

# 2018 정보보호 R&D 데이터 챌린지

---

- AI기반 악성코드 탐지 트랙 [대학(원)생] -

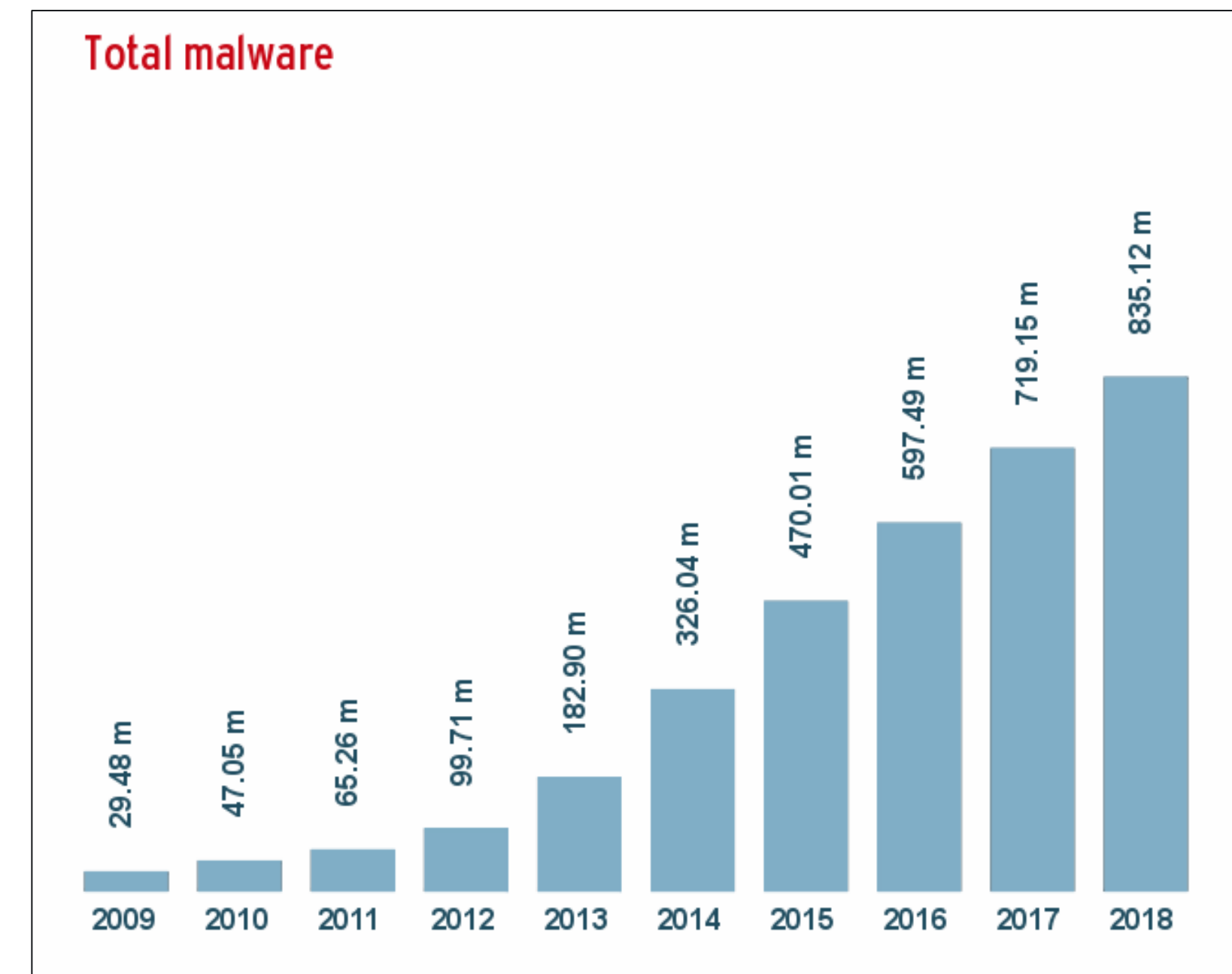
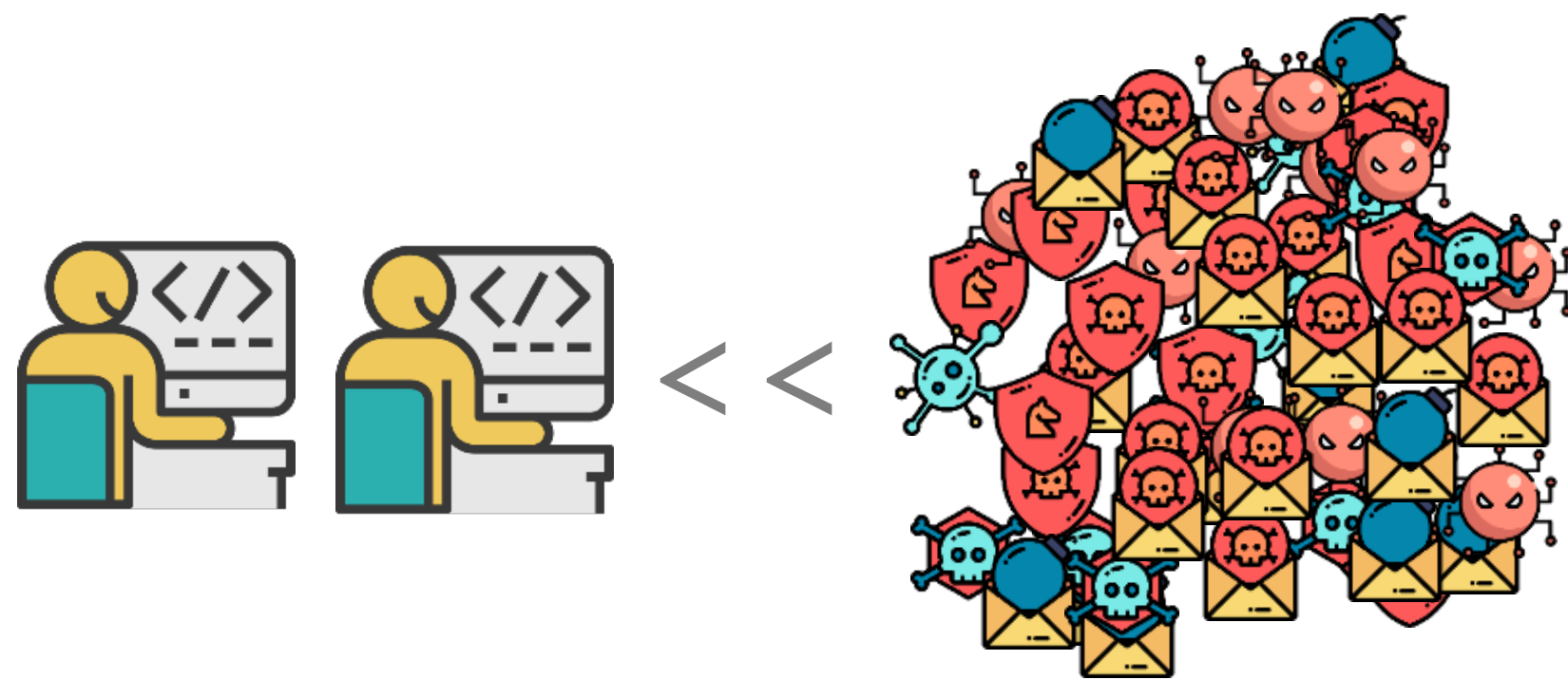
Dept. of Computer Science, Kookmin University  
Information Security Lab.

TEAM : KMU INFOSEC (발표자: 정성민)



- 서론
- 악성코드 탐지 시스템
  - 딥러닝 기반 악성코드 탐지 모델
  - 엘라스틱서치 기반 최대 유사 악성코드 검색 모델
  - 머신러닝 기반 앙상블 모델
- 특징 정보
  - 정적 특징
  - 동적 특징
- 데이터셋
- 결과

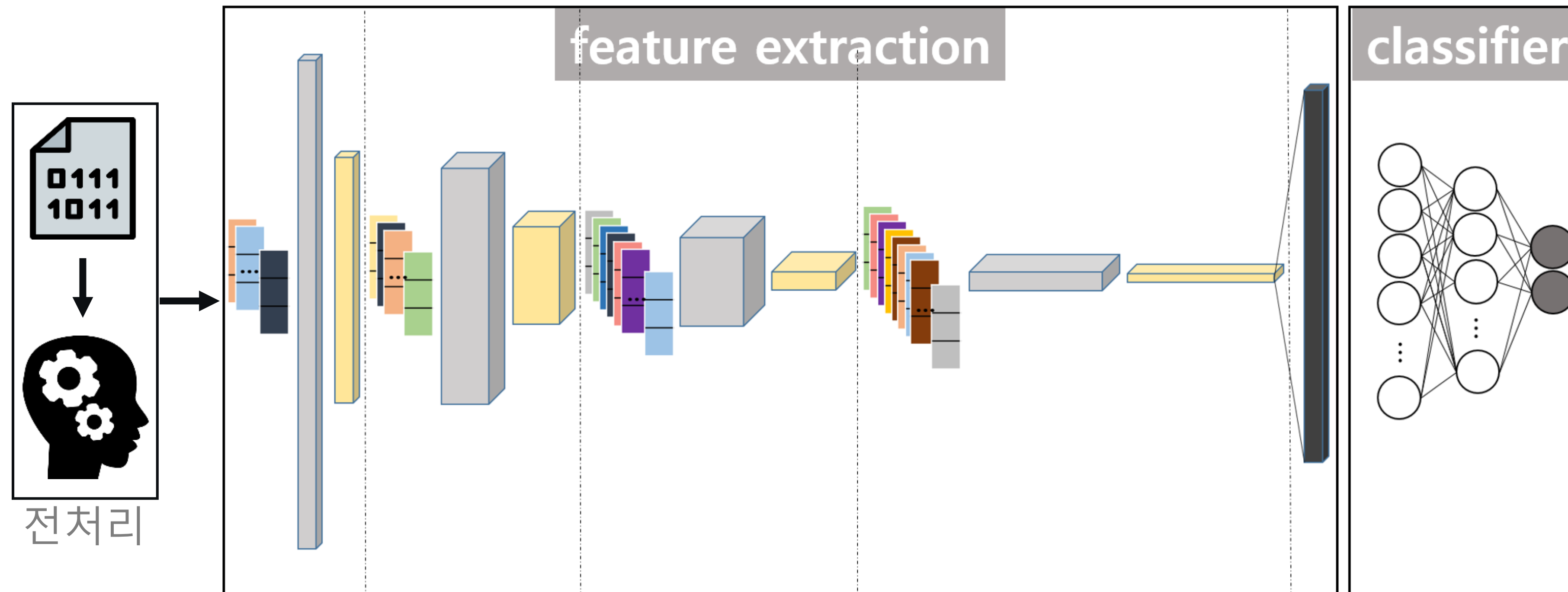
- 악성코드의 기하급수적인 증가, 변종/신종 악성코드 출현
  - More than 430 Million malware were discovered in 2015 (increasing 36% in 2014)
    - Internet Security Threat Report (Symantec, 2016)
  - AV-TEST Statistics
- 그에 비해 현저히 적은 악성코드 분석가 수



<https://www.av-test.org/en/statistics/malware/>

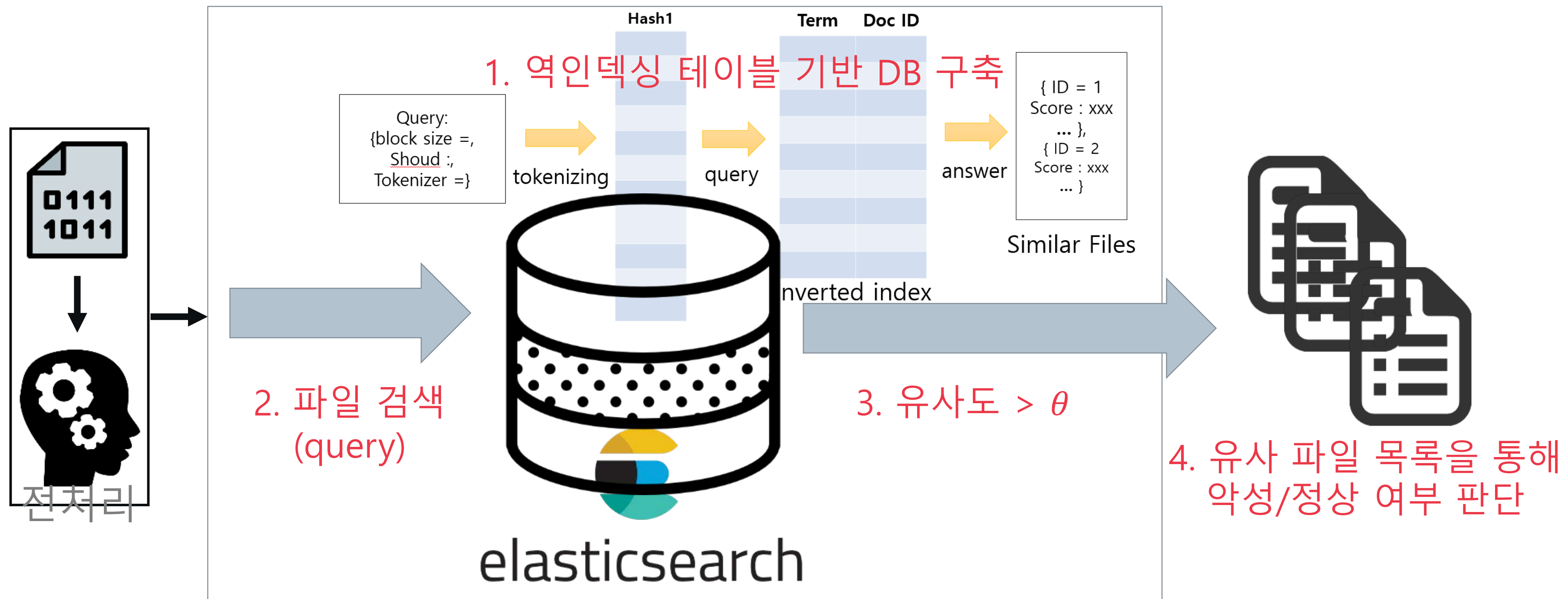
- 전통적인 악성 코드의 분석 방법
  - 정적 분석: w/o running (코드 정보, 엔트로피, 헤더 정보, 파일 크기 등)
    - 특징: fast, simple
  - 동적 분석: with running (실행 정보, API call, 네트워크 통신 정보 등)
    - 특징: slow, compact(critical), robust to obfuscation
- Hand-crafted feature → Machine Learning-based approach 로의 변화
  - 자동화
  - 양질의 특징 추출

- 딥러닝 기반 탐지 모델
  - 전처리 (preprocessing)
  - CNN 기반의 딥러닝 학습 모델링



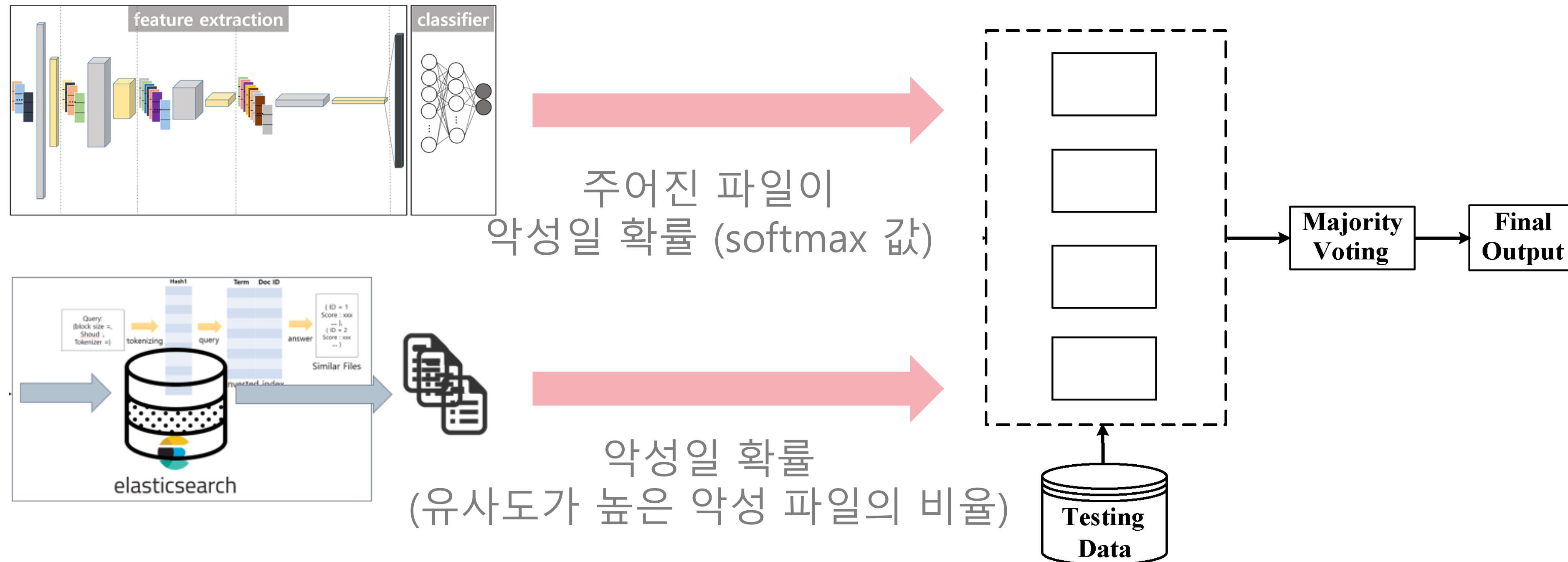
[HYPER PARAMETER]  
Objective function: cross entropy  
Optimizer: Adam  
Activation Function: Leaky ReLU  
Conv filter → Batch Normalization  
No dropout  
Learning rate = 0.0001

- 엘라스틱서치 기반 최대 유사 악성코드 검색 모델
  - 전처리 (preprocessing) → 특징 추출
  - 특징 정보를 많이 가질 수록 유사 점수(score) 상승 → 점수에 따라 악성/정상 여부 판단



- 머신 러닝 기반 앙상블 모델

- 다양한 특징 정보에 대한 두 모델(딥러닝, 엘라스틱서치)의 결과를 상호보완함
- 사용 앙상블 모델: Random Forest






- 정적 정보 [IDA Pro]

- structured opcode sequence
- disassemble information
- strings
- ssdeep
- Import information



- 동적 정보 [Cuckoo Sandbox]
- API call sequence



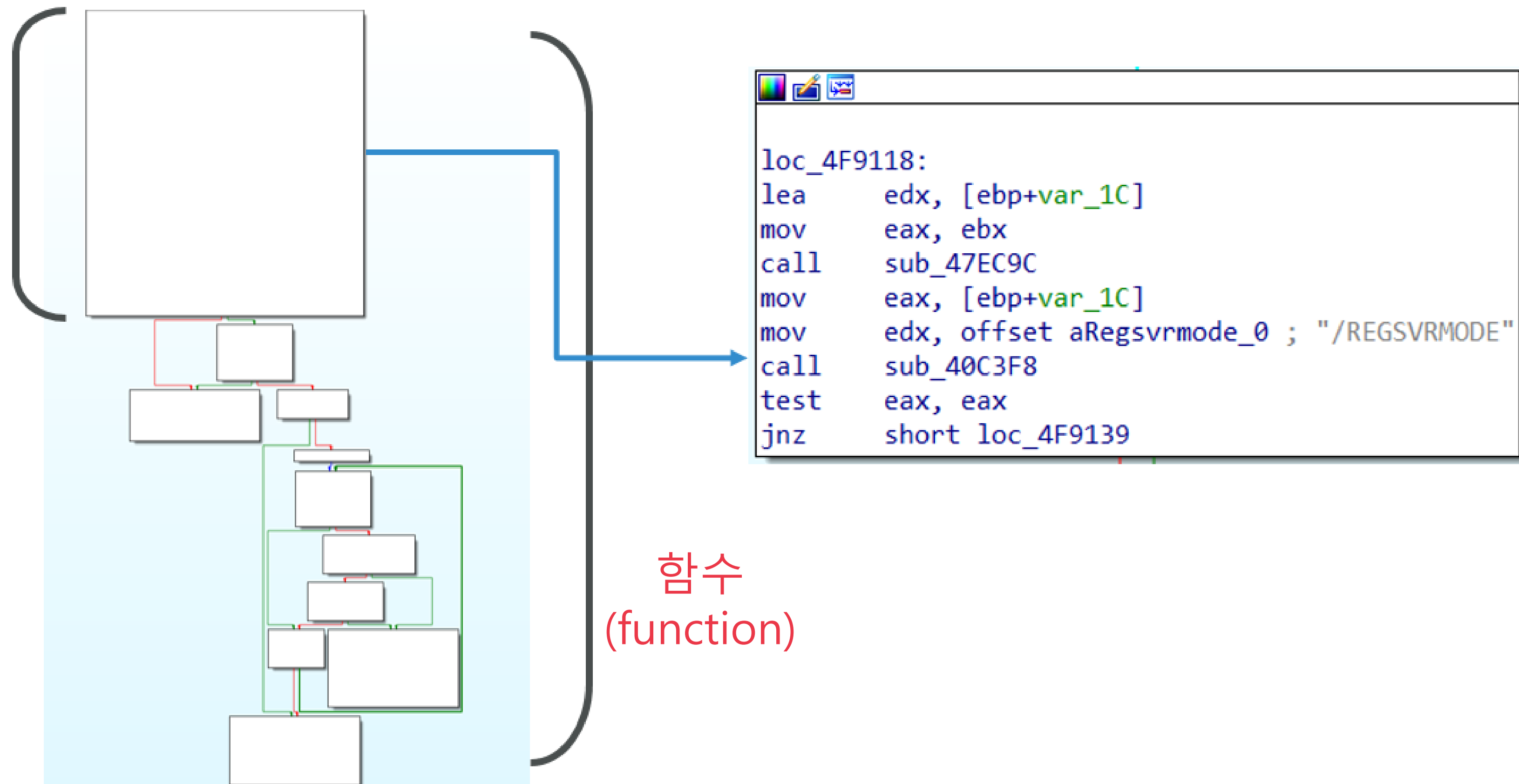
- 정적 정보 [IDA Pro]

- structured opcode sequence
- disassemble information
- strings
- ssdeep
- Import information

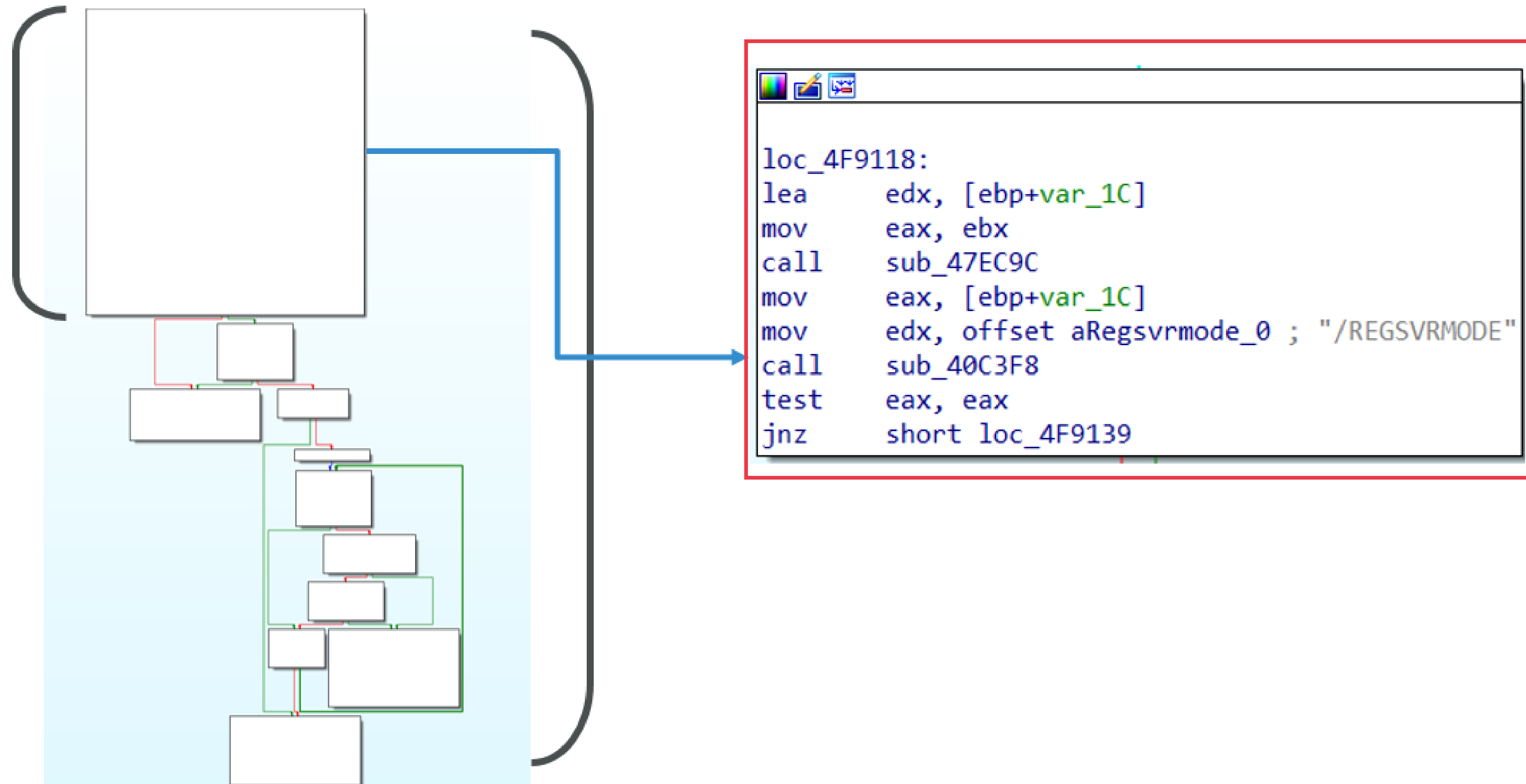
기본 블록  
(basic block)

- 동적 정보 [Cuckoo Sandbox]

- API call sequence



- 정적 정보 [IDA Pro]
  - structured opcode sequence
  - **disassemble information**
  - strings
  - ssdeep
  - Import information
- 동적 정보 [Cuckoo Sandbox]
  - API call sequence



## Binary file 에 대해 정규식을 이용하여 추출

- structured opcode sequence
  - disassemble information
  - **strings**
  - ssdeep
  - Import information
- 
- 동적 정보 [Cuckoo Sandbox]
    - API call sequence

```
[Wx1f !"#$%&w'()*+,-  
./0123456789;<=>?@ABCDEFGHIJKLMNOPQRSTUVWXYZ[Ww]^_`abcdefghijklmnopqrstuvwxyz{|}~',  
'Wx1f !"#$%&w'()*+,-./0123456789;<=>?@abcdefghijklmnopqrstuvwxyz[Ww]^_`abcdefghijklmnopqrstuvwxyz{|}~',  
'Wx1fgm]zs', 'Wx1ftuVS;%', ' ', ' ', ' ', ' ', 'H', " <requestedExecutionLevel  
level='asInvoker' uiAccess='false' />", ' ', ' ', ' ', ' ', '</requestedPrivileges>', '  
<requestedPrivileges>', ' ', '</security>', ' ', '<security>', ' ', '</trustInfo>', ' ', '<trustInfo xmlns="urn:schemas-  
microsoft-com:asm.v3">', ' ', '4P/qD', " Base Class Array'", ' Base Class Descriptor at (' , " Class Hierarchy Descriptor"  
" Complete Object Locator"', " Type Descriptor"', ' delete', ' delete[]', ' new[]', '!25jK94', '!O/.E', '!E>^~A', ""B <1=',  
""hLFGp', ""~):+nN', '#49sOD', '#9Wx1fRM6', '#IsMPQ', "$cdAV", '$5PiU&n', '$f[2@D', '$h[#XC', '%;s%^,',  
'%S#[kwx1f=', "'! 2XVa", ""'+T^=^$", ""')/f3.", ""<v(TW", ""HtkeQ", '(/+uxrT', '(Y-b0', '(null)', '@Ag^#|', ')lkWWp)',  
')))]}|, ""*'3IU32LR", '*HgW&66v', "*VP[yC:", '+0WWW{Oh', '+6@m_', "",0{x8'vZ", "-tfmH", '.00cfg', '.3MY^^',  
'.?AVbad_alloc@std@@', '?AVbad_array_new_length@std@@', '?AVbad_exception@std@@',  
'.?AVexception@std@@', '?AVlength_error@std@@', '?AVlogic_error@std@@', '?AVout_of_range@std@@',  
'.?AVtype_info@@', '.CRT$XCA', '.CRT$XCAA', '.CRT$XCZ', '.CRT$XIA', '.CRT$XIAA', '.CRT$XIAC', '.CRT$XIC',  
'.CRT$XIZ', '.CRT$XLA', '.CRT$XLZ', '.CRT$XPA', '.CRT$XPX', '.CRT$XPXA', '.CRT$XPZ', '.CRT$XTA', '.CRT$XTZ', '.PF bi',  
'data$r', '.gfids', '.gfids$x', '.gfids$y', '.idata$2', '.idata$3', '.idata$4', '.idata$5', '.idata$6', '.rdata', '.rdata$T',  
'rdata$r', '.rdata$sxdata', '.rdata$zzzdbg', '.rsrc$01', '.rsrc$02', '.rtc$IAA', '.rtc$IZZ', '.rtc$TAA', '.rtc$TZZ', '.text$mnn',  
'text$x', '.tls$ZZZ', '.xdata$x', '/#8_o8_o8_o', "/f1X1'|p5", '0Wx1f020j0r0', '0 0(00080@0H0P0X0`0h0p0x0', '0#0-0=0',  
'0$0,040<0D0LOTOwwwd0l0t0|0', '0$1C1'1{1', "0'090K0|0o0", '0*0/0@0F0Q0Y0d0j0u0{0', '0,161<1B1', '0.LyJo',  
'0=xak/', ' '
```

- 정적 정보 [IDA Pro]

- structured opcode sequence
- disassemble information
- strings
- **ssdeep**
- Import information

- 동적 정보 [Cuckoo Sandbox]

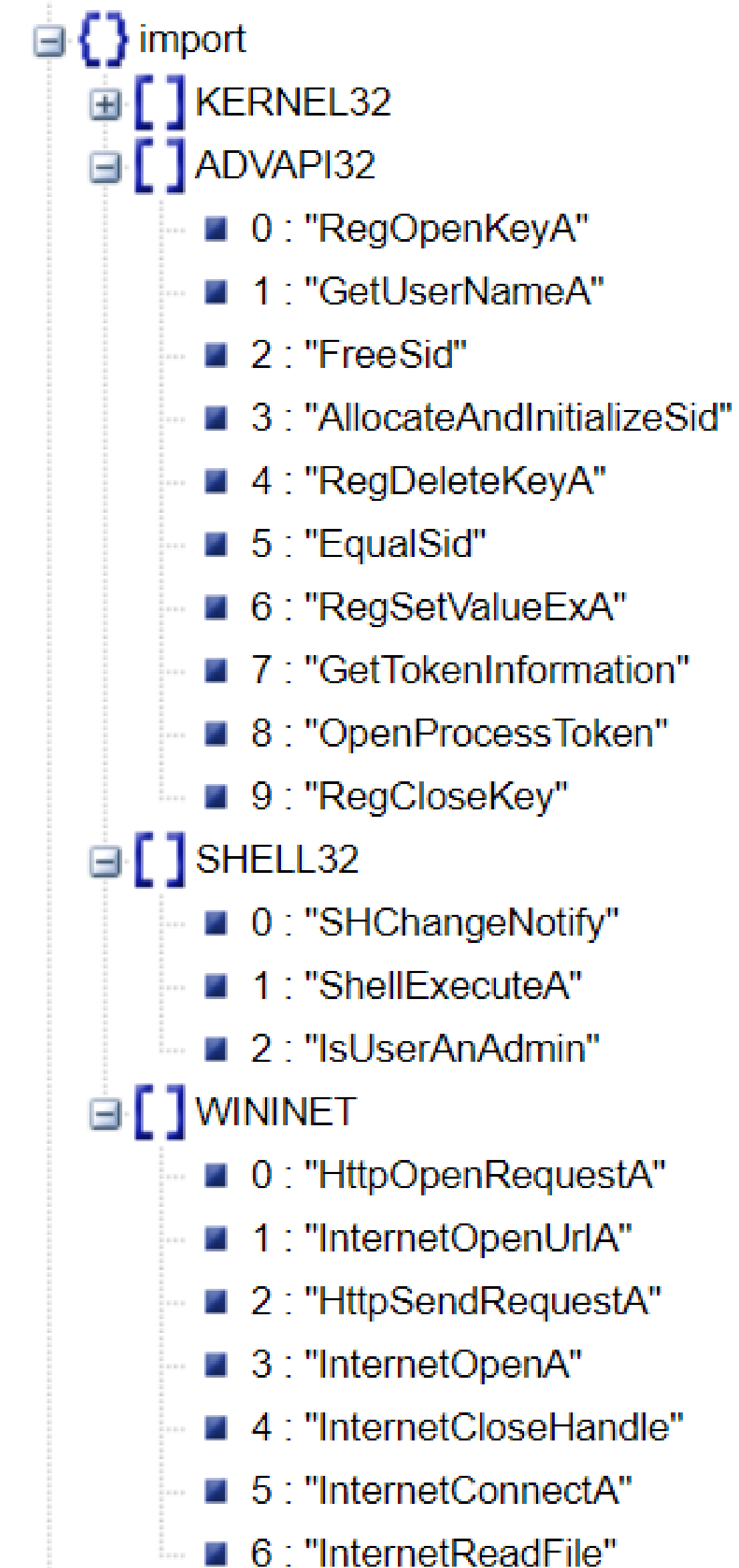
- API call sequence

- **ssdeep**

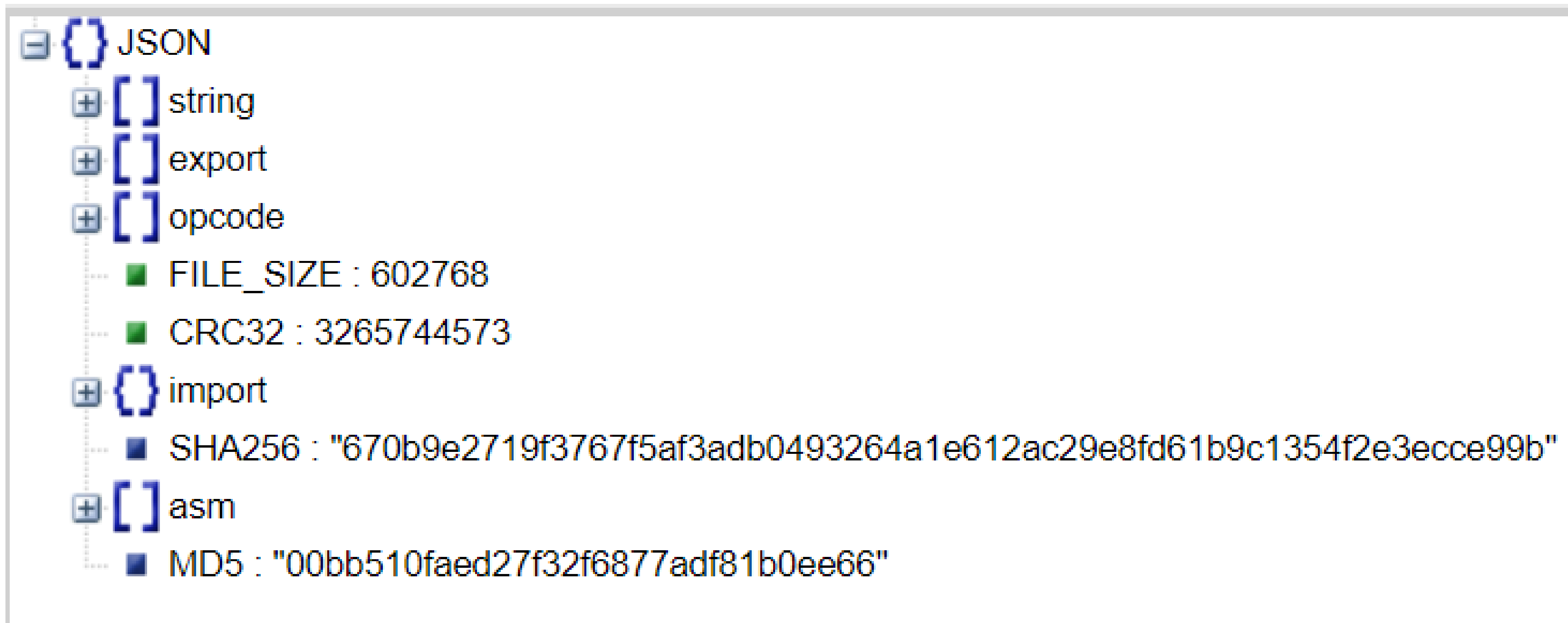
- Fuzzy 해시 기반의 파일 유사도 측정 프로그램
- **Fuzzy 해시**: 일정 크기 단위로 구분, 각 블록의 해시 값을 생성하는 방식

```
{ -
  "_index": "kisa_only_our_data_ssdeep",
  "_type": "type_ssdeep",
  "_id": "b4263c6c73b32ff1a5b2b462f2c9192f",
  "_version": 1,
  "found": true,
  "_source": { -
    "chunk_size": "768",
    "chunk": "zXjvSeUvK9U1AhqCf51yySKFsn6g8fDrrD0mjPJ0Mw80xGSTwGKn23+zjGk",
    "double_chunk": "zuHQJf3jFC6gBYJ0MwLxFTX2Gk"
  }
}
```

- 정적 정보 [IDA Pro]
  - structured opcode sequence
  - disassemble information
  - strings
  - ssdeep
  - **Import information: IAT 에서 라이브러리 정보 추출**
- 동적 정보 [Cuckoo Sandbox]
  - API call sequence



- 정적 정보 [IDA Pro]
  - IDAPython 기반 정적 정보 추출 자동화



- 정적 정보 [IDA Pro]

- structured opcode sequence
- disassemble information
- strings
- ssdeep
- Import information

- 동적 정보 [Cuckoo Sandbox]

- API call sequence

## Cuckoo Sandbox 를 통해 프로세스 별 API 호출 시퀀스 추출

```
▼ 3440 {81}
  RegCreateKeyExW : 7
  NtDuplicateObject : 7
  NtOpenSection : 4
  NtFreeVirtualMemory : 2
  RegCloseKey : 90
  NtQueryKey : 2
  NtReadFile : 4
  LdrUnloadDll : 12
  HttpOpenRequestA : 11
  GetSystemInfo : 1
  RegQueryValueExA : 64
  getaddrinfo : 2
  InternetOpenA : 11
  OpenServiceW : 1
  CopyFileA : 1
  NtTerminateProcess : 3
  NtClose : 118
  RegCreateKeyExA : 5
  GetFileSize : 5
  InternetCloseHandle : 33
  SetFileTime : 3
  NtDelayExecution : 13
  InternetConnectA : 11
  InternetQueryOptionA : 1
  NtDeviceIoControlFile : 24
  NtAllocateVirtualMemory : 22
  ReadProcessMemory : 5
  RegOpenKeyExA : 23
  DeleteFileW : 1
  NtWriteFile : 2
  LdrGetDllHandle : 11
  OpenServiceA : 2
  NtWriteVirtualMemory : 1
  SetFilePointer : 4
  NtQueryValueKey : 1
  NtResumeThread : 2
  CryptAcquireContextA : 1
```



- 목적

- 빅데이터 기반의 신뢰도 높은 모델 생성
- 다양한 백신사의 라벨 활용

- 데이터셋

- 악성: 자체 수집 약 500,000개
- 정상: 자체 수집 약 100,000개
- \* 금년 KISA 학습/예선 데이터 적극 활용

- 본선 대회 대비 정확도
  - KISA 예선 대회 데이터를 통한 예상 정확도 : 97~99%

- 본선 실제 대회 정확도

## AI기반 악성코드 탐지\_대학(원)부문

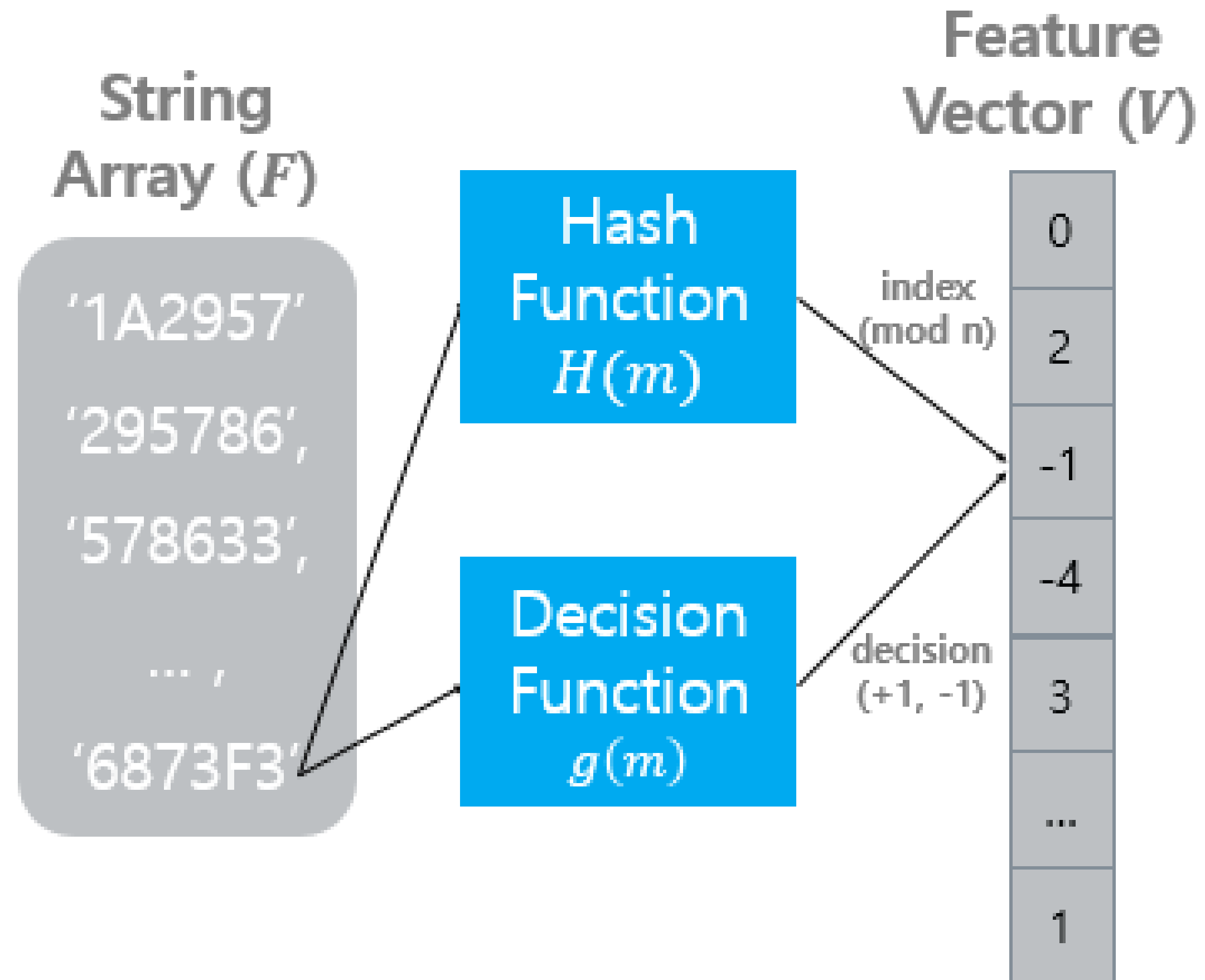
순위	팀명	탐지율									
		1차 1회	1차 2회	1차 3회	1차 4회	1차 5회	2차 1회	2차 2회	2차 3회	2차 4회	2차 5회
1위	KMU InfoSec	14:21:40 94.14%	17:31:16 94.64%	17:43:43 95.91%	18:25:36 94.72%	18:26:29 95.91%	12:55:10 93.87%	14:09:12 95.43%	14:47:07 95.63%	15:22:53 96.28%	

Thank you  
&  
Q n A

# API call sequence (for Deep Learning)

- N-gram chunking (N=4)
  - Feature Hashing
    - **Frequency**-based feature vector

```
def hashing_vectorizer(features: array of string, N: integer):  
    x := new vector[N]  
    for f in features:  
        idx := hash(f) mod N  
        if g(f) == 1:  
            x[idx] += 1  
        else:  
            x[idx] -= 1  
    return x
```

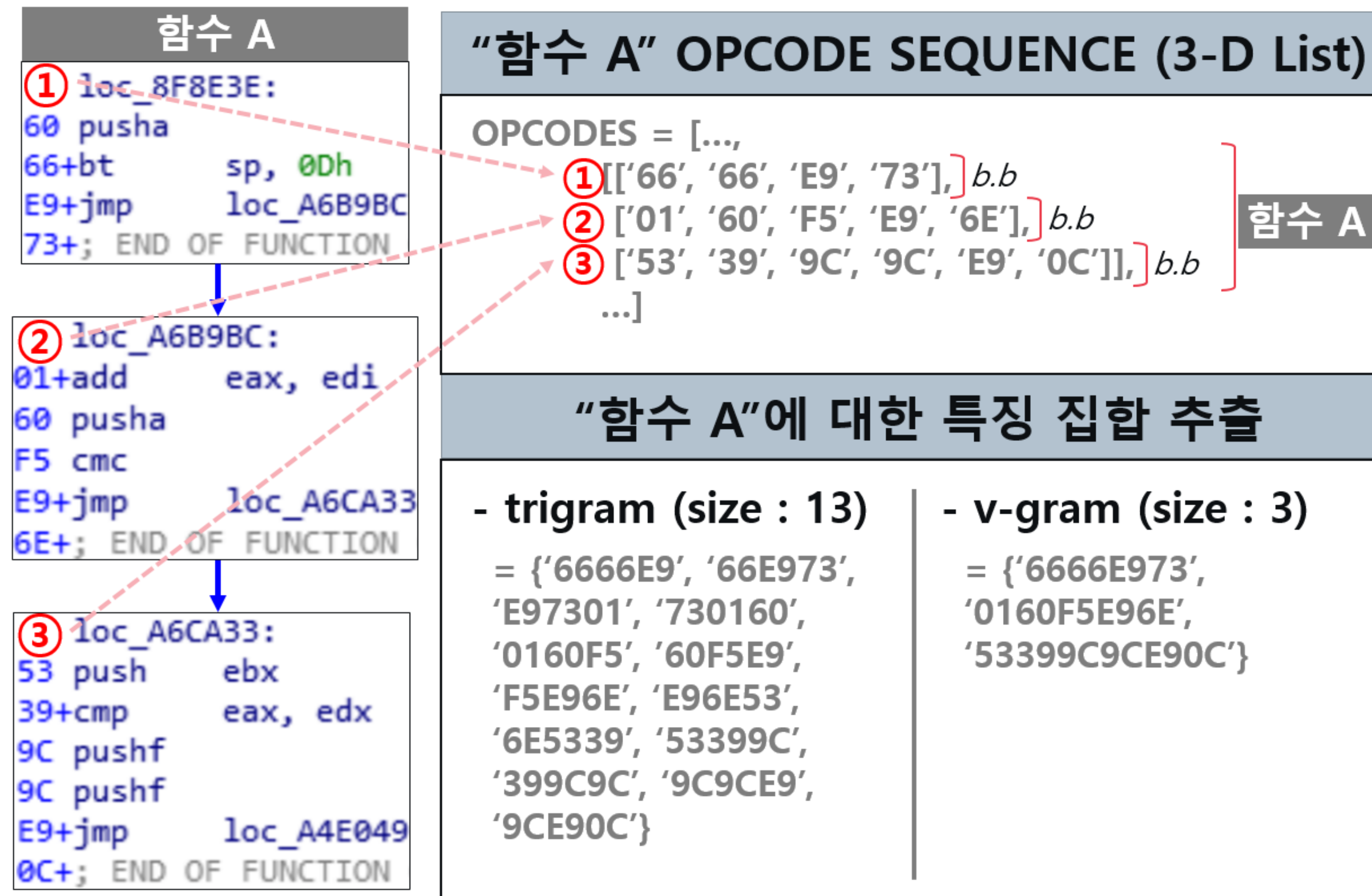


# APPENDIX



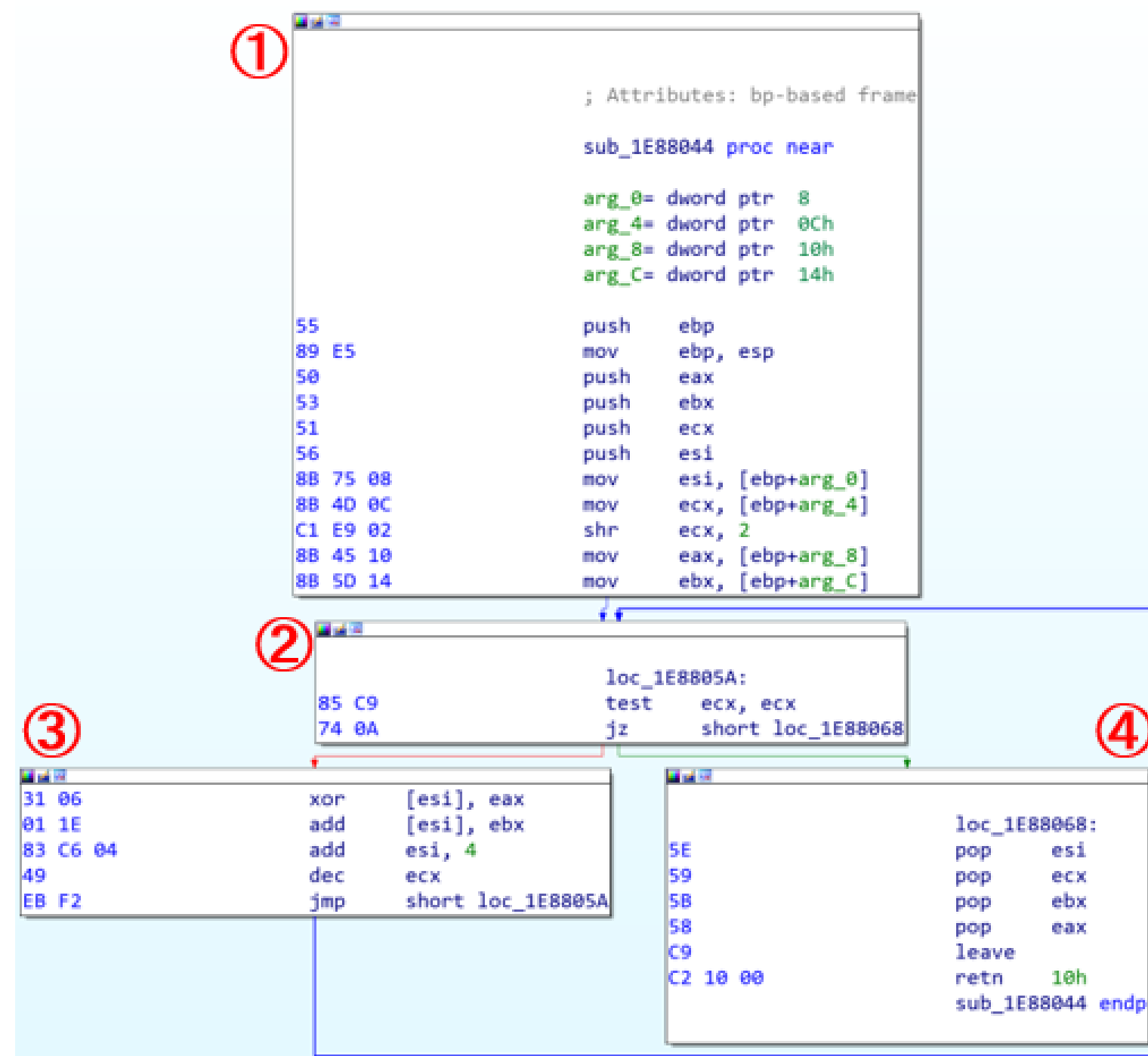
# structured opcode sequence (for Deep Learning)

- N-gram ? V-gram !



# structured opcode sequence (for Deep Learning)

- N-gram ? V-gram !



```
; Attributes: bp-based frame
sub_1E88044 proc near
    arg_0= dword ptr 8
    arg_4= dword ptr 0Ch
    arg_8= dword ptr 10h
    arg_C= dword ptr 14h
55      push    ebp
89 E5   mov     ebp, esp
50      push    eax
53      push    ebx
51      push    ecx
56      push    esi
8B 75 08 mov     esi, [ebp+arg_0]
8B 4D 0C mov     ecx, [ebp+arg_4]
C1 E9 02 shr     ecx, 2
8B 45 10 mov     eax, [ebp+arg_8]
8B 5D 14 mov     ebx, [ebp+arg_C]

loc_1E8805A:
85 C9   test    ecx, ecx
74 0A   jz      short loc_1E88068

31 06   xor     [esi], eax
01 1E   add     [esi], ebx
83 C6 04 add     esi, 4
49      dec     ecx
EB F2   jmp     short loc_1E8805A

loc_1E88068:
5E      pop     esi
59      pop     ecx
5B      pop     ebx
58      pop     eax
C9      leave
C2 10 00 retn    10h
sub_1E88044 endp
```

3-D LIST

FILE = [..., [['55', '89', ..., '8B', '8B'], ['85', '74'], ['31', '01', ..., 'EB'], ['5E', '59', ..., 'C2']], ...]

function

- Basic block

	md5 hash
① ['55', '89', '50', '53', '51', '56', '8B', '8B', 'C1', '8B', '8B']	→ C6A9FBE75628012CC69FD1FC415B20F1
② ['85', '74']	→ E345FAC6BC5C868F0222430C733FA26E
③ ['31', '01', '83', '49', 'EB']	→ AC2D461042F560151E439BE4F0F2A6B4
④ ['5E', '59', '5B', '58', 'C9', 'C2']	→ C6A9FBE75628012CC69FD1FC415B20F1

- Function: XOR operation result of basic block representative value

→ C6A9FBE75628012CC69FD1FC415B20F1 XOR E345FAC6BC5C868F0222430C733FA26E XOR AC2D461042F560151E439BE4F0F2A6B4 XOR C6A9FBE75628012CC69FD1FC415B20F1

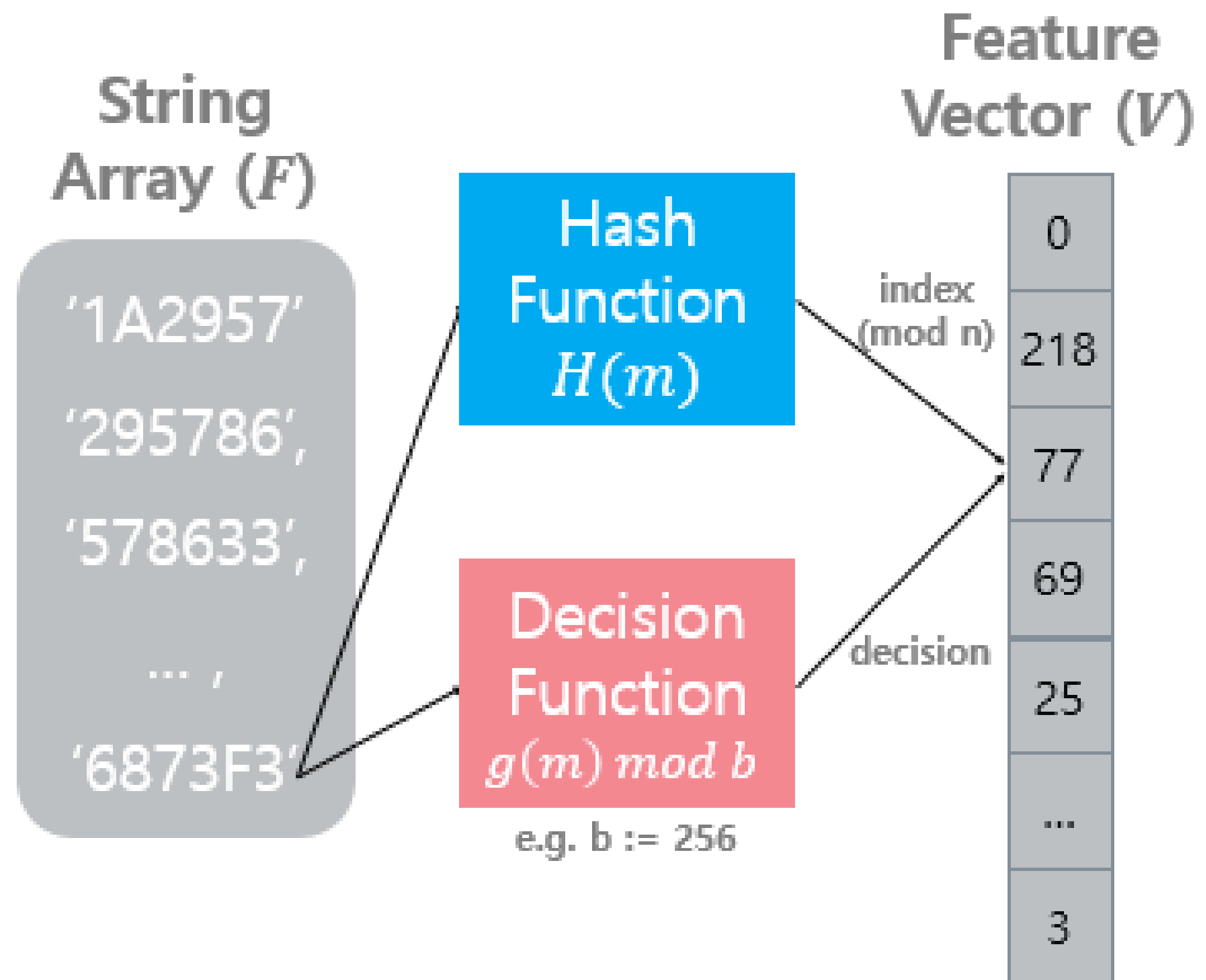
= 4F68BCD6FEA9E69A1C61D8E883CD04DA



# structured opcode sequence (for Deep Learning)

- V-gram chunking
  - Feature Hashing
    - **Content**-based feature vector
      - overwrite max value

```
def hashing_vectorizer(features: array of string, N: integer,  
                      B: integer):  
    x := new vector[N]  
    for f in features:  
        idx := hash(f) mod N  
        value := G(f) mod B  
        x[idx] = value > x[idx] ? value : x[idx]  
  
    return x
```



- disassemble, ssdeep, strings, import
  - 엘라스틱서치의 value 로 가공하지 않고 사용