



Lecture 38: Reverse Engineering

Professor Adam Bates
CS 46I / ECE 422
Fall 2019

Goals for Today



- Learning Objectives:
 - Understand the objectives and methods of reverse engineering
 - Review several examples of important RE results in the literature
- Announcements, etc:
 - Forensics CP2 due **December 6th**
 - Final Exam — 7pm December 13th

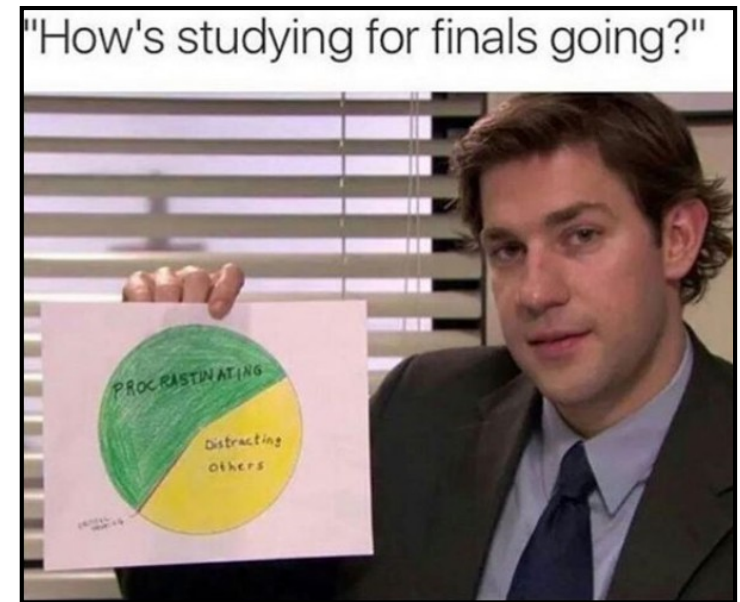


Reminder: Please put away devices at the start of class

Final Details



- December 13th, 7-10pm
- Here, 1404 Siebel
- Multiple choice + short answer
- **Closed book.**
- No electronic devices permitted (or necessary)!
- **Content:** All lectures, and MP3, MP4, MP5.
- Sample exams are not available for the final; feel free to re-review midterm sample exams



Final Details



- Changes from midterm:
 - Multiple Choice will now be scantron. Only one answer is correct per question. No points deducted for guessing.
 - More multiple choice questions, ~30+ as opposed to 20.
 - More short answer questions
 - MP questions will account for a smaller percentage of overall score.



- *“The process of analyzing a subject system to identify the system's components and their interrelationships, and to create representations of the system in another form or at a higher level of abstraction.”*
- Term originates in hardware, attempting to decipher design documents from finished products.
- Term can be applied to a variety of sociotechnical products today.

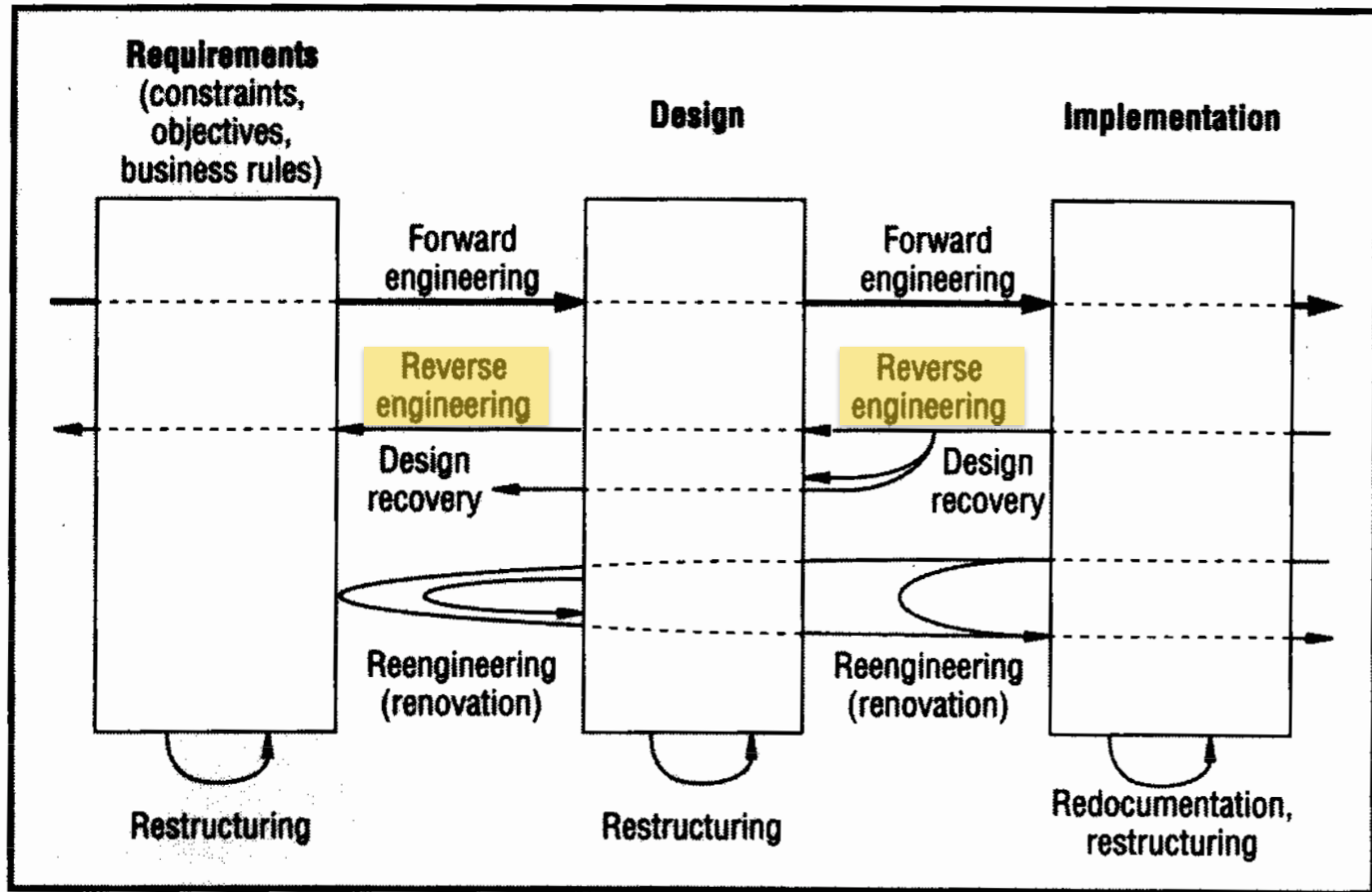
[Chikofsky and Cross, IEEE'90]

This semester...



- Arguably, we've actually already covered reverse-engineering a lot:
 - Binary Exploitation almost always requires RE
 - Software Testing lecture discussed RE
 - Malware analysis to determine malware payloads
 - During cloud lecture, we discussed the reverse engineering of an Amazon data center policy
 - ...
- But if you want more, I'll give you more! :)

Reverse Engineering



[Chikofsky and Cross, IEEE'90]



- Reverse Engineering (RC), Reverse Code Engineering (RCE)
- reverse engineering -- process of discovering the technological principles of a [insert noun] through analysis of its structure, function, and operation.
- The development cycle ... backwards

Why Reverse Engineer?

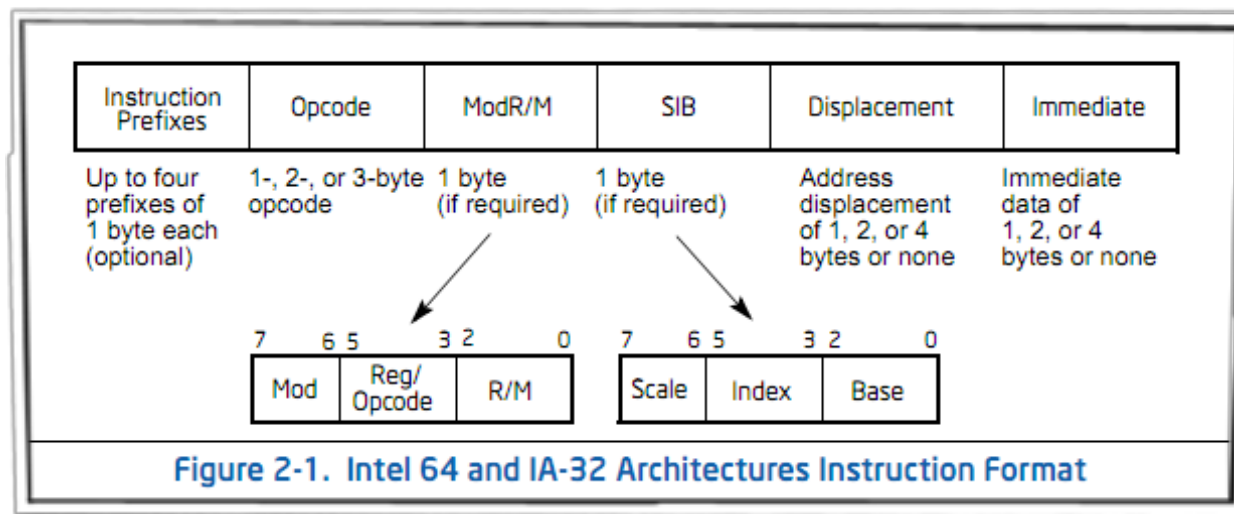


- Malware analysis
- Vulnerability or exploit research
- Check for copyright/patent violations
- Interoperability (e.g. understanding a file/protocol format)
- Copy protection removal



- Static Code Analysis (structure)
 - Disassemblers
- Dynamic Code Analysis (operation)
 - Tracing / Hooking
 - Debuggers

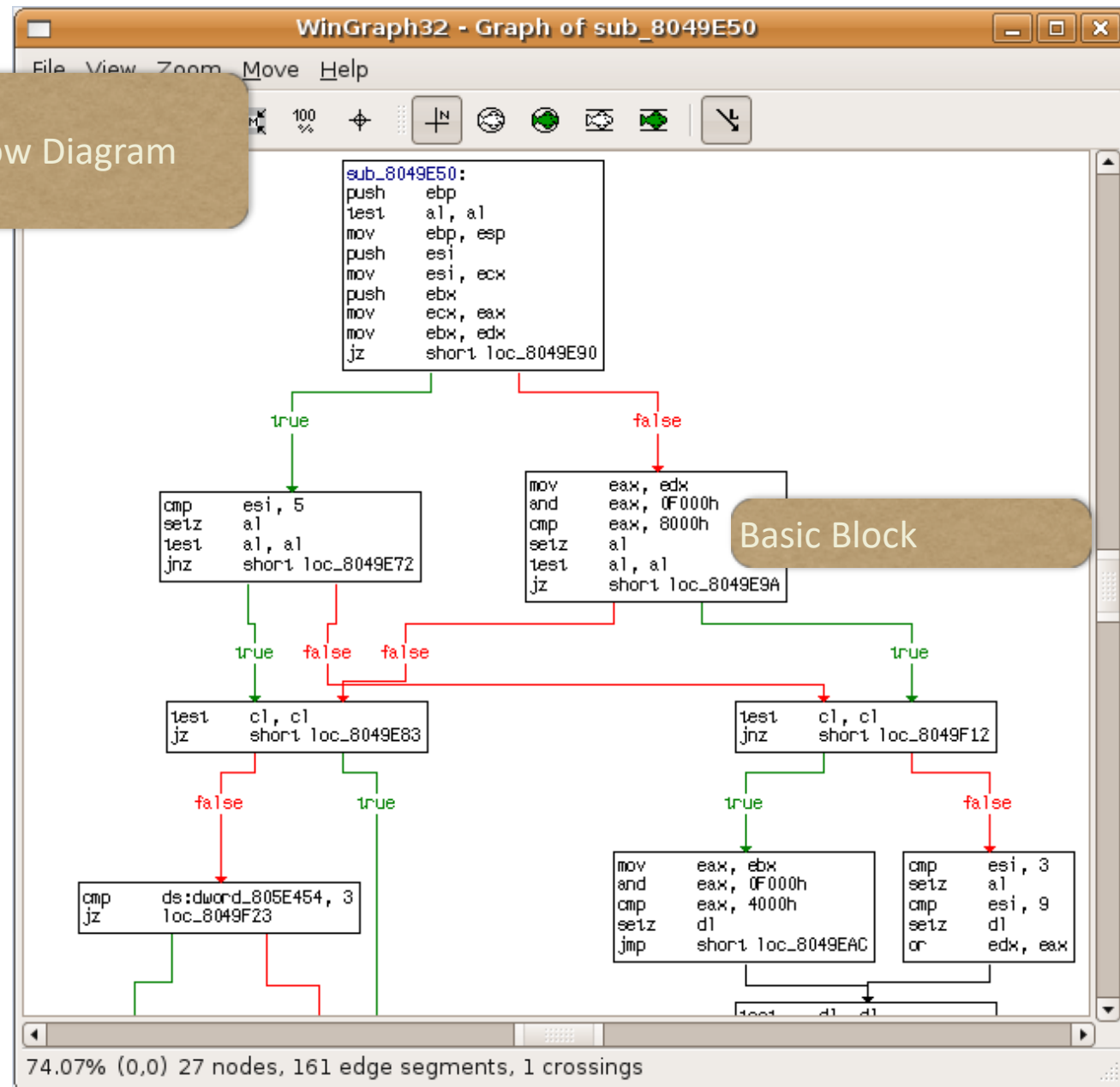
Disassembly



Disassembly



Control Flow Diagram



Challenges to Disassembly



- Static analysis isn't perfect at recovering original code
- “Benign” Optimizations
 - Constant folding
 - Dead code elimination
 - Inline expansion
 - etc...
- Intentional Obfuscation (note: not just for attackers!)
 - Packing
 - No-op instructions

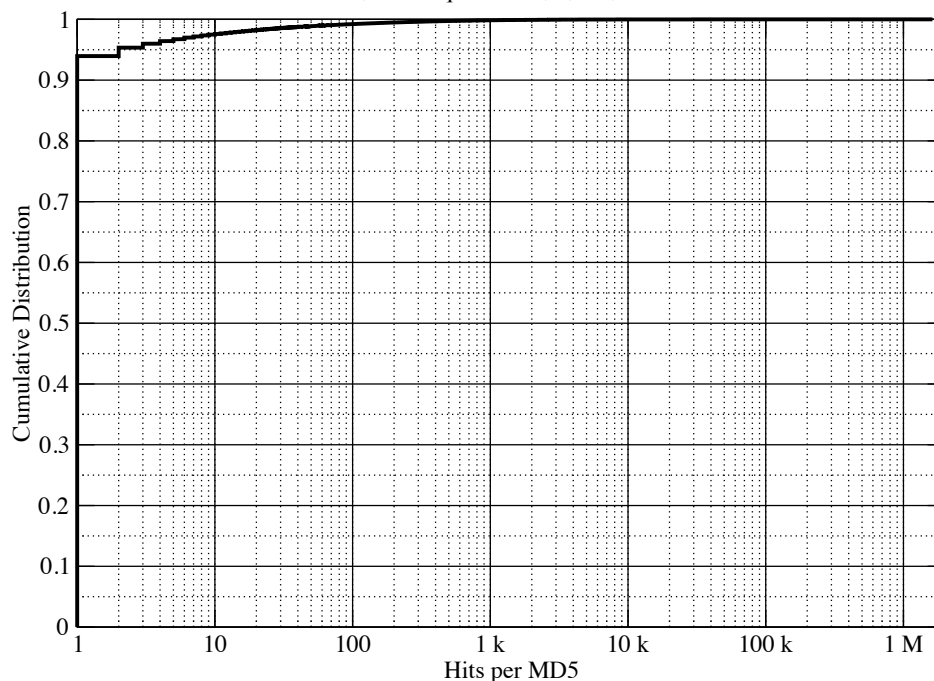
E.G.: Malware Packing



- Packing is used by software developers specifically to avoid reverse engineering... just ask malware authors!

Cumulative Distribution of Hits per MD5

246,952 unique MD5s, 5,772,891 Hits



Packer identification
98,801 malware samples

PEiD	Count
UPX	11244
Upack	6079
PECompact	4672
Nullsoft	2295
Themida	1688
FSG	1633
tElock	1398
NsPack	1375
ASpack	1283
WinUpack	1234

Identified: 59,070 (60%)

Top 10: 33.3%

SigBuster	Count
Allapple	22050
UPX	11324
PECompact	5278
FSG	5080
Upack	3639
Themida	1679
NsPack	1645
ASpack	1505
tElock	1332
Nullsoft	1058



- Trace every instruction a program executes -- single step
- Or, let program execute normally until an exception
- At every step or exception, can observe / modify:
- Instructions, stack, heap, and register set
- May inject exceptions at arbitrary code locations
- INT 3 instruction generates a breakpoint exception

Debugging



C CPU - main thread, module ollydbg

Address	Hex dump	Command	Comment
00401020	• 6A 00	PUSH 0	
00401022	• E8 85C60E00	CALL <JMP.&KERNEL32.Get	[Module = KERNEL
00401027	• 8BD0	MOV EDX, EAX	
00401029	• E8 C6E20D00	CALL 004DF2F4	
0040102E	• 5A	POP EDX	
0040102F	• E8 24E20D00	CALL 004DF258	
00401034	• E8 FBE20D00	CALL 004DF334	
00401039	• 6A 00	PUSH 0	
0040103B	• E8 14F80D00	CALL 004E0854	ollydbg [Arg1 = ollydbg
00401040	• 59	POP ECX	
00401041	• 60 00F14F00	PUSH 00F14F00	

Dest=ollydbg.004E0854

Address	Hex dump
004EE080	DC 88 4E 00 00 03 C8 8A 4E 00 00 06 00 90 4E 00
004EE090	00 01 00 91 4E 00 00 01 B0 93 4E 00 00 00 34 95
004EE0A0	4E 00 00 00 4C 95 4E 00 00 20 D1 BA 4E 00 00 1F
004EE0B0	84 D4 4E 00 00 1E 4C EB 4D 00 00 1E B8 12 4D 00
004EE0C0	00 1E 08 70 4D 00 00 1E C0 73 4D 00 00 1E 10 71
004EE0D0	4D 00 00 1E EC A8 4D 00 00 1E A0 70 4D 00 00 20
004EE0E0	F4 EE 4D 00 00 00 3C F2 4D 00 00 1F 78 D4 4E 00
004EE0F0	00 20 D4 F3 4D 00 00 20 00 F5 4D 00 00 01 07 FE
004EE100	4D 00 00 00 28 0A 4E 00 00 00 74 32 4E 00 00 03

Registers (FPU)
EAX 00000000
ECX 0012FFB4
EDX 7C90EB94 ntdll.
EBX 7FFDA000
ESP 0012FFC0
EBP 0012FFF0
ESI 00000000
EDI 00000000
EIP 0040103B ollydbg
C 0 ES 0023
P 1 CS 001B
A 0 SS 0023
Z 1 DS 0023
S 0 FS 003B 32bit
T 0 GS 0000 NULL

Address	Value
0012FFC0	00000000
0012FFC4	7C816D4F
0012FFC8	00000000
0012FFCC	00000000
0012FFD0	7FFDA000
0012FFD4	8054A6ED
0012FFD8	0012FFC8
0012FFDC	892FFAA8
0012FFE0	FFFFFFFF

L Log data

Address	Message
773D0000	Module C:\WINDOWS\WinSxS\x86_Microsoft.Windows.Common-Controls_6595b641
00401000	Entry point of main module
0040103B	INT3: EAX = 0 EBX = 7FFDA000 (2147328000.) [DWO]
0040103B	Breakpoint

OllyDbg
Debugger

Debugging Benefits



- Sometimes easier to just see what code does
- A workaround to packing routines:
 - just let the code unpack itself and debug as normal
 - ... unless the code tests for a debugging/VM environment!
- Most debuggers have in-built disassemblers anyway
- Can always combine static and dynamic analysis

Debugging Challenges



- We are now executing potentially malicious code
 - use an isolated virtual machine
- Anti-Debugging
 - detect debugger and [exit | crash | modify behavior]
 - IsDebuggerPresent(), INT3 scanning, timing, VM-detection, pop ss trick, etc., etc., etc.
 - Anti-Anti-Debugging can be tedious



- angr: a Python framework for binary analysis
- Based on Valgrind VEX intermediate representation (IR)
 - In Valgrind, used to transform machine code to IR, instrument the IR, then compile back to machine code.
- Uses Symbolic Execution to simulate the effects of statements, expressions, operations, etc. in the IR code.
- Symbolic Execution allows us to learn constraints of the program without actually executing!

Symbolic Execution



```
x = int(input())
if x >= 10:
    if x < 100:
        print "You win!"
    else:
        print "You lose!"
else:
    print "You lose!"
```

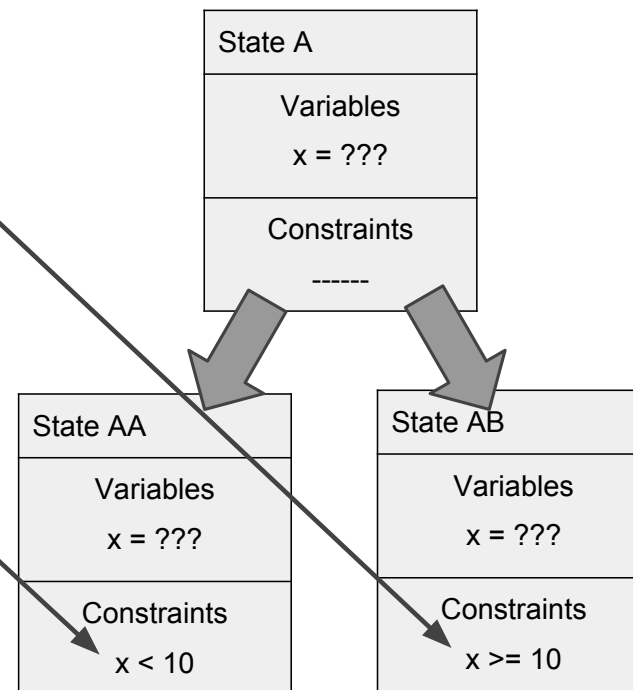
State A
Variables x = ???
Constraints -----

[Shellphish Lab, UCSB]

Symbolic Execution



```
x = int(input())
if x >= 10:
    if x < 100:
        print "You win!"
    else:
        print "You lose!"
else:
    print "You lose!"
```



[Shellphish Lab, UCSB]

Symbolic Execution



```
x = int(input())
if x >= 10:
    if x < 100:
        print "You win!"
    else:
        print "You lose!"
else:
    print "You lose!"
```

State AA
Variables x = ???
Constraints x < 10

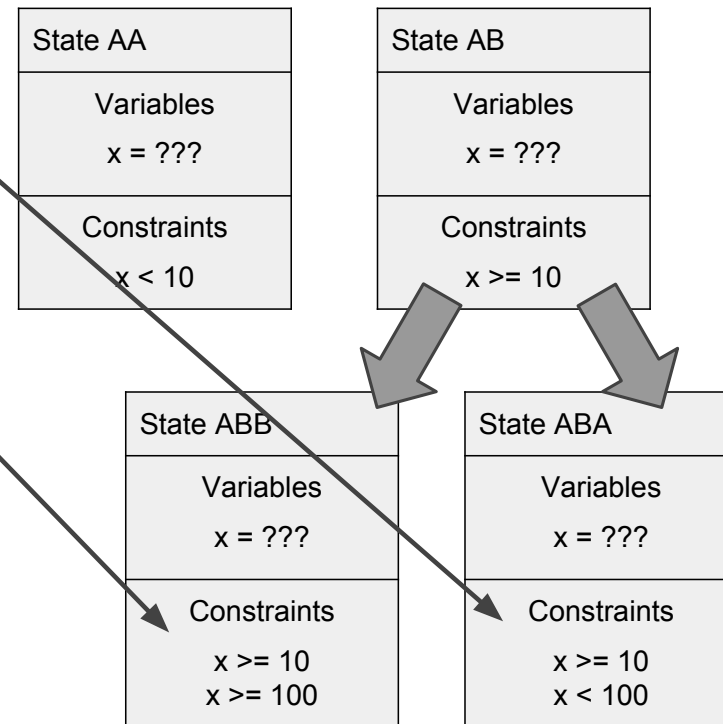
State AB
Variables x = ???
Constraints x >= 10

[Shellphish Lab, UCSB]

Symbolic Execution



```
x = int(input())
if x >= 10:
    if x < 100:
        print "You win!"
    else:
        print "You lose!"
else:
    print "You lose!"
```

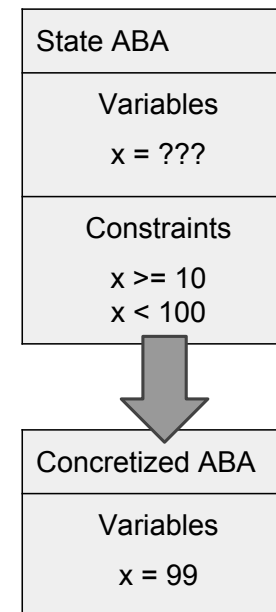
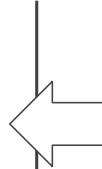


[Shellphish Lab, UCSB]

Symbolic Execution



```
x = int(input())
if x >= 10:
    if x < 100:
        print "You win!"
    else:
        print "You lose!"
else:
    print "You lose!"
```



[Shellphish Lab, UCSB]



- Advantages
 - Semantically Aware
 - Targetable
- Disadvantages
 - path explosion
 - constraint solving

angr analyses

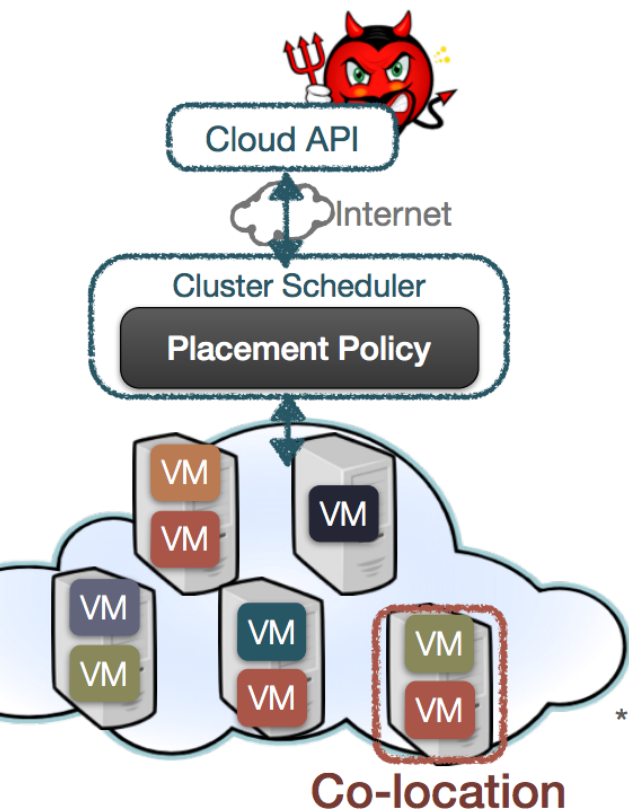


Analysis	Description
CFGFast/CFGAccurate	Control flow graph recovery
Disassembly	Linear disassembly rendering routine
DDG	Data dependence analysis
VFG	Value-flow graph recovery, performs
BackwardSlicing	Backward program slicing based on
BoyScout*	Determines architecture of binary blobs
GirlScout*	Determines base addresses of binary



- What else can we reverse engineer?
 - A policy?
 - A machine learning model??
 - Other stuff???

Ex: Cloud Co-Residency Attacks



If a truly random placement policy was used...

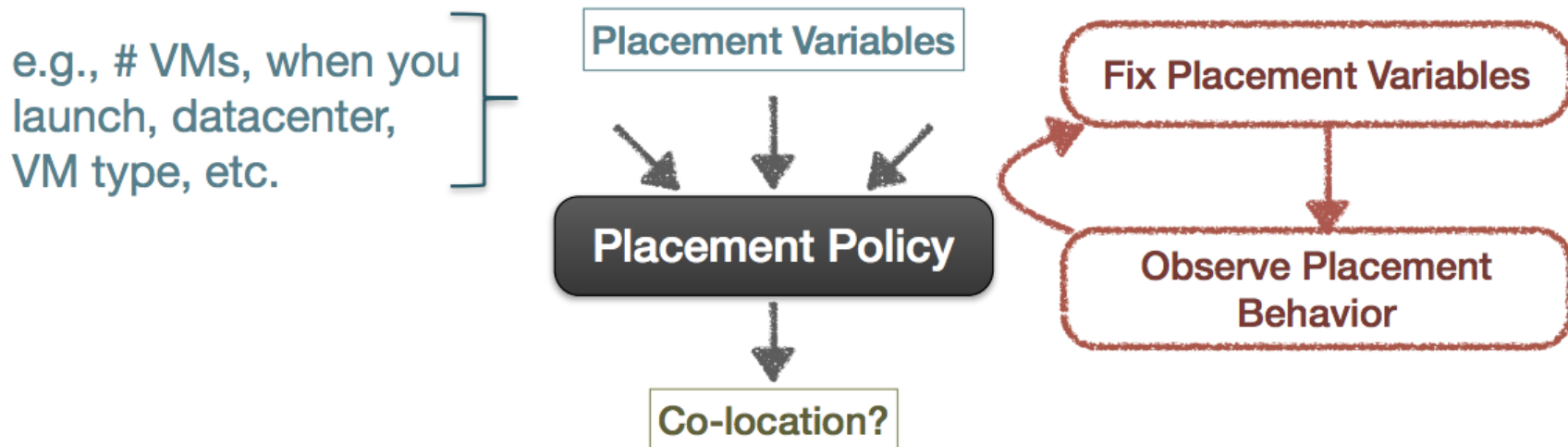
- $N = 50,000$ machines
- v victim VMs and a attacker VMs
- Probability of Collision:

$$P_c = 1 - \left(1 - \frac{v}{N}\right)^a$$

v	$a = \ln(1 - P_c) / \ln(1 - v/N); P_c = 0.5$
10	3466
20	1733
30	1155

[Varadarajan et al., Security'15]

Reverse Engineering Amazon's VM Placement Policy

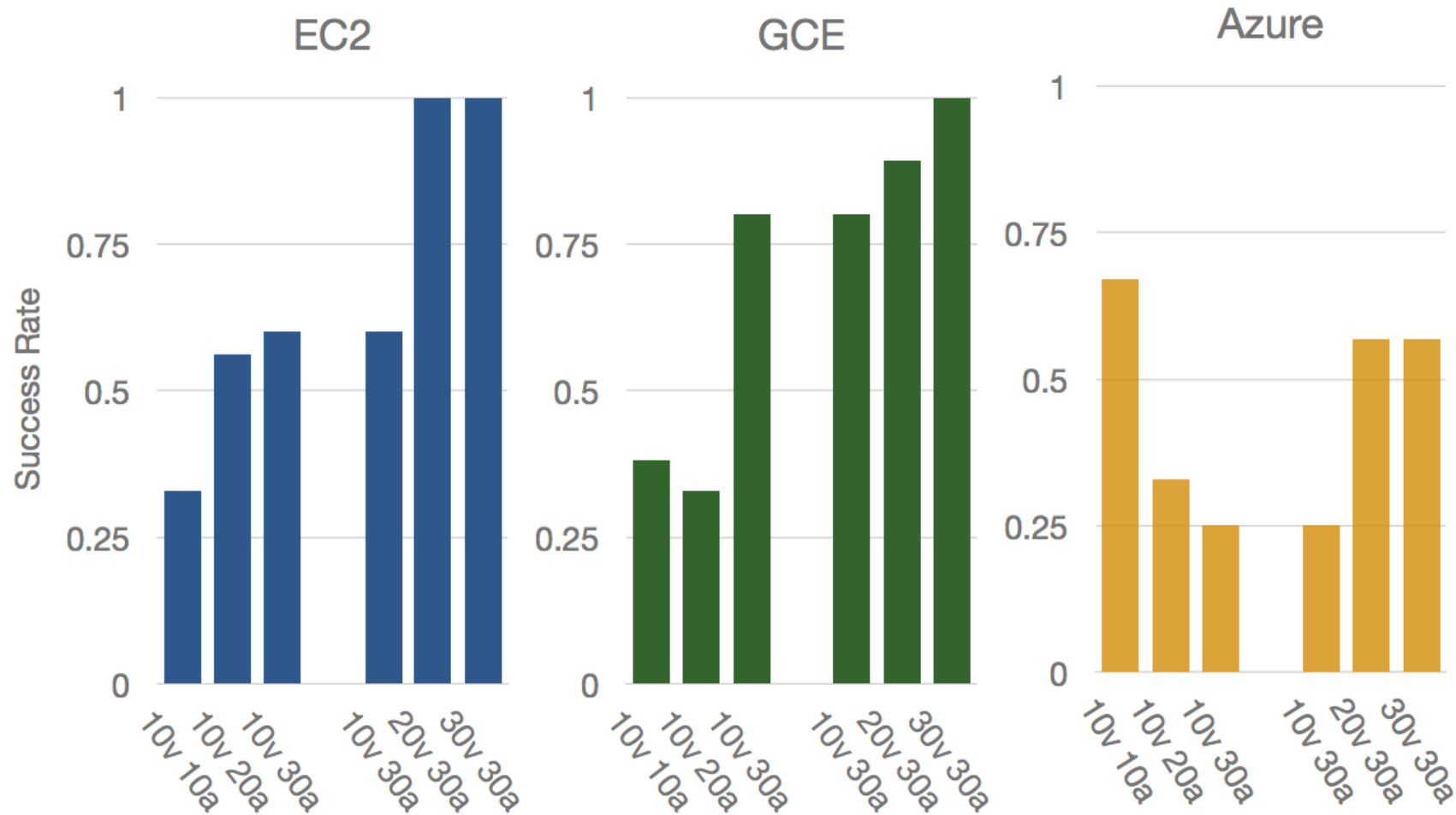


- 6 placement variables: # victim & attacker VMs, delay b/w launches, time of day, day of week, datacenter, cloud provider Small instance type
- 9 samples per strategy with 3 runs per time of day and 2 days of week (weekday/weekend).



[Varadarajan et al., Security'15]

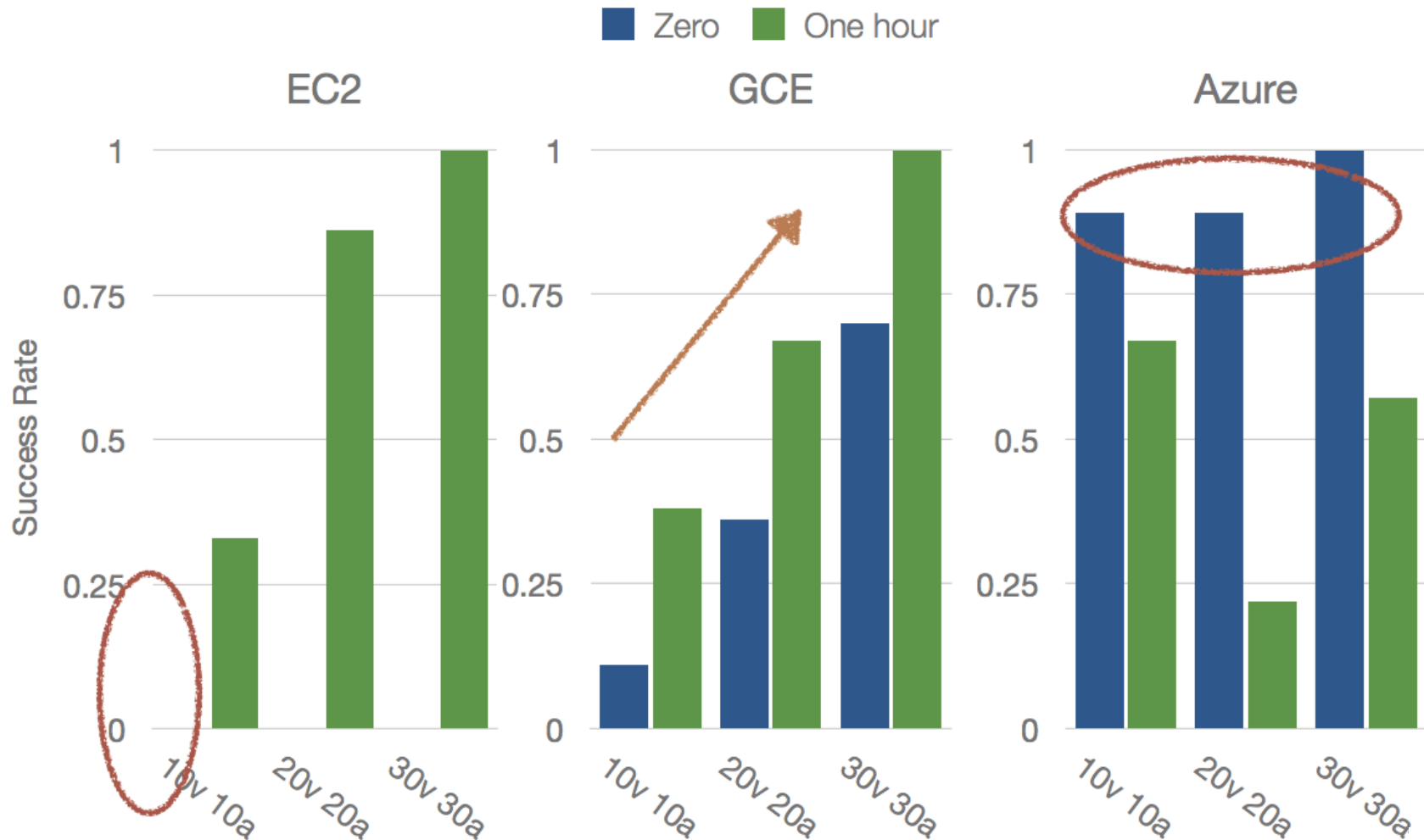
Reverse Engineering Amazon's VM Placement Policy



***Co-location is possible with as low as 10 VMs and
always achieves co-location with 30 VMs***

[Varadarajan et al., Security'15]

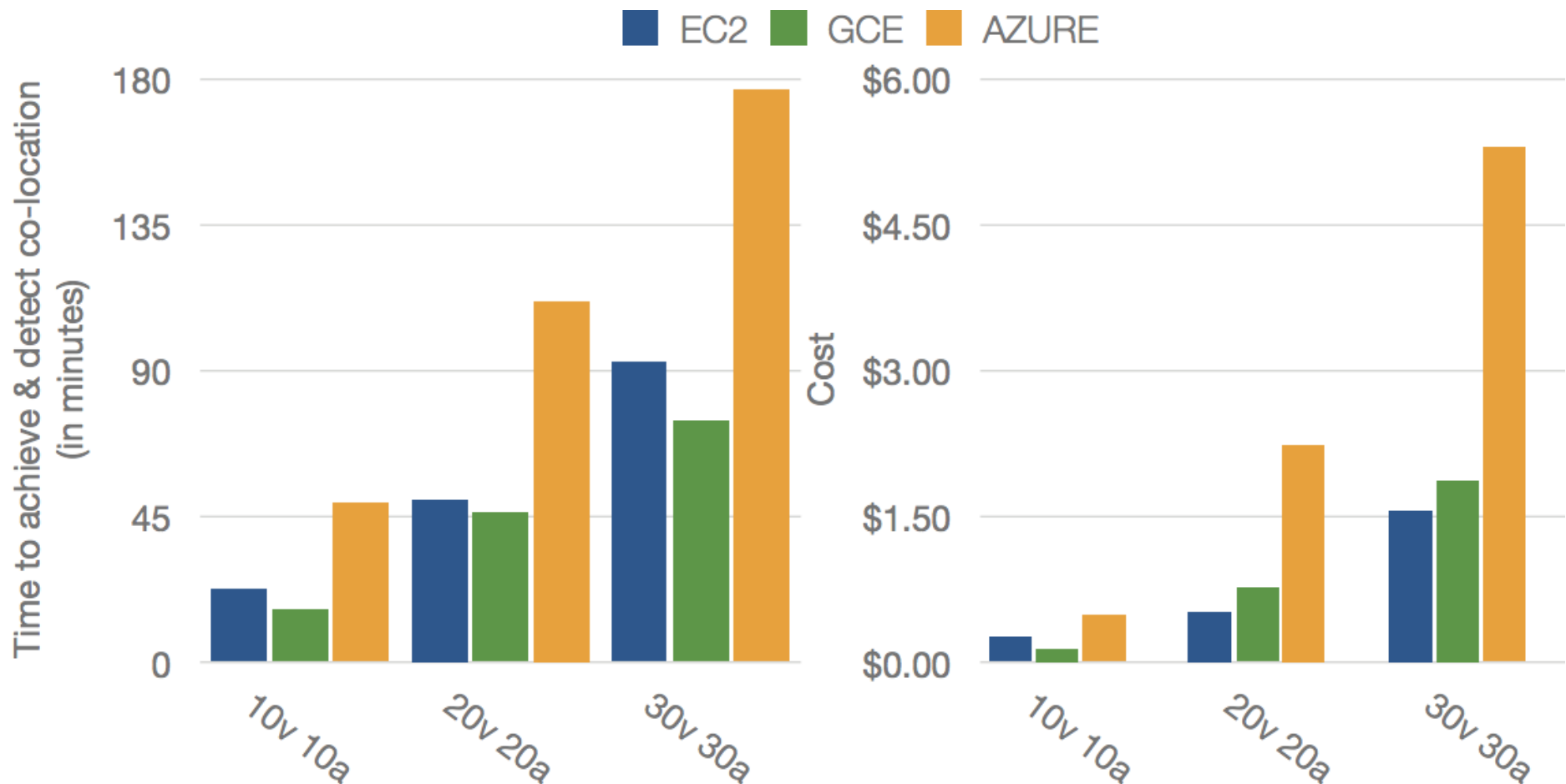
Reverse Engineering Amazon's VM Placement Policy



***Different clouds have wildly different
temporal placement strategies***

[Varadarajan et al., Security'15]

Reverse Engineering Amazon's VM Placement Policy



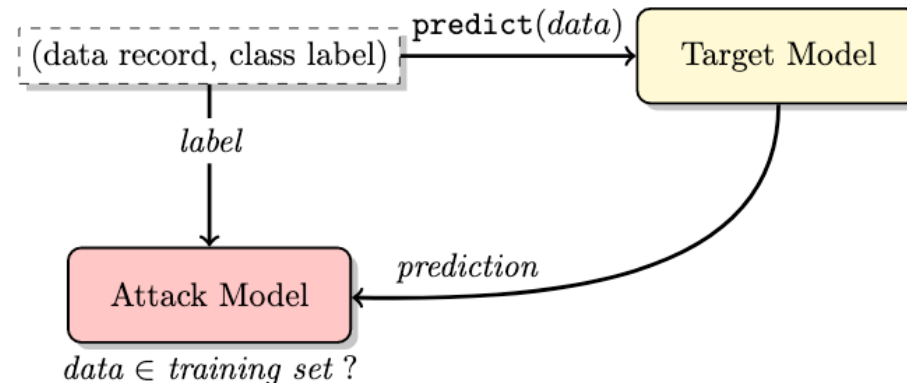
Successful co-location as affordable as 14 cents.

[Varadarajan et al., Security'15]

Ex: Membership Inference Attacks



- Attacker queries the target model with a data record and obtains the model's prediction on that record.
- The prediction is a vector of probabilities, one per class, that the record belongs to a certain class.
- This prediction vector, along with the label of the target record, is passed to the attack model, which infers whether the record was in or out of the target model's training dataset.



[Shokri et al., Oakland'17]

Ex: Model Inversion Attacks



- Also exploits confidence values revealed alongside predictions.
- Perturb features of some input image (e.g., transform pixel values in facial recognition), record change in confidence values.
- Use gradient descent to identify local minima (in this case, highest obtained confidence values).
- Can be used to launch a reconstruction attack on the original training data



Training Data



Inverted Model

[Fredrikson et al., CCS'15]

ICES Evaluation



Please write “**A. Bates CS 461 AL4**” at the top of your ICES form.

You can doodle it in just under “INSTRUCTOR AND COURSE
EVALUATION SYSTEM.”