# Benchmarking for HTTP/2

Kit Chan (kichan@yahoo-inc.com)
Kenny Peng (kennyp@yahoo-inc.com )
Nora Shoemaker (nshoemak@andrew.cmu.edu)

# HTTP/2 is Great!!!

# So...

- Which server software should I use? (Hopefully ATS)
- How many machines I need to buy?

## Perspective of a System Engineer

# Benchmarking Tool Needed

- h2load to the rescue
- part of nghttp2
- latest version is 1.4.0

# Some examples

- Basic

*h2load -n100 -c10 -m10 https://www.google.com/*

- Adding/Changing header

*h2load -n100 -c10 -m10 --header="accept-encoding: gzip" https://www.google.com/*

# More Examples

## Multi-threading support

*h2load -t2 -n100 -c10 -m10 --header="accept-encoding: gzip" [https://www.google.com/](https://www.google.com/)*

## Timeout

*h2load -t2 -n100 -c10 -m10 --header="accept-encoding: gzip" --connection-active-timeout=3 --connection-inactivity-timeout=3 [https://www.google.com/](https://www.google.com/)*

# Protocols & Ciphers

- --npn-list
  - Allows you to define preferences of protocol to be used
  - Allows you to load test with h2, spdy/3.1 or http/1.1
- --ciphers
  - Allows you to define the ciphers to be used

# Rate Mode / Timing Script

- Rate Mode - you can control # of connections per seconds
- Timing Script - you can define a list of URLs for each connection to cycle through. (each line is a time in millisecond, a tab and then the url) e.g

100.0    https://screen.yahoo.com/
200.0    /__rapid-worker-1.1.js
300.0    /__test.css

- Example Usage

*h2load -c100 -r10 -t5 --header="accept-encoding: gzip" --timing-script-file=/tmp/myscript.txt*

# h2load output

```
[kichan@loadtest3 ~]$ time h2load -n80 -c10 -r1 -t1 --header=":authority: screen.yahoo.com" --header="user-agent: Mozilla/5.0 (Macintosh; Intel Mac OS X 10_10_3
) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/42.0.2311.135 Safari/537.36" --header="accept-encoding: gzip"  --input-file=/home/kichan/exp/v2test.nds.h2.test
starting benchmark...
spawning thread #0: 10 total client(s). Up to 1 client(s) will be created every 1s 80 total requests
TLS Protocol: TLSv1.2
Cipher: ECDHE-RSA-AES128-GCM-SHA256

finished in 9.04s, 8.84872 req/s, 617.84KB/s
requests: 80 total, 80 started, 80 done, 78 succeeded, 2 failed, 0 errored, 0 timeout
status codes: 78 2xx, 0 3xx, 2 4xx, 0 5xx
traffic: 5719906 bytes total, 75302 bytes headers, 5635540 bytes data
                   min        max       mean        sd      +/- sd
time for request:  16.67ms   525.94ms   33.59ms   76.51ms   97.50%
time for connect:  10.32ms    16.07ms   13.09ms    1.92ms   50.00%
time to 1st byte:  28.17ms    36.50ms   32.28ms    2.61ms   60.00%
```
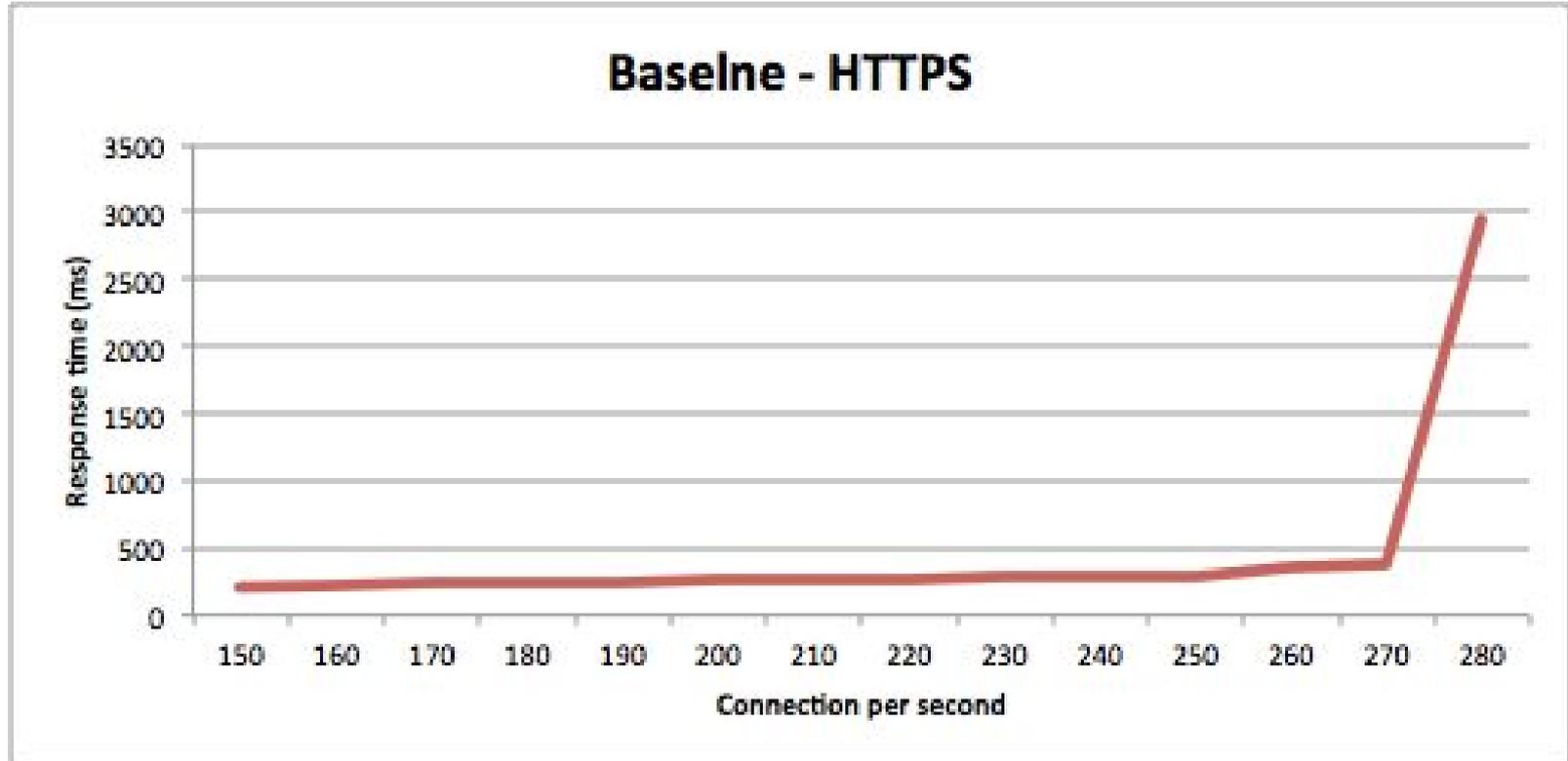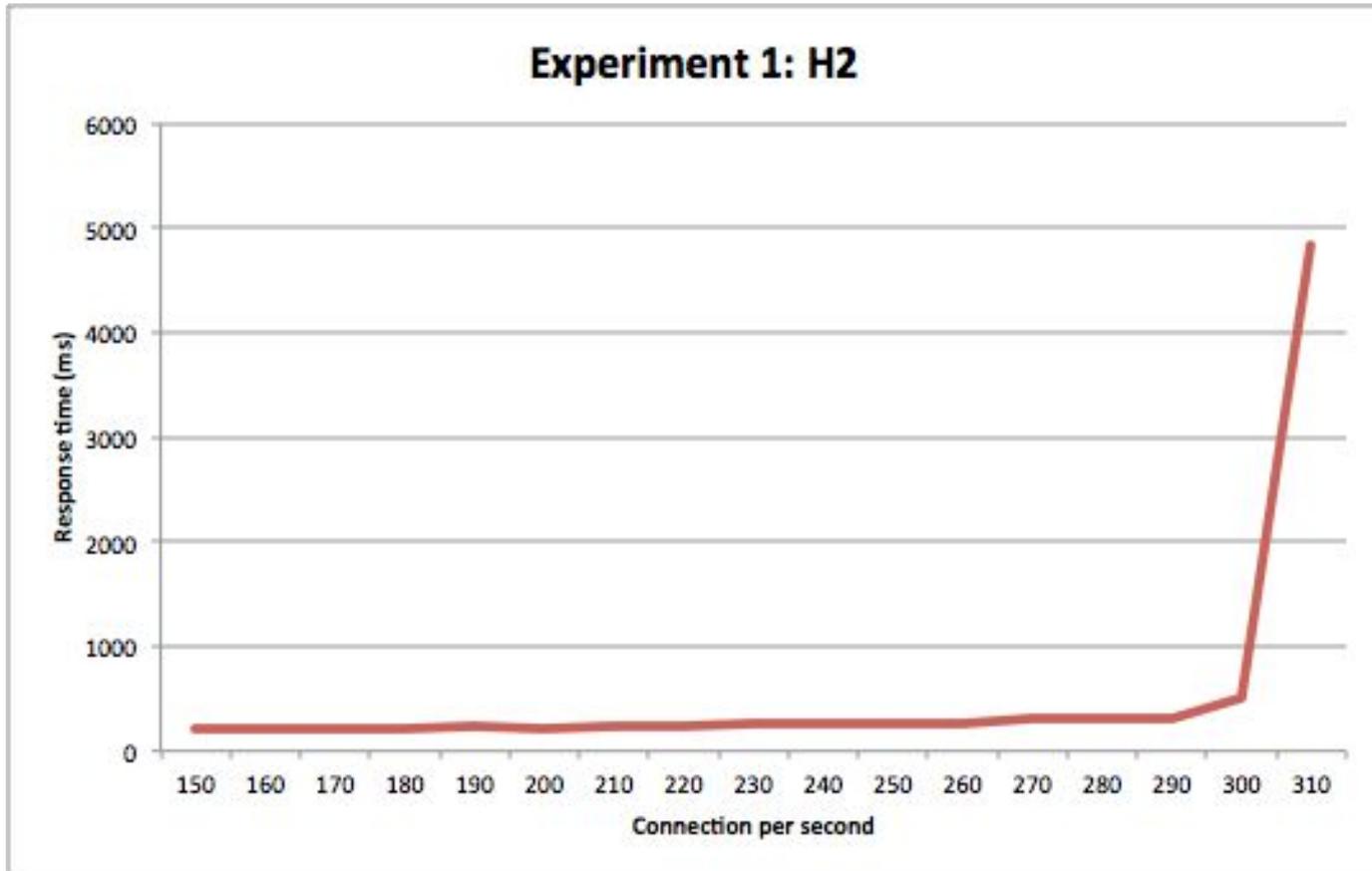
# Experiment 0: Baseline

Background

- One page
- Heavy use of ESI - i.e. ATS is doing page assembly
- Other plugins use to validate cookie, finding out locations, determining buckets for testing + other stuff
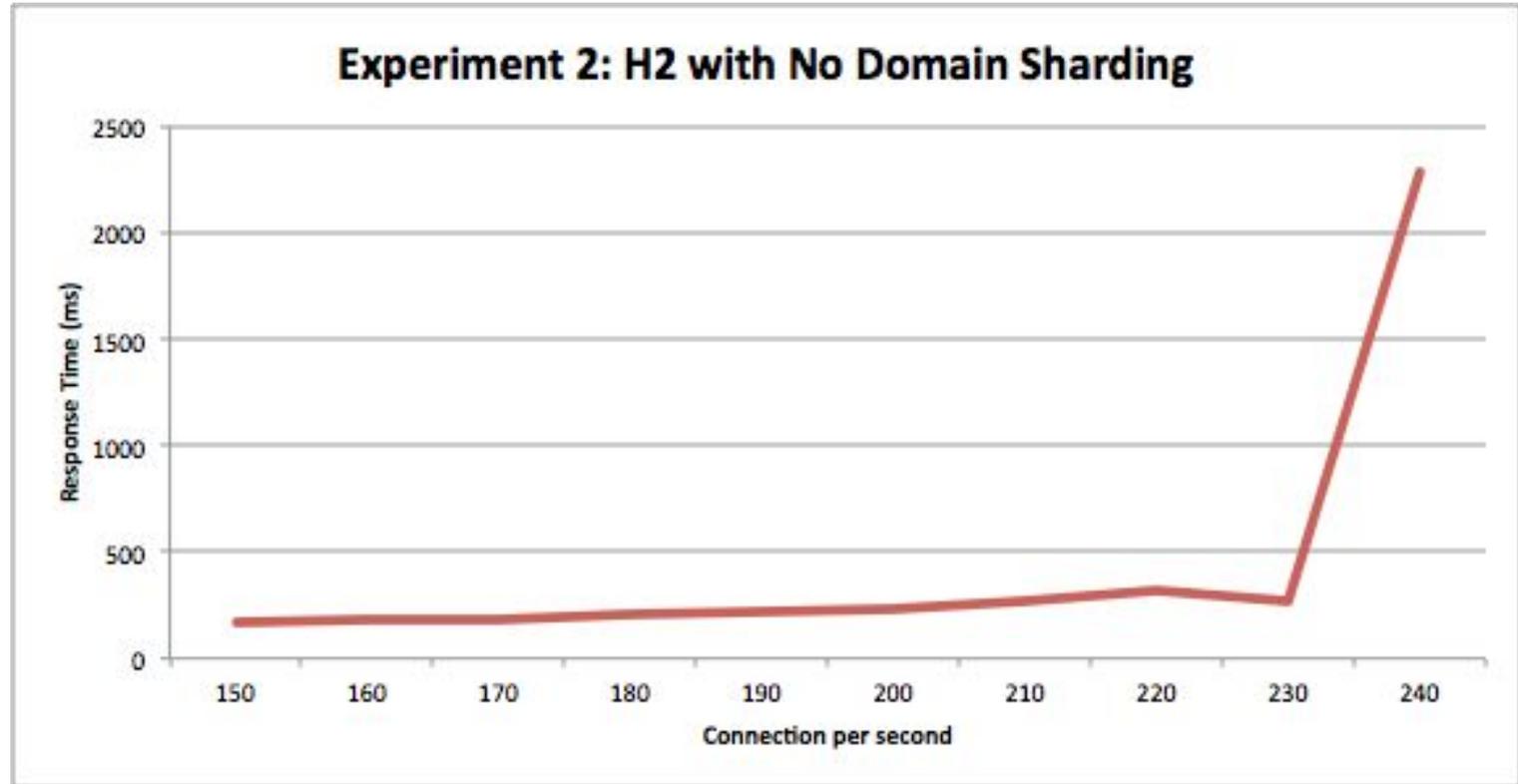
# Experiment 0: Baseline



Baselne - HTTPS

# Experiment 1: Just turn H2 on

# Experiment 2: No Domain Sharding

- Same page + 8 assets in one connection (e.g. CSS/JS/SWF/WOFF etc)
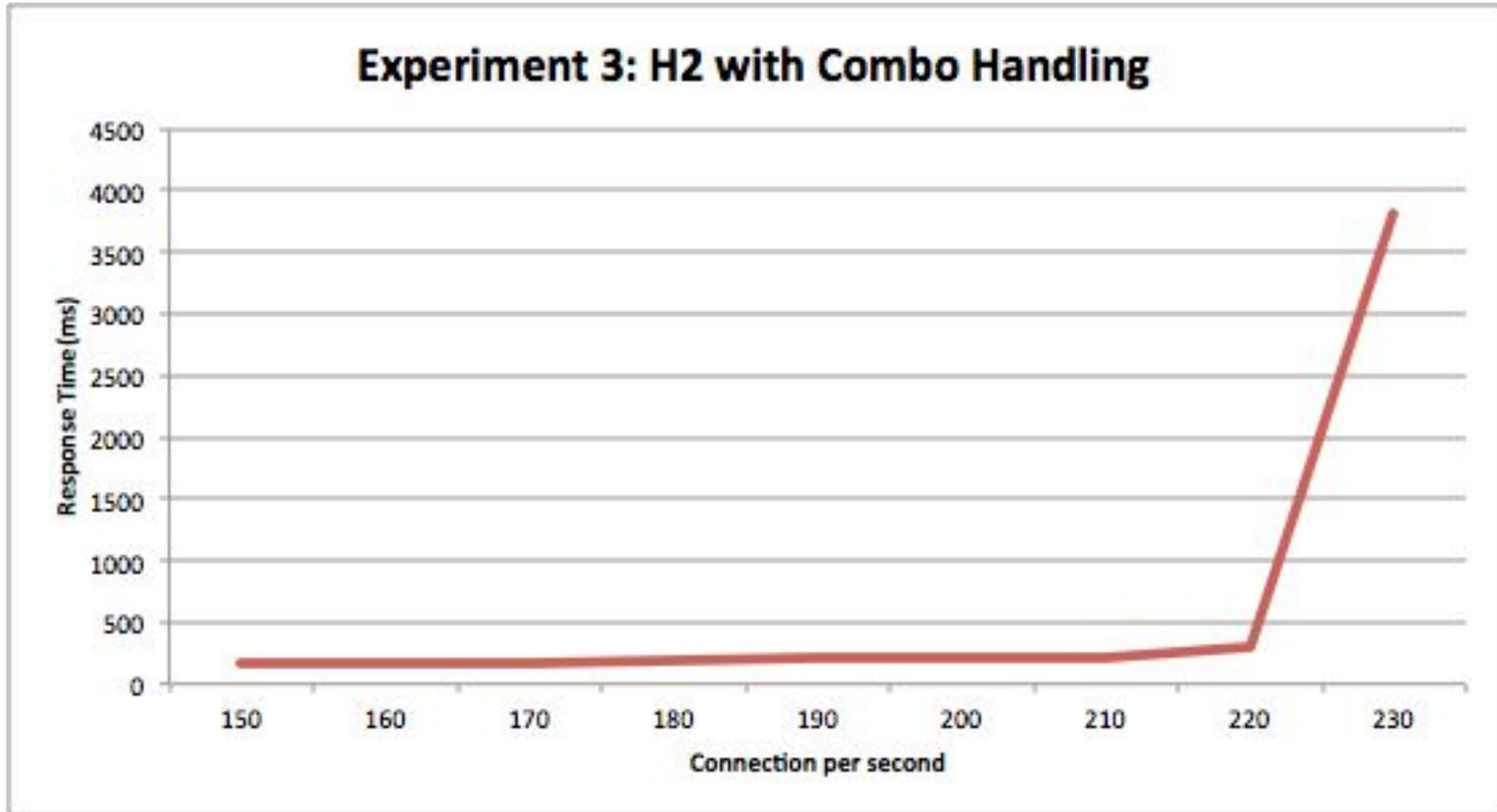
# Experiment 2: no domain sharding

# Experiment 3 & 4: combo handling or not

- Combo Handling of assets. e.g.

*https://s.yimg.com/zz/combo?/os/stencil/2.0.26/styles.css&/os/mit/td/lasso-1.2.197 /cinematron-simple-dark/cinematron-simple-dark-min.css*

- First experiment retrieving one page and 5 combo assets URLs
- Second experiment retriveing one page and the "un-combo-ed" URLs of the above 5 combo assets URLs  (80+ URLs)

# Experiment 3: Combo Handling

# Experiment 4: No combo handling

- Latency/Response time is too high even with low CPS

# Experiment 5: ATS vs nghttpx

- Settings - access/error log turned on, no ocsp stapling, no cache for ATS, same # of execution theads
- Just do http/2 termination and proxy the requests
- Requests - 3 large image objects (60K to 200K) per connection

# Experiment 5: Results

- Process CPU utilization during idle - ATS: 0.5%, nghttpx: 0%
- System CPU utilization during idle - ATS: 7%, nghttpx: 6.5%
- Peak System CPU utilization under same traffic load - ATS: 18%, nghttpx: 15%


- Imply nothing! Simply a comparison worth investigating further for a very particular scenario.

# Final Words

- We need to consider server capacity for H2 and related deployment
- h2load far from perfect
    - Contribution opportunities!!!

# Credits/Shoutouts

Nora - patches for timeout, rate mode and its multithread support

Kenny - patches for header, running most of the experiments

Tatsuhiro Tsujikawa - Owner of the nghttp2 project - https://nghttp2.org/

## h2load: Don't DOS our server!

⅄ master    ⬙ v1.4.0

Browse files

🟫 **tatsuhiro-t** committed 29 days ago

1 parent 5594e3d    commit 11cb4ea214f3432f9117f4e95663aeffcfff9c0f

📄 Showing **1 changed file** with **7 additions** and **0 deletions**.

Unified    Split

7 🟩🟩🟩🟩🟩 src/h2load.cc

View

```
       @@ -1918,6 +1918,13 @@ int main(int argc, char **argv) {
1918  1918      config.nv.push_back(std::move(cva));
1919  1919    }
1920  1920
      1921  +  // Don't DOS our server!
      1922  +  if (config.host == "nghttp2.org") {
      1923  +    std::cerr << "Using h2load against public server " << config.host
      1924  +              << " should be prohibited." << std::endl;
      1925  +    exit(EXIT_FAILURE);
      1926  +  }
      1927  +
1921  1928    resolve_host();
1922  1929
1923  1930    std::cout << "starting benchmark..." << std::endl;
```

Thanks

# Bonus - Generating HAR

e.g. -

*nghttp -nv --har=/tmp/sample.out https://www.google.com/*

# Bonus - HAR Viewer