

# Very Short Intermittent DDoS Attacks in an Unsaturated System

**Huasong Shan**<sup>+</sup>, Qingyang Wang<sup>+</sup>, Qiben Yan<sup>\*</sup>

<sup>+</sup>*Louisiana State University*

<sup>\*</sup> University of Nebraska-Lincoln

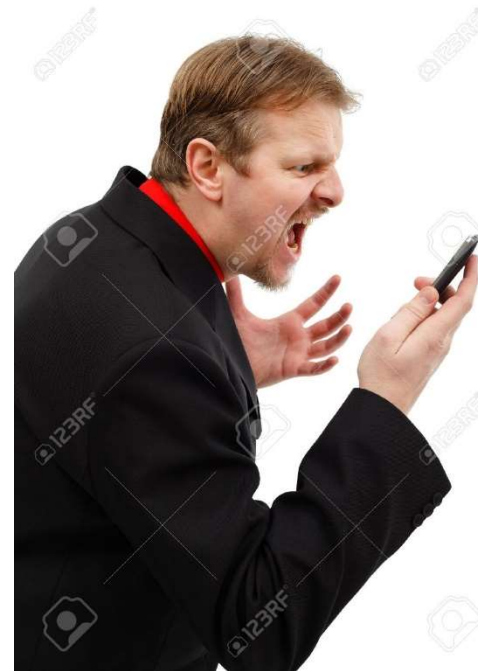


LSU

N

UNIVERSITY OF  
Nebraska<sup>®</sup>  
Lincoln

# Responsiveness and Patience



- Systems that respond to user actions quickly (within **100ms**) feel more fluid and natural

–[Card et al. *SIGCHI Conference on Human factors in computing systems* '91]

# Responsiveness of Web Applications

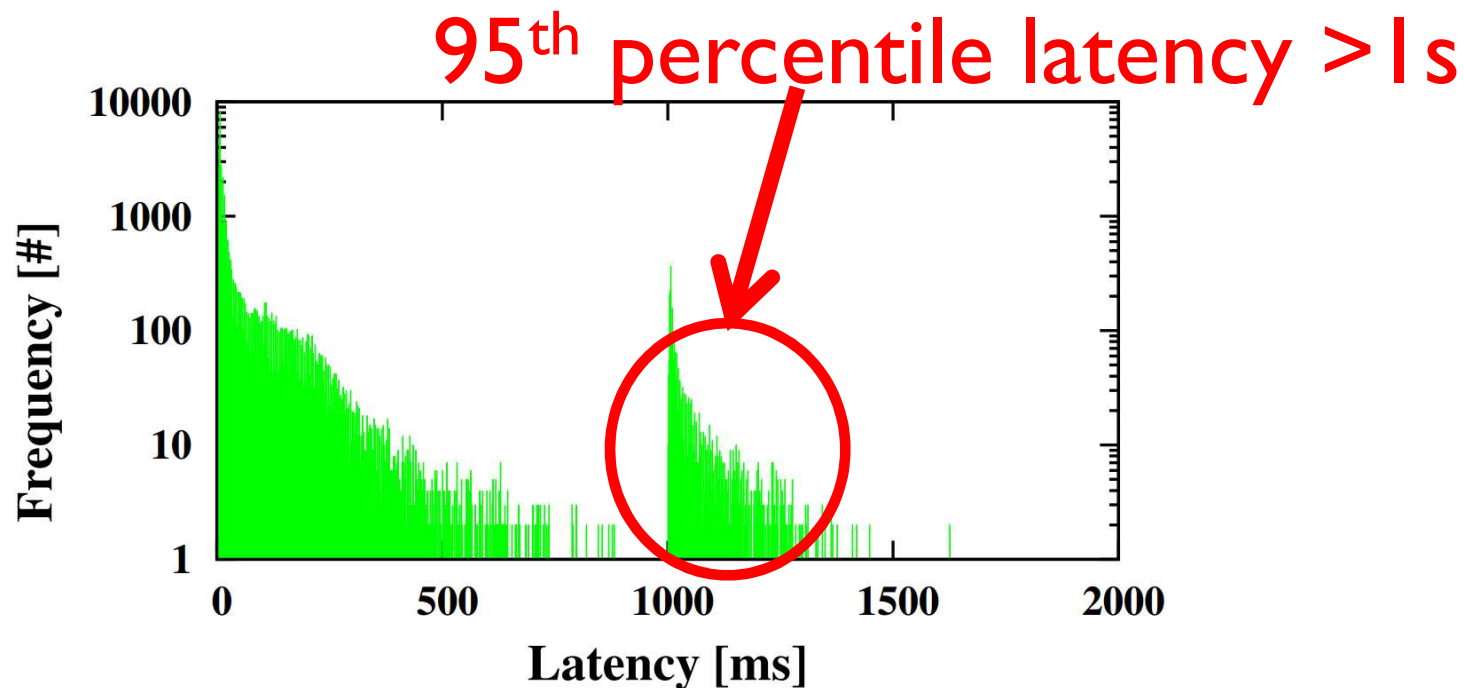


- Google: 99%-ile latency in tens of milliseconds  
—[Lo et al. *ISCA'15*]
- Amazon: 100ms latency increase -> 1% sales decrease  
—[Kohavi et al. *Computer'07*]
- Websites Service Level Objective(SLO):  
tail latency  
—[Beset. *Operating Systems Review'12*]

# Very Short Intermittent DDoS Attacks (VSI-DDoS)



- Hurt **the responsiveness** of web services (Long tail latency problem)



# Very Short Bottleneck (VSB)



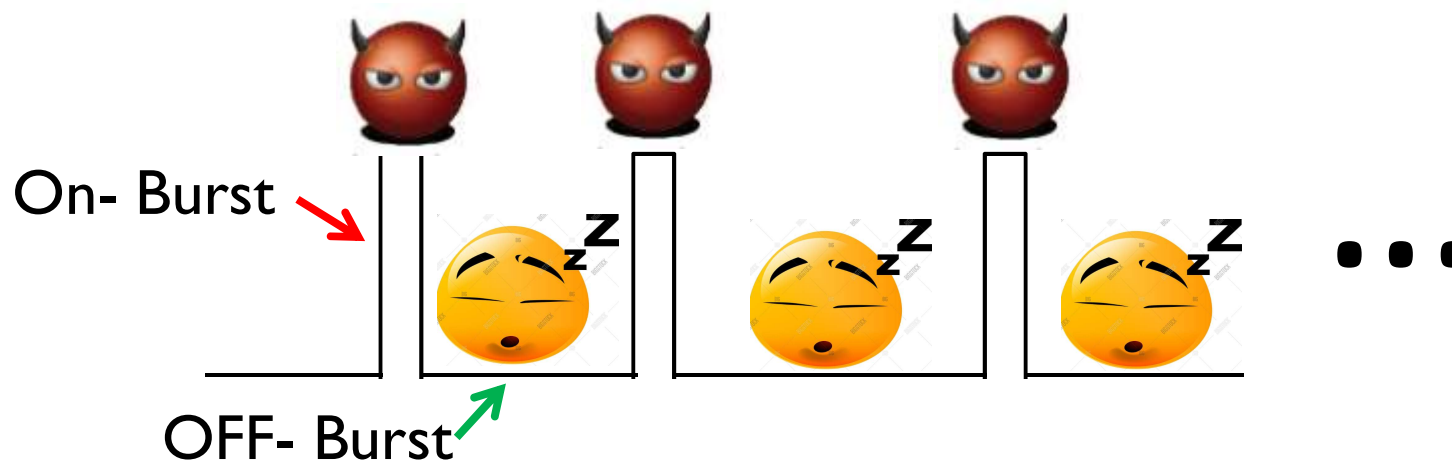
- ❑ A performance vulnerability in n-tier systems
  - ▶ **Very short bottlenecks** (tens or hundreds of milli-seconds)  
—[Wang et al. ICDCS'17, TRIOS'14]
  - ▶ Very long response time (seconds)
- ❑ VSI-DDoS Attack Approach
  - ▶ Create **very short bottleneck (VSBs)**, causing long tail latency

# VSI-DDoS Attacks Scenario



□ ON: a burst of HTTP requests

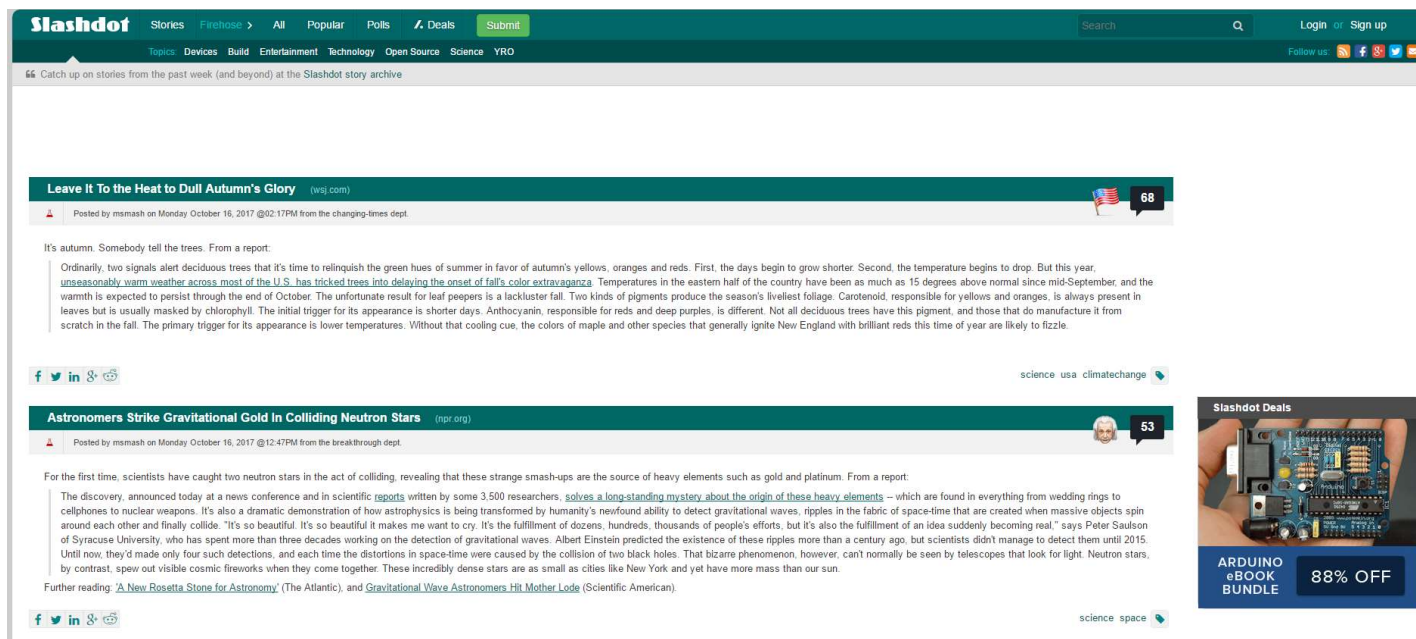
□ OFF: null



# Benchmark Application

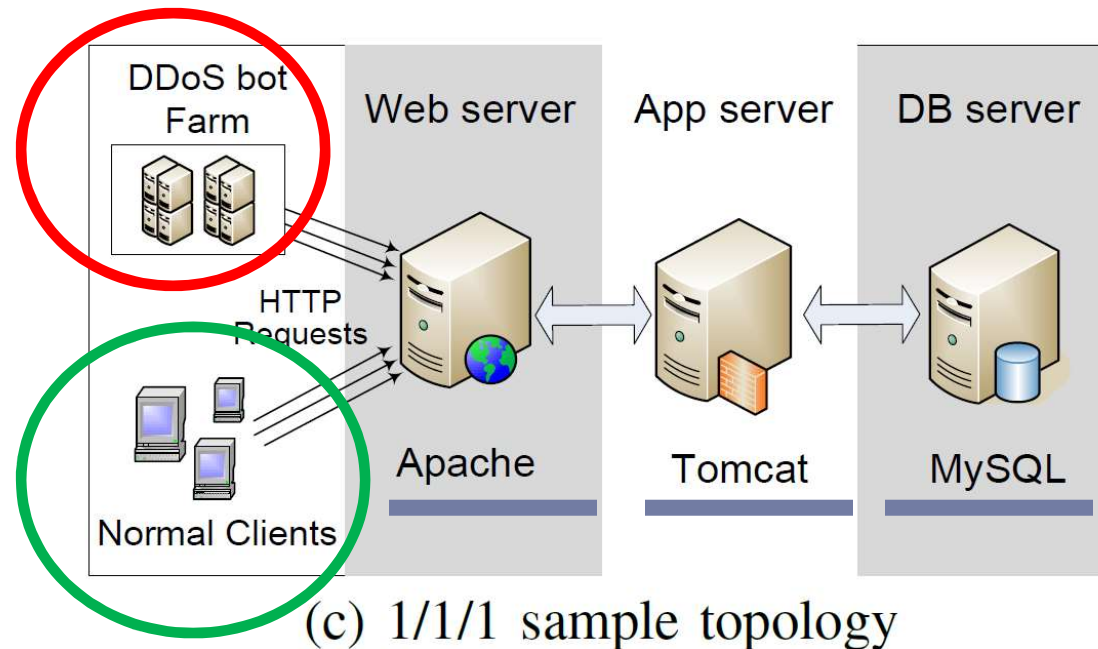
## ▣ RUBBoS benchmark

- Bulletin board system like Slashdot ([www.slashdot.org](http://www.slashdot.org))
- N-tier architecture
- 24 web interactions



# Experimental Sample Topology

## □ Attackers' workloads: Apache Bench

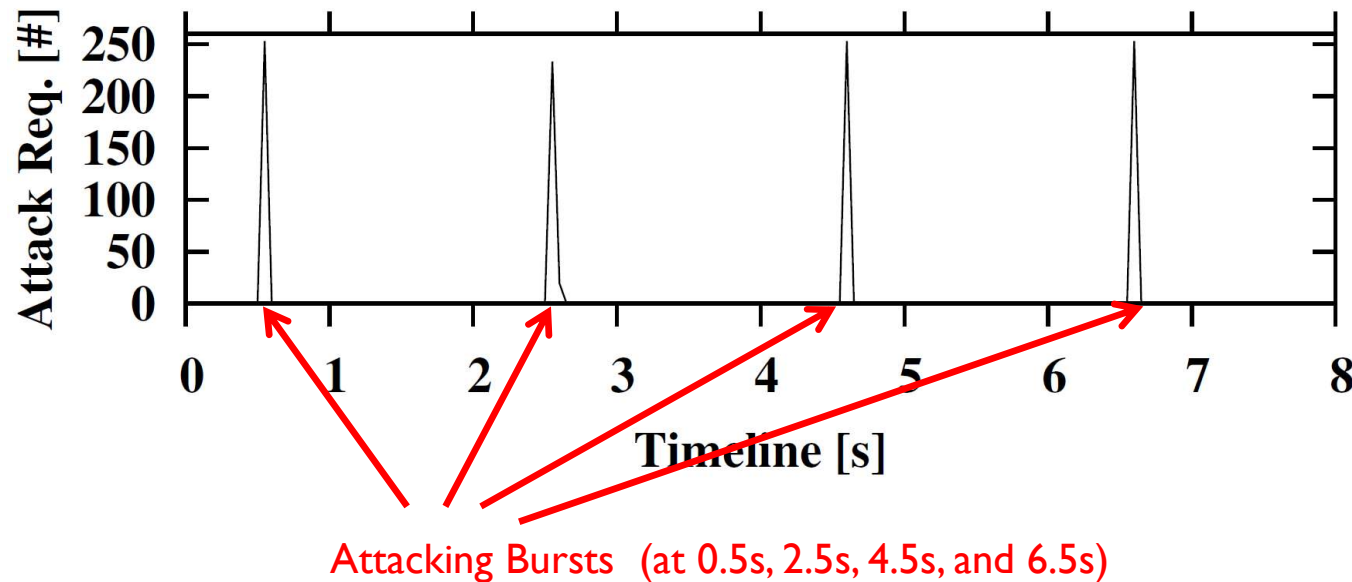


## □ Normal clients' workloads: RUBBoS Clients

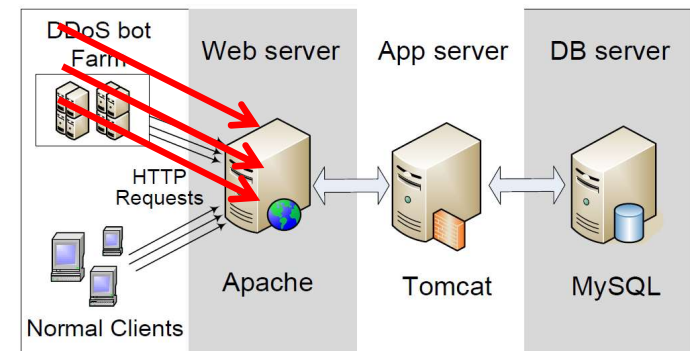
- Concurrent users (e.g., 3000 )
- An average 7-second think time



# (1) A Concrete Attack Scenario

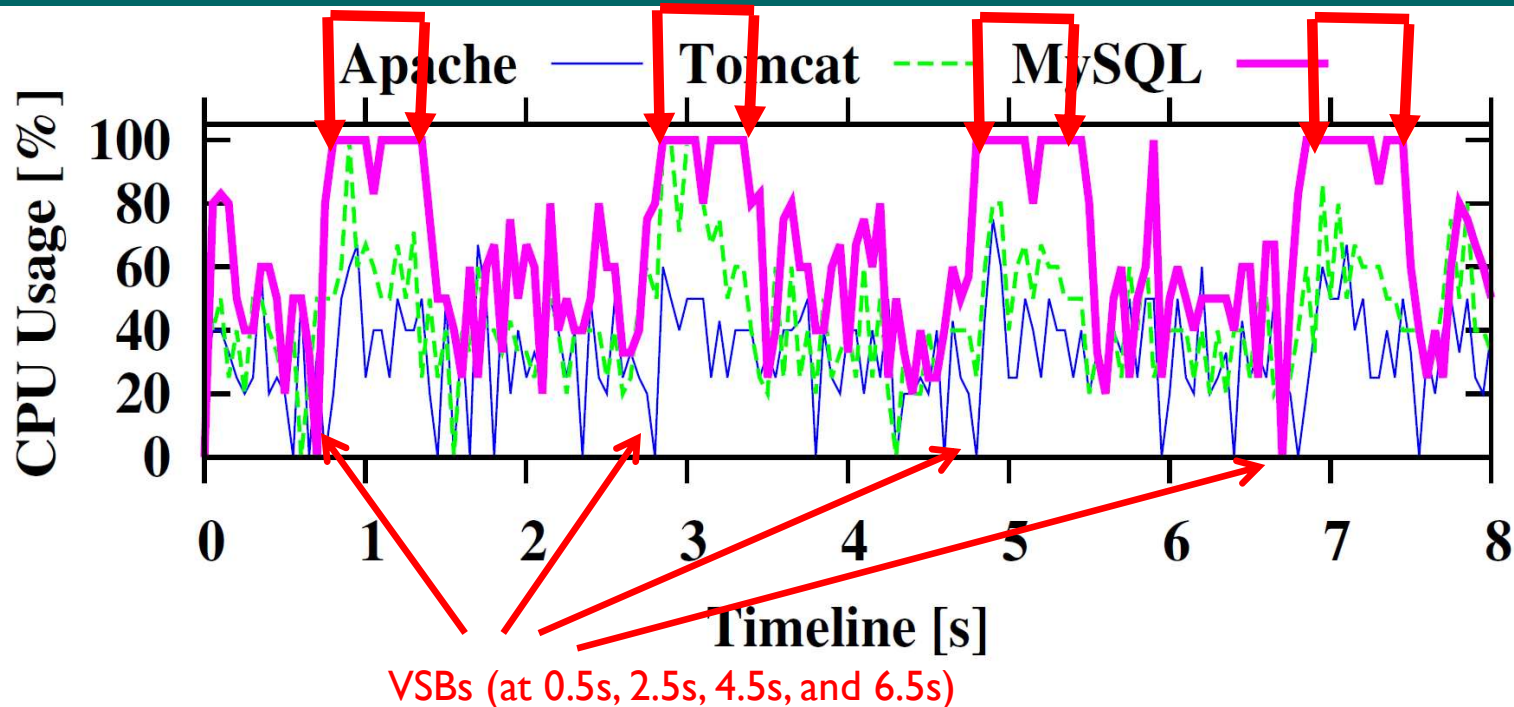


- **Attackers:** send a burst of 250 HTTP requests with 50ms, repeat in every 2 seconds

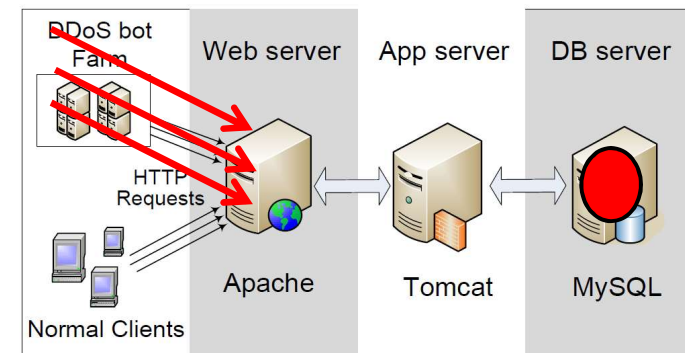


(c) 1/1/1 sample topology

## (2) Milli-Bottlenecks in MySQL

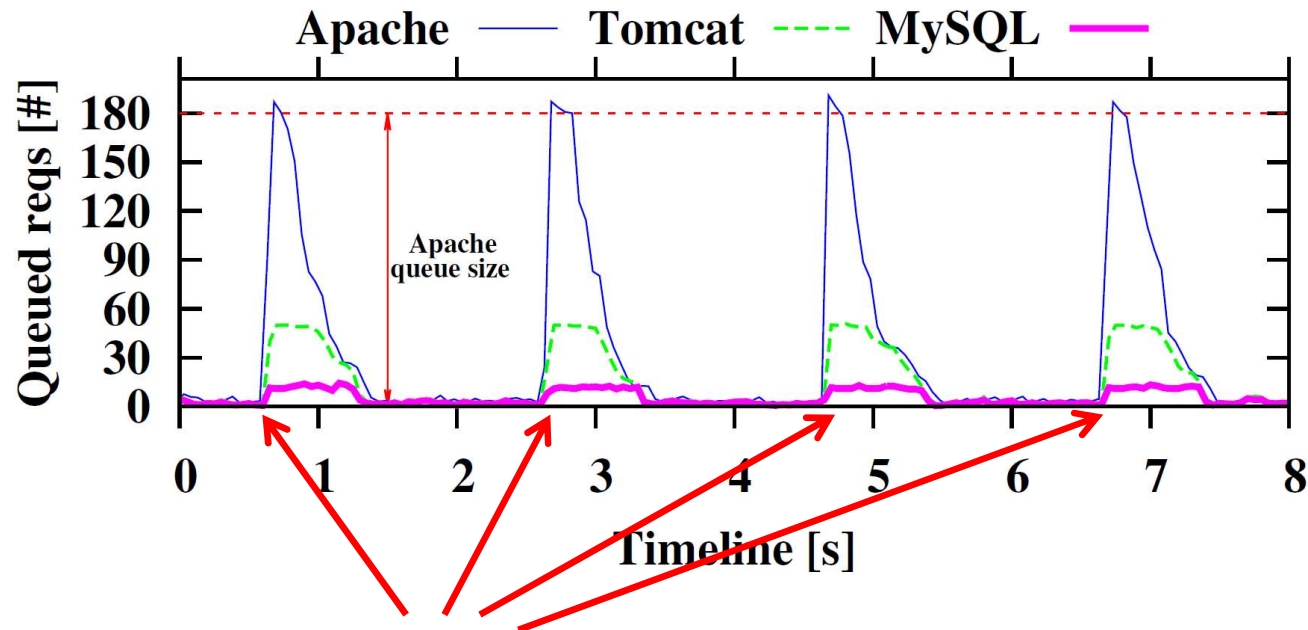


- MySQL: very short bottlenecks in MySQL (less than 500ms)

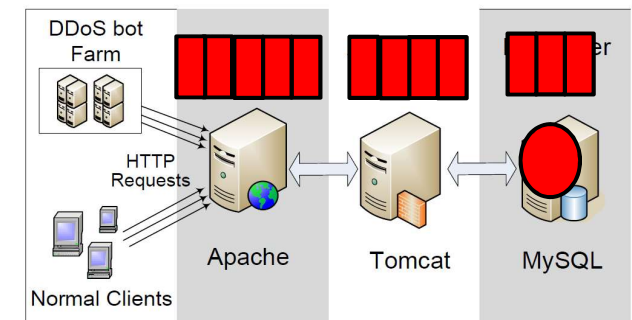


(c) 1/1/1 sample topology

### (3) Queue Overflow Propagation

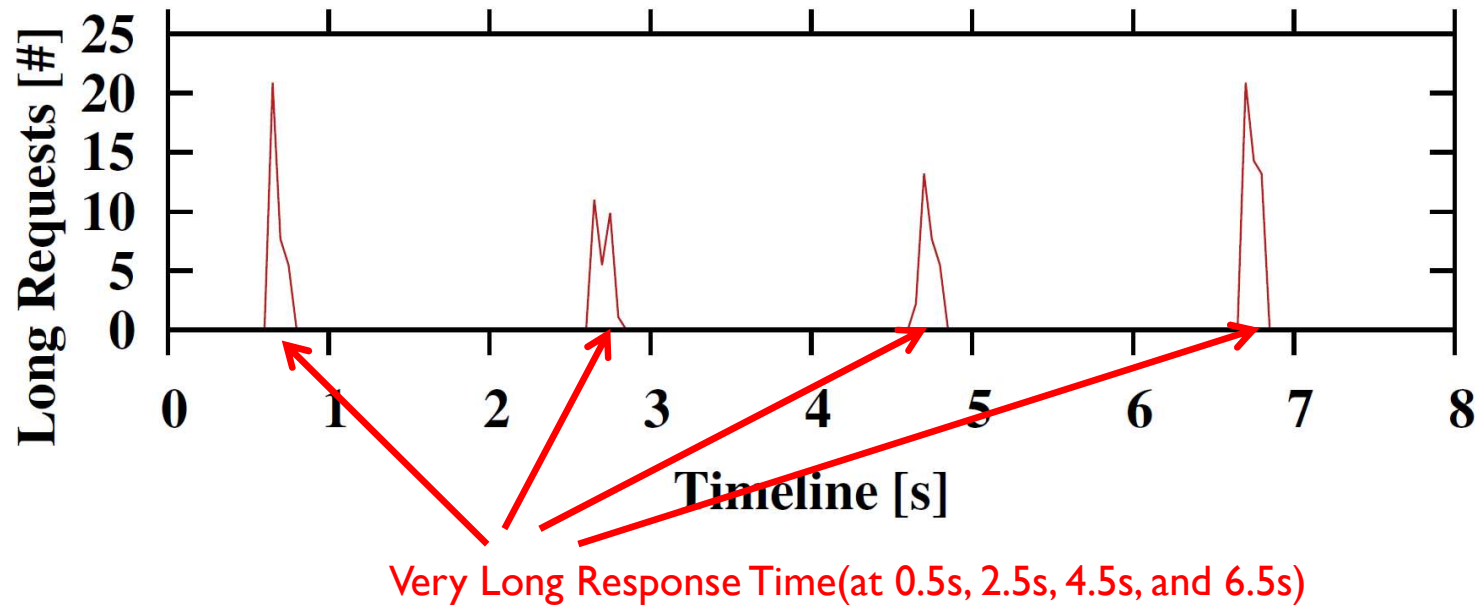


- 3-tier System: queue overflow from MySQL to Tomcat, Apache



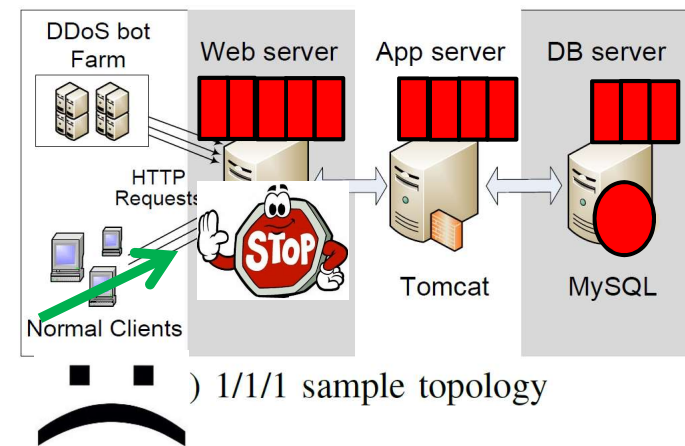
(c) 1/1/1 sample topology

## (4) Very Long Response Time



### □ Legitimate Users:

- Queue full in Apache
- Drop new Req. by Apache
- TCP retransmission (min: 1s)
- Long response time(> 1s)



# Damage of VSI-DDoS Attacks

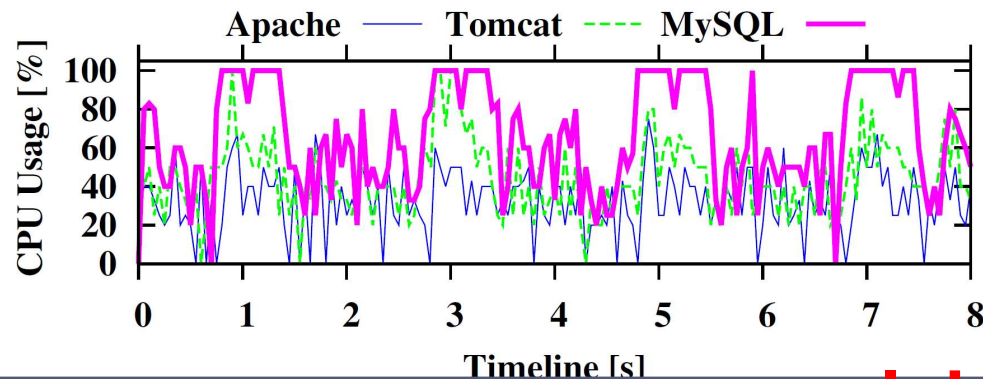
- 95%-ile latency of the target system  $> 1$  second



VSI-DDoS Attacks is **harmful** to tail-latency sensitive web applications (e.g., e-payment, web search, online gaming)

Response time [s]

# Stealthiness of VSI-DDoS Attacks



VSI-DDoS Attacks is **stealthy** to defense tools and Cloud scaling

## ❑ Defense Tools:

Cisco Adaptive Security Appliance (1s)  
Snort (1s)

## ❑ Cloud Scaling:

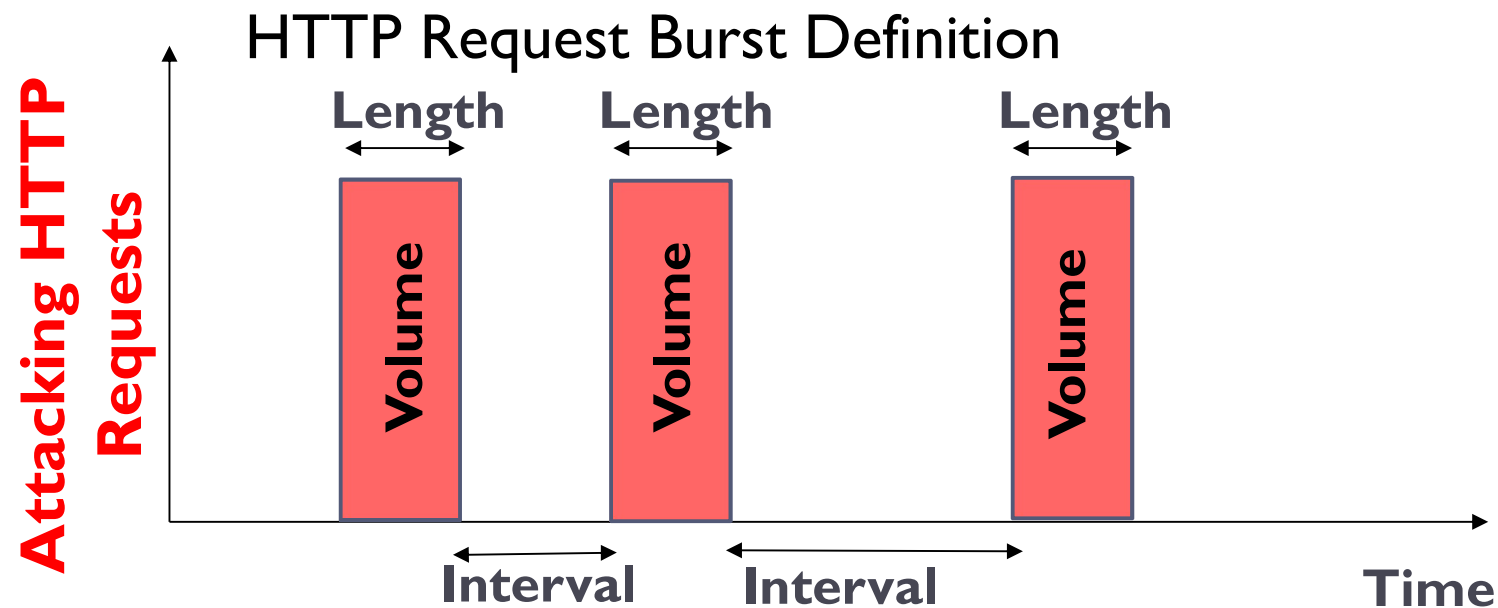
Amazon CloudWatch (1min)  
Microsoft Azure Application Insights (1min)

# Challenges to Launch Effective VSI-DDoS Attacks



❑ How to trigger Very Short Bottleneck (VSB)?

➤ HTTP Burst

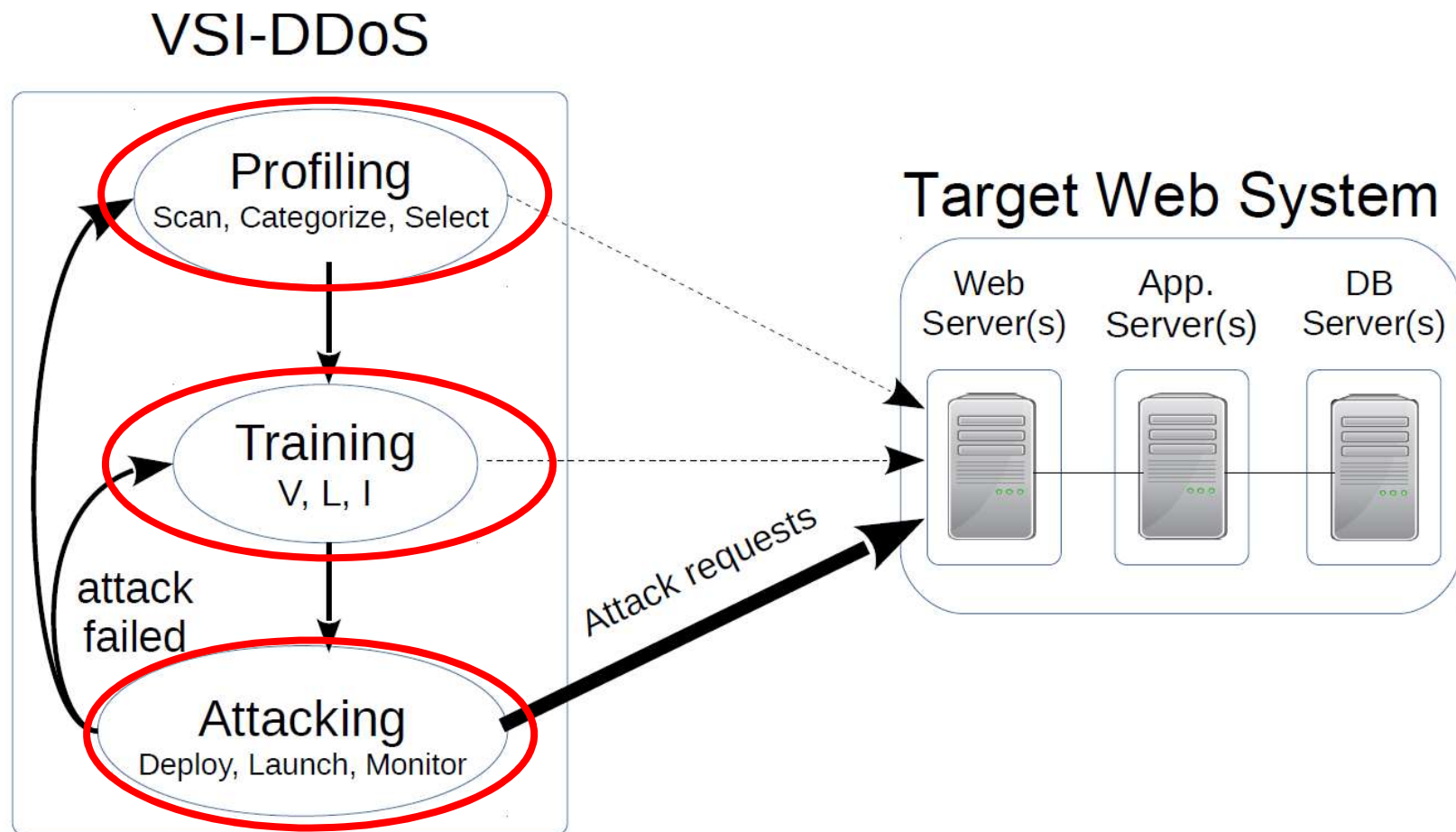


❑ How to quantify the damage of VSI-DDoS Attacks?

➤ Tail latency (percentile response time)



# VSI-DDoS Attacks Control Framework



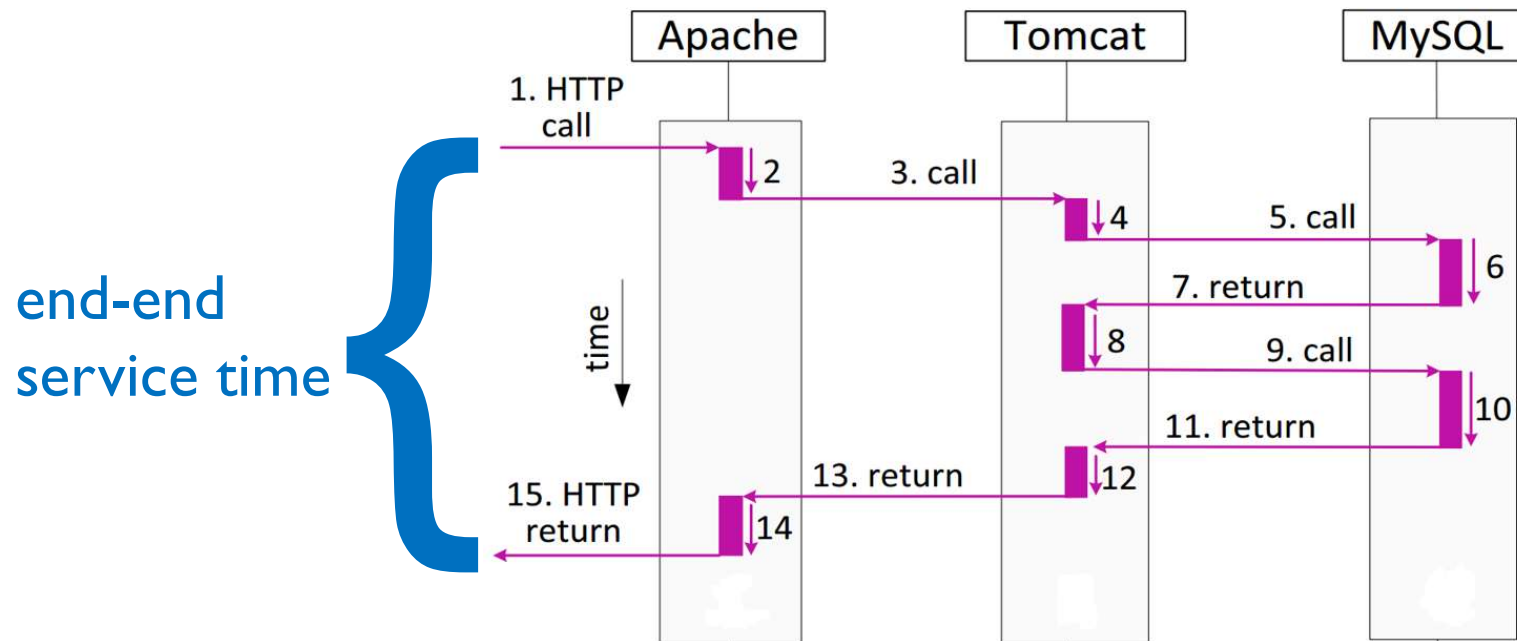


# (1)Profiling

Profiling  
Scan, Categorize, Select

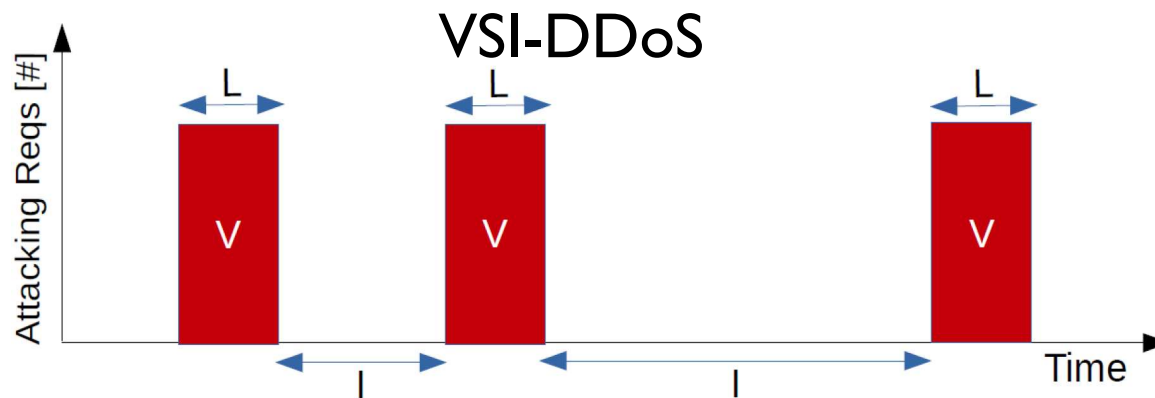
## □ Select attack requests

- Profile **end-end service time** of HTTP requests
- Req. with **long service time** as candidate attack req.



## (2) Training

Training  
V, L, I

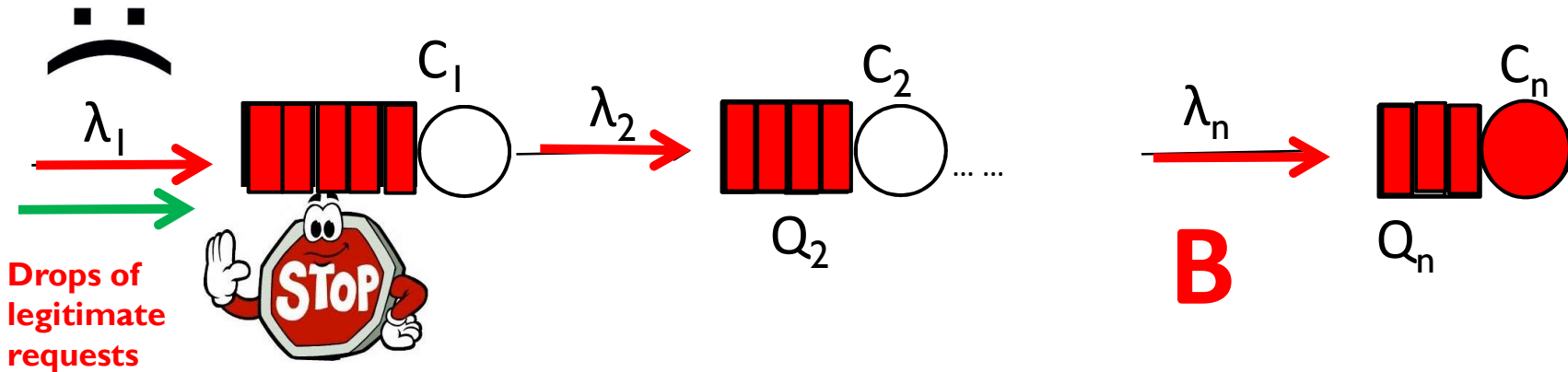


### □ Optimal attacking parameters

- Effective VSBs + Long tail latency + Moderate average utilization

# Training: optimal burst volume $V$

- Optimal  $V$  to create effective VSBs to cause long tail latency

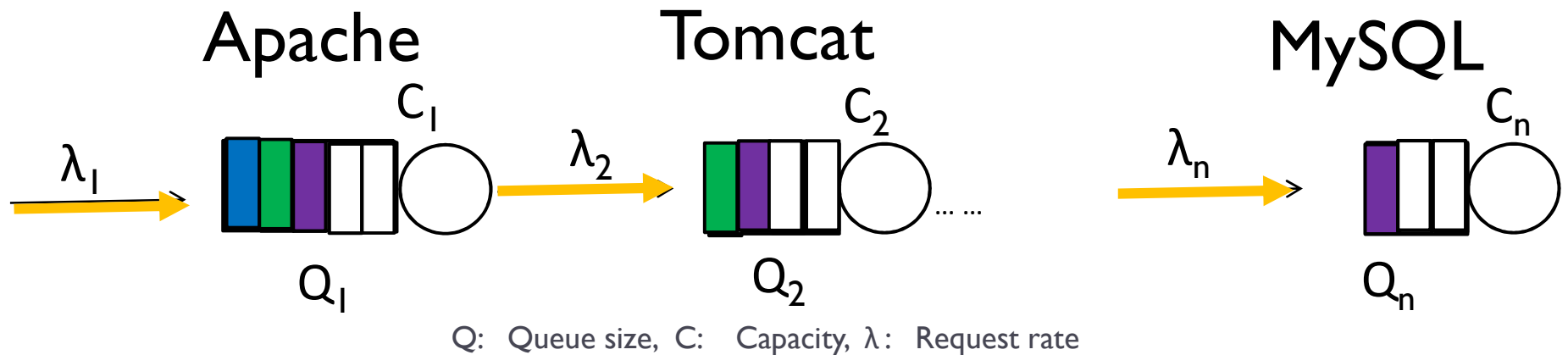


attack burst (B) -> Very short bottlenecks (milliseconds) ->  
 fill up  $n$ -th queue -> fill up all queues ->  
 drop new requests ->  
 long response time (TCP retr.Timeout: seconds)

- Approach: increase  $V$  step by step until occurrence of long response time

# Training: optimal burst length $L$

- Optimal  $L$  to occupy queue as long as possible

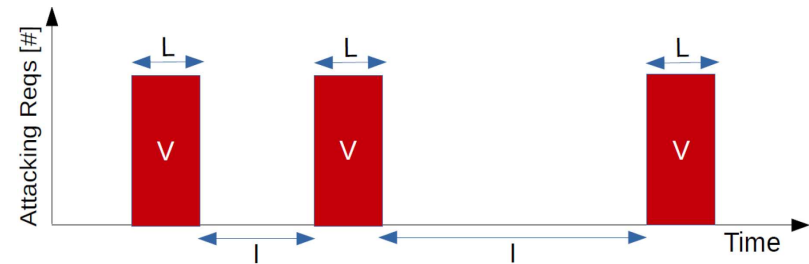


- Approach:  $L$  = end-end service time of attacking requests

# Training: optimal burst Interval $I$

## □ Optimal $I$

- too long, attack fail
- too short, flooding DDoS

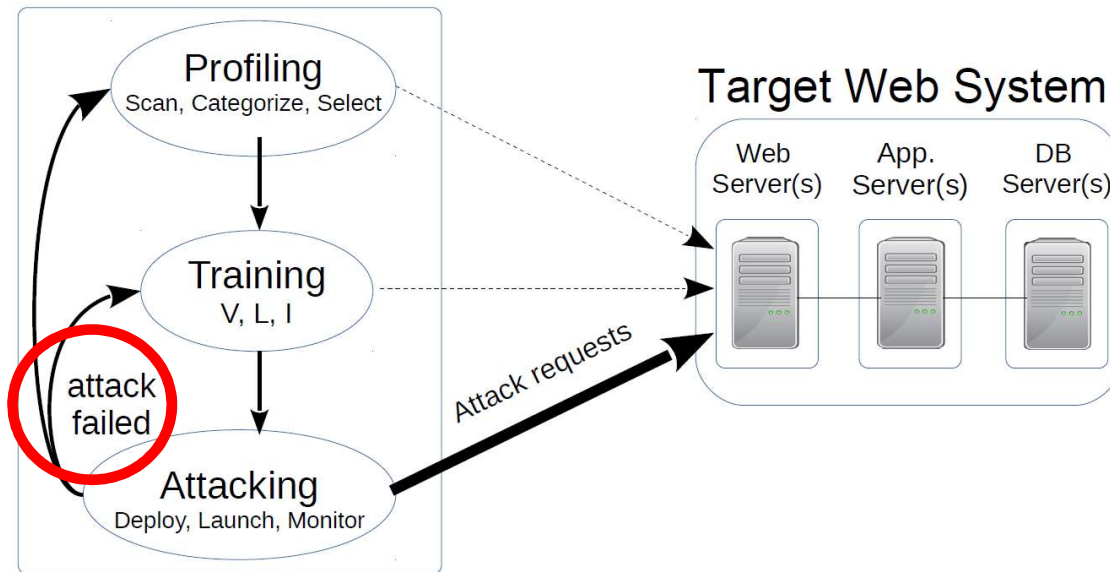


- Approach: Increase/Reduce  $I$  step by step until achieving the attack goal

## (3)Attacking

Attacking  
Deploy, Launch, Monitor

VSI-DDoS



❑ Redo if attack failed

➤ Variation of background workload or system state

❑ Further detailed control refers to our ACM CCS'17 Paper  
“Tail Attacks on Web Applications”

# Possible Countermeasures

---

## ❑ Fine-Grained VSBs Detection

- High overhead

## ❑ Threshold-Based Monitoring and Detection

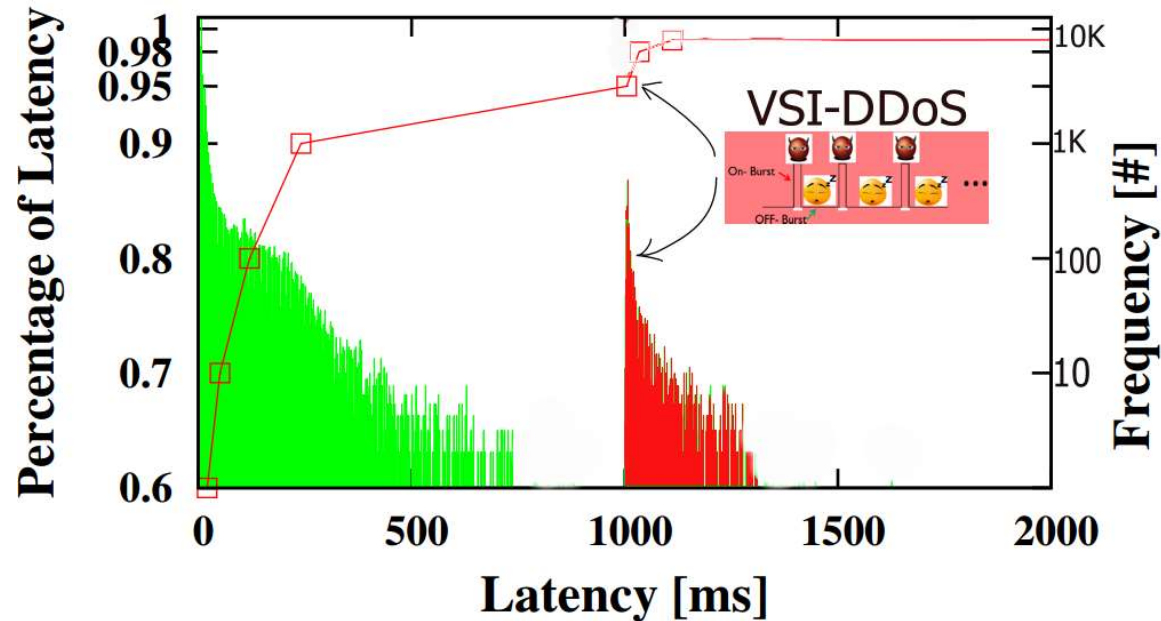
- Too high threshold, can not detect anomalies
- Too low threshold, high false positive error

## ❑ User Behavior Model Validation

- Distinguish humans and bots
- Bots learn from those model

# Conclusion

- A low-volume application layer DDoS attacks: VSI-DDoS Attacks
- Very Short Bottlenecks + Long-tail latency + Moderate average utilization







Thank You.  
Any Questions?



**LSU**



UNIVERSITY OF  
**Nebraska**  
Lincoln®

# Backup

---



## Extras



# Coordinate and Synchronize Bots

## □ Centralized

–[Guirguis et al. ICNP'04, Zhang et al. NDSS'07, Ramamurthy et al. ATC'08]

## □ Decentralized

–[Ke et al. AsiaCCS'16]

M. Guirguis, A. Bestavros, and I. Matta. Exploiting the transients of adaptation for DoS attacks on internet resources. In IEEE ICNP, 2004

P. Ramamurthy, V. Sekar, A. Akella, B. Krishnamurthy, and A. Shaikh. Remote profiling of resource constraints of web servers using mini-flash crowds. In USENIX ATC, 2008

Y. Zhang, Z. M. Mao, and J. Wang. Low-rate tcp-targeted DoS attack disrupts internet routing. In NDSS, 2007

Y.-M. Ke, C.-W. Chen, H.-C. Hsiao, A. Perrig, and V. Sekar. Cicadas: Congesting the internet with coordinated and decentralized pulsating attacks. In AsiaCCS, 2016.

# Related Works: Pulsating DDoS Attacks



## ❑ Low-rate network-layer pulsating DDoS attacks

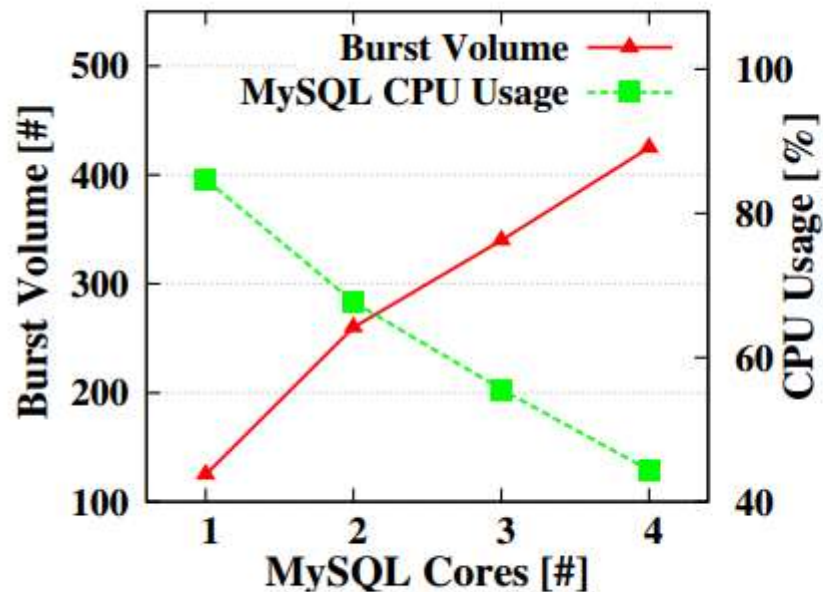
–[Kuzmanovic et al. SIGCOMM'03, Guirguis et al. ICNP'04]

- ▶ Temporarily saturate the network bandwidth in TCP layer

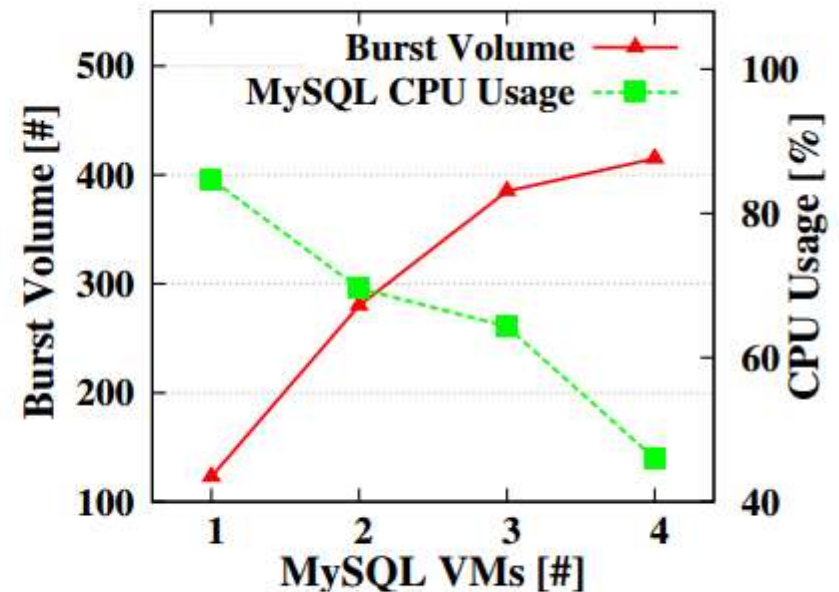
A. Kuzmanovic and E.W. Knightly. Low-rate tcp-targeted denial of service attacks: the **shrew** vs. the mice and elephants. In ACM SIGCOMM, 2003.

M. Guirguis, A. Bestavros, and I. Matta. Exploiting the transients of adaptation for **roq** attacks on internet resources. In IEEE ICNP, 2004

# VSI-DDoS Attacks under Cloud Scaling



(a) Scaling CPU cores of MySQL.



(b) Scaling Virtual Machines of MySQL.

# VSI-DDoS Attacks under Defense Tools

Table 1: Measured HTTP traffic in the cases of 95th, 98th and 99th percentile response time ( $>1s$ ) as candidate attacking goals. All of measured metrics are less than the predefined thresholds set based on system capacity when the corresponding attacking goal is achieved.

Metrics	Threshold	2000 low load				4000 high load			
		95th	98th	99th	B/L	95th	98th	99th	B/L
In. packets(#/min)	299K	158K	119K	111K	99K	224K	214K	208K	201K
Out. packets(#/min)	349K	171K	134K	127K	116K	259K	249K	241K	233K
In. speed(MB/sec)	9.32	4.68	3.96	3.62	3.11	7.08	6.76	6.45	6.23
Out. speed(MB/sec)	17.83	7.62	6.83	6.48	5.94	12.78	12.46	12.12	11.89

**In.:** HTTP Incoming, **Out.:** HTTP Outgoing, **B/L:** Baseline

# (1)Profiling

Profiling  
Scan, Categorize, Select

## ❑ Choose attack requests

### (a) Scanning the supported HTTP requests

- GET requests: crawling tools (e.g., scrapy)
- POST requests: browser-based tools (e.g., PhantomJS)

### (b) Identifying candidate attack requests

- Req. with long service time as candidate attack req.
  - ❑ Consume more resources, low attack cost

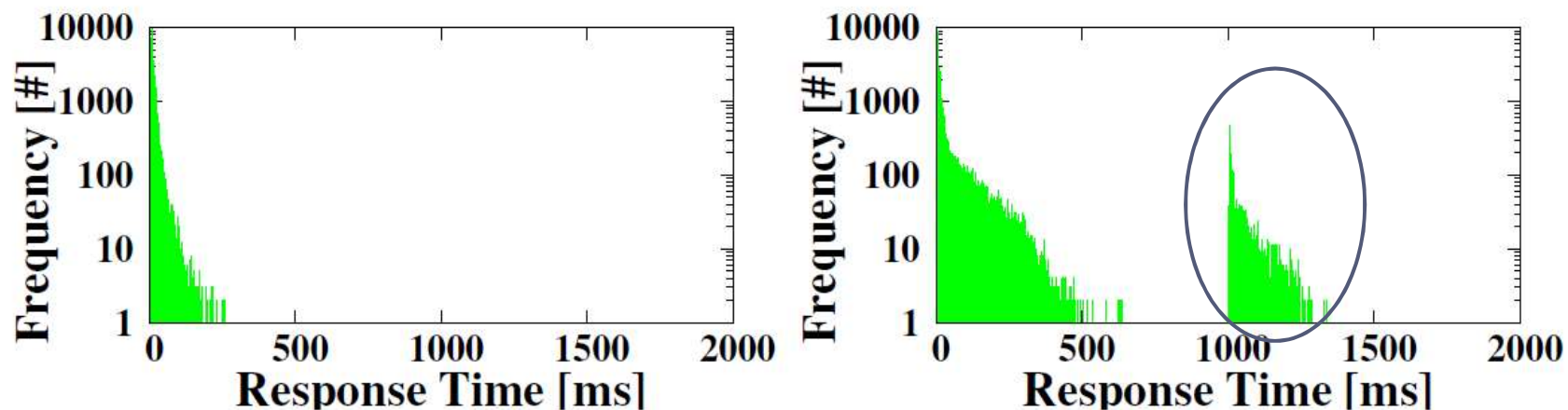
### (c) Selecting attack requests

- Reasonable request flow to cater to user behavior model



# Training: optimal burst volume $V$

- **Optimal  $V$**  to create effective VSBs to cause long tail latency

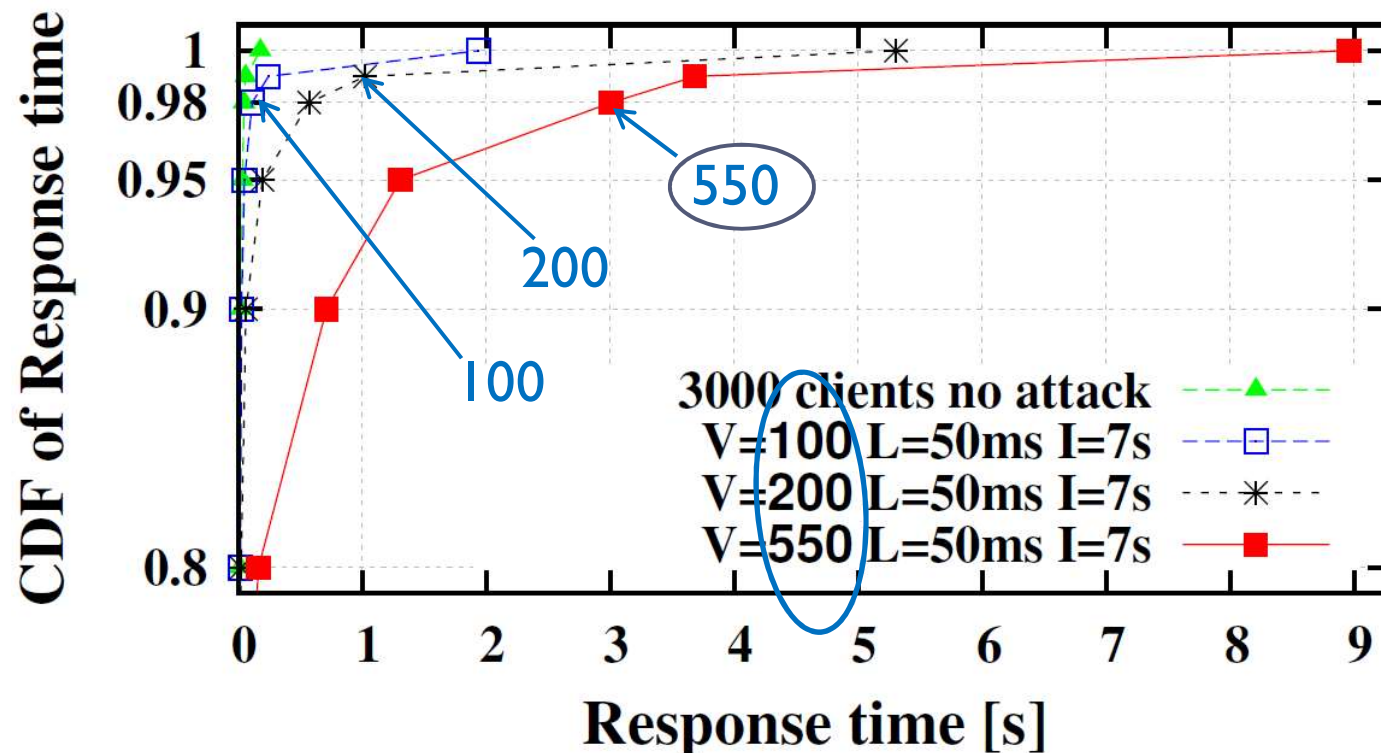


(a) 20 requests per burst case. (b) 100 requests per burst case.

- Approach: increase  $V$  step by step until occurrence of long response time

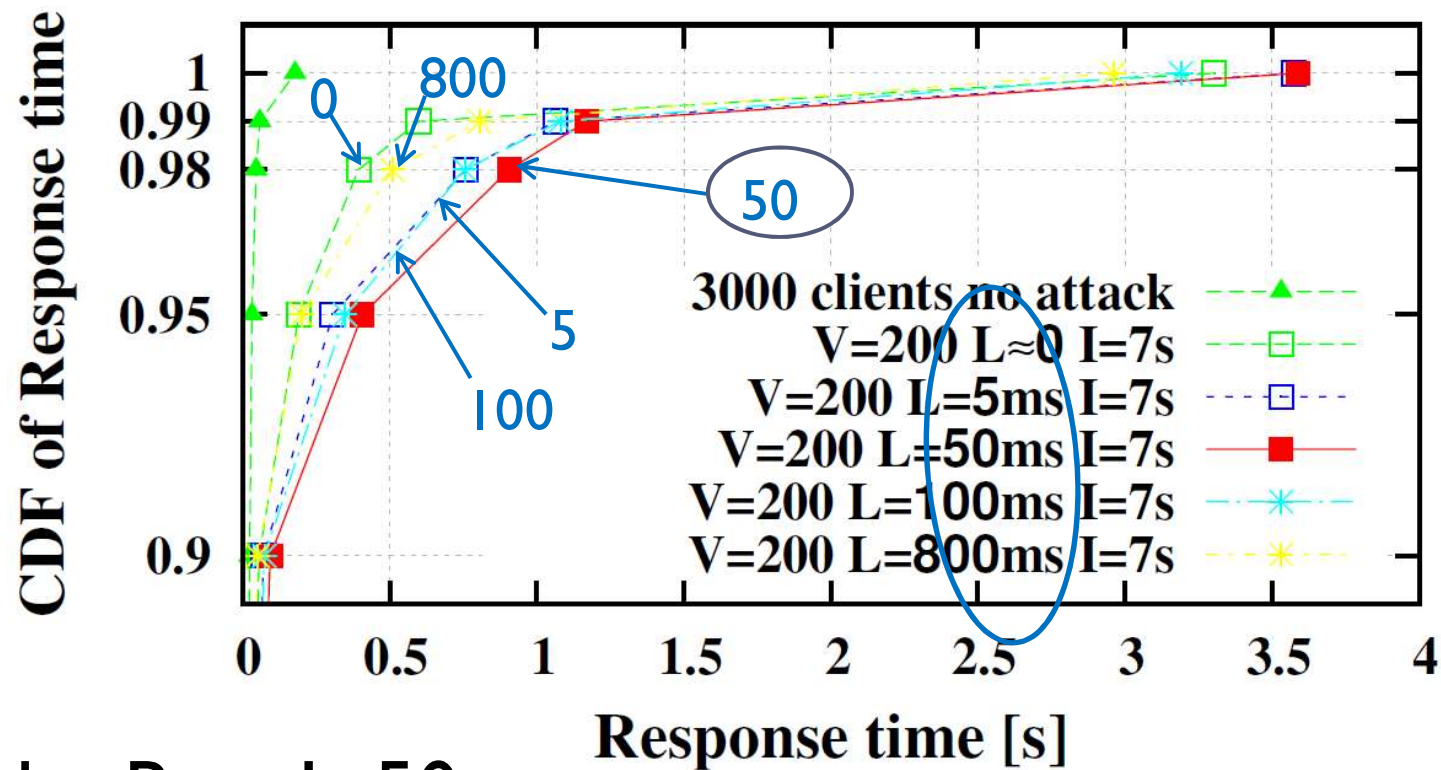


# Training: upper bound of attack volume $V$



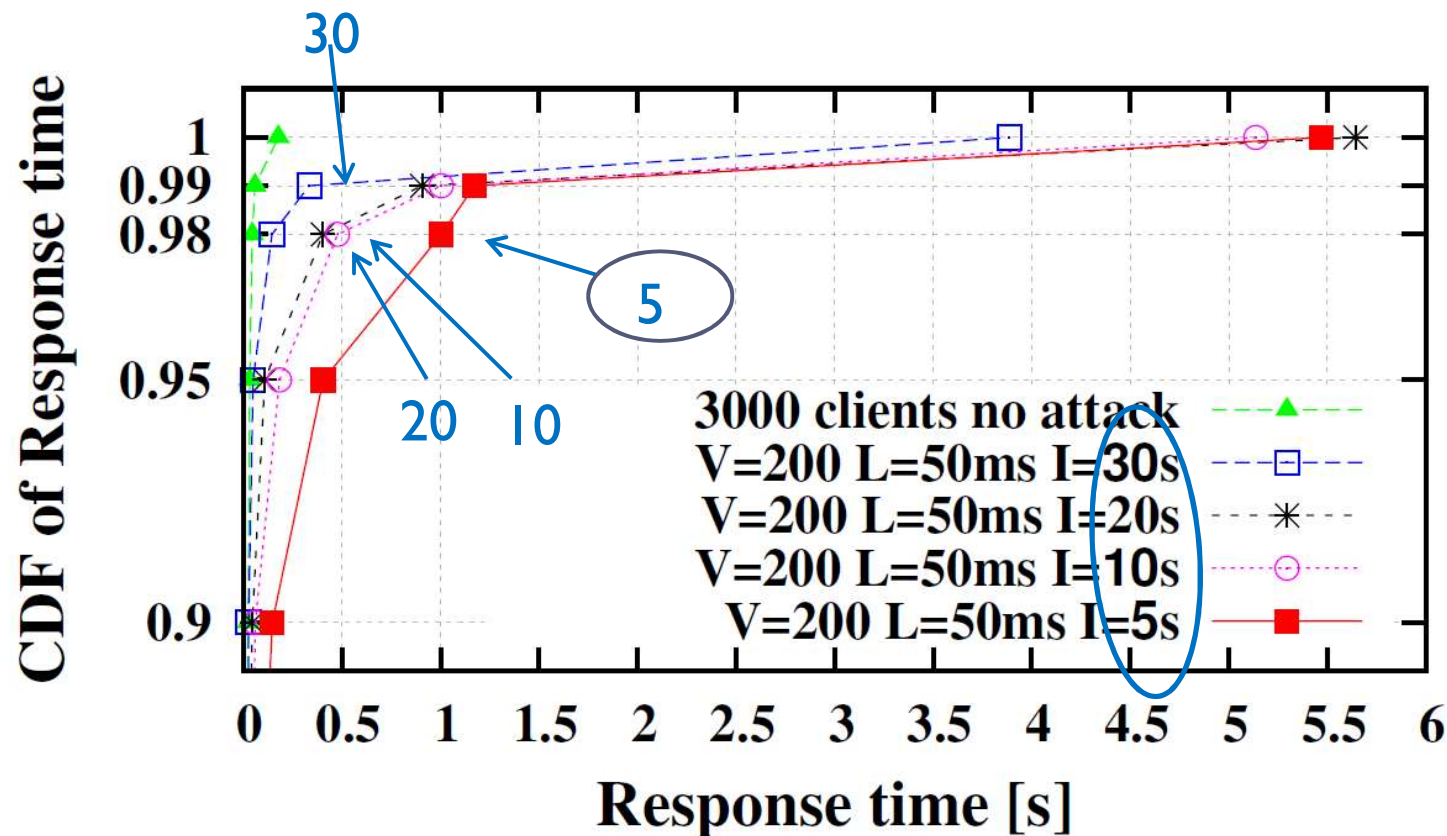
- too high  $V$  triggers target systems to alarm users
- whether this volume will trigger the alarm as a signal

# Training: optimal burst length L



□ the Best L: 50ms

# Training: optimal burst Interval I

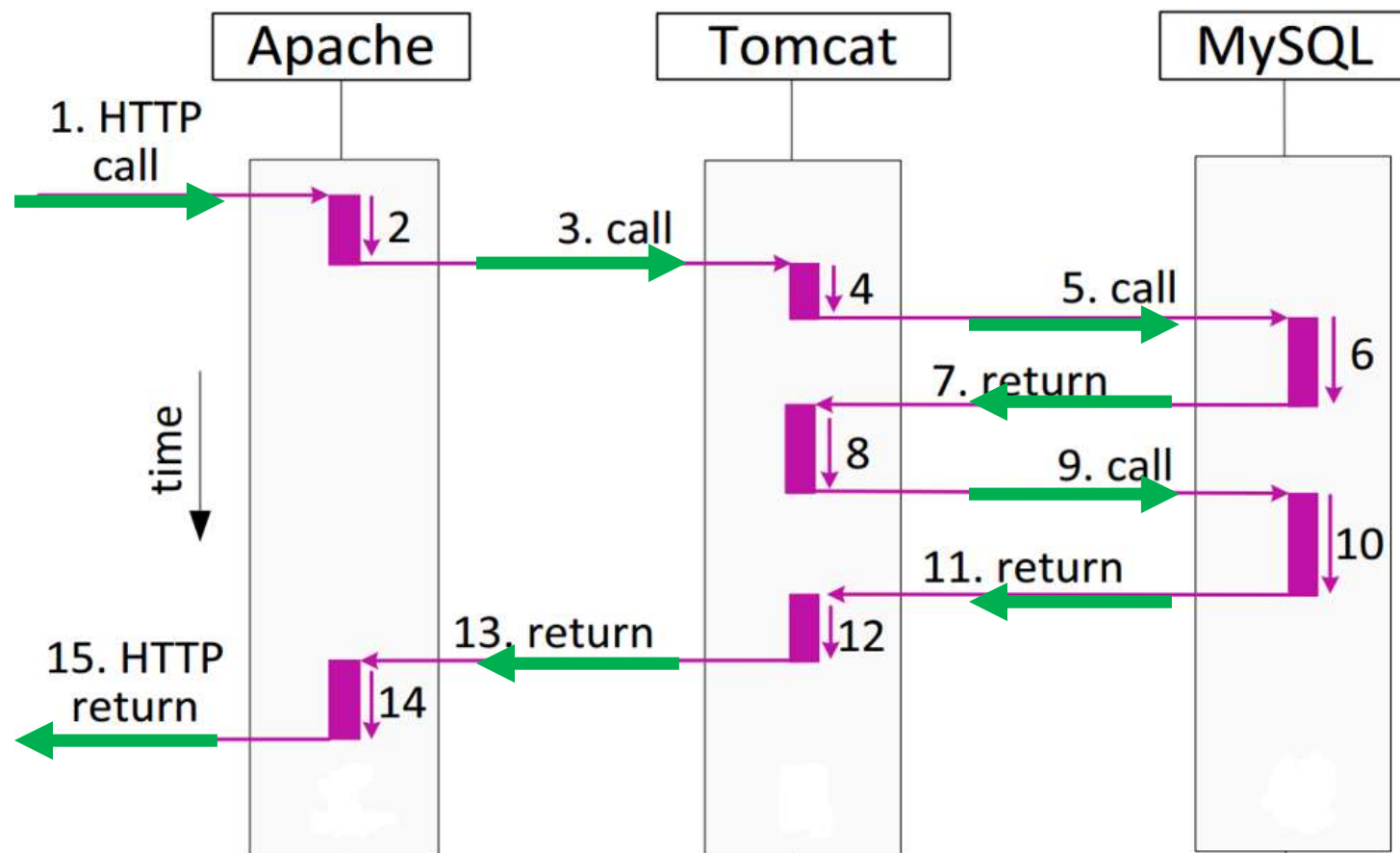


- Reduce the interval step by step (e.g., 5s) until the percentile RT meets the attacking goal

# A Life-time of a Request in the N-tier System (How to Process a Request)



N-tier: the call/response RPC (Remote Procedure Call) style synchronous communication





# Asynchronization

- ❑ Can reduce TCP retransmission
- ❑ Can not reduce the queued time
  - ▶ Overly large buffers result in longer queues and higher latency



# Buff Amplification

- ❑ Queue size in User space, overhead
- ❑ Bufferbloat: TCP backlog buff size in Kernel space
  - ▶ Overly large buffers result in longer queues and higher latency