# Web Load Balancing on a Budget

# Pain

- Hosting 60+ websites
- Single web server
    - Redundant subsystems (disk, power)
- SPOF
- Inconvenient maintenance windows
    - Clients
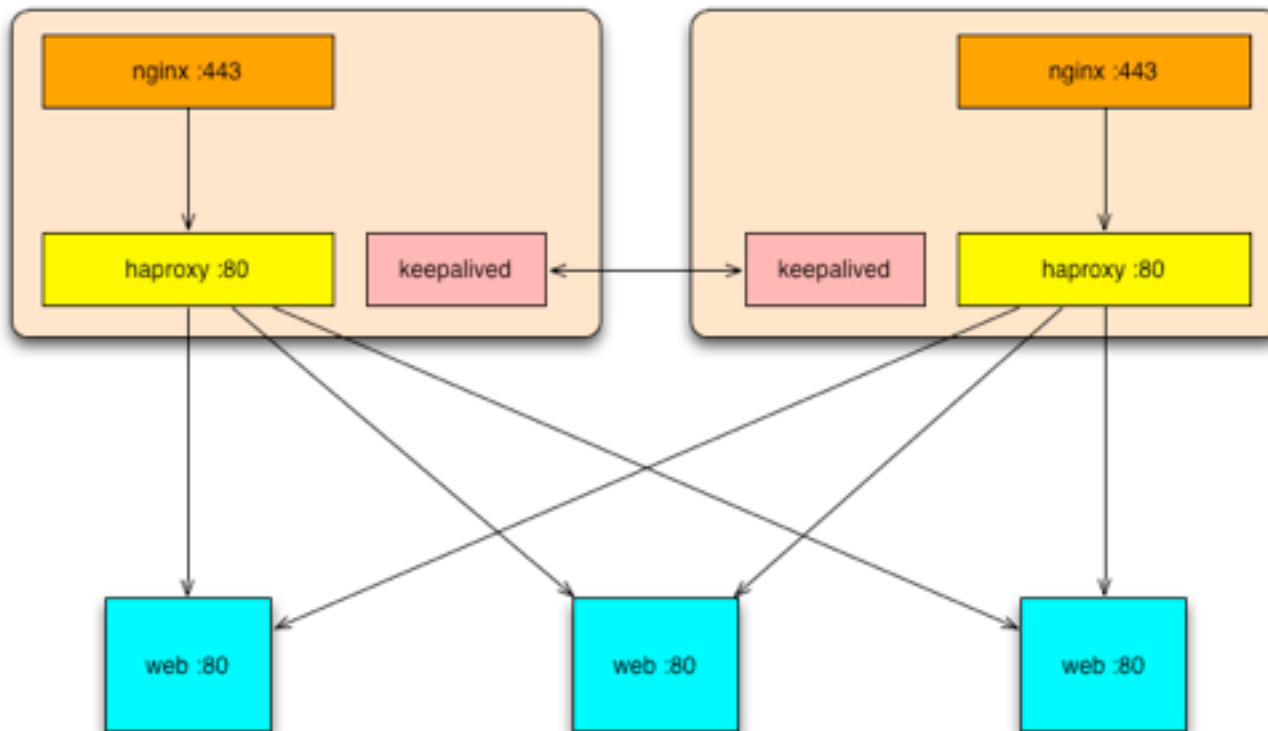    - MY TEAM!

# Scope

- Simple. Availability.
  - Minimize/mitigate downtime
    - Outages
    - Planned maintenance

- Session state persistence (for failed web backends) was not required

# Stack

- Vmware ESXi (all hosts)
- CentOS (load balancers)
  - Nginx (SSL termination)
  - Haproxy (web load balancing)
  - Keepalived (VRRP)
- Windows Server 2008 (IIS)

# Design

```
 1  # Let us bind to addresses and ports that may not be "real" (nginx and haproxy need this)
 2  net.ipv4.ip_nonlocal_bind = 1
 3
 4  # Reuse sockets that are in a TIME_WAIT state so we don't exhaust resources
 5  net.ipv4.tcp_tw_reuse = 1
 6
 7  # Give us a larger useable port range (default: 32768 61000)
 8  net.ipv4.ip_local_port_range = 1024 65535
 9
10  # Close TCP FIN connections faster to help lower resources used by the network stack
11  net.ipv4.tcp_fin_timeout = 30
12
13  # Increase the number of slots that iptables has for tracking connections
14  net.ipv4.netfilter.ip_conntrack_max = 131072
15
16  # Decrease the time that iptables waits to close sockets in TIME_WAIT (def: 120)
17  net.ipv4.netfilter.ip_conntrack_tcp_timeout_time_wait = 30
```

```
 1 http {
 2     server_tokens off;              # make the hackers work for it
 3
 4     server_name_in_redirect off;   # use requested Host header
 5
 6     proxy_read_timeout 1500;        # 25 min; required for long-running reports
 7
 8     upstream www.example.com {
 9         server 192.168.1.2:80;
10     }
11
12     server {
13         listen 192.168.1.2:443;
14         access_log /var/log/nginx/access.log;
15         ssl on; # go!
16         ssl_certificate /usr/local/ssl/wildcard.example.com.cert;
17         ssl_certificate_key /usr/local/ssl/wildcard.example.com.key;
18         ssl_prefer_server_ciphers on;        # prefer SSLv3 and TLSv1 ciphers
19         ssl_ciphers HIGH:+MEDIUM:!ADH:!MD5; # use 128-bit and higher ciphers; exclude ADH, MD5
20         ssl_protocols TLSv1 SSLv3;          # guarantee only TLSv1 and SSLv3 protocols
21
22         location / {    # match all HTTP requests
23             proxy_set_header X-Forwarded-For $proxy_add_x_forwarded_for;
24             proxy_set_header secureCookie YES;
25             proxy_pass http://www.example.com;
26         }
27     }
28
29 }
```

```
 1 global
 2     maxconn 32500
 3     nbproc 1
 4
 5 defaults
 6     log global          # log all proxy instances
 7     option httplog      # verbose HTTP logging
 8     mode http           # layer 7 proxy
 9
10     retries 2
11     timeout connect 2s  # 2 retries x 2s timeout = 4s until haproxy looks for new backend
12     option redispatch   # redispatch requests originally sent to downed backend servers
13
14     timeout client 5s
15     timeout server 1500s # 25 minutes, for long-running reports (matches nginx config)
16
17     balance roundrobin  # also leastconn
18
19     option forwardfor   # X-Forward-For header
20     option httpclose    # force HTTP connections closed so all connections are logged
21
22 listen statsWWW 192.168.1.1:8080
23     stats enable        # enable the haproxy stats page (uri-stem = /haproxy?stats)
24     stats auth god:sex
25     stats refresh 120s
26
27 listen www.example.com
28     bind 192.168.1.2:80
29     acl acl_port_80 dst_port eq 80
30     acl acl_secure hdr(secureCookie) YES
31     acl acl_sourceLocal src 192.168.1.2
32     acl acl_excludeCSS url_dir -i CSS
33     redirect location https://www.example.com/sslRedirect.asp code 301 if acl_port_80 !acl_secure !
acl_sourceLocal !acl_excludeCSS
34     cookie SERVERID insert indirect nocache
35     server www1 192.168.1.3 cookie www1 weight 1 check inter 1s fall 3 rise 2
35     server www2 192.168.1.4 cookie www2 weight 1 check inter 1s fall 3 rise 2
35     server www3 192.168.1.5 cookie www3 weight 1 check inter 1s fall 3 rise 2
38     option httpchk GET /ping.htm HTTP/1.1\r\nHost:www.example.com   # health check
39     capture request header X-Forwarded-For len 15
```

```
 1 global_defs {
 2     lvs_id lb1
 3 }
 4
 5 vrrp_sync_group virtualGroup1 {
 6     group {
 7         virtualInstance1
 8     }
 9 }
10
11 vrrp_instance virtualInstance1 {
12     interface eth0
13     advert_int 1          # advertise interval in seconds
14     state MASTER          # either MASTER or BACKUP
15     virtual_router_id 1 # MUST be between 1 and 255
16     priority 101          # 101 on MASTER, 100 on BACKUP
17     authentication {
18         auth_type AH
19         auth_pass "sexdrugsandrockandroll"
20     }
21     virtual_ipaddress { # only 20 IPs per block, blah
22         192.168.1.1/24 brd 192.168.1.255 dev eth0
23         192.168.1.2/24 brd 192.168.1.255 dev eth0
24         192.168.1.3/24 brd 192.168.1.255 dev eth0
25         192.168.1.4/24 brd 192.168.1.255 dev eth0
26         ...
27         192.168.1.19/24 brd 192.168.1.255 dev eth0
28     }
29 }
```
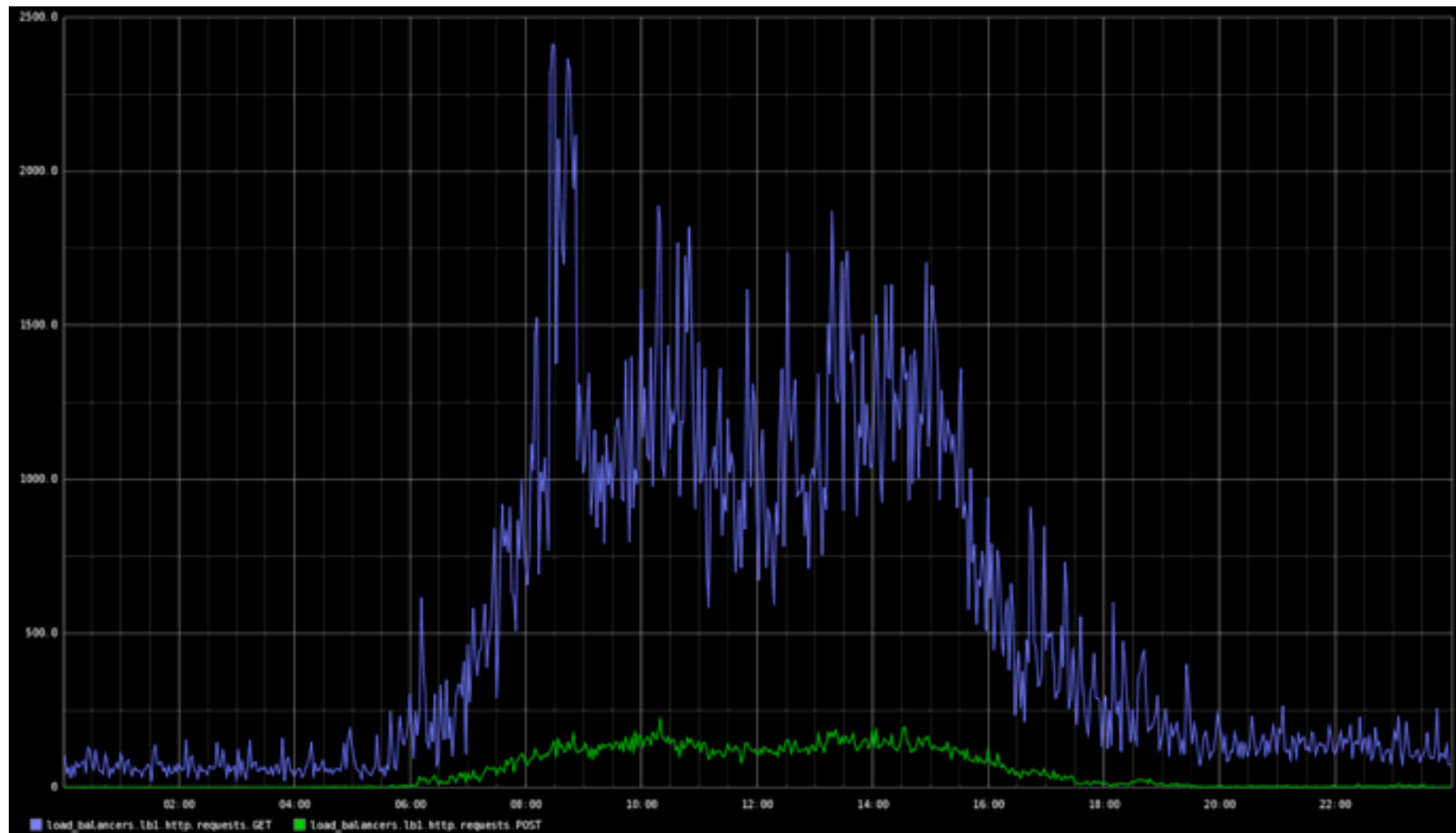
# Pitfalls

- VMware and TSO
  - When jumbo frames go bad
  - HTTP 500 errors
  - Web servers "offline" for 20-30 minutes
  - Disable TSO
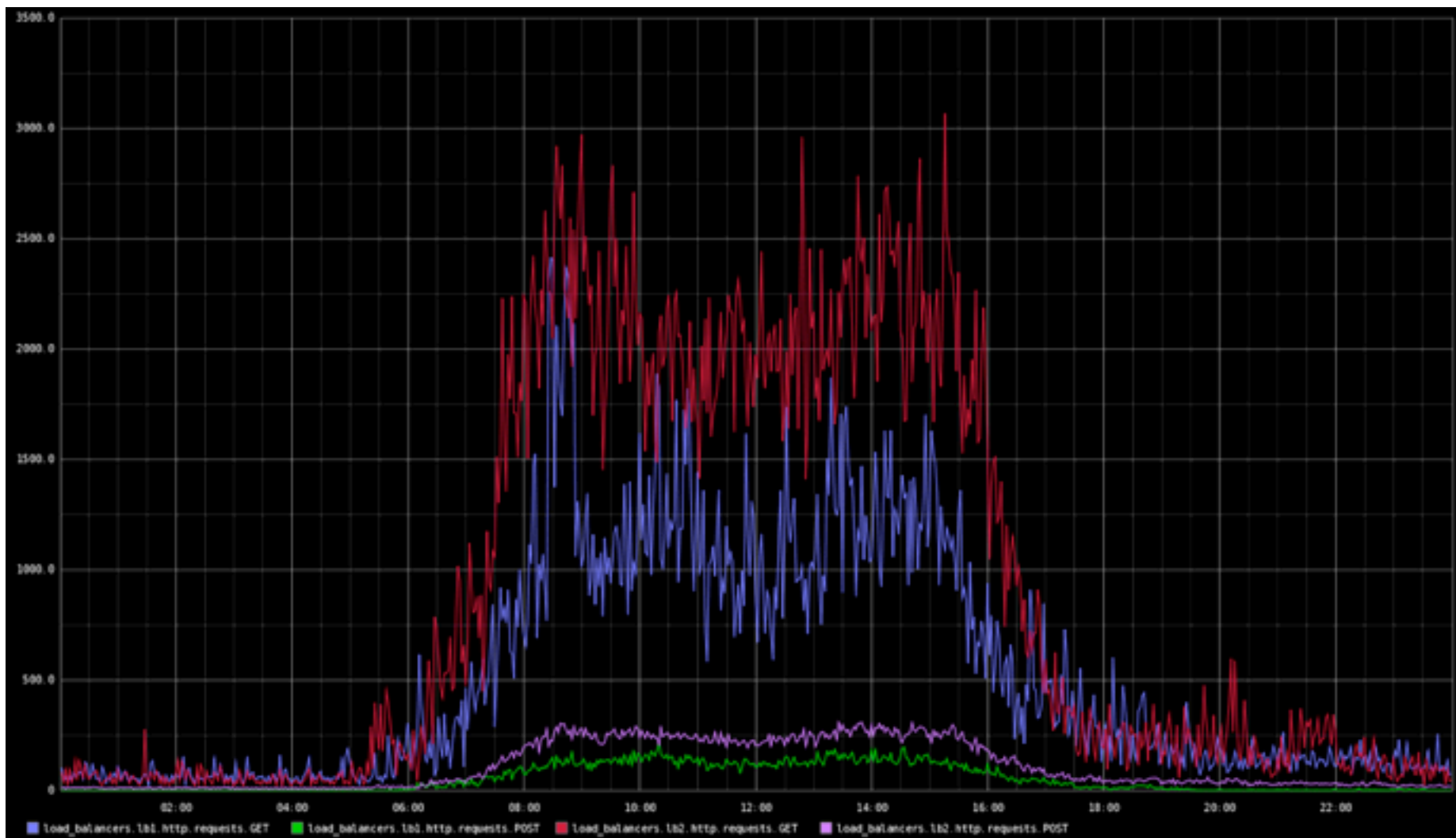    - http://www.ryanfrantz.com/2011/02/03/tcp-segmentation-offload/

# What's Happening?

- No visibility into performance
  - Request volume (nginx/haproxy/IIS)
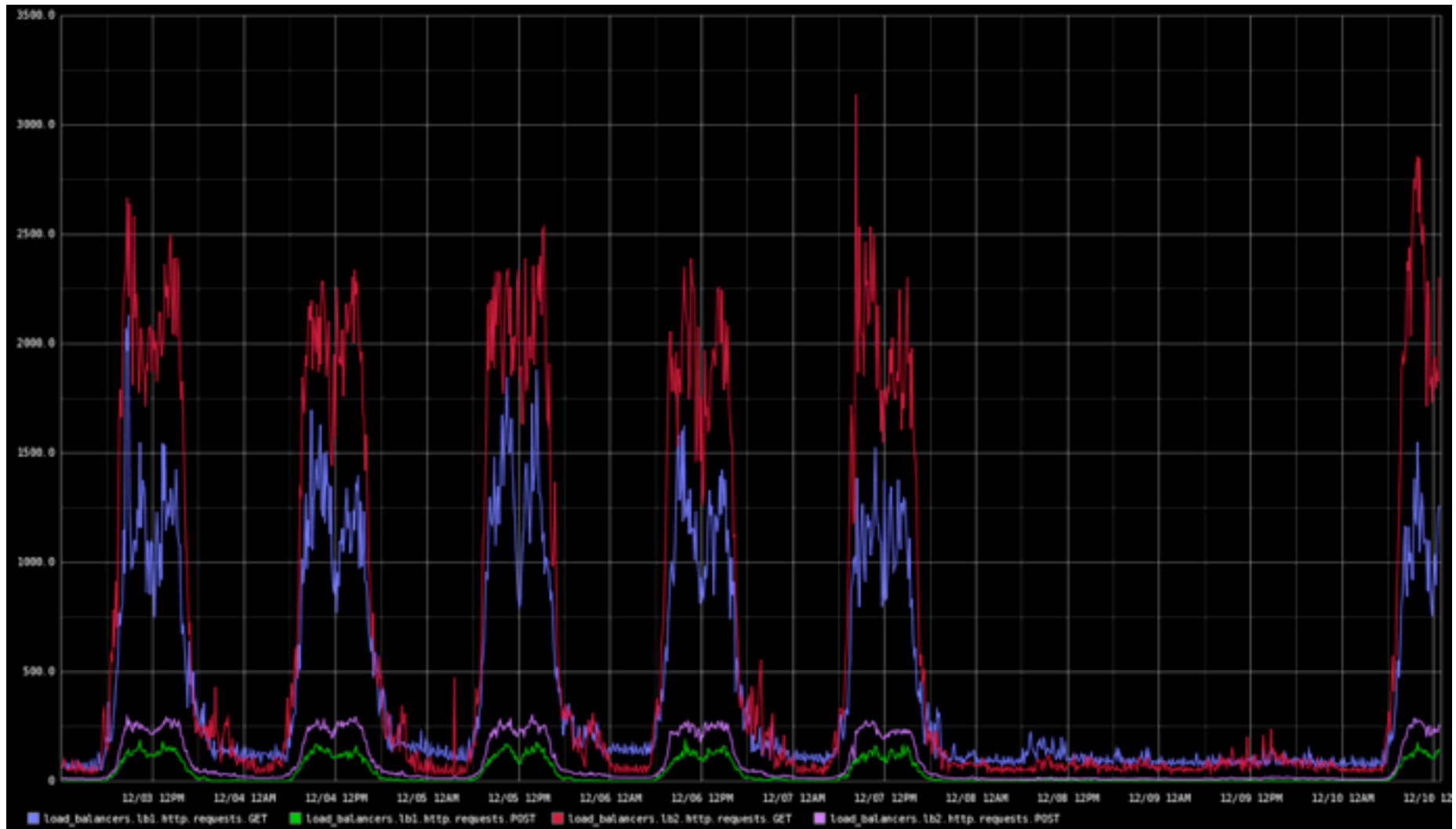  - Seasonality
- Health (beyond up/down)
  - keepalived

# HTTP Traffic: One Load Balancer



load_balancers.lb1.http.requests.GET   load_balancers.lb1.http.requests.POST

# HTTP Traffic: Two Load Balancers



load_balancers.lb1.http.requests.GET  load_balancers.lb1.http.requests.POST  load_balancers.lb2.http.requests.GET  load_balancers.lb2.http.requests.POST

# HTTP Traffic: One Week



load_balancers.lb1.http.requests.GET　　load_balancers.lb1.http.requests.POST　　load_balancers.lb2.http.requests.GET　　load_balancers.lb2.http.requests.POST

# TCP Connections



load_balancers.lbl.connections.tcp.established   load_balancers.lbl.connections.tcp.close_wait   load_balancers.lbl.connections.tcp.close

# TCP Connections: + TIME_WAIT



load_balancers.lb1.connections.tcp.established   load_balancers.lb1.connections.tcp.close_wait   load_balancers.lb1.connections.tcp.close   load_balancers.lb1.connections.tcp.time_wait

# Thank You for Attending LOPSA-East '13

**Please fill out the Trainer Evaluation**
**http://lopsa-east.org/2013/training-survey**

**Rate LOPSA-East '13**
**http://www.lopsa-east.org/2013/rate-lopsa-east-13**

LOPSA
the league of professional system administrators

@Ryan_Frantz
rfrantz@etsy.com
github.com/RyanFrantz
www.ryanfrantz.com