

# Control Flow Analysis



*Edgar Barbosa*

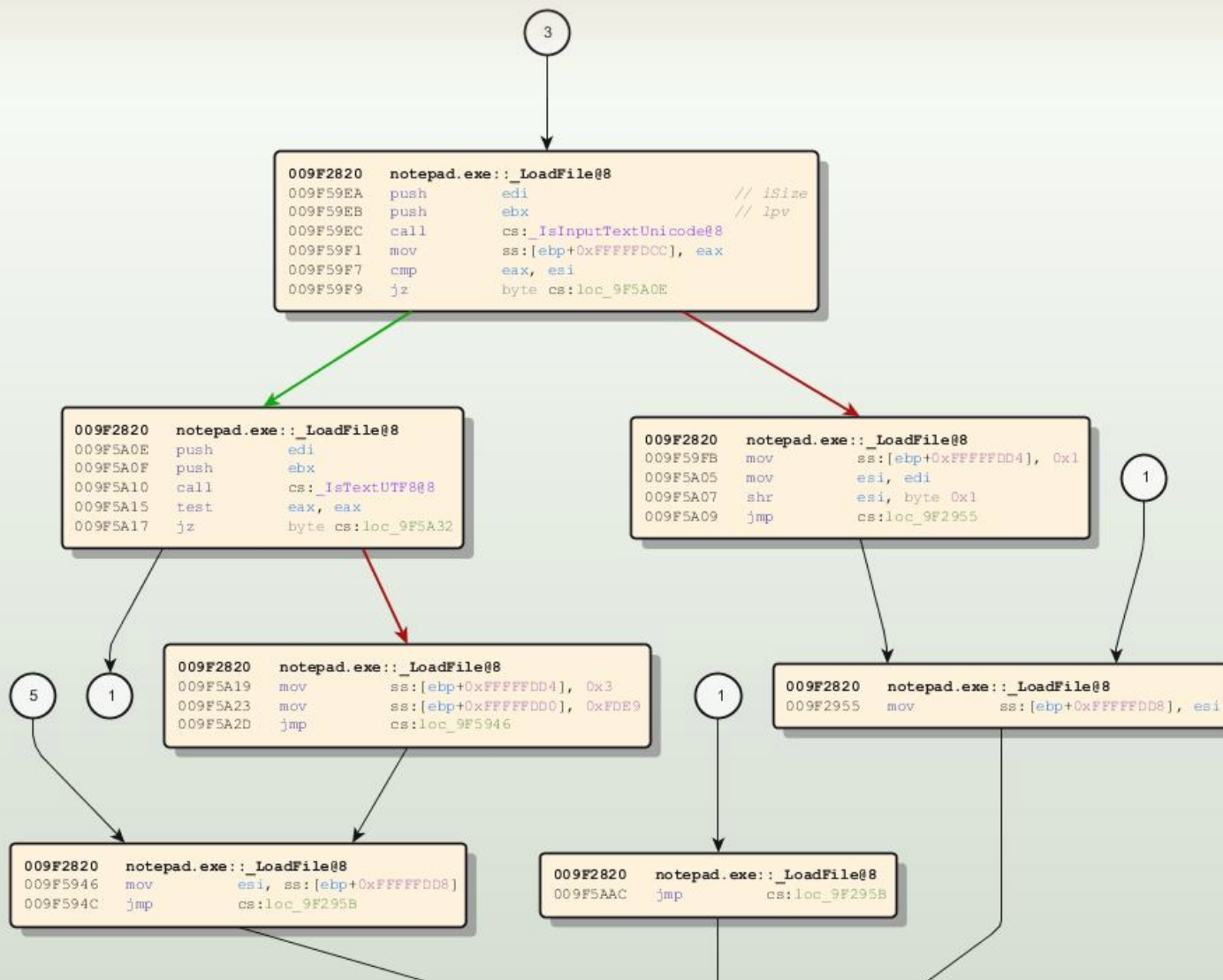
*H2HC 2011*

*São Paulo - Brazil*

# Who am I?



- ✧ Edgar Barbosa
- ✧ Senior Security Researcher at COSEINC (Singapore)
- ✧ One of the developers of Blue Pill, a hardware-based virtualization rootkit. Also presented a way to detect this type of rootkit.
- ✧ Discovered the Windows kernel KdVersionBlock data structure used for some forensic tools.
- ✧ Focus: RCE, Windows Internals, Virtualization and Program Analysis.
- ✧ Currently working on the COSEINC SMT Project, which aims to automate the bug finding process with the help of SMT solvers. The current presentation is part of the research done for the SMT project.



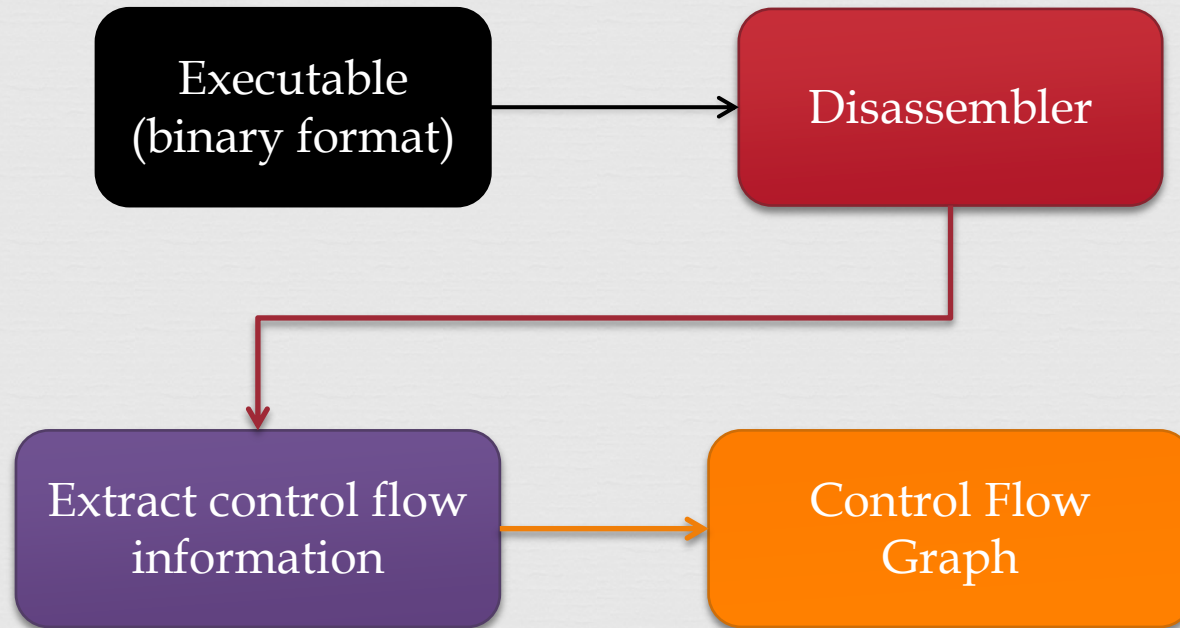
# Control Flow Analysis

# Control Flow Analysis



- ❧ Control Flow Analysis (CFA)
- ❧ Static analysis technique to discover the *hierarchical flow of control* within a procedure (function).
- ❧ Analysis of all possible *execution paths* inside a program or procedure.
- ❧ Represents the control structure of the procedure using *Control Flow Graphs*.
- ❧ Compiler theory - optimization
- ❧ The focus of this presentation is to demonstrate CFA for Reverse Code Engineering, where the source code isn't available.

# RCE and CFA

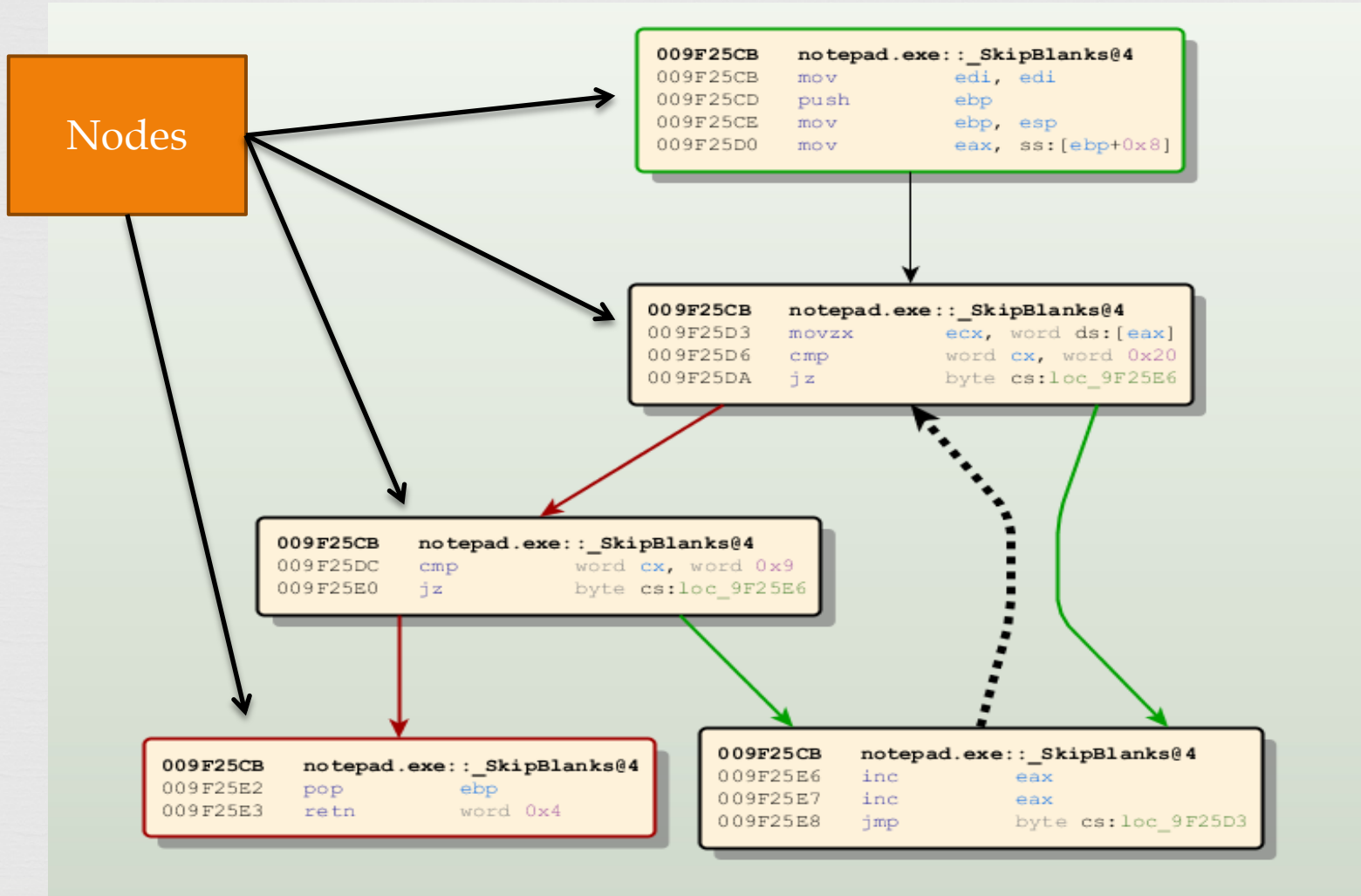


# What is a CFG?



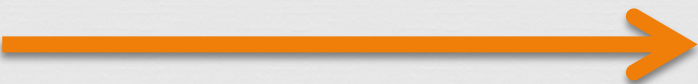
- ❧ A Control Flow Graph (CFG) is a directed graph  $G(V;E)$  which consists of a set of vertices (nodes)  $V$ , and a set of edges  $E$ , which indicate possible flow of control between nodes
- ❧ Or, is a **directed graph** that represents a superset of *all possible execution paths* of a procedure.
- ❧ Graph nodes represents objects called *Basic Blocks* (BB)

# CFG



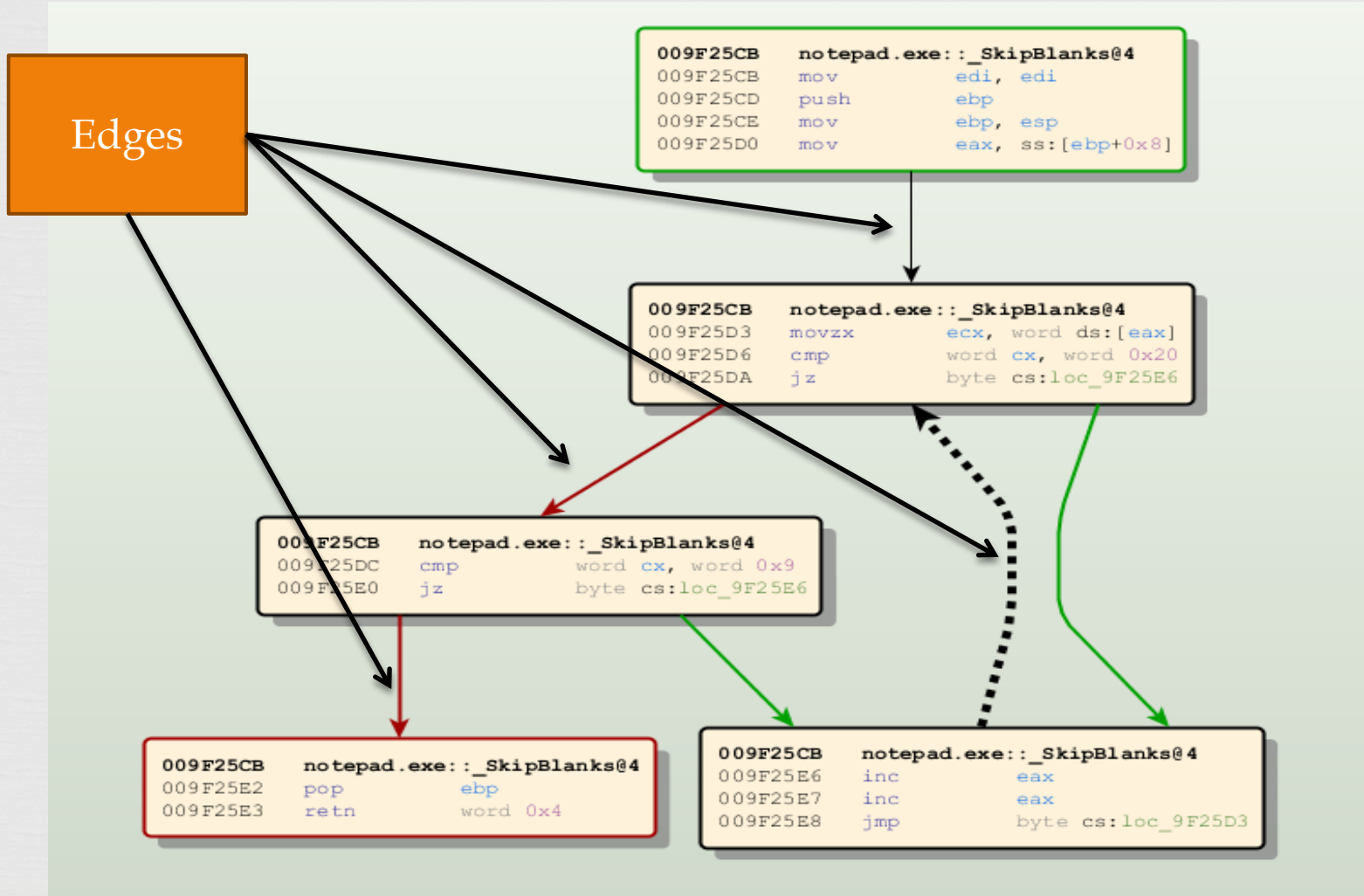
# Edges



**tail**  **head**



# CFG



# BinNavi



☞ Views

☞ Nodes

☞ Edges

# CFG properties



- ∞ In the CFA literature the algorithms assume the following CFG properties:
  - ∞ Unique Start node (Entry node)
  - ∞ All the nodes of must be reachable from the START node.
  - ∞ Unique Exit node
- ∞ Real-world:
  - ∞ Easy to find multiple exit nodes (RETURN) on the disassembly of a function
- ∞ Create a new exit node, add it to the graph and modify the return instructions to jump to the new node.

# BB identification



- ✧ In general, the problem of discovering all the possible execution paths of a code is *undecidable*. (cf. *Halting problem*).
- ✧ First step for CFG reconstruction is to identify all the basic blocks.
- ✧ A basic block is a *maximal sequence* of instructions that can be entered only at the first of them and exited only from the last of them

# Basic Block (BB)



# Basic Blocks



- ❧ First instruction of a BB (the *leader* instruction):
  1. The entry point of the routine
  2. The target of a branch instruction
  3. The instruction immediately following a branch
- ❧ Although CALL is a branch instruction, the target function is assumed to *always* return and therefore it is allowed in the middle of a BB.
- ❧ To build the BB's we need to identify all the *leader* instructions. This requires the disassembly of the instructions.
- ❧ Two disassembly algorithms

# 1. Linear Sweep



⌘ A linear sweep algorithm starts with the first byte in the code section and proceeds by decoding each byte until an illegal instruction is encountered <sup>[a]</sup>

>> 8B FF 55 8B EC 8B 45 08

8B FF	mov	edi, edi
55	push	ebp
8B EC	mov	ebp, esp
8B 45 08	mov	eax, [ebp+8]

# 2. Recursive Traversal



⌘ Linear sweep algorithm doesn't take into account the control flow behaviour of some instructions.

```
>> EB 01 FF 8B 45 FC
```

```
EB 01    jmp short 0x401020
```

```
FF      ???    ;invalid
```

⌘ Recursive traversal disassemblers interpret branch instructions in the program to translate only those bytes which can actually be *reached* by control flow. <sup>[b]</sup>



# 2. Recursive Traversal



EB 01      jmp short 0x401020

FF            ???      (*UNREACHABLE*)

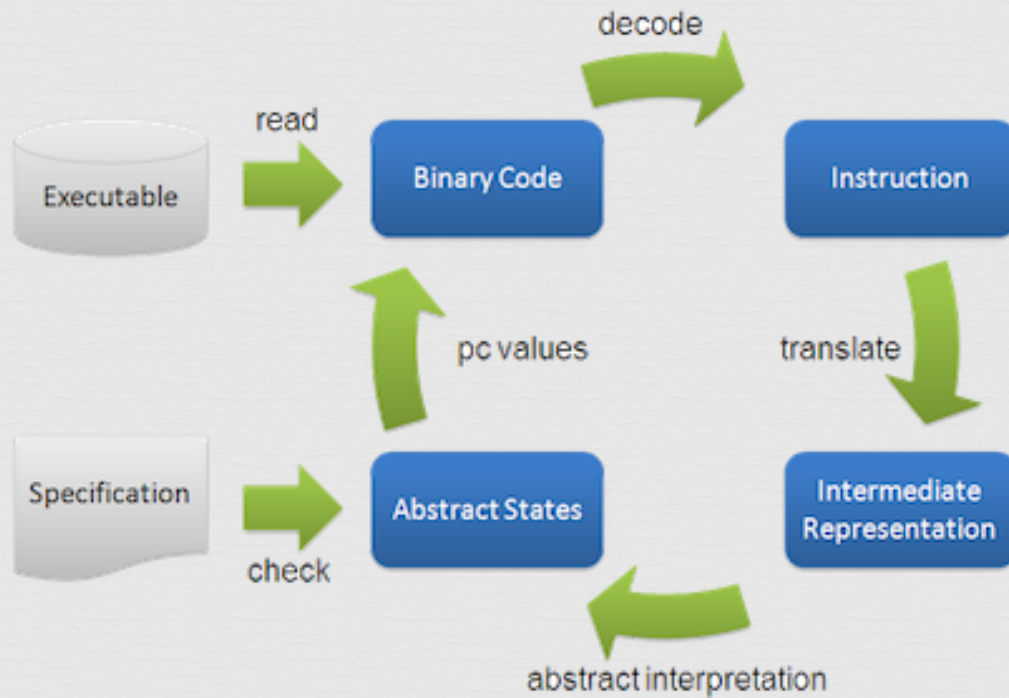
8B 45 FC    mov eax, dword ptr ss:[ebp-4]

# State-of-art CFG reconstruction

---

- ☞ Once identified the basic blocks, the CFG construction is done after the addition of the edges.
- ☞ CFG construction is especially difficult when the code includes *indirect calls*. (`call dword ptr[eax]`)
- ☞ State-of-art CFG construction available is the open-source Jakstab tool (Java Toolkit for Static Analysis of Binaries) from Johannes Kinder.
- ☞ Provides better results than IDAPro.

# Jakstab<sup>[d]</sup>



# Self-modifying code



Control Flow Analysis

# Self-modifying code

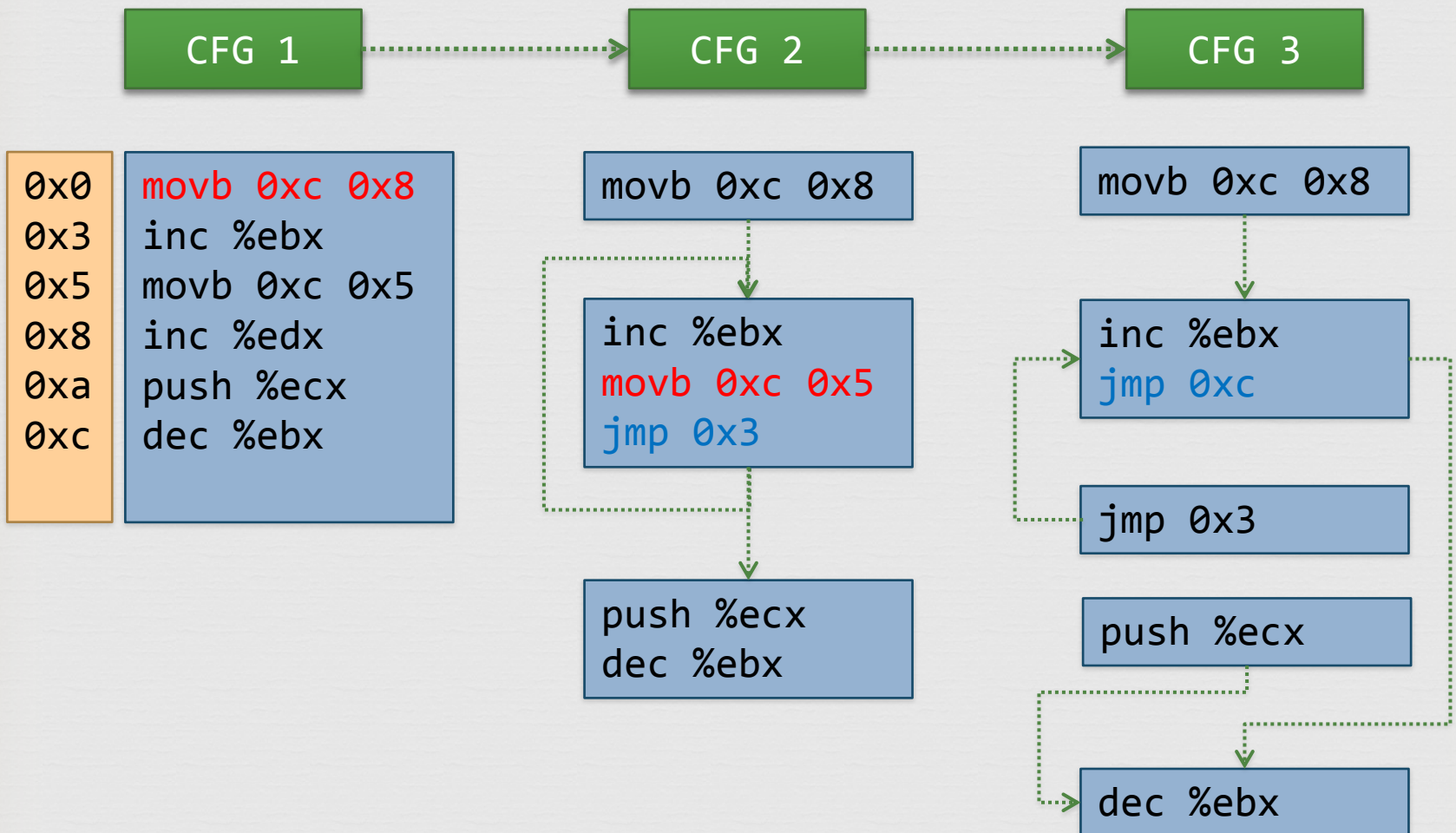


☞ Consider the following example<sub>[c]</sub> (*not real x86 opcodes*)

Address	Assembly	Binary
0x0	movb 0xc 0x8	c6 0c 08
0x3	inc %ebx	40 01
0x5	movb 0xc 0x5	c6 0c 05
0x8	inc %edx	40 03
0xa	push %ecx	ff 02
0xc	dec %ebx	48 01

☞ A linear sweep or recursive traversal algorithm execution on the above code would result in a single Basic Block (single entry/single exit/no branches)

# SMC

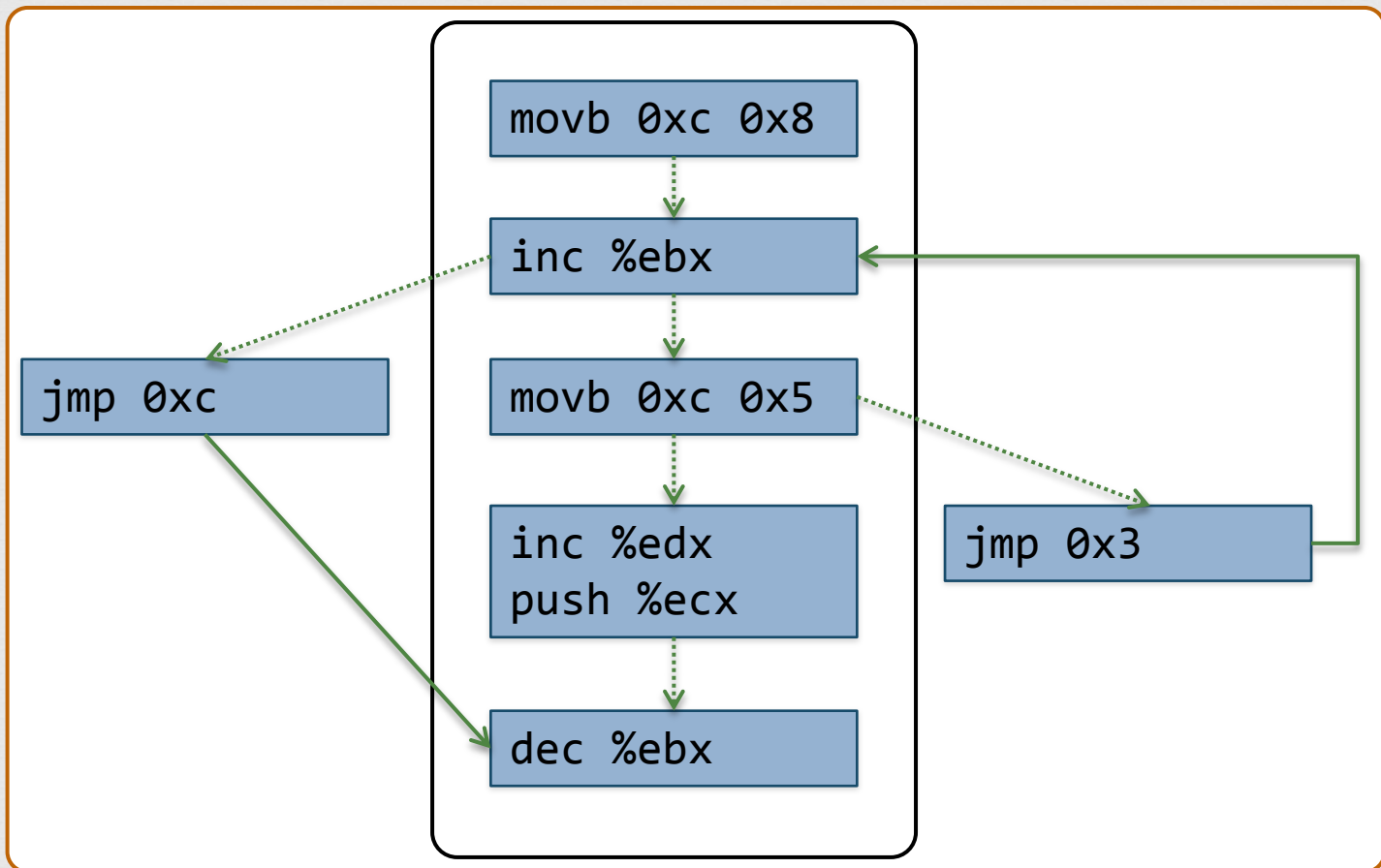


# SE-CFG

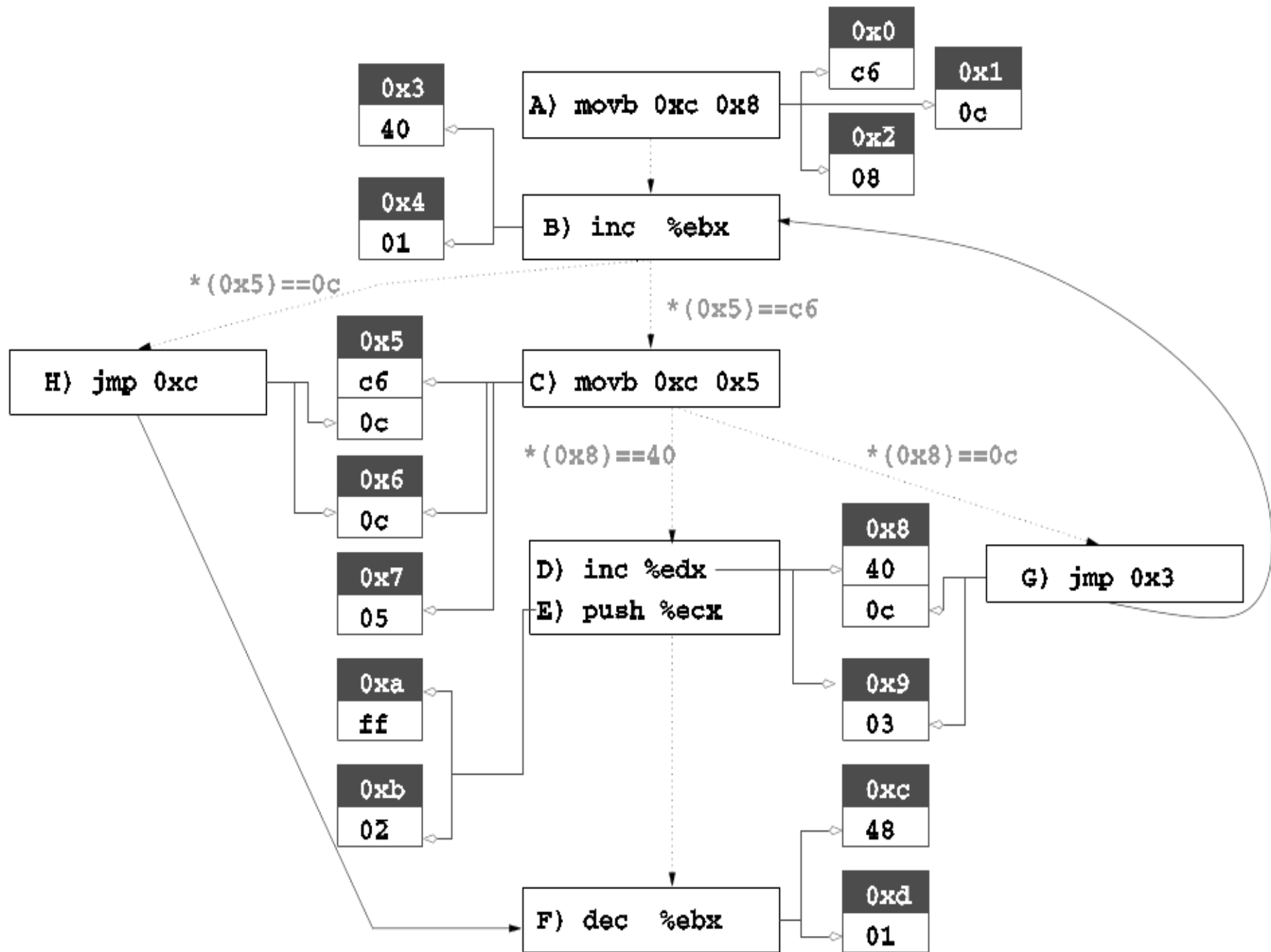


- ❧ State-Enhanced Control Flow Graph (SE-CFG)
- ❧ CFG augmented with extensions to support SMC.
- ❧ Allows the use of control flow analysis algorithms for SMC.
- ❧ “A Model for Self-Modifying Code”
- ❧ Codebyte extensions – Codebyte conditional edges
- ❧ Implemented in a link-time binary rewriter: Diablo.
- ❧ It can be downloaded from
  - ❧ <http://www.elis.ugent.be/diablo>

# SMC - CFG







# Dominators



Control Flow Analysis

# Dominance relation



- Relation about the nodes of a control flow graph.
- “Node A **dominates** Node B if *every path* from the **entry node** to B includes A”.
- Representation:  $A \text{ dom } B$
- Properties:
  - Antisymmetric (either  $A \text{ dom } B$  or  $B \text{ dom } A$ )
  - Reflexive ( $A \text{ dom } A$ )
  - Transitive (If  $A \text{ dom } B$  and  $B \text{ dom } C$  then  $A \text{ dom } C$ )
- Can be represented by a tree, the Dominator Tree.

# Control Flow Graph

Entry Node

```
009F25CB notepad.exe::_SkipBlanks@4
009F25CB mov     edi, edi
009F25CD push   ebp
009F25CE mov     ebp, esp
009F25D0 mov     eax, ss:[ebp+0000]
```

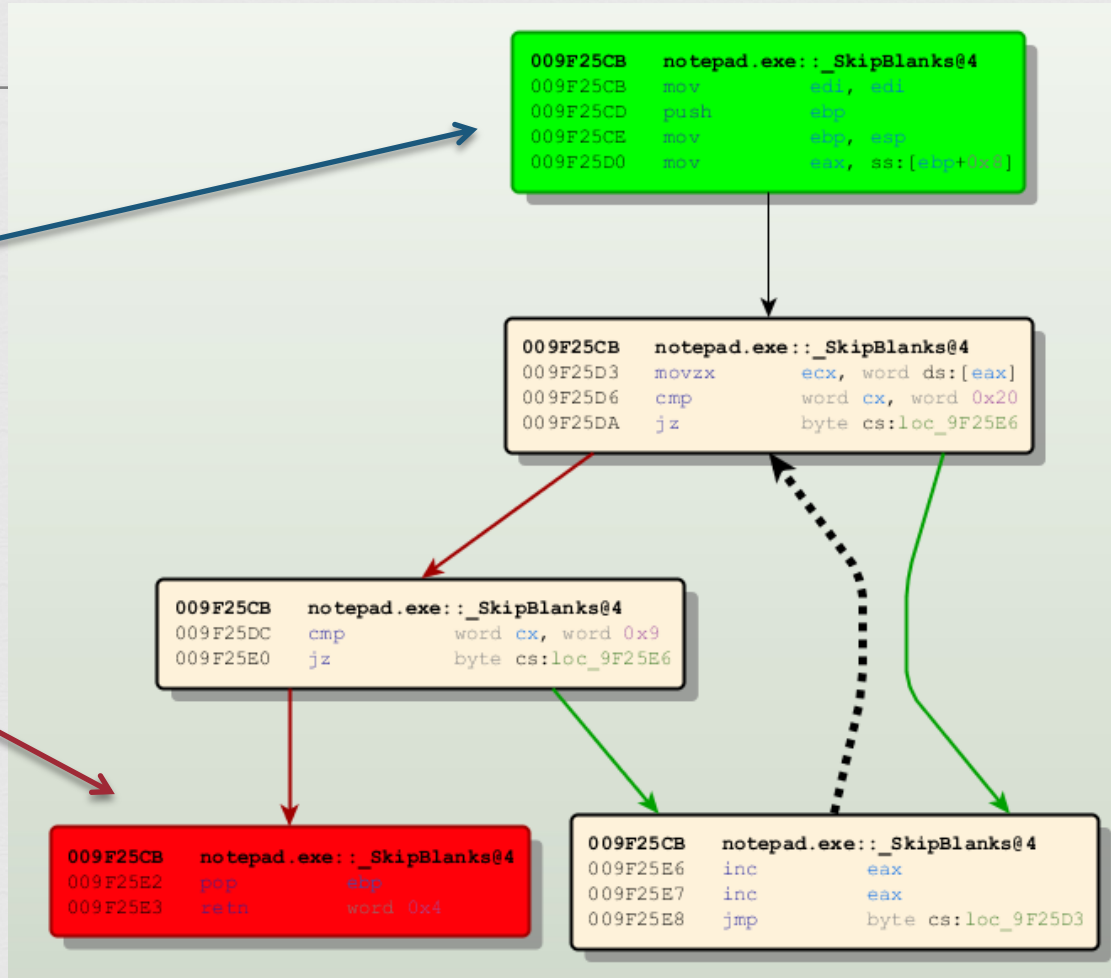
```
009F25CB notepad.exe::_SkipBlanks@4
009F25D3 movzx  ecx, word ds:[eax]
009F25D6 cmp    word cx, word 0x20
009F25DA jz     byte cs:loc_9F25E6
```

```
009F25CB notepad.exe::_SkipBlanks@4
009F25DC cmp    word cx, word 0x9
009F25E0 jz     byte cs:loc_9F25E6
```

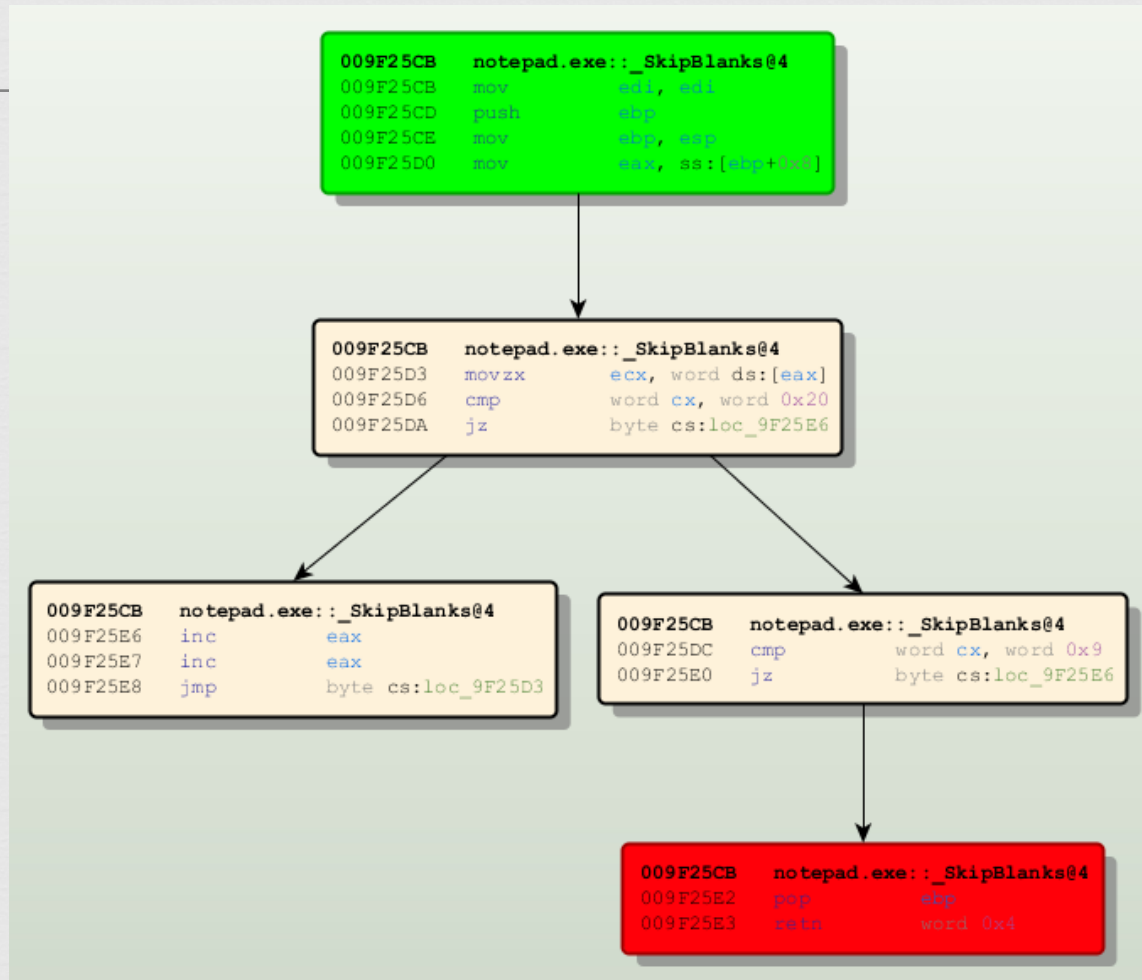
Exit node

```
009F25CB notepad.exe::_SkipBlanks@4
009F25E2 pop    ebp
009F25E3 retn  word 0x4
```

```
009F25CB notepad.exe::_SkipBlanks@4
009F25E6 inc    eax
009F25E7 inc    eax
009F25E8 jmp    byte cs:loc_9F25D3
```



# Dominator Tree



# Implementations



- ❧ Classic reference:
  - ❧ Lengauer-Tarjan algorithm
- ❧ Boost C++ library
- ❧ Immunity Debugger
  - ❧ libcontrolflow.py
    - ❧ Class DominatorTree
- ❧ BinNavi API
  - ❧ GraphAlgorithms *getDominatorTree()*
  - ❧ DEMO: Gui plugin

# Natural loops



- ⌘ We can use the Dominator Tree to identify loops.
- ⌘ Locate the back edges
- ⌘ Back edge:
  - ⌘ An edge whose **head** dominates its **tail**.
- ⌘ A loop consists:
  - ⌘ of **all nodes** dominated by its entry node (head of the back edge) from which the entry node can be reached
- ⌘ These loops are named *Natural Loops*.

Loop Header

```

009F25CB notepad.exe::_SkipBlanks@4
009F25CB mov     edi, edi
009F25CD push   ebp
009F25CE mov     ebp, esp
009F25D0 mov     eax, ss:[ebp+0x8]

```

```

009F25CB notepad.exe::_SkipBlanks@4
009F25D3 movzx   ecx, word ds:[eax]
009F25D6 cmp     word cx, word 0x20
009F25DA jz      byte cs:loc_9F25E6

```

```

009F25CB notepad.exe::_SkipBlanks@4
009F25DC cmp     word cx, word 0x9
009F25E0 jz      byte cs:loc_9F25E6

```

```

009F25CB notepad.exe::_SkipBlanks@4
009F25E2 pop     ebp
009F25E3 retn   word 0x4

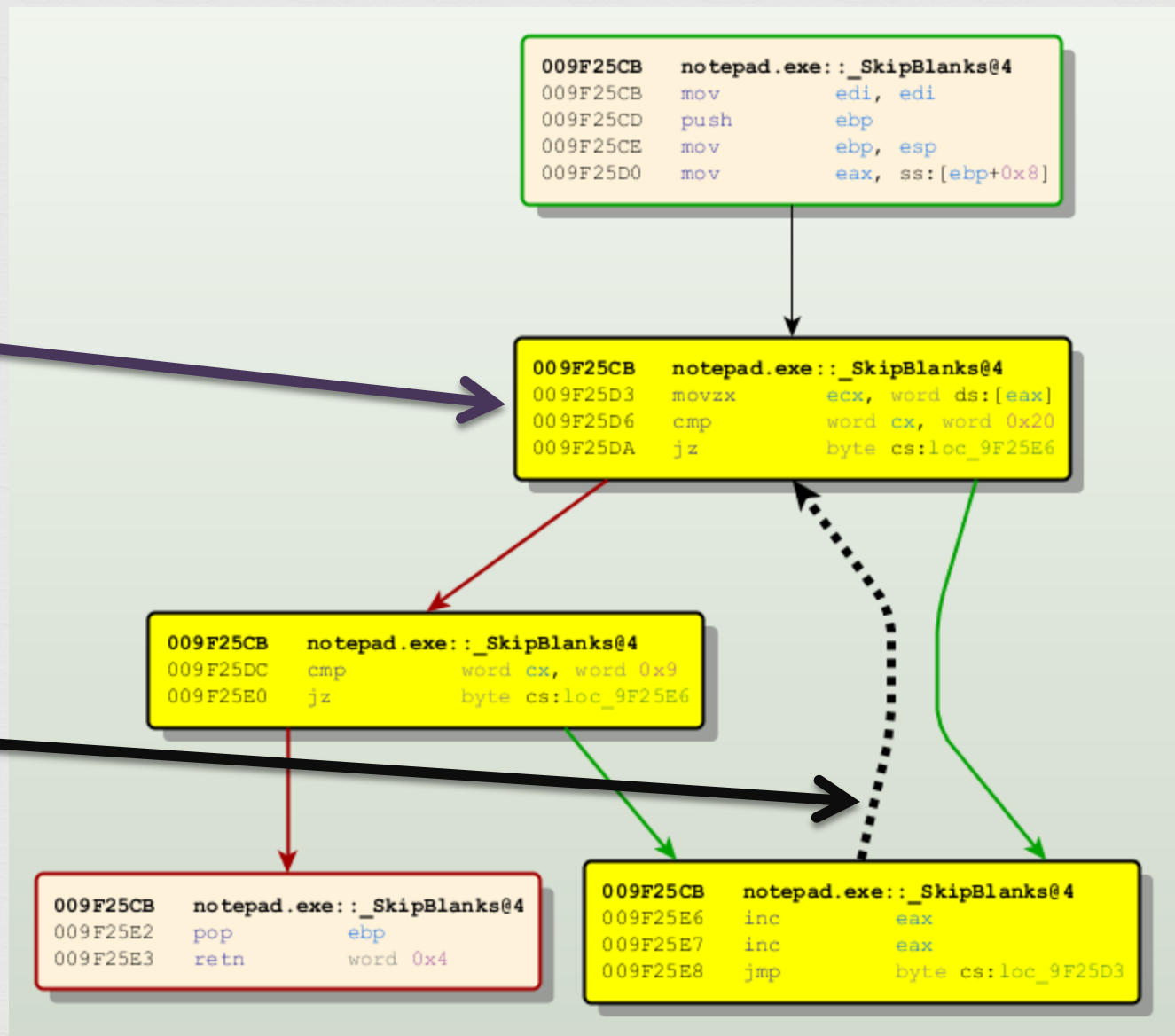
```

```

009F25CB notepad.exe::_SkipBlanks@4
009F25E6 inc     eax
009F25E7 inc     eax
009F25E8 jmp     byte cs:loc_9F25D3

```

Back Edge





# ImmunityDbg !findloop

```
00B: 00B225D3 LOOP! from:0x00b225d3, to:0x00b225e6
00B: 00B225DC Loop node:0x00b225dc
00B: 00B225D3 Loop node:0x00b225d3
00B: 00B225E6 Loop node:0x00b225e6
00B: Done!
ImmDbg\PyCommands\findloop.py

00B225CA 90 NOP
00B225CB $ 8BFF MOV EDI,EDI
00B225CD . 55 PUSH EBP
00B225CE . 8BEC MOV EBP,ESP
00B225D0 . 8B45 08 MOV EAX,DWORD PTR SS:[EBP+8]
00B225D3 > 0FB708 MOUZX ECX,WORD PTR DS:[EAX] \ Loop 0x00B225D3 Node
00B225D6 . 66:83F9 20 CMP CX,20 |
00B225DA . 74 0A JE SHORT notepad.00B225E6 |
00B225DC . 66:83F9 09 CMP CX,9 | Loop 0x00B225D3 Node
00B225E0 . 74 04 JE SHORT notepad.00B225E6 |
00B225E2 . 5D POP EBP
00B225E3 . C2 0400 RETN 4
00B225E6 > 40 INC EAX | Loop 0x00B225D3 Node
00B225E7 . 40 INC EAX |
00B225E8 ^ EB E9 JMP SHORT notepad.00B225D3 /
00B225EA 90 NOP
00B225EB 90 NOP
```

Address	Hex dump	ASCII
00B2C000	00 00 00 00 78 00 00 00	....x...
00B2C008	01 00 00 00 FF FF FF FF	...ÿÿÿÿ
00B2C010	4E E6 40 BB B1 19 BF 44	Næ@>± ±;D
00B2C018	00 00 00 00 00 00 00 00	.....
00B2C020	00 00 00 00 00 00 00 00	.....
00B2C028	00 00 00 00 00 00 00 00	.....
00B2C030	00 00 00 00 00 00 00 00	.....

```
!findloop -a 0xb225cb
```

# Strongly connected components



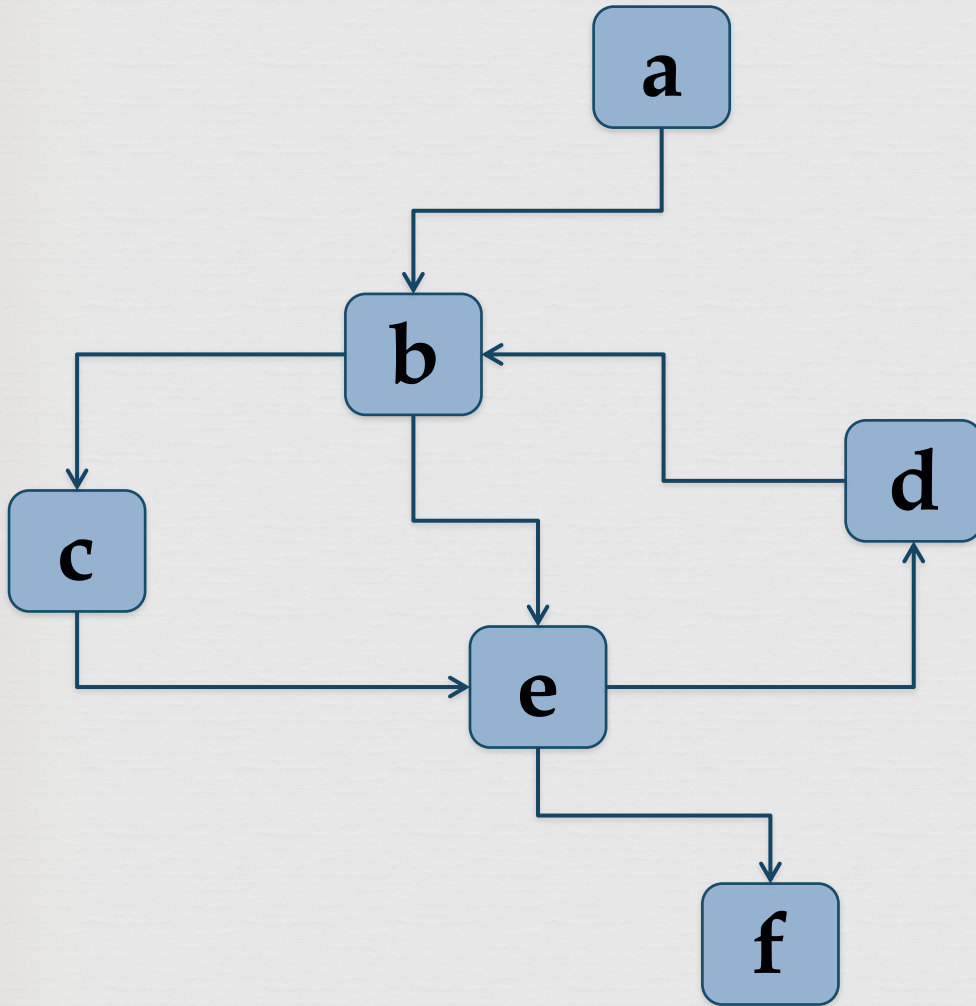
Control Flow Analysis

# SCC



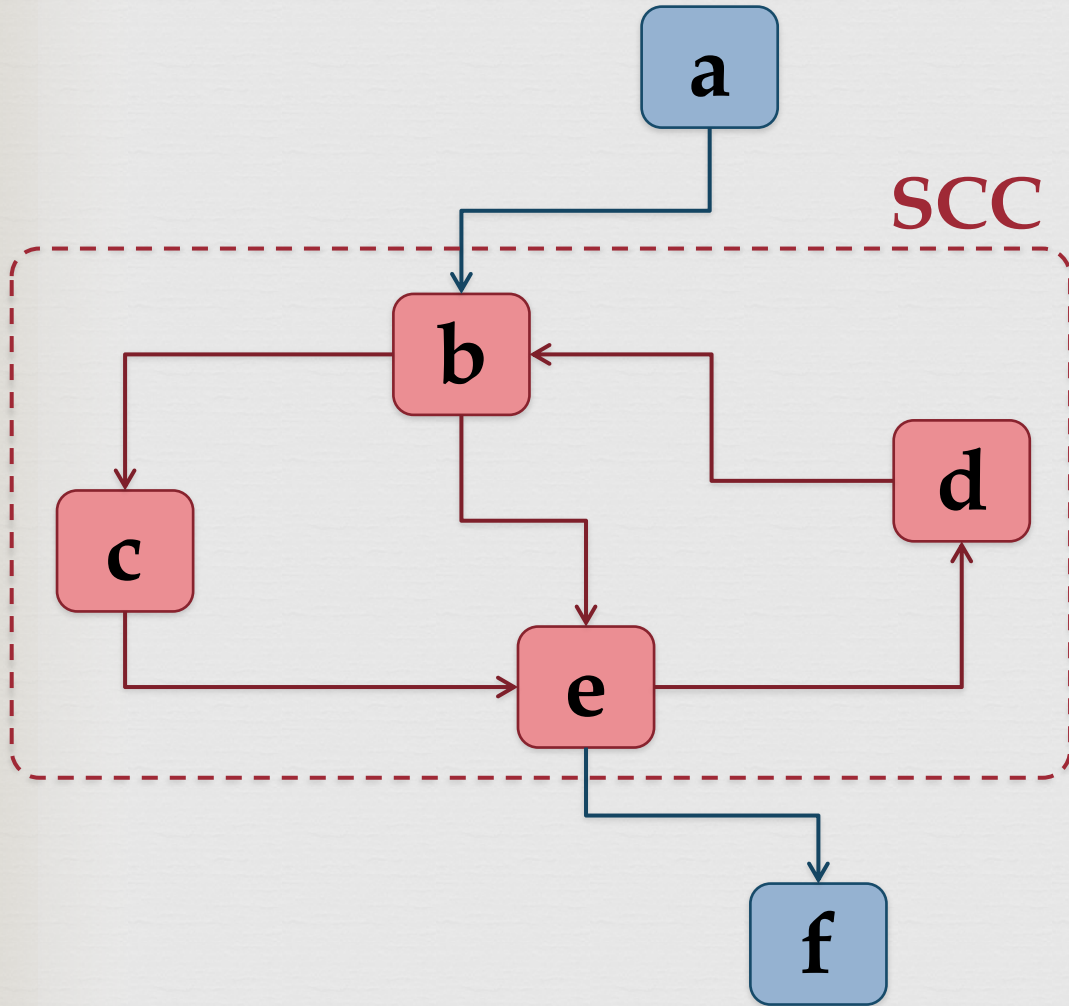
- ∞ SCC → Strongly connected components
- ∞ A graph (directed/undirected) is called *strongly connected* if there is a *path* from each *vertex* to every other vertex
- ∞ Any loop is a strongly connected component

# SCC



*This graph is **not** strongly connected.*

# SCC



*But it contains  
a subgraph  
which is  
strongly-  
connected.*

# SCC - algorithms



- ❧ Tarjan algorithm
  - ❧ fast algorithm - complex
- ❧ Kosaraju-Sharir algorithm
  - ❧ simple, but slower than Tarjan's algorithm
- ❧ Implementations available for all languages:
  - ❧ C#/Python/Lua/Ruby/Java

```
009F360C notepad.exe::_SkipProgramName@4
009F360C mov     edi, edi
009F360E push  ebp
009F360F mov     ebp, esp
009F3611 mov     eax, eax [ebp+0x0]
009F3614 mov     ecx, word ptr ds:[eax]
009F3617 push  edi
009F3618 xor     esi, esi
009F361A test   word ptr word ptr
009F361D ja     byte ptr loc_9F361F
```

```
009F360C notepad.exe::_SkipProgramName@4
009F361F mov     ecx, word ptr
```

```
009F360C notepad.exe::_SkipProgramName@4
009F3622 cmp     word ptr word ptr
009F3626 ja     byte ptr loc_9F362E
```

```
009F360C notepad.exe::_SkipProgramName@4
009F362E cmp     word ptr word ptr
009F3630 ja     byte ptr loc_9F3630
```

```
009F360C notepad.exe::_SkipProgramName@4
009F3635 test   esi, esi
009F3637 ja     byte ptr loc_9F3637
```

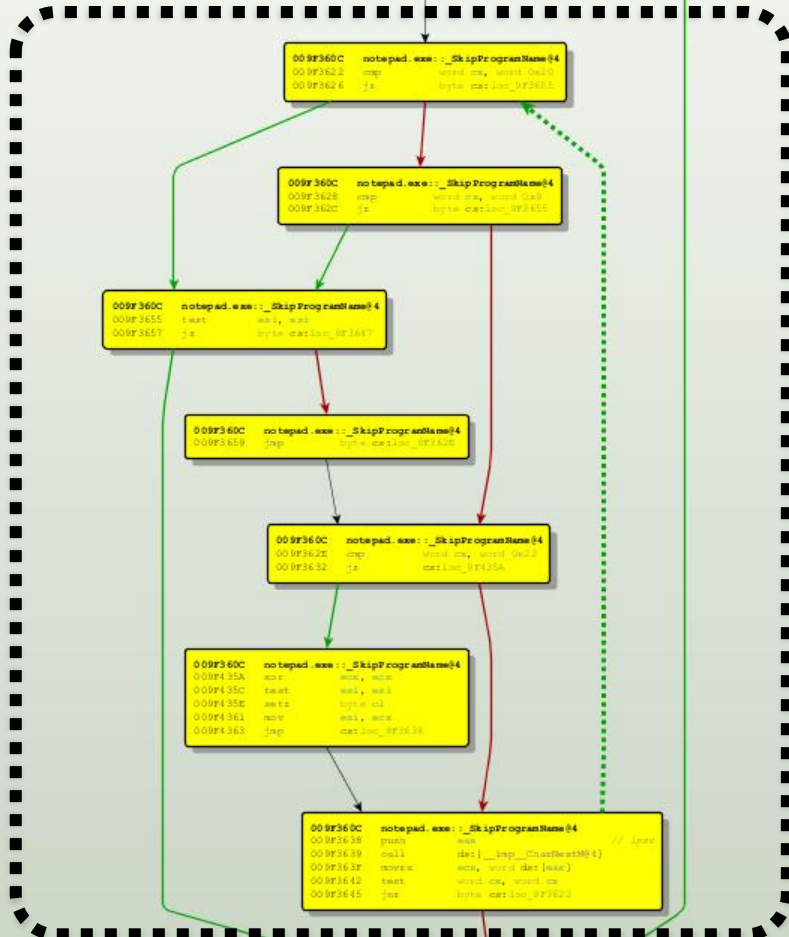
```
009F360C notepad.exe::_SkipProgramName@4
009F3639 jmp     byte ptr loc_9F3639
```

```
009F360C notepad.exe::_SkipProgramName@4
009F363E cmp     word ptr word ptr
009F3640 ja     byte ptr loc_9F3640
```

```
009F360C notepad.exe::_SkipProgramName@4
009F364A xor     esi, esi
009F364C test   esi, esi
009F364E scasd  byte ptr
009F3651 mov     esi, esi
009F3653 jmp     byte ptr loc_9F3653
```

```
009F360C notepad.exe::_SkipProgramName@4
009F3658 push  esi
009F365B call   dword ptr [ebp+0x0] // jmp
009F365F mov     ecx, word ptr ds:[eax]
009F3662 test   word ptr word ptr
009F3665 jmp     byte ptr loc_9F3662
```

```
009F360C notepad.exe::_SkipProgramName@4
009F3667 pop     esi
```



# Interval Analysis



Control Flow Analysis

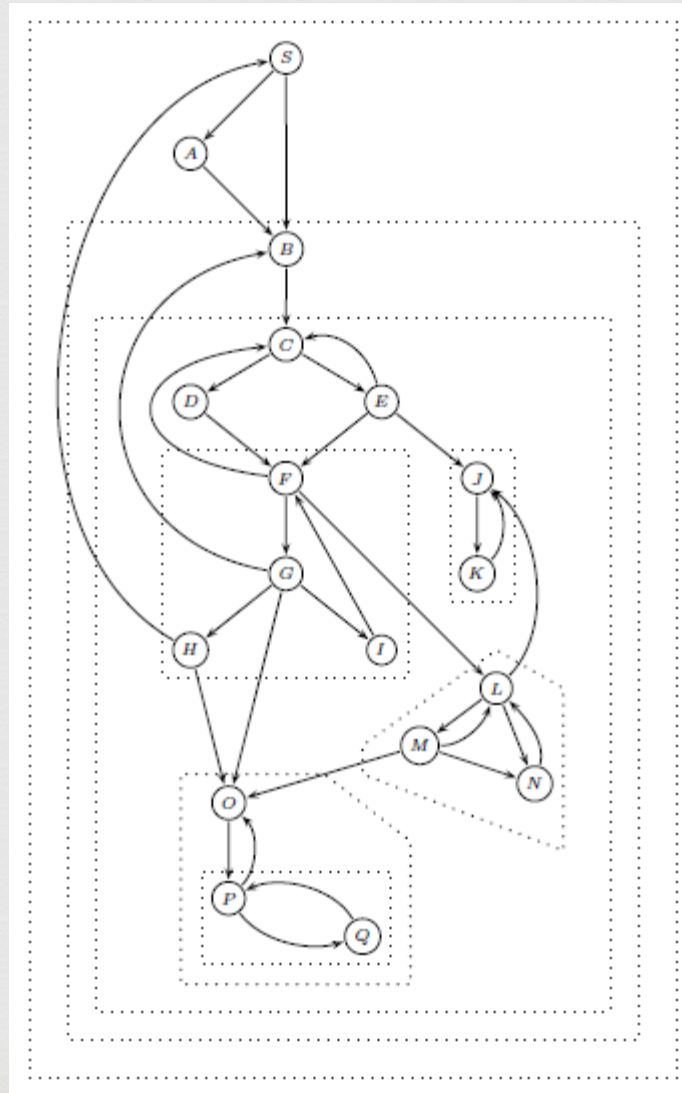


# Regions and intervals



- ❧ Unfortunately SCC isn't able to identify *nested* loops
- ❧ Interval Analysis
  - ❧ Divides the CFG into regions and consolidate them into new nodes (*abstract nodes*) resulting in an *abstract flowgraph*.
- ❧ We need to identify *regions* and *pre-intervals*
- ❧ Region:
  - ❧ A *region* in a flow graph is a sub graph  $H$  with an unique entry node  $h$
- ❧ Pre-Interval:
  - ❧ A *pre-interval* in a flow graph is a region  $\langle H, h \rangle$  such that every cycle (loop) in  $H$  includes the header  $h$ .
- ❧ Similar to a unique entry SCC.

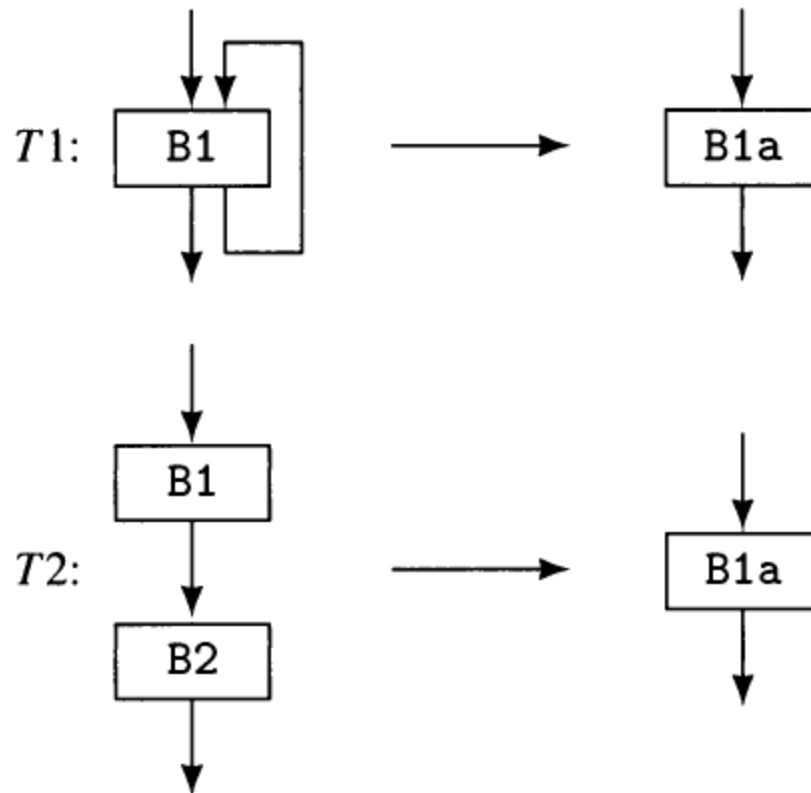
# Nested Intervals



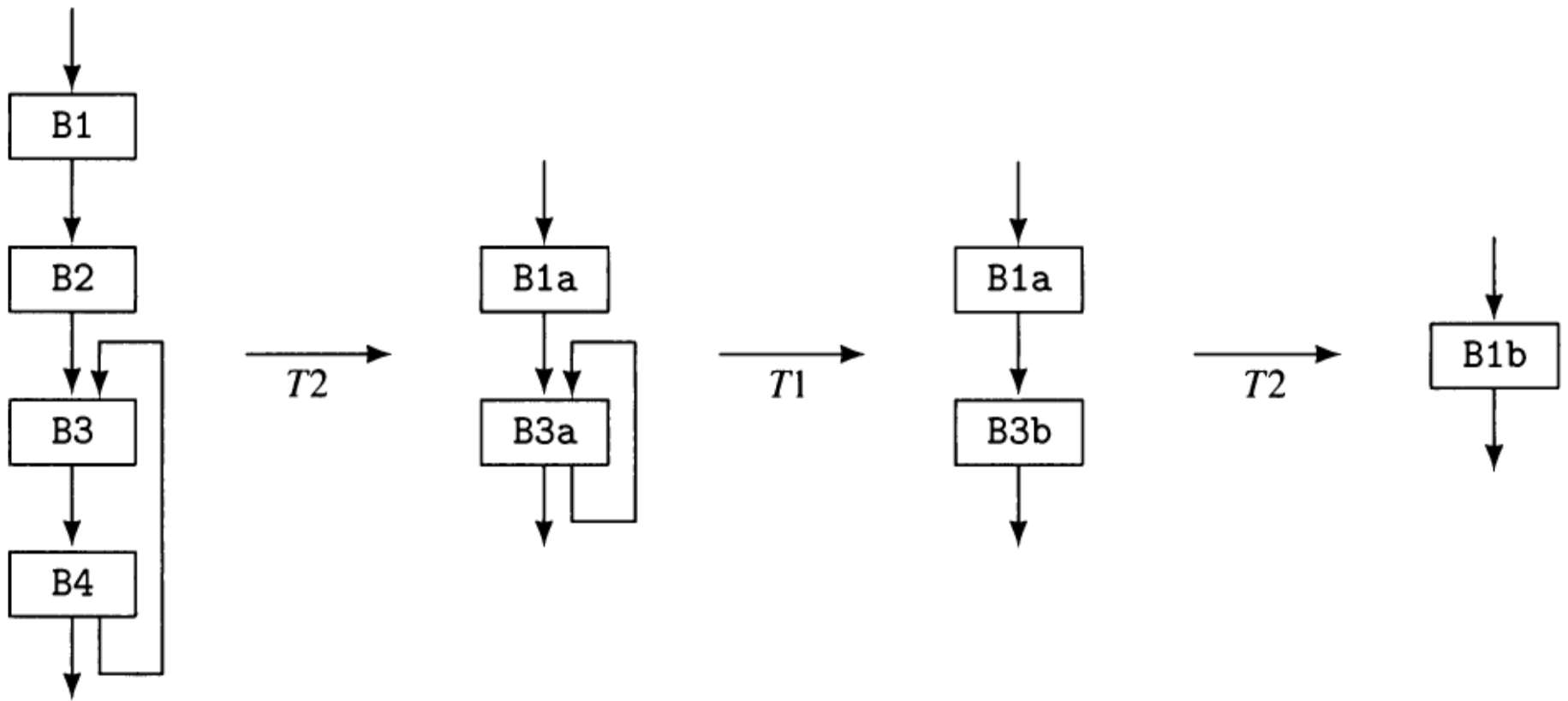
# T1/T2 transformations



- ✧ Reduction of graphs
- ✧ We can collapse nodes from a region to a single node. This is called t1/t2 transformation. If we apply it to all loops, the graph becomes a *cycle-free* one.
- ✧ Cycle-free graphs are easier to analyze.



*T1-T2 transformations.*



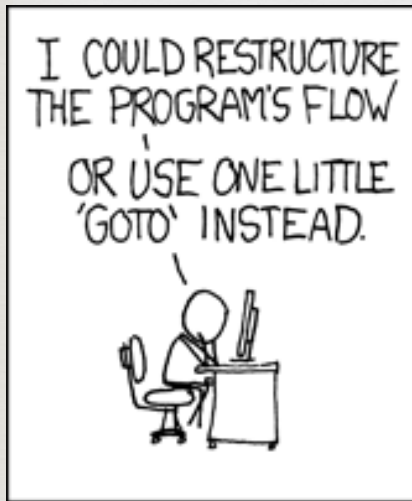
Example of  $T1$ - $T2$  transformations.

# Interval analysis



 DEMO

# GOTO considered harmful...



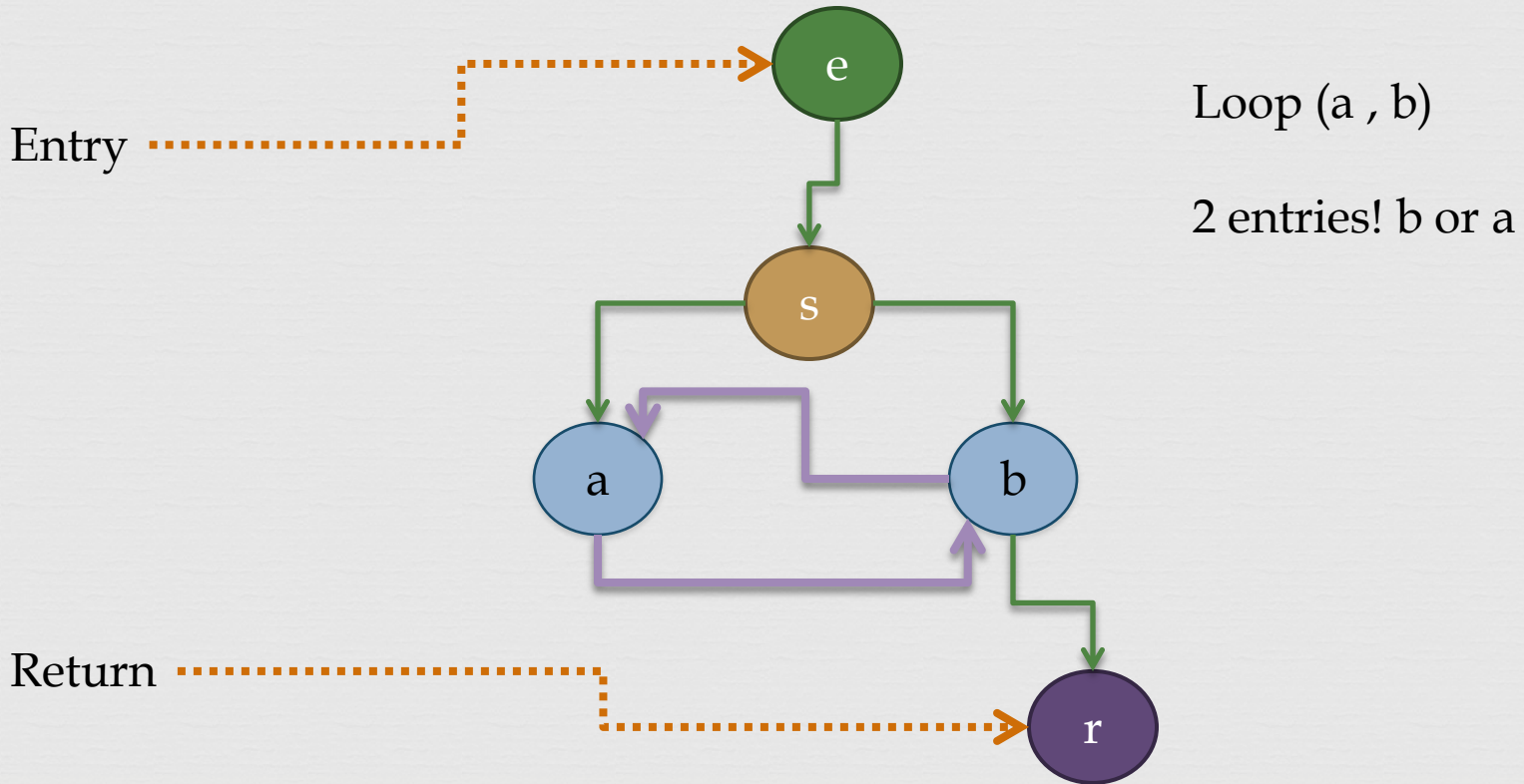
# Irreducible graphs



- ☞ All the loops identified by the previous methods (dominance tree/interval analysis) are called *natural loops*.
- ☞ They are *unique entry* loops.
- ☞ There another type of loop:
  - ☞ irreducible graphs or improper regions



# Irreducible graph



# Irreducible graphs



- ❧ Who codes like that?
  - ❧ Anyone who uses GOTO
  - ❧ It is rare, but it does exist
    - ❧ notepad.exe
    - ❧ ntoskrnl.exe (Windows Kernel)
- ❧ What's the problem?
  - ❧ Most of the algorithms are unable to handle irreducible graphs!!! Including Interval analysis.
  - ❧ Can't apply T1/T2

# translateString



```
int *__stdcall TranslateString(int a1)
{
    wchar_t v1; // cx@1
    ...
    if ( v1 )
    {
        while ( 1 )
        {
            v5 = &v22 + v26;
            ...
            LABEL_49:
            v1 = *(_WORD *)v7;
            ...
        }
        goto LABEL_49;
    }
}
```

Jump inside the  
WHILE statement

# Solutions



- There are 2 main solutions to handle *irreducible graphs*:
  - Structural Analysis
  - DJ-Graphs

# Structural Analysis



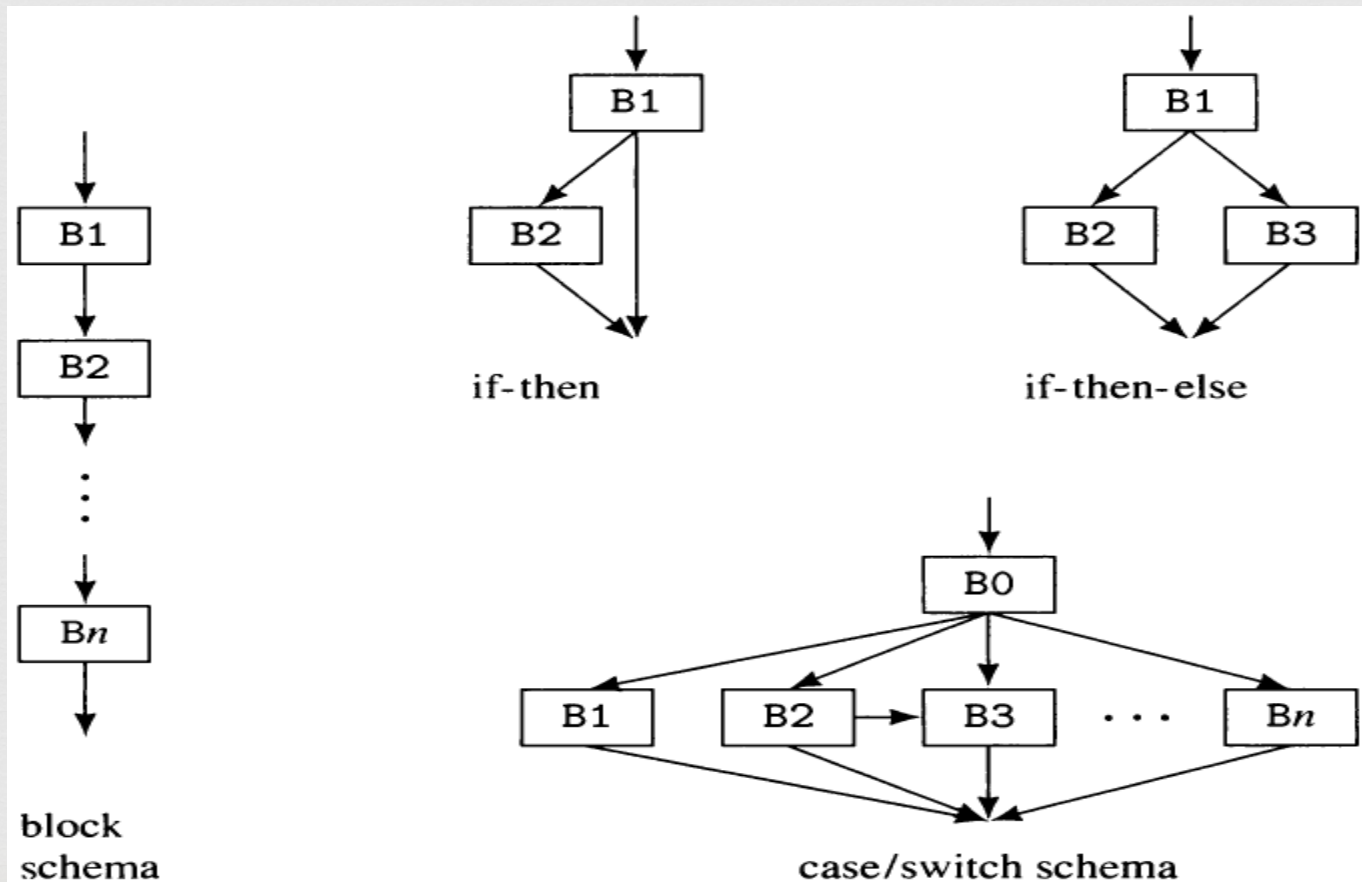
Control Flow Analysis

# Structural Analysis

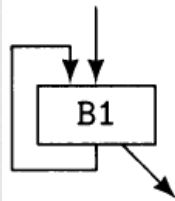


- ✧ Structural analysis will identify the main language constructs inside a flow graph using *region schemas*.
- ✧ Do you want to build your own decompiler?
  - ✧ Hex-Rays decompiler internally uses Structural Analysis
- ✧ Created by Micha Sharir
- ✧ Reference paper:
  - ✧ Structural analysis: a new approach to flow analysis in optimizing compilers (1979)

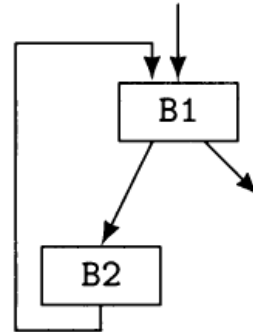
# Acyclic schemas



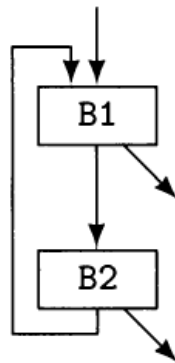
# Cyclic schemas



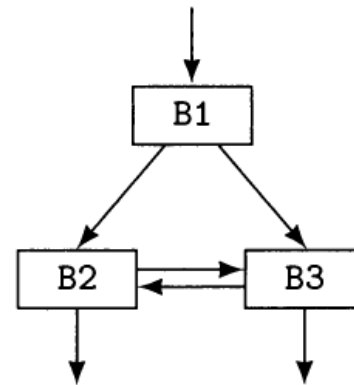
self loop



while loop



natural loop schema



improper interval schema

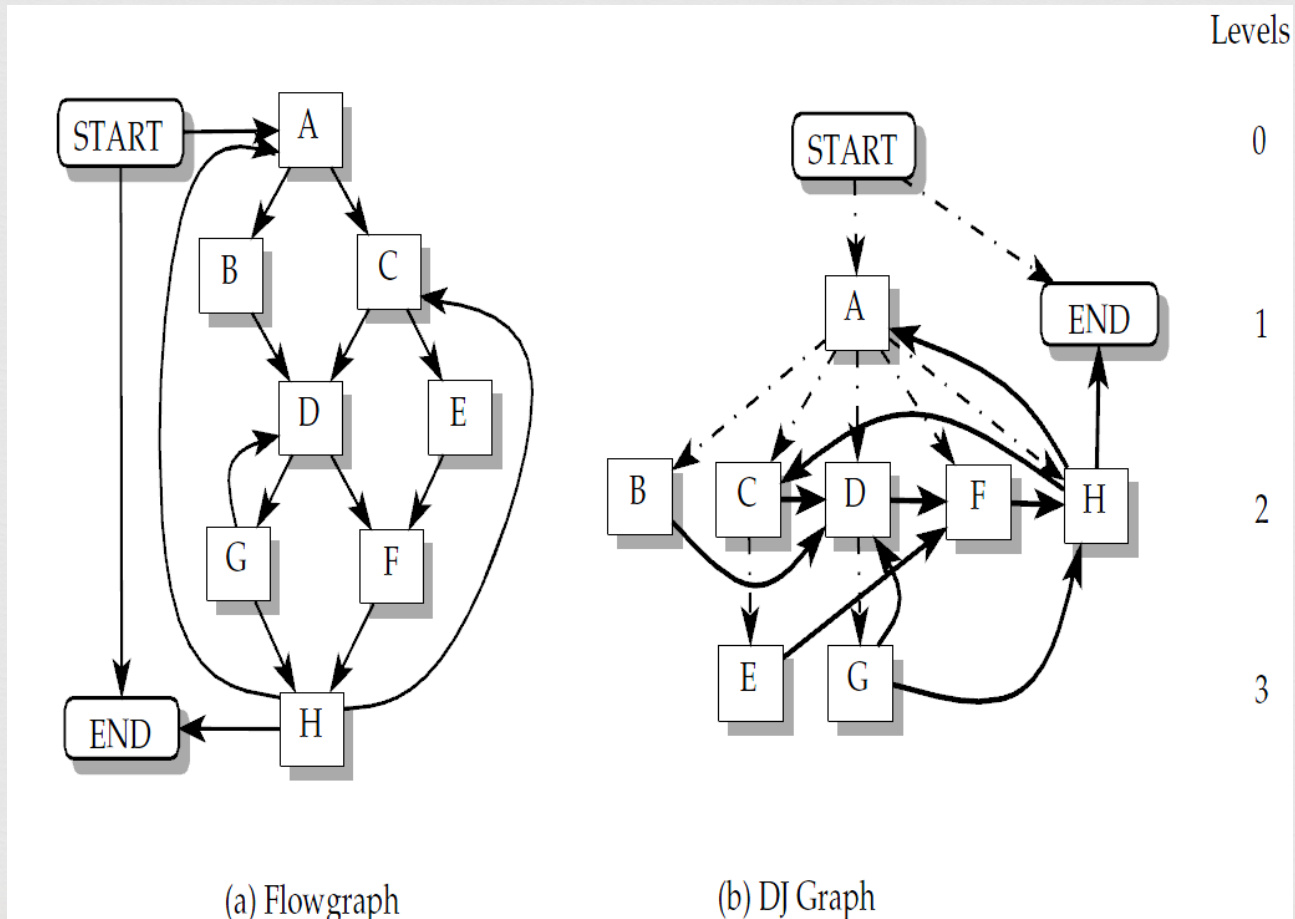


# DJ-Graphs



- ✧ Another way to handle *irreducible graphs*.
- ✧ It is also able to identify all types of structures, including improper regions and nested structures.
- ✧ Uses a combination of the dominance tree and the original flowgraph with two additional types of edges:
  - ✧ the D edge (*Dominator*)
  - ✧ the J edges
- ✧ Paper: *Identifying loops using DJ graphs*.<sup>[e]</sup>

# DJ-Graphs



# Applications



- ❧ Taint analysis
  - ❧ Control dependency (dominators, post-dominators)
- ❧ Diff Slicing
  - ❧ Execution Indexing (view the CFG as a grammar)
    - ❧ Execution alignment
  - ❧ Identification of root causes of software crashes
- ❧ Decompilation
- ❧ Code coverage
- ❧ Bug finding

# References



- ⌘ a - [http://www.usenix.org/event/usenix03/tech/full\\_papers/prasad/prasad\\_html/node5.html](http://www.usenix.org/event/usenix03/tech/full_papers/prasad/prasad_html/node5.html)
- ⌘ b - An Abstract Interpretation-Based Framework for Control Flow Reconstruction from Binaries
- ⌘ c - Bertrand Anckaert, Matias Madou, and Koen De Bosschere. 2006. A model for self-modifying code. In *Proceedings of the 8th international conference on Information hiding (IH'06)*
- ⌘ d - <http://www.jakstab.org/>
- ⌘ e - Vugranam C. Sreedhar, Guang R. Gao, and Yong-Fong Lee. 1996. Identifying loops using DJ graphs. *ACM Trans. Program. Lang. Syst.* 18, 6 (November 1996), 649-658.
- ⌘ f - Advanced compiler implementation - Steven Muchnick
- ⌘ g - Notes on Graph Algorithms Used in Optimizing Compilers - Carl D. Offner

# Questions?



- ☞ Contact: [edgarmb@gmail.com](mailto:edgarmb@gmail.com)  
[edgar@research.coseinc.com](mailto:edgar@research.coseinc.com)
- ☞ twitter: @embarbosa