004	.9500	.9604	.9552
VC17,32,none(Od)	2	.0463	.9517	.0002	.0007	.0001	.0001	.0000	.0000	.0000	.0000	.0001	.0000	.0000	.0000	.0000	.0002	.0000	.0000	.0006	.9625	.9517	.9570
VC03,32,max(Ox)	3	.0026	.0001	.9875	.0061	.0001	.0003	.0000	.0002	.0000	.0000	.0000	.0002	.0000	.0000	.0000	.0022	.0000	.0000	.0006	.9786	.9875	.9830
VC17,32,max(Ox)	4	.0004	.0010	.0144	.9774	.0001	.0005	.0000	.0001	.0000	.0001	.0001	.0008	.0000	.0001	.0000	.0044	.0000	.0000	.0007	.9887	.9774	.9830
VC17,64,none(Od)	5	.0002	.0001	.0002	.0001	.9978	.0008	.0000	.0000	.0000	.0000	.0000	.0000	.0000	.0001	.0000	.0001	.0001	.0001	.0004	.9977	.9978	.9977
VC17,64,max(Ox)	6	.0002	.0001	.0004	.0004	.0013	.9931	.0000	.0000	.0000	.0005	.0000	.0000	.0001	.0006	.0000	.0001	.0000	.0023	.0008	.9955	.9931	.9943
GCC,32,none(O0)	7	.0002	.0000	.0001	.0000	.0000	.0000	.9973	.0009	.0004	.0001	.0003	.0004	.0000	.0000	.0002	.0000	.0000	.0000	.0000	.9972	.9973	.9973
GCC,32,max(O3)	8	.0001	.0000	.0005	.0002	.0000	.0000	.0011	.9921	.0000	.0003	.0002	.0051	.0000	.0000	.0000	.0004	.0000	.0000	.0000	.9946	.9921	.9933
GCC,64,none(O0)	9	.0001	.0001	.0000	.0000	.0000	.0000	.0008	.0000	.9970	.0006	.0000	.0000	.0010	.0002	.0000	.0000	.0001	.0000	.0000	.9979	.9970	.9975
GCC,64,max(O3)	10	.0000	.0000	.0000	.0000	.0001	.0002	.0001	.0006	.0003	.9800	.0000	.0003	.0003	.0168	.0000	.0001	.0001	.0012	.0000	.9819	.9800	.9809
Clang,32,none(O0)	11	.0003	.0002	.0000	.0000	.0000	.0000	.0004	.0004	.0000	.0001	.9972	.0006	.0007	.0000	.0001	.0000	.0000	.0000	.0000	.9959	.9972	.9965
Clang,32,max(O3)	12	.0001	.0000	.0006	.0011	.0000	.0001	.0004	.0068	.0000	.0003	.0006	.9869	.0000	.0019	.0000	.0012	.0000	.0000	.0000	.9790	.9869	.9829
Clang,64,none(O0)	13	.0000	.0000	.0000	.0000	.0000	.0001	.0000	.0000	.0009	.0003	.0010	.0000	.9973	.0003	.0000	.0000	.0000	.0001	.0000	.9979	.9973	.9976
Clang,64,max(O3)	14	.0000	.0000	.0001	.0001	.0001	.0006	.0000	.0001	.0002	.0147	.0000	.0015	.0002	.9810	.0000	.0001	.0001	.0013	.0001	.9796	.9810	.9803
ICC,32,none(Od)	15	.0001	.0000	.0001	.0000	.0000	.0000	.0002	.0000	.0000	.0000	.0000	.0000	.0000	.0000	.9994	.0002	.0000	.0000	.0000	.9995	.9994	.9994
ICC,32,max(Ox)	16	.0002	.0001	.0029	.0028	.0000	.0001	.0000	.0003	.0000	.0001	.0000	.0005	.0000	.0000	.0001	.9930	.0000	.0001	.0000	.9916	.9930	.9923
ICC,64,none(Od)	17	.0000	.0000	.0000	.0000	.0001	.0000	.0000	.0000	.0000	.0000	.0000	.0000	.0000	.0000	.0000	.0000	.9996	.0001	.0000	.9995	.9996	.9995
ICC,64,max(Ox)	18	.0000	.0000	.0001	.0000	.0001	.0016	.0000	.0001	.0000	.0012	.0000	.0000	.0000	.0013	.0000	.0001	.0002	.9952	.0000	.9948	.9952	.9950
Others	19	.0000	.0000	.0000	.0001	.0000	.0000	.0000	.0000	.0000	.0000	.0000	.0000	.0000	.0000	.0000	.0000	.0000	.0000	.9997	.9964	.9997	.9980



# Comparison of Performance of Related Work

		o-glassesX		o-glasses[23]		Rosenblum's[29]	ORIGIN[26]		BinComp[25]	
		(L=64)	(L=16)	(L=64)	(L=16)		(SVM)	(CRF)		
Accuracy	All components	.9884 (19)	.9555 (19)	.9421 (19)	.9323 (19)	.924 (3)	.604 (18)	.918 (18)	.801 (6)	
	Compiler family	.9886 (4)	.9670 (4)	.9140 (4)	.9271 (4)	.924 (3)	.983 (3)	.999 (3)	-	
	Optimization	.9989 (2)	.9943 (2)	.9830 (2)	.9864 (2)	-	.971 (2)	.999 (2)	.917 (2)	
	Architecture	.9997 (2)	.9985 (2)	.9959 (2)	.9940 (2)	-	-	-	-	
	Code/Non-code	.9999 (2)	.9995 (2)	.9991 (2)	.9987 (2)	-	-	-	-	
ML model		Attention	Attention	CNN	CNN	CRF	SVM	CRF	(k-means)	
Dataset	Features	64 Instructions	16 Instructions	64 Instructions	16 Instructions	Byte Seq.	1 Function	Function Seq.	1 File	
	#Samples	1,793,478	1,900,000	471,124	1,886,521	81,886,169	955,000	955,000	1,177	
	#Binaries	28,074	28,074	28,074	28,074	1,119	2,686	2,686	1,177	
	#Variations	VC	6	6	6	6	1	6	6	2
		GCC	4	4	4	4	1	8	8	2
		Clang	4	4	4	4	-	-	-	-
		ICC	4	4	4	4	1	4	4	2
		Non-code	1	1	1	1	-	-	-	-
K-fold cross-validation		4	4	4	4	(Anomalous)	(Anomalous)	(Anomalous)	10	

# Calculating ‘Why’ with the Attention Mechanism

	ENTER 0x558b, 0x8	.222
	MOV [EDX+0xc], ECX	.115
	CMP DWORD [EBP+0x10], 0x80	.029
	JNZ 0x40001b7	.016
	MOV EAX, [EBP+0x8]	.019
	MOV ECX, [EAX+0xc]	.079
	MOV [EBP-0x8], ECX	.036
Input	MOV EDX, [EBP-0x8]	.041
	SHR EDX, 0x10	.063
	AND EDX, 0xff	.255
	MOV EAX, [EDX*4+0x0]	.153
	AND EAX, 0xff000000	.044
	MOV ECX, [EBP+0x8]	.062
	XOR EAX, [ECX]	.207
	MOV EDX, [EBP-0x8]	.039
	SHR EDX, 0x8	.063

## Algorithm 1 Automatic feature extraction procedure

```
1: InputFile ← Compile(SourceFile, Optimization).pullCode()
2: Offset ← 0
3: while Offset < InputFile.EndOffset do
4:   Input ← InputFile.pull(Offset, L)
5:   Result ← model.predictor(Input)
6:   if Result.Label = RealLabel then
7:     if Result.Confidence > 0.99 then
8:       i ← argmax(L2(Result.AttentionWeight))
9:       Count[input[i]] ++
10:    end if
11:  end if
12:  Offset ++
13: end while
14: Print(Ranking(Count))
```

# Typical Instructions for each compiler

e.g.,  
when focusing on the NOP instruction...

<pre>VC03,32,none(Od) 1 XOR ECX, [EAX*4+0x0] 2 AND EAX, 0xff000000 3 XOR EDX, [EAX*4+0x0] 4 AND ECX, 0xff0000 5 XOR EAX, [EDX*4+0x0]</pre>	<pre>VC17,32,none(Od) 1 IMUL ECX, EAX, 0x0 2 IMUL ECX, EAX, 0x9 3 IMUL EDX, ECX, 0xc 4 XOR ECX, EBP 5 IMUL ECX, EAX, 0x3</pre>	<pre>GCC,32,none(O0) 1 AND EAX, 0xff000000 2 MOVZX EAX, AL 3 LEA EDX, [EAX*4+0x0] 4 ADD DWORD [EBP-0x8], 0x1 5 AND EAX, 0xff0000</pre>	<pre>Clang,32,none(O0) 1 AND EAX, 0xff000000 2 NOP WORD [EAX+EAX+0x0] 3 AND ECX, 0xff0000 4 XOR EAX, [ECX*4+0x1028] 5 NOP DWORD [EAX+0x0]</pre>	<pre>ICC,32,none(Od) 1 IMUL EAX, EAX, 0x4 2 AND EAX, 0xff000000 3 IMUL EDX, EDX, 0x4 4 CALL 0x4003c76 5 MOV EAX, 0xffffffff</pre>
<pre>VC03,32,max(Ox) 1 AND ESI, 0xff0000 2 AND ECX, 0xff0000 3 AND EDI, 0xff 4 XOR ECX, EBP 5 AND ECX, 0xff</pre>	<pre>VC17,32,max(Ox) 1 NOP 2 MOVZX ECX, BYTE [EAX+EBP] 3 PUSH DWORD [ESP+0x24] 4 SUB DWORD [ESP+0x2c], 0x1 5 XOR EAX, ESP</pre>	<pre>GCC,32,max(O3) 1 LEA ESI, [ESI+0x0] 2 LEA EDI, [EDI+0x0] 3 AND EAX, 0xff000000 4 AND EAX, 0xff0000 5 AND EBX, 0xff0000</pre>	<pre>Clang,32,max(O3) 1 NOP WORD [CS:EAX+EAX+0x0] 2 MOV EBP, 0xff000000 3 MOV EDX, 0xff00 4 MOV DL, [ESP+0x15] 5 MOV ESI, 0xff00</pre>	<pre>ICC,32,max(Ox) 1 NOP DWORD [EAX+0x0] 2 LEA EBX, [EDX+EDX] 3 LEA ECX, [ESI+ESI] 4 LEA EDX, [ESI+ESI] 5 PXOR XMM0, [ESP+0x20]</pre>

Fig. 5. Typical instructions for each compiler against aes.c

# Case Study: Various Optimization Levels

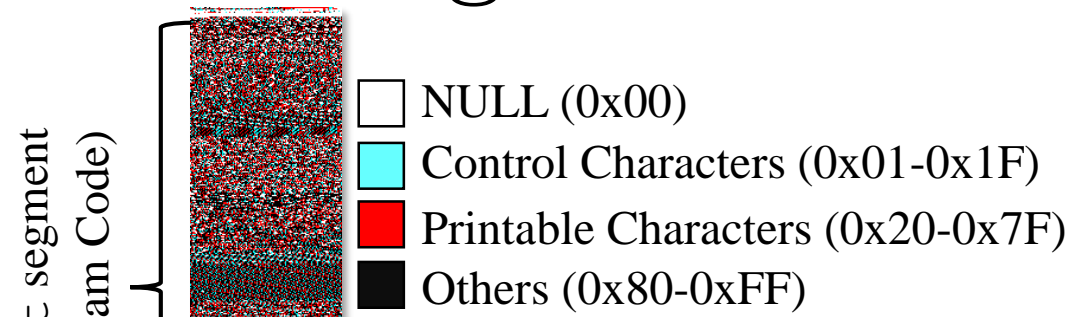
Function name	Segment	Start	Length
f sub_401010	.text	00401010	0000014C
f sub_401160	.text	00401160	0000023A
f sub_4013A0	.text	004013A0	00000147
f sub_4014F0	.text	004014F0	000000BF
f sub_4015B0	.text	004015B0	0000039D
f <b>_main</b>	.text	<b>00401950</b>	<b>00004129</b>
f sub_405AA0	.text	00405AA0	00000046
f sub_405AF0	.text	00405AF0	00000085
f sub_405B80	.text	00405B80	0000042D
f sub_405FB0	.text	00405FB0	000001E7
f sub_4061A0	.text	004061A0	00000015
f sub_4061C0	.text	004061C0	0000002D
f MakeMnemonic	.text	004061EE	00000006
f InstrDasm	.text	004061F4	00000006
f InstrDecode	.text	004061FA	00000006
f ___crtCorExitProcess	.text	00406200	0000002B
f ___crtExitProcess	.text	0040622B	00000017
f __lockexit	.text	00406243	00000009
f __unlockexit	.text	0040624C	00000009
f __init_pointers	.text	00406255	00000033
f __initterm_e	.text	00406288	00000024
f __cinit	.text	004062AC	00000097
f _doexit	.text	00406343	00000140
f <b>_exit</b>	.text	<b>00406483</b>	<b>00000016</b>
f <b>__exit</b>	.text	<b>00406499</b>	<b>00000016</b>
f <b>__cexit</b>	.text	<b>004064AF</b>	<b>0000000F</b>
f <b>__c_exit</b>	.text	<b>004064BE</b>	<b>0000000F</b>
f __amsg_exit	.text	004064CD	0000001D
f <b>_printf</b>	.text	<b>004064EB</b>	<b>000000A7</b>

OfficeMalScanner.exe (OMS)

Original Code  
(no optimization)

Static Link Library  
(maximum optimization)

# Bit-Image of OMS



Stirling<sup>[10]</sup> : a hex editor

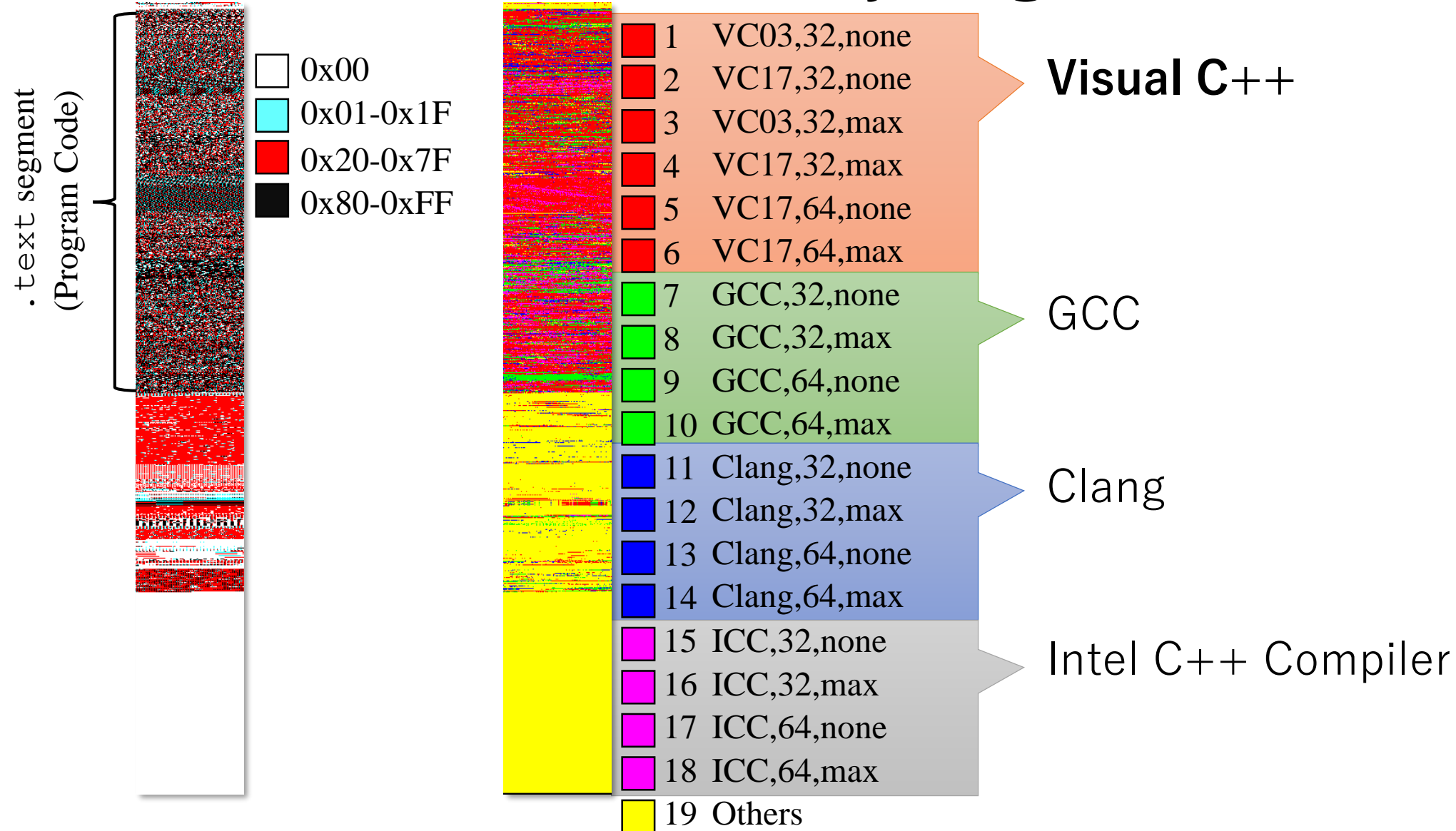
ADDRESS	00	01	02	03	04	05	06	07	08	09	0A	0B	0C	0D	0E	0F	0123456789ABCDEF
00000000	4D	5A	90	00	03	00	00	00	04	00	00	00	FF	FF	00	00	MZ.....
00000010	B8	00	00	00	00	00	00	00	40	00	00	00	00	00	00	00	ク.....@.....
00000020	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	00	.....り...
00000030	00	00	00	00	00	00	00	00	00	00	00	D8	00	00	00	00	..コ..エ..!ク..L..!Th
00000040	0E	1F	BA	0E	00	B4	09	CD	21	B8	01	4C	CD	21	54	68	is program canno
00000050	69	73	20	70	72	6F	67	72	61	6D	20	63	61	6E	6E	6F	t be run in DOS
00000060	74	20	62	65	20	72	75	6E	20	69	6E	20	44	4F	53	20	mode....\$......
00000070	6D	6F	64	65	2E	0D	0D	0A	24	00	00	00	00	00	00	00	3載tw./謡./謡./謡
00000080	33	9D	41	C4	77	FC	2F	97	77	FC	2F	97	77	FC	2F	97	./預./裂.!色./痢
00000090	77	FC	2F	97	61	FC	2F	97	F4	E0	21	97	57	FC	2F	97	..%.../謡..羅./.
000000A0	9F	E3	25	97	09	FC	2F	97	77	FC	2E	97	85	FD	2F	97	..<庸./痢.\$録./倫
000000B0	15	E3	3C	97	66	FC	2F	97	9F	E3	24	97	CE	FC	2F	97	.)要./由ichw./.
000000C0	CF	FA	29	97	76	FC	2F	97	52	69	63	68	77	FC	2F	97	.....PE..L...
000000D0	00	00	00	00	00	00	00	00	50	45	00	00	4C	01	04	00	

Bit-Image View

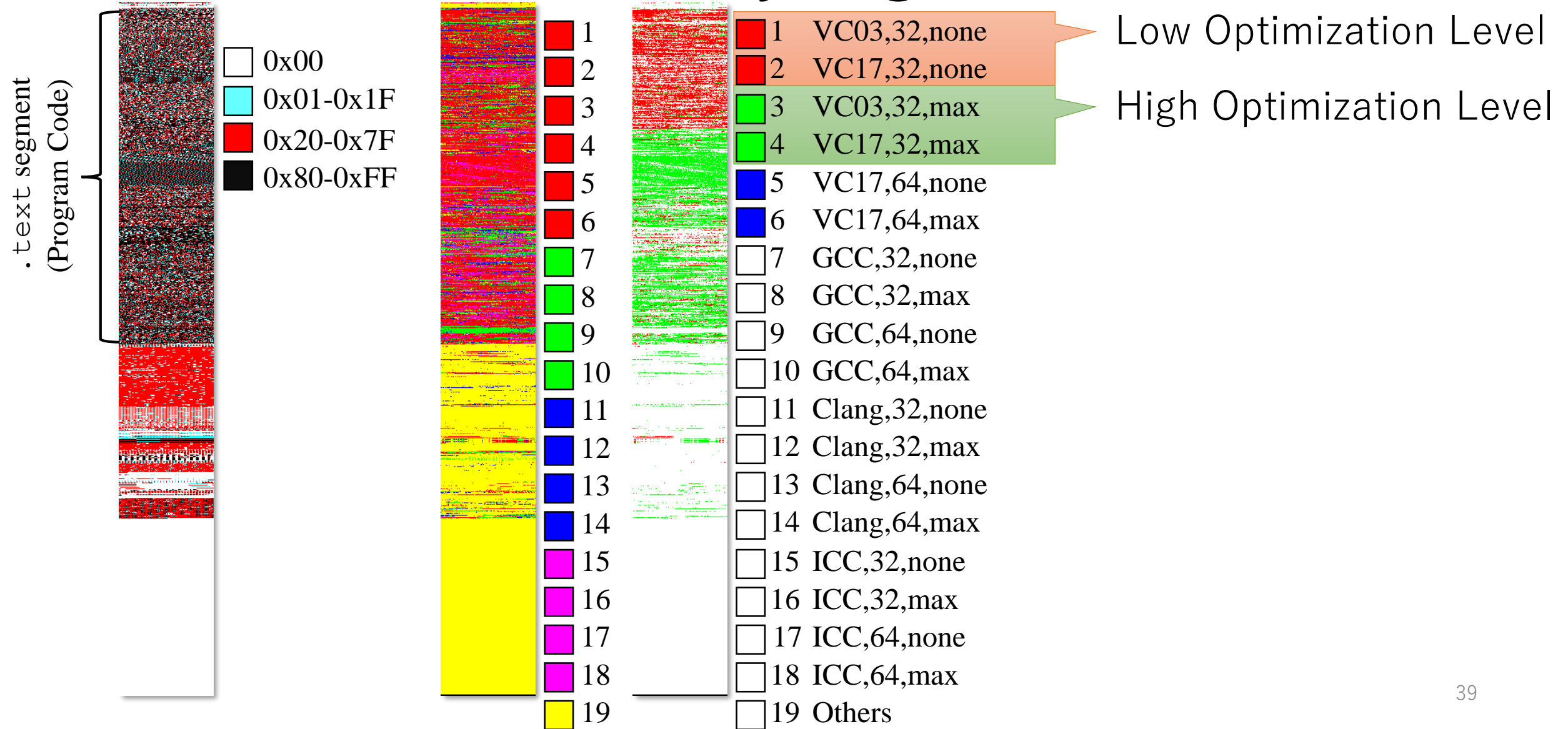
Hex View

[10] K. Goto, "Stirling," <https://www.vector.co.jp/soft/win95/util/se079072.html>, 1998.

# Visualization of OMS by o-glassesX



# Visualization of OMS by o-glassesX



# Case Study: Tracking Emdivi RATs in Dev. Env.

**EMDIVI** malware family is used in targeted email attacks against Japanese organizations. It allows machines to be remotely controlled by attackers for malicious commands and other activities.

Japan pension system hacked, 1.25 million cases of personal data leaked

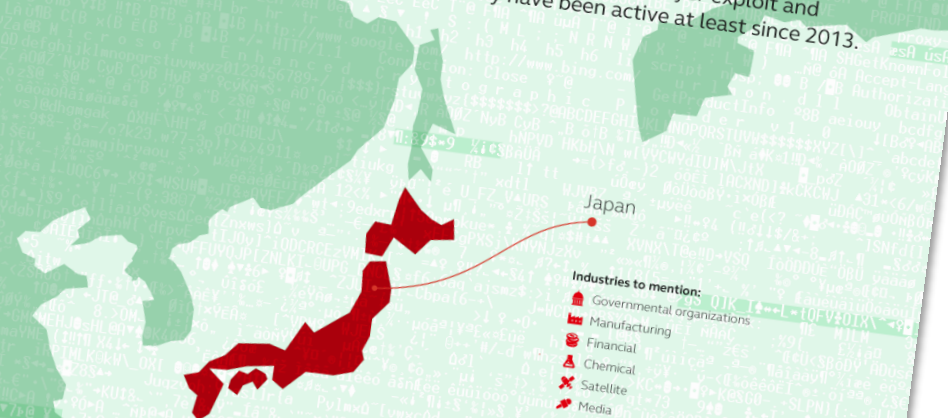
2 MIN READ



年金情報125万件流出

Blue Termite cyber-espionage campaign targets hundreds of organizations in Japan

The attackers hunting confidential information utilizing a zero-day Flash player exploit and a sophisticated backdoor, customized to each victim. They have been active at least since 2013.



- Industries to mention:
- Governmental organizations
  - Manufacturing
  - Financial
  - Chemical
  - Satellite
  - Media
  - Educational



# The Version of Emdivi

```
sub_4282BD proc near
push    offset aT17_08_26_kenp ; "t17.08.26.KENP00202"
mov     ecx, offset unk_437590
call    sub_401F54
push    offset sub_42905B
call    _atexit
pop     ecx
retn
sub_4282BD endp
```

- Almost attached malware is compiled just before used  
So, the sender and the developer may be in the same group
- Frequently Updated
  - t17 : For initial compromise
  - t19,t20 : For expanding the intrusion (High stealth performance)

# Emdivi dataset

Analysis Report made by Macnica Networks

[https://www.macnica.net/security/report\\_01.html/](https://www.macnica.net/security/report_01.html/)

## 163 MD5 Hashes of the Emdivi Family



Appendix Emdivi RAT ハッシュ値		
MD5	Version	Compile Time (UTC+8)
7fa87d1adc06bb19dde13689afe8f8ef	t9_4_sender	2012/05/28 10:02:45
2f210e5e55eb90880c12019e358c43fb	t9_5_system	2012/05/30 12:49:25
66680364d2f006db747dd640b044efe3	t9_5_system	2012/05/30 11:51:39
d953cadc4be2ab27219ef87a6a1aad87	t9_5_system	2012/05/30 14:47:22
3a68b60202787c4c779f8534ea186c75	t9_5_system	2012/05/30 12:49:25
b1f967dfe09603844a2354977356165f	t9_5_system	2012/05/30 14:47:22
4aa0d9c2b300d627c1f5abd048331597	t9_6	2012/06/01 22:07:15
094d87782555477fdc6325c56c28ff30	t9_6	2012/06/01 22:07:15
a219e2c31784bec4fc159400b229f4e0	t9_6	2012/06/01 18:27:44
dfb0ad1e22d60716512855602d47392d	t9_6	2012/06/01 19:24:13
e01e34660211bb8c7c746a6819f81c2b	t9_6	2012/06/01 19:24:13
d2f46428e1651ab6555d6f5ee87b04e9	t11.05	2012/09/13 12:58:32
1c462660b33130f5e9c2ad664eedb40	t15.07	2013/04/22 13:17:43
402406d85bccc25b14e2baa1c7b9584c	t16.19	2013/07/25 16:21:30

# Difference of Emdivi Rats in Dev. Env.

## All Sample

Architecture : 32-bit (x86)  
Compiler family : Visual C++  
Optimization level : max

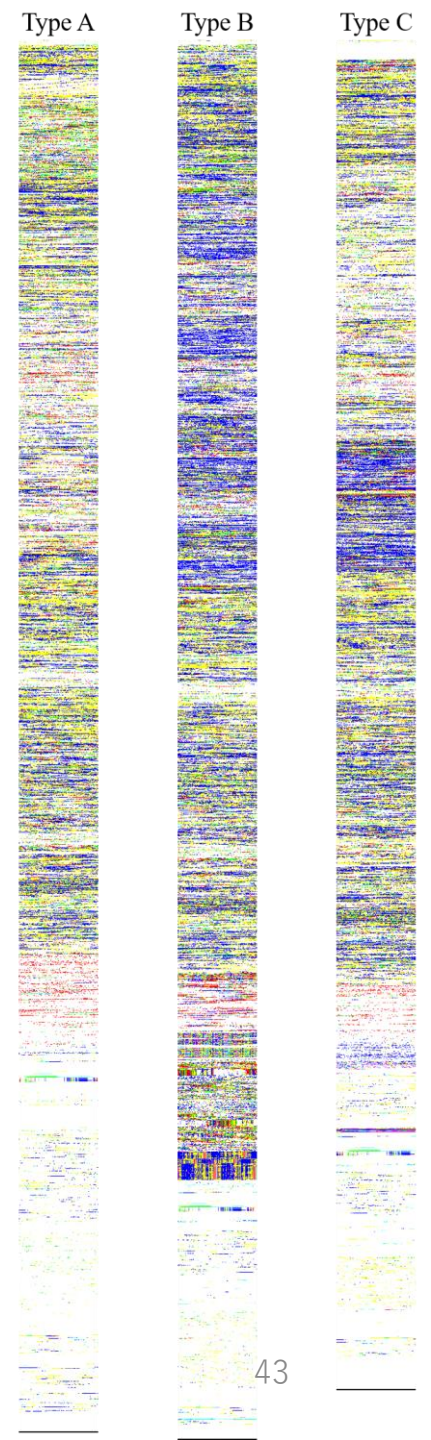
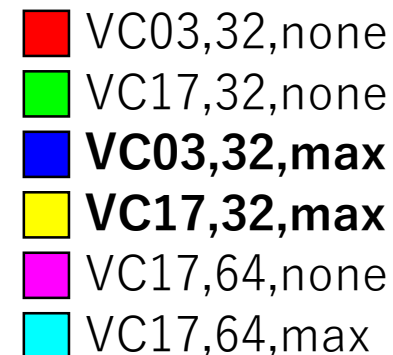
## Focusing on Compiler Version

Yellow : relatively new compiler  
Blue : relatively old compiler

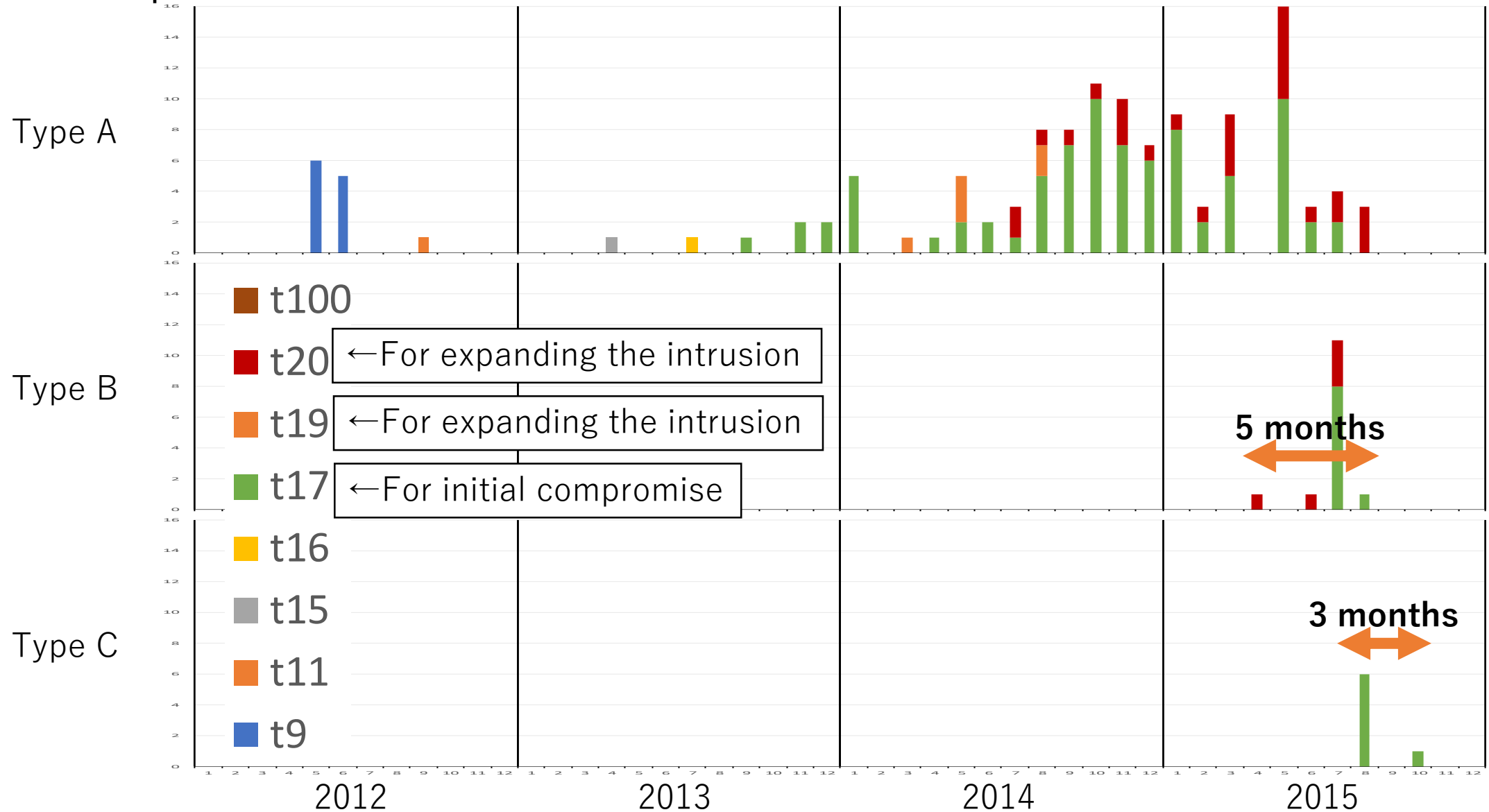
Type A: Yellow

Type B: Blue -> Yellow

Type C: Yellow -> Blue -> Yellow



# Compile-time and version of Emdivi RATs



# Limitation

## **Obfuscated Code**

Input machine code need to be already de-obfuscated.

## **Multi CPU Architectures**

This method may be applied to many CPU architecture besides x86.

Splitting binary by instruction is difficult in two more CPU architectures inputs at the same time.

This limitation will be resolved, if new bin2vec method supporting multi CPU Arc. is released...

# Conclusion

## High Recognition Rate for Stripped Machine Code

16-instruction input: .956 accuracy

64-instruction input: .988 accuracy

## Solution to Black Box Problem

o-glassesX can calculate how much input data contributes to output in units of instructions

## Case Study: Emdivi

It has been revealed that there are three attackers in the same attack group.

o-glassesX and our dataset are available at

<https://github.com/yotsubo/o-glassesX>

