

# Correlating Low-Level Events To Identify High-Level Bot Behaviors

Liz Stinson  
John Mitchell  
Stanford University

Matt Fredrikson  
Somesh Jha  
University of Wisconsin

Lorenzo Martignoni  
University of Milan

# Our anti-inspirations

- “Pe **Too ambiguous** n app is  
con connecting to the network
- Host-level methods that inundate us with  
in **Too noisy; devoid of meaning**,  
file accesses/changes) without providing a  
higher-level assessment of what’s going on

# Problem Statement

- >5M “distinct, active” bot-infected machines detected between January - June, 2007
  - “active”: carried out at least one attack
  - Symantec Threat Report, Volume XII
- The \*best\* anti-virus signature scanners fail to detect anywhere from 30% to 50% of malware samples *seen in the wild*
  - NB: The best AV scanners may not be who you think they are...

# Problematic Asymmetry

work Malware writers know they have the advantage here and they exploit it. (nt )

- AV mal exc Tens of thousands of novel malware variants created annually d must

# Existing behavior-based detection

- May identify *incidental*, rather than **fundamental** behaviors
  - which dir(s) does app live in: write to: App shadow?
  - App survives reboot? Spawns/terminates other
- For ML-based approaches, may be other ways to achieve same end (i.e. ways not included in model) detect
- More general characterizations
  - Abstract: spyware monitors/reports user actions
  - Concrete: rootkits that load kernel modules

# Broad spectrum. How to evaluate?

- How *effectively* does this method distinguish malicious behavior from benign?
- How *thoroughly* is target behavior captured?
- How *complex* is the identified behavior?
- How *fundamental* is the behavior to the malware's purpose?

# Goals

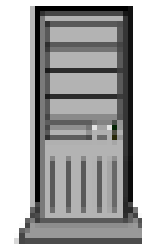
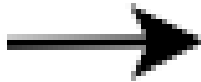
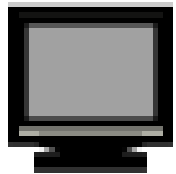
- We want to sample bot commands  
– “downloading and executing a program”

## Sample bot commands

- `http.execute <URL> <local_path>`
- `harvest.registry <reg_key>`
- `redirect <lport> <rhost> <rport>`
- `startkeylogger`
- Via monitoring process execution
- Distinguish malicious from benign instances of above by identifying if *remotely initiated*

# Example: Acting like a proxy

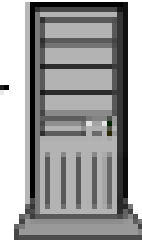
web browser



web proxy

**tcp connection**

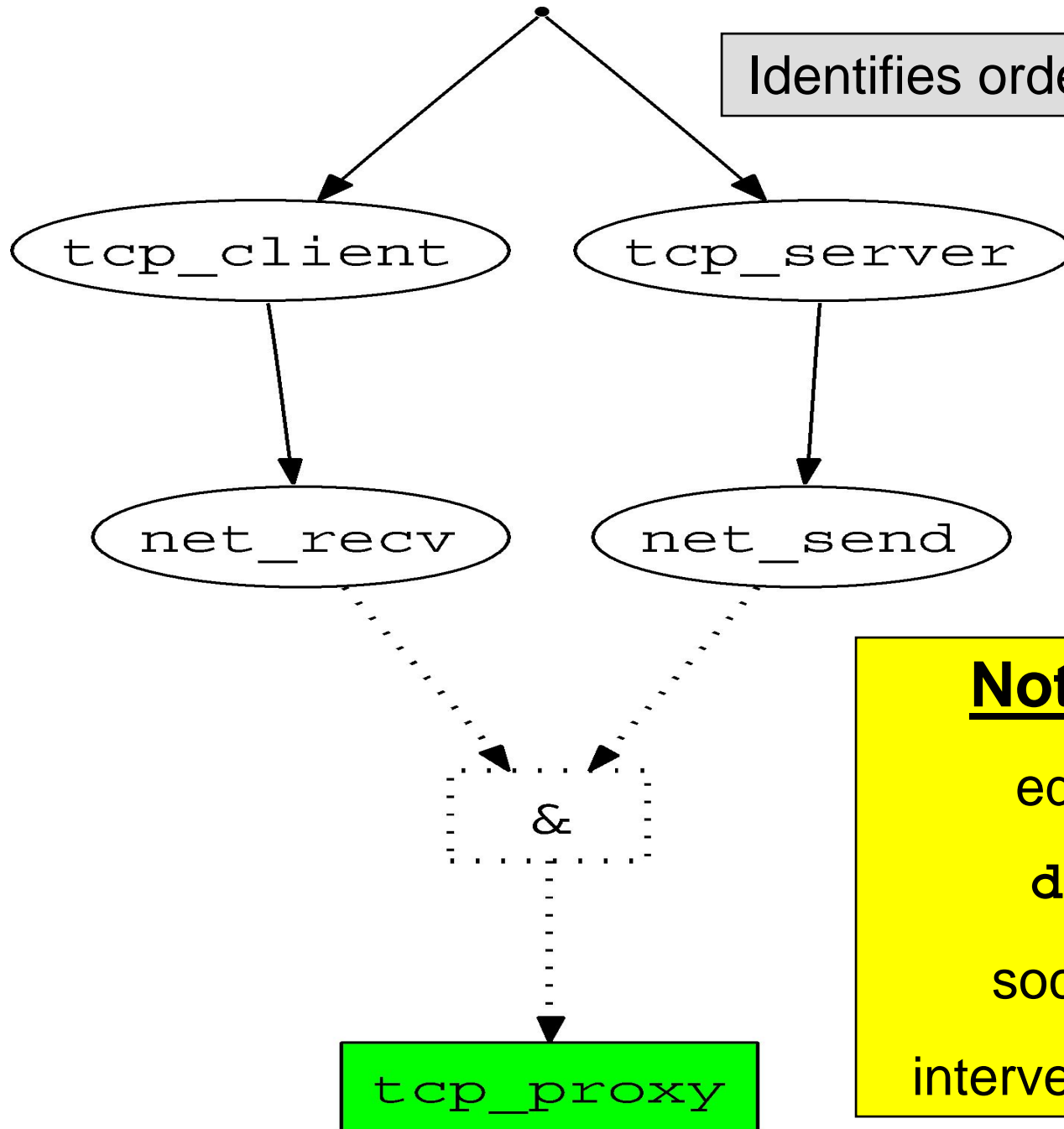
web server



**tcp connection**



Identifies ordering dependencies



**Not shown here**

edge constraints

**die** operations

socket duplication

intervening irrelevant ops

Including parameters and constraints

`tcp_client(sd0, rem_ip, rem_port)`

`tcp_server(sd1, loc_ip, loc_port, cli_ip, cli_port)`

`(sd2 == sd0)`

`(sd3 == sd1)`

`net_recv(sd2, recv_buf) -> recv_len`

`net_send(sd3, send_buf) -> send_len`

`(recv_len > 0)`

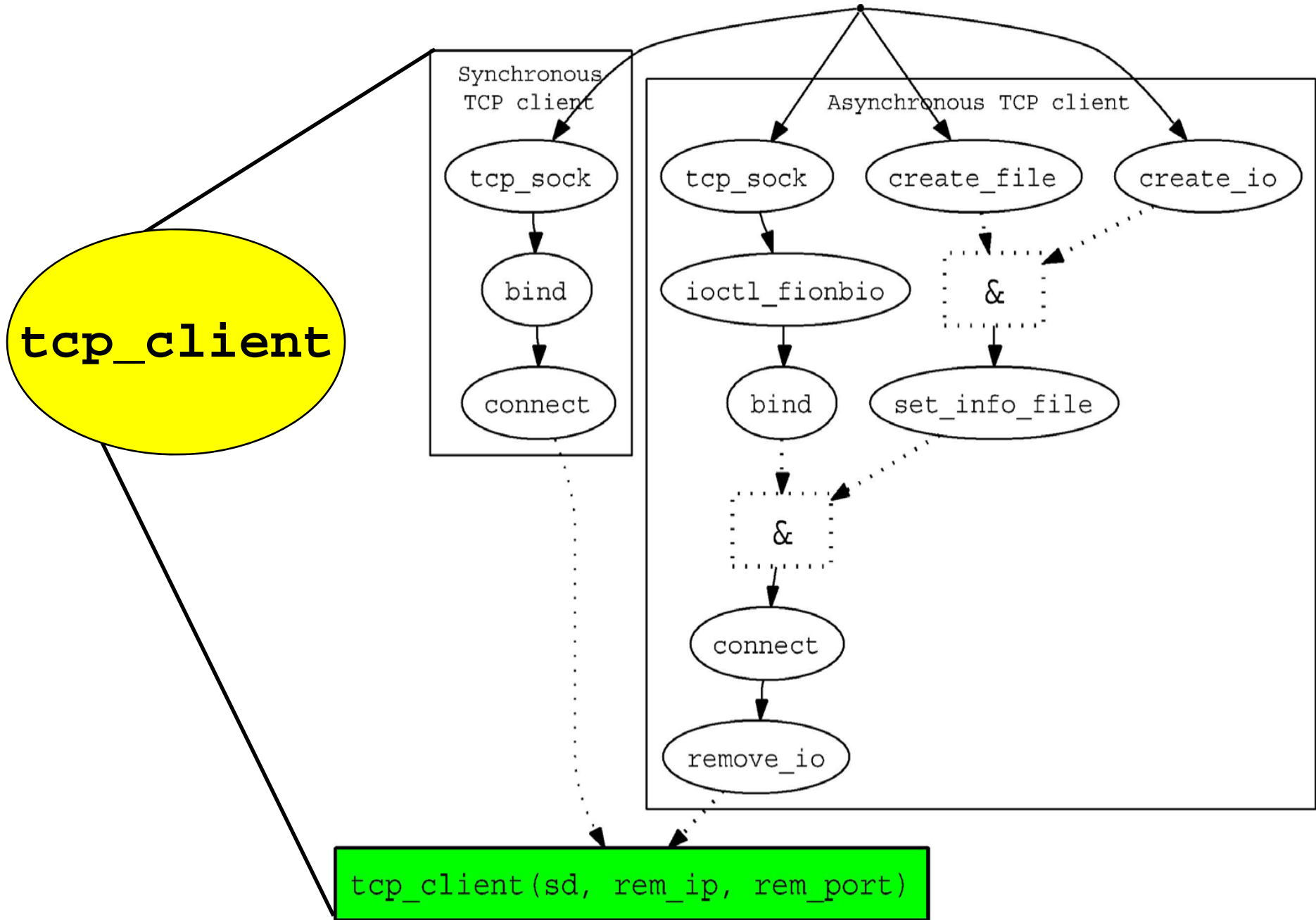
`(send_len > 0)`

Constraints can be pre-conditions or post-conditions

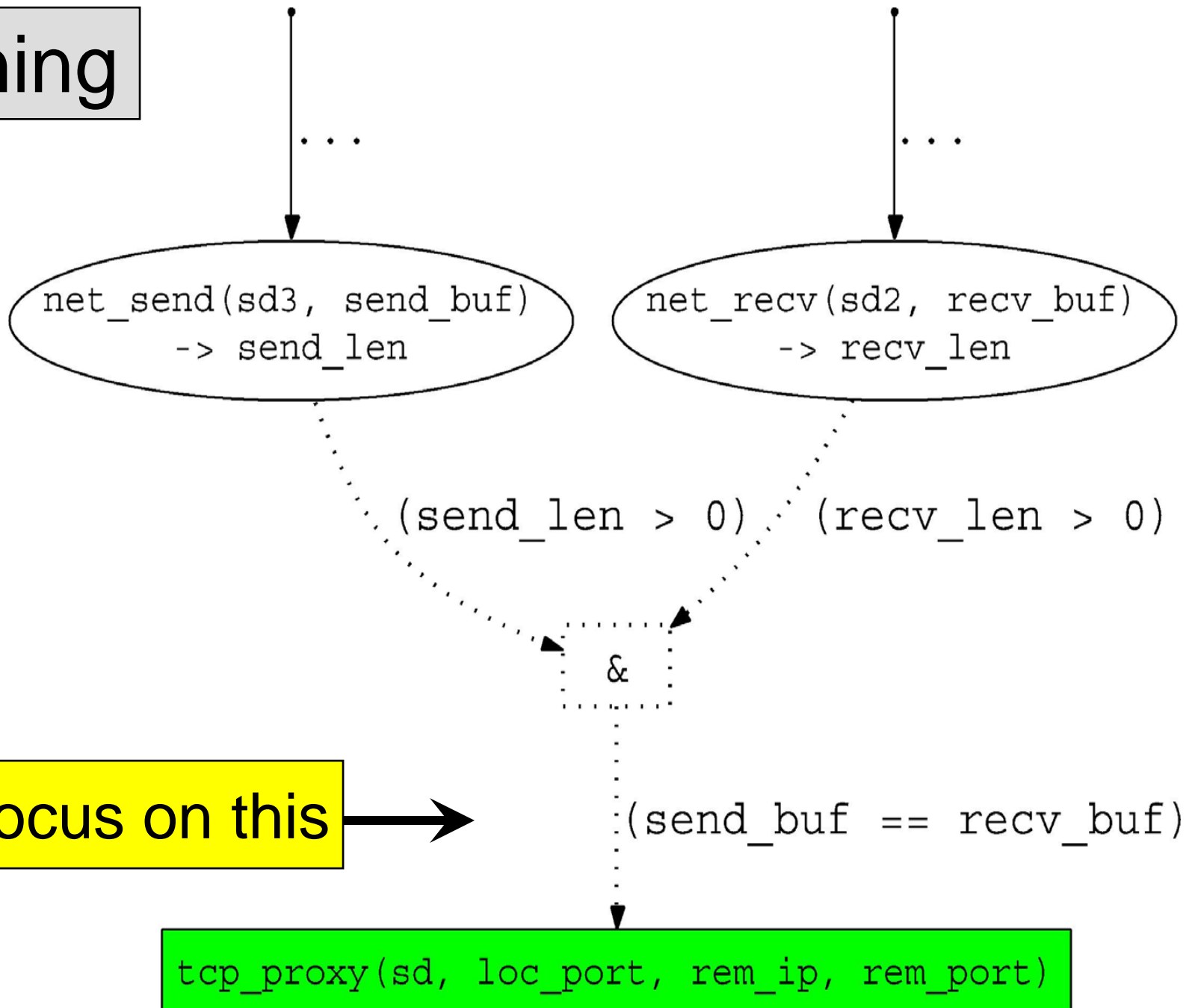
`&`

`(send_buf == recv_buf)`

`tcp_proxy(sd, loc_port, rem_ip, rem_port)`



# Refining



```
( send_buf == recv_buf )
```

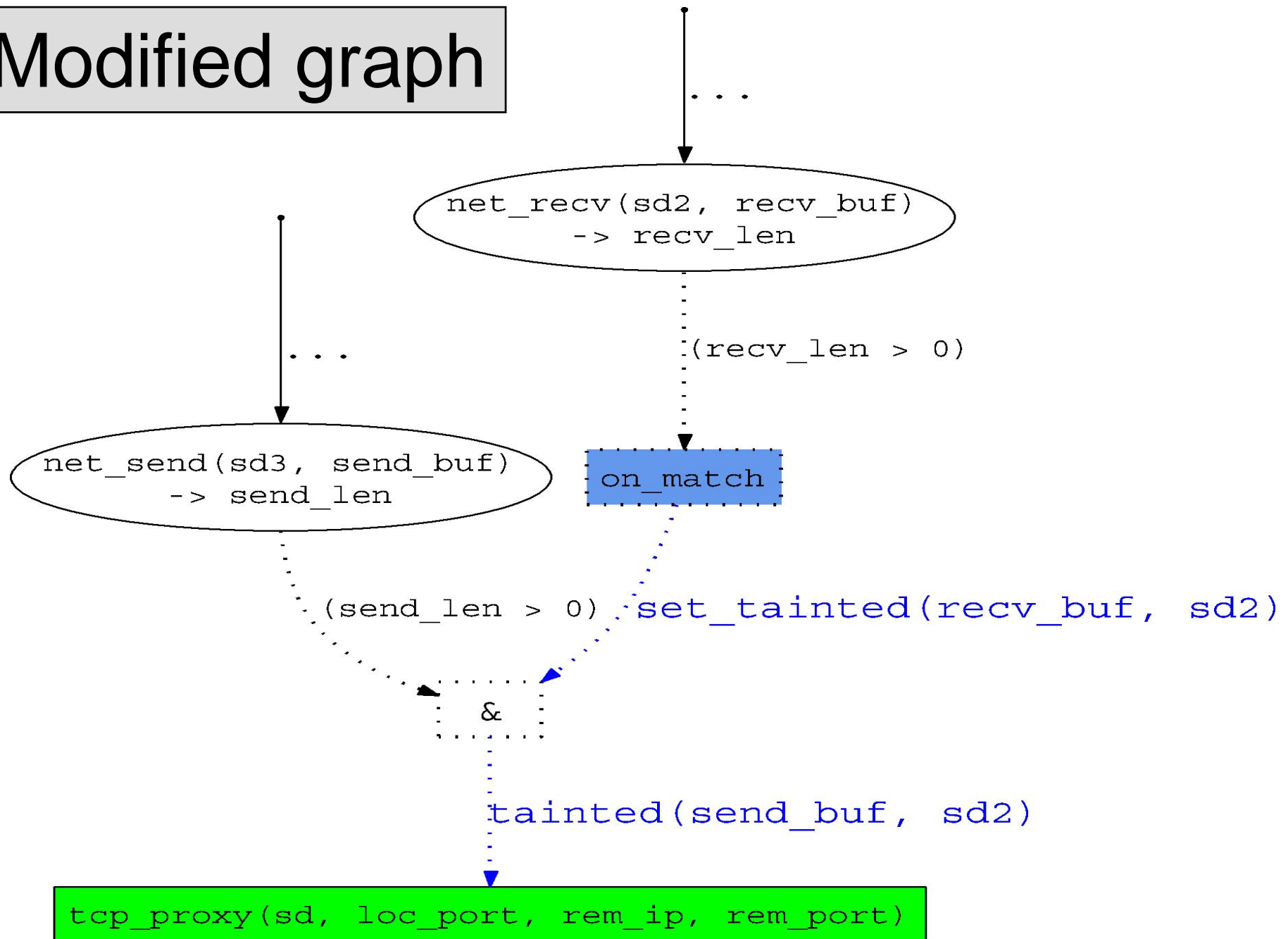
- Too constrained; really want to express: the buffer that is sent is *derived from* a buffer that is received
- Augment (add action to): `on_match` of `net_recv`

```
set_tainted( recv_buf, sd2 /*taint label*/ )
```

- Change condition to:

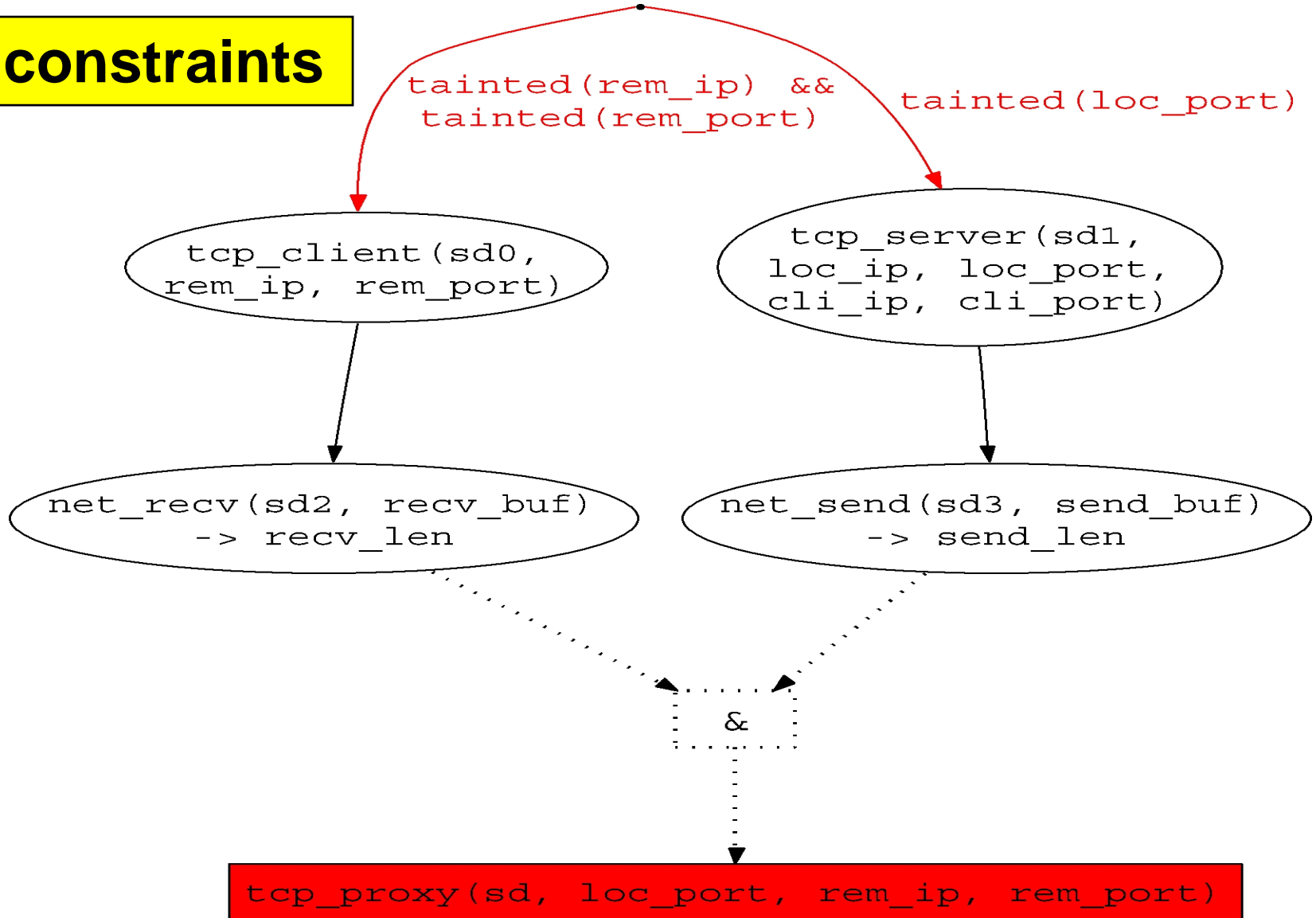
```
tainted( send_buf, sd2 /*taint label*/ )
```

# Modified graph



**.redirect <loc\_port> <rem\_host> <rem\_port>**

**Add constraints**



# “Language” our system exports

- Set of high-level primitives that can be combined to describe interesting behaviors
  - `tcp_client`, `tcp_server`, `net_send`, `net_recv`, `create_exec_file`, ...
- Using these, we can detect:
  - Leak private data (reg key values, file contents, system info, ...)
  - Download and execute a program
  - Send email
  - Proxy
  - Keystroke logging



# Challenges

- Posed by proprietary-OS environment
  - Opacity; identifying operations & constraints
  - Replicating OS semantics
- Posed by syscall interposition generally
- Posed by hypothetical attempts to evade
  - Split behavior across processes or across runs of the same application
  - Expropriate kernel functionality
    - e.g. raw sockets

# Summary

- Target the behaviors that make bots useful
- Identify the essential ops in those behaviors
- Use data-flow analysis info variously
- Good initial results against bots
  - Including: rbot, agobot, dsnxbot, spybot, ...
  - Use bot commands as inspiration
  - Resilient to encryption of bot communications
- Good initial results against benign progs
  - When testing against specifications that encode remote-control requirement
  - Performing user-input tracking