# Varnish http accelerator

- A 2006 software design

Poul-Henning Kamp

phk@FreeBSD.org

# Varnish Cheat-Sheet

- Web-accelleration for slow CMS systems
- Narrow focus on server side speedup
  - No FTP etc.
  - Content provider features
- High Performance
  - 32 & 64bit, large RAM, sendfile, accept filters
  - SMP/Multicore friendly architecture
  - 2006 software design
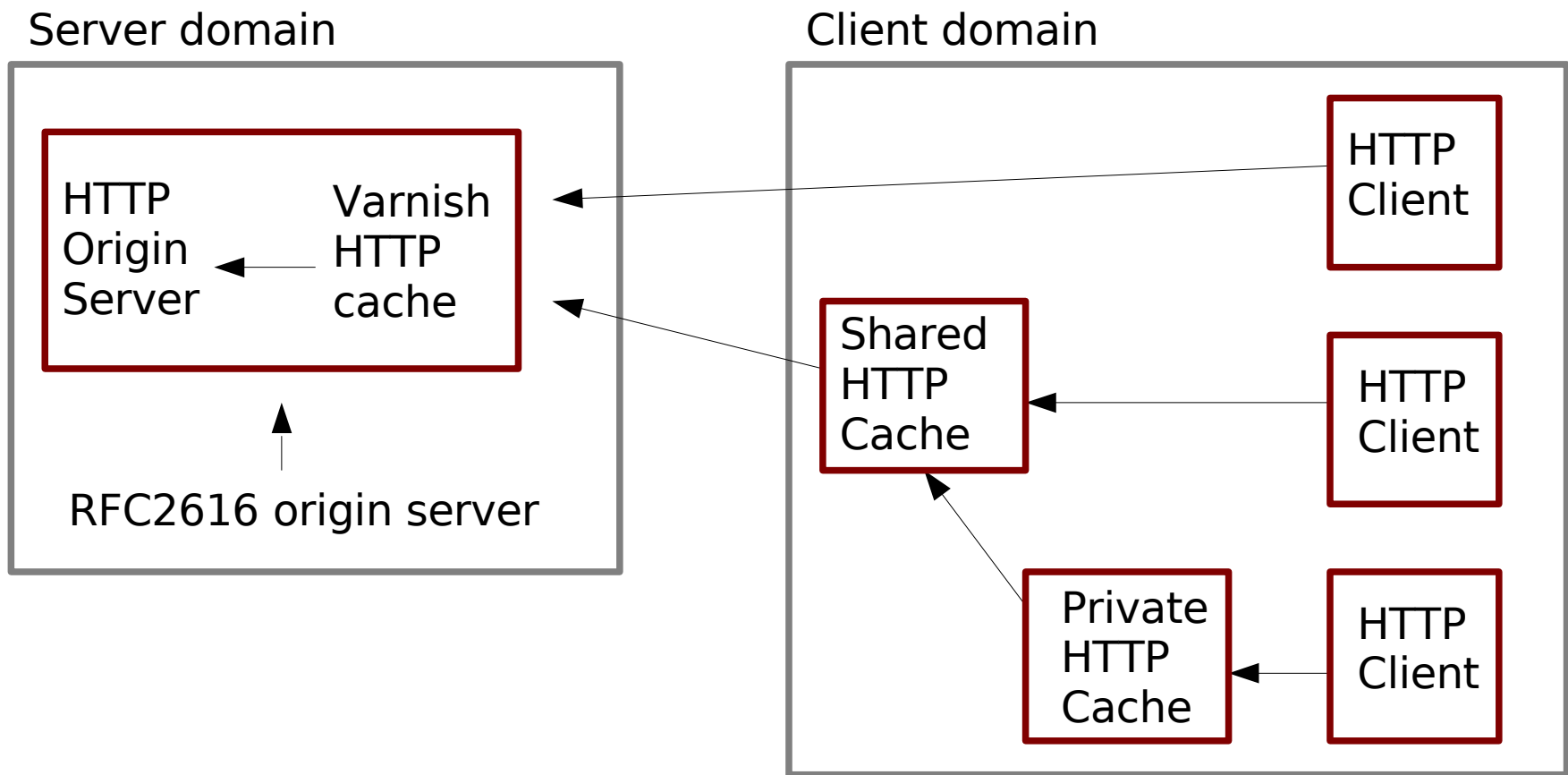  - 11 syscall + 7 locks per cache hit

# dictionary:Varnish

- tr. v. var·nished, var·nish·ing, var·nish·es

  - 1. To cover with varnish.

  - 2. To give a smooth and glossy finish to.

  - 3. **To give a deceptively attractive appearance to.**

# RFC2616 and Varnish

**Server domain**

HTTP Origin Server    Varnish HTTP cache

RFC2616 origin server

**Client domain**

HTTP Client

Shared HTTP Cache

HTTP Client

Private HTTP Cache

HTTP Client

# Client Cache Situation

- Origin servers are adversarial.
- Anything the origin server says is law
  - … if we can make sense of it.
- If in doubt: don't cache.
- Be semantically transparent at any cost.
- If origin server does not reply: error.

# Server Cache Situation

- Backend (origin server) is on our side
  - More precisely: We are on its side.
- We might be <u>responsible</u> for modifying the origin servers instructions.
  - Change TTL, rewrite URLs etc.
- Whatever happens: protect the backend.
- If backend does not reply: do something!

# Content Provider Features

- Instant URL invalidation
  - Regexp matching
  - Easy Integration to CMS system

- Each object checked max 1 time
  - When used next time
  - Many objects will expire without being checked.

# Content Provider Features

- Varnish Configuration Language
  - Full expiry time control
  - Load/Situation mitigation
  - Content substitution
  - URL editing
- Prefetching (v2 feature)
  - Inspect object usage stats
  - Compression for bandwidth savings

# Varnish Config Language

- Simple domain specific language
  - Compiled via C language to binary
    - Transparantly handled by manager process.
  - Dynamically loaded into cache process
  - Multiple VCLs can be loaded concurrently
- Instant switch from one VCL to another.
  - From CLI or VCL

# VCL example

```
acl journalists {
        10.0.0.0/8;
}

if (client.ip ~ journalists) {
        pass;
}

if (req.url.host ~ "cnn.no$") {
        rewrite req.url.host "cnn.no" "vg.no";
}

if (!backend.up) {
        if (obj.exist) {
                set obj.ttl += 10m;
                deliver;
        }
        switch_config "ohhshit";
}
```
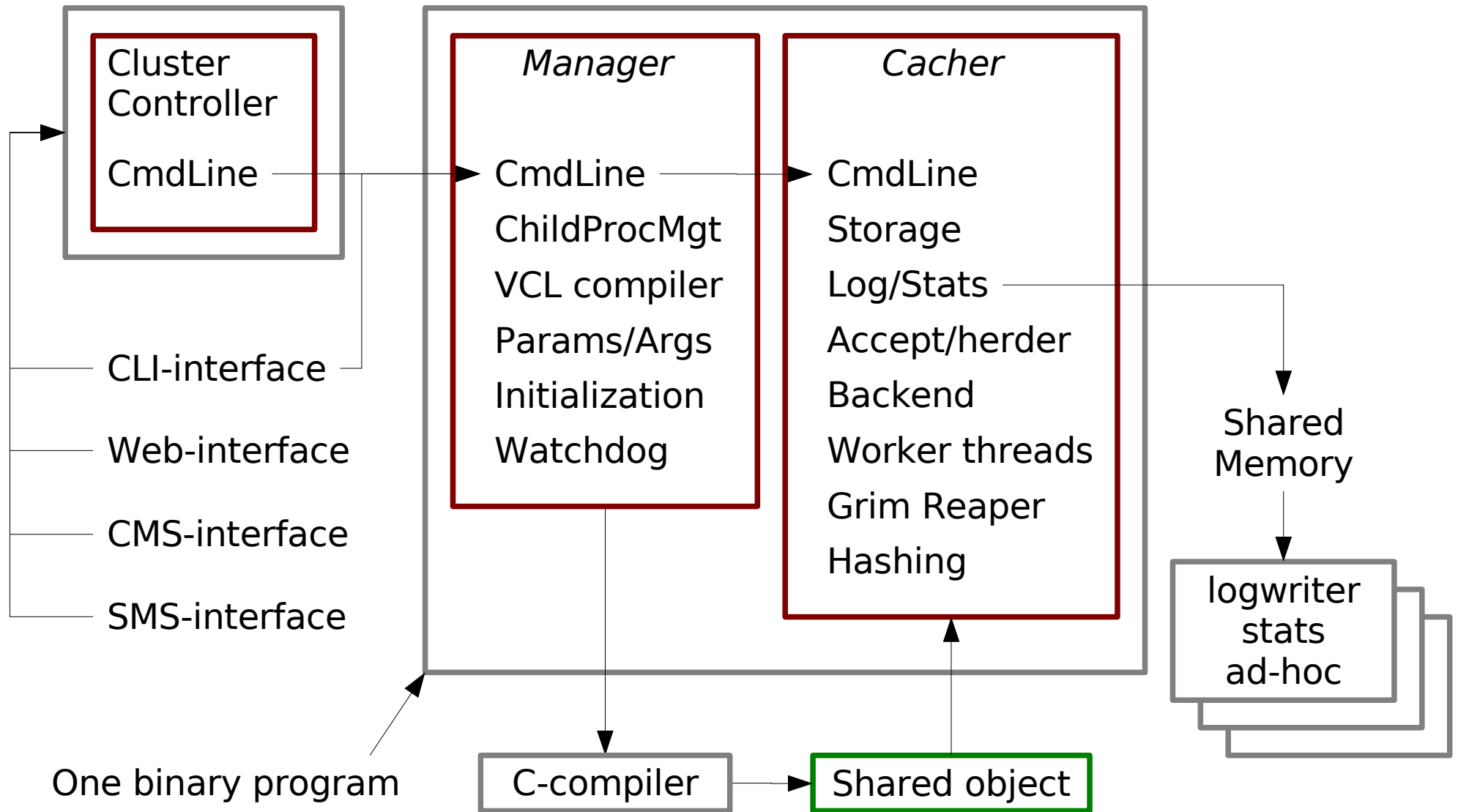
# Operational Features

- Two process design
  - Automatic restarts on crash
- One binary program
  - Easy installs
- Command Line Interface
  - No need for browser with java
  - Change parameters on the fly
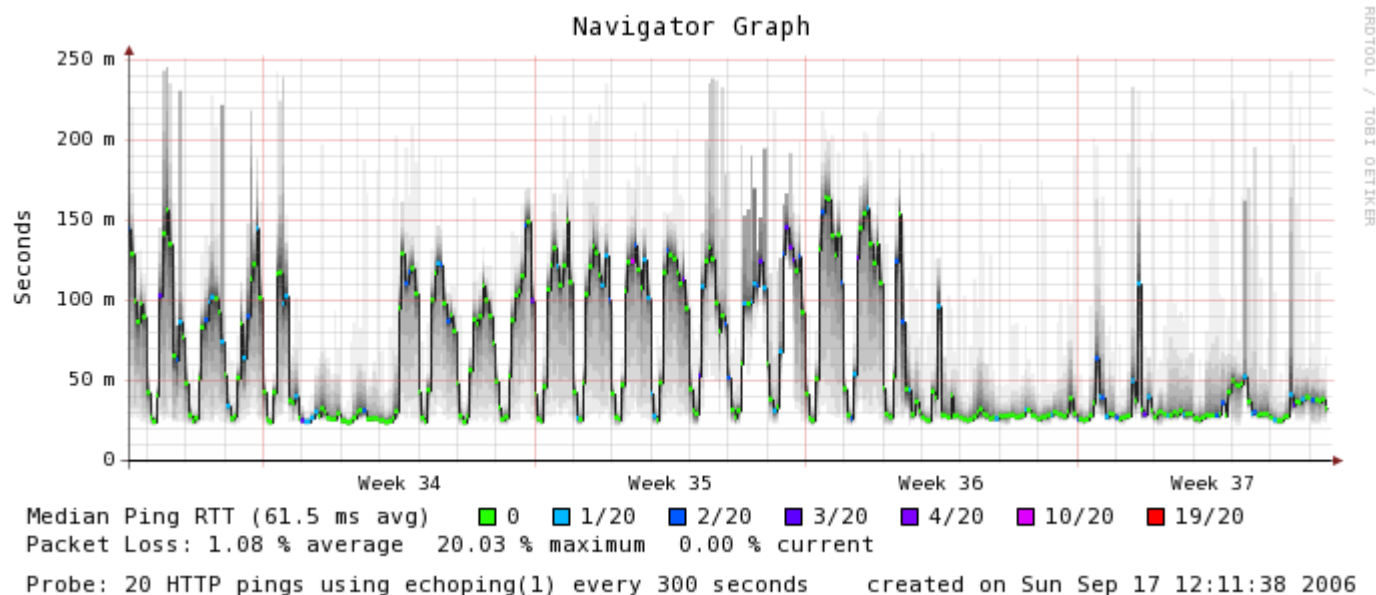
# Varnish Architecture

**Cluster Controller**

CmdLine

CLI-interface

Web-interface

CMS-interface

SMS-interface

*Manager*

CmdLine
ChildProcMgt
VCL compiler
Params/Args
Initialization
Watchdog

*Cacher*

CmdLine
Storage
Log/Stats
Accept/herder
Backend
Worker threads
Grim Reaper
Hashing

Shared Memory

logwriter
stats
ad-hoc

One binary program

C-compiler

Shared object

# Varnish is fast!

- 1000 cache hits @ 100 Mbit/sec:

| Sample | Min | Max | Median | Average | Stddev |
|---|---|---|---|---|---|
| Full | 12,9 | 3001 | 16,2 | 71,1 | 338,5 |
| 90% fractile | 12,9 | 26 | 15,9 | 16,3 | 1,7 |

(all times are in microseconds)



Navigator Graph

Median Ping RTT (61.5 ms avg) ■ 0 ■ 1/20 ■ 2/20 ■ 3/20 ■ 4/20 ■ 10/20 ■ 19/20
Packet Loss: 1.08 % average   20.03 % maximum   0.00 % current
Probe: 20 HTTP pings using echoping(1) every 300 seconds   created on Sun Sep 17 12:11:38 2006

# Varnish Coordinates

- **WWW.varnish-cache.org**
  - Yes, we need to work on the web-page.
- Coded by: Poul-Henning Kamp
  - phk@FreeBSD.org
- Runs on: Any resonable modern UNIX
  - FreeBSD, Linux etc.
- Open Source
  - BSD license

# Why I wrote Varnish

- I was sceptical when VG approached me with the Varnish Project.
    - Not really my specialty
    - Not really my preferred area
- On the other hand...
    - I'm tired of people writing lousy programs and blaming it on "my" kernel
    - A fine chance to educate by example
- Also, I could use the money.

# Performance Programming
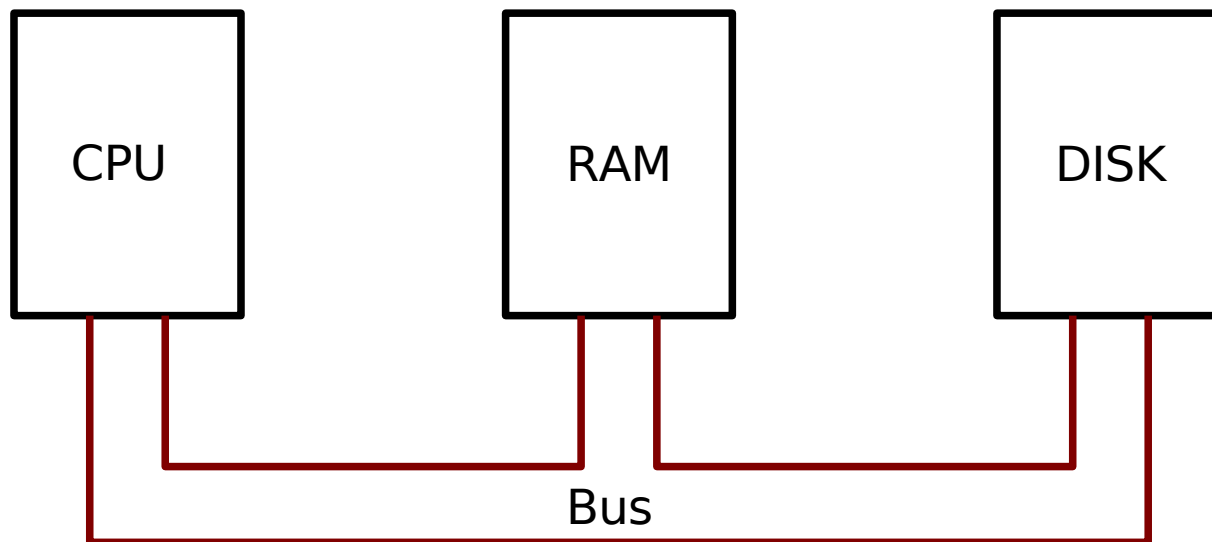
- Understand and exploit the envelope of your requirements

- Avoid unecessary work

- Use few cheap operations

- Use even fewer expensive operations

- Don't fight the kernel

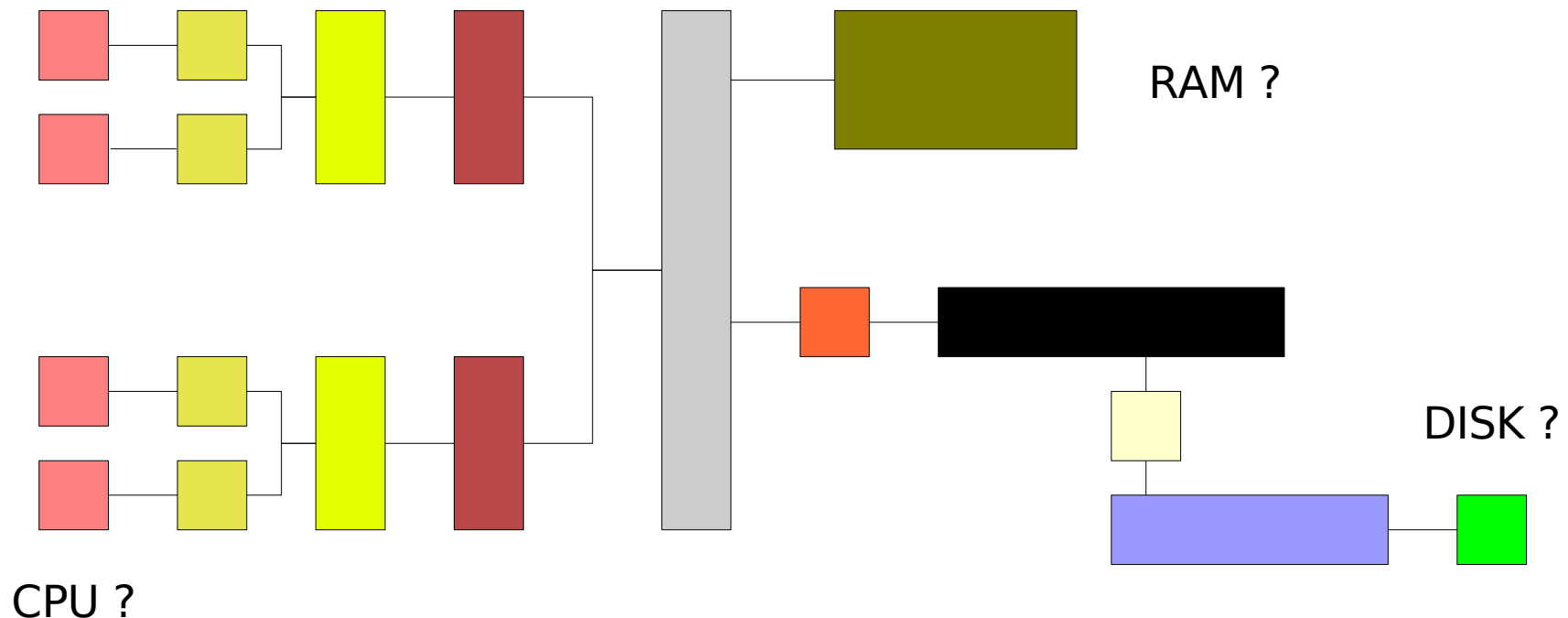- Use the features the nice kernel guys have given you

# People program like 1970

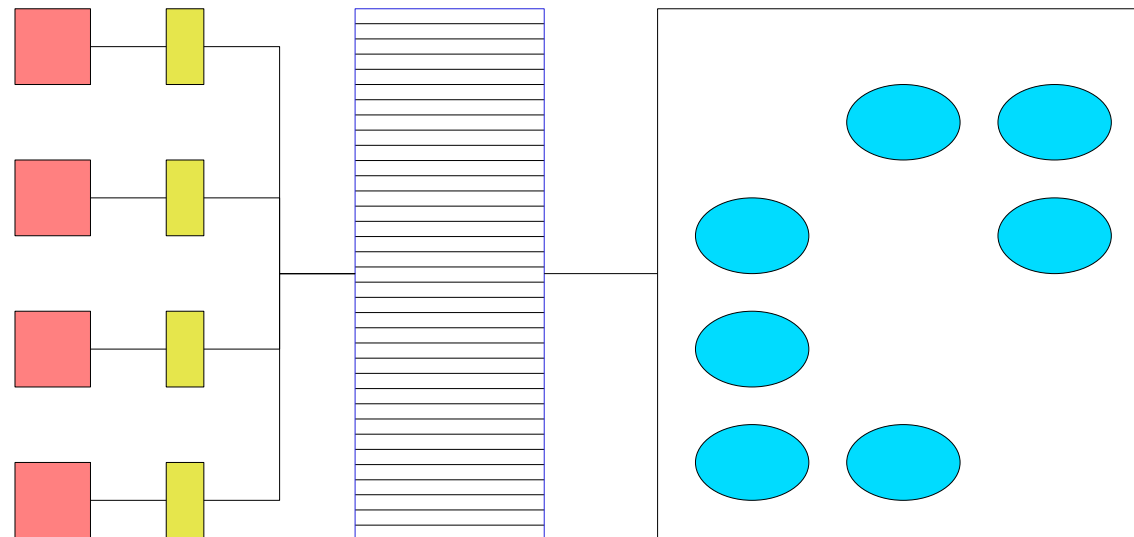- "A computer consists of a CPU, internal storage and external storage"

# 2006 reality

- "A Computer consists of a number of processing units ("cores") connected via a number of caches and busses to a virtualized storage."

RAM ?

DISK ?

CPU ?

# 2006 conceptual model

- ”A computer consists of a number of cores with private caches, a shared, page granularity, cache for external storage objects.”

# Squid is 1970 model

- Some objects "in Memory"

- Other objects "on disk"

- Explicit transfer between the two states:

  - Unused memory object gets written to disk

  - Disk object gets read to memory for transmission.

# Fight the kernel (1)

- Squid creates a memory object
- Gets used some, then falls into disuse
- Kernel runs short of RAM
- Kernel pages unused pages out to disk

# Fight the kernel (2)

- Squid decides memory object is unused and should be moved to disk

- Creates file, issues write(2)

- Page is not in RAM, generates Page-Fault

- Kernel reads data back from disk

kernel may need to page something out to make space

- Kernel fixes pagetable, continues squid

# Fight the kernel (3)

- Squids write(2) can write object to disk

- Squid can reuse memory for other object

# Fight the kernel (4)

- Squid needs first object back

- Push something to disk




- Read in first object back in

# Varnish

- Varnish puts object in virtual memory

- Object used some, then falls into disuse

- Kernel pages object out to free RAM for better use.

- Varnish access object in virtual memory

- Kernel page-faults, reads object from disk.
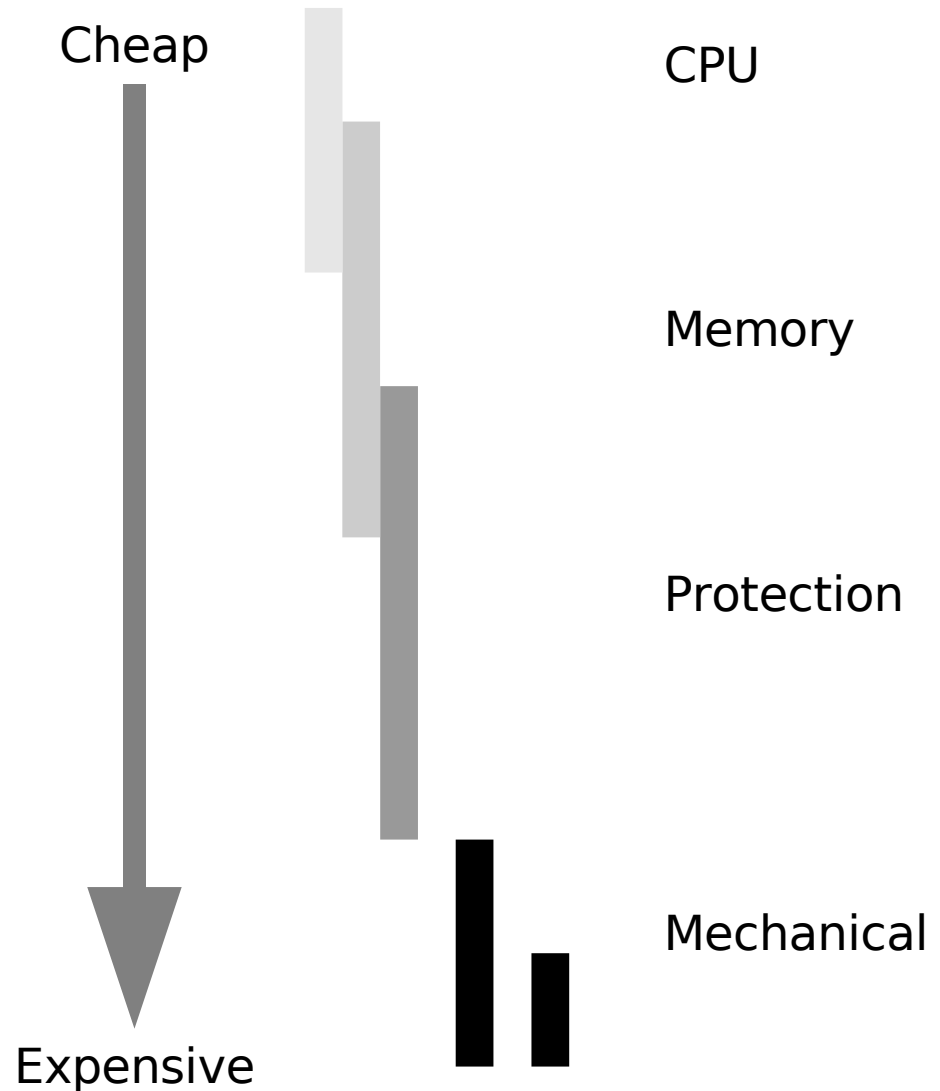
# Virtual memory shortage

- On 32 bit systems VM is limited to ~2GB

- If you have a webserver with a larger working set of content, you can afford to buy a new 64 bit CPU.

# Performance price List

- char *p += 5;
- strlen(p);
- memcpy(p, q, l);
- Locking
- System Call
- Context Switch
- Disk Access
- Filesystem operation

Cheap

Expensive

CPU

Memory

Protection

Mechanical

# How Cheap or Expensive ?

- CPU clock: 4 Ghz

- Multiple instructions per clock cycle

- Several layers of cache

- Conservative rule of thumb:

  **1 Billion instructions per second**

# Cost of a system call ?

- Context switch to kernel is expensive.
  - Thousands of clocks to save/restore registers
  - Argument validation
  - Page Table modifications

- 1 μsec = 0.000001s = 1,000 instructions

# Cost of disk I/O

- Enter kernel

- Navigate disk layout

  - Filesystem, partitions etc.

- Via device driver to hardware

- Busses, handshaking etc.

- Mechanical movement.

- 1 msec= 0.001s = 1,000,000 instructions

# Varnish logging

- Logging to shared memory
  - Practically Free (as in peanuts)
  - No slowdown for the <u>real</u> workload.

- Generate "real" output in separate processes, which "tail" the shmlog:
  - Apache format
  - Custom format
  - Realtime views

# Logging

- Logging to a file is expensive:

```
FILE *flog;                              Filesystem operation

flog = fopen("/var/log/mylog", "a");              Once only

[...]

fprintf(flog, "%s Something went wrong with %s\n",
    timestamp(), myrepresentation(object));
fflush(flog);
                        Disk I/O



                  Millons of calls
```

# Cheap logging

- Shared memory logging is cheap:

```
char *logp, *loge;

fd = open(...);
logp = mmap(..., size);
loge = logp + size;

[...]
LOCK(&shmlog);
logp[1] = LOG_ERROR;
logp[2] = sprintf(logp + 3,
    "Something went bad with %s",
    myrepresentation(obj));
logp[3 + logp[2]] = LOG_END;
logp[0] = LOG_ENTRY;
logp += 3 + logp[2];
UNLOCK(&shmlog);
```

Filesystem ops

Once only

Memory and arithmetic

Millions of calls.

# But it can be even cheaper

- strlen(3) and ASCIIZ format is not cheap
- Keep track of both ends of strings
- Use memcpy(3) instead of strcpy(3)
  - memcpy(3) can move data in 32 or even 64 bit chunks.
  - strcpy(3) **has** to use 8 bit chunks.

# But it can be **even** cheaper

- Add a shmlog buffer to worker threads
- Logging to this buffer is lockless
- Batch "commit" local buffer to "real" shmlog at synchronization points.
  - Cost: One lock + one memcpy().

- => Practically no mutex contention.

# …and

- Coredumps contains the most recent log records

- We can afford to log much more detail, it does not take up disk-space

- Approx 30 records per request
  - =>100.000's records/sec

# "varnishtop"

- Example of ad-hoc shmlog reader:

- ./varnishtop -i sessionclose

- Why do sessions close ?

```
132232.12 timeout
 39002.59 EOF
  5232.97 not HTTP/1.1
  3069.28 Connection: close
  1383.40 remote closed
   630.93 silent
    15.49 pipe
     9.15 Not in cache.
     6.36 Purged.
     5.41 no request
     0.98 Bad Request
```

# Varnishtop

- What is my most popular URL ?

```
1304.86 /tmv11.js
 989.08 /sistenytt.html
 495.05 /include/global/art.js
 491.01 /css/hoved.css
 490.05 /gfk/ann/n.gif
 480.08 /gfk/ann/ng.gif
 468.12 /gfk/front/tipsvg.png
 352.66 /css/ufront.css
 317.75 /t.gif
 306.79 /gfk/plu2.gif
 298.84 /css/front.css
 292.84 /gfk/min2.gif
 280.94 /css/blog.css
 279.84 /
```

# Where does traffic come from ?

```
33913.74 Referer: http://www.vg.no/
 4730.72 Referer: http://vg.no/
  925.62 Referer: http://www.vg.no
  510.10 Referer: http://www.vg.no/pub/vgart.hbs?art
  434.37 Referer: http://www.vg.no/export/Transact/m
  349.55 Referer: http://www.vg.no/pub/vgart.hbs?art
  344.66 Referer: http://www.vg.no/pub/vgart.hbs?art
  324.06 Referer: http://www.vg.no/export/Transact/t
  297.25 Referer: http://www.nettby.no/user/
  263.82 Referer: http://www.vg.no/sport/fotball/
  242.55 Referer: http://www.vg.no/pub/vgart.hbs?art
```

# Statistics

- Same story

- Stored in shared memory

- Programs can monitor & present data

- No system calls necessary to get up to date numbers

# Real-time statistics

```
16:23:13
Hitrate ratio:            9         9         9
Hitrate avg:         0.9986    0.9986    0.9986

   17772105         435.55          301.26 Client connections accepted
  130213161        3623.22         2207.26 Client requests received
  129898315        3617.23         2201.93 Cache hits
      85043           0.00            1.44 Cache hits for pass
     227180           4.99            3.85 Cache misses
     313630           4.99            5.32 Backend connections initiated
        439           0.00            0.01 Backend connections recyles
         54           0.00            0.00 Backend connections unused
       6196           1.00            0.11 N struct srcaddr
       1656         -24.97            0.03 N active struct srcaddr
       3222           0.00            0.05 N struct sess_mem
       2258         -51.95            0.04 N struct sess
      65685           5.99            1.11 N struct object
      65686           5.99            1.11 N struct objecthead
```

# Memory Strategy

- How long is a HTTP header ?
  - 80 char ? 256 char ?  2048 char ?

- Classic strategy:
  - Allocate small, extend with realloc(3)

- Realloc(3) to more memory is expensive
  - Needs memory copy 99% of the times.

# Memory Strategy

Special offer!

Virtually **FREE** memory

offer good from 1980 onwards

# Memory Strategy

- Allocate enough memory for 99% of the cases up front

- Unused & untouched memory is free
  - as in "peanuts"

- If data is long lived, trim back with realloc(3)

# Memory management

- During the handling of a request we need various bits of memory

- Use a private pool
  - Keeps things close together (good caching)
  - Avoids locking against other CPUs

- Free everything with a single pointer assignment at the end
  - As cheap as can be

# Memory/Lock interaction

- Do minimum work while holding lock

- Allocate memory before grabbing lock

- If you didn't need it anyway:
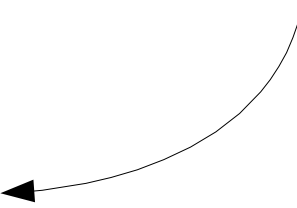
  - cache it until next time

# Preallocation

```
LOCK(storage_lock);
TAILQ_FOREACH(...) {
        if (...)
                break;  /* found */
}
if (ptr == NULL) {
        ptr = malloc(sizeof *wp->store);
        assert(wp->store != NULL);
        [...] // fill ptr->stuff
        TAILQ_INSERT_TAIL(...)
}
UNLOCK(storage_lock)
```

Malloc(3)

May need to synchronize between CPUs with locks

May need to get memory from kernel with syscall
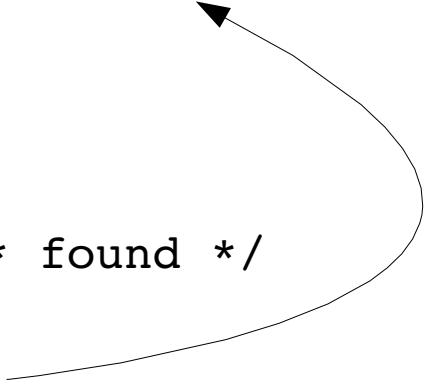
# Preallocation

```
if (wp->store == NULL) {
        wp->store = malloc(sizeof *wp->store);
        assert(wp->store != NULL);
}

LOCK(storage_lock);
TAILQ_FOREACH(...) {
        if (...)
                break;  /* found */
}
if (ptr == NULL) {
        ptr = wp->store;
        wp->store = NULL
        [...] // fill ptr->stuff
        TAILQ_INSERT_TAIL(...)
}
UNLOCK(storage_lock)
```

malloc call moved out
of locked code section.

Lock held shorter time.

# Avoid unnecessary work

- Don't copy data when we can avoid it
  - Receive headers into one chunk of memory and leave them there

- Don't text-process all HTTP headers
  - We only care about few of them

- Transport headers by reference
  - Use scatter/gather I/O (see writev(2) syscall).

# Avoid strlen(3)

- `find_hdr(sp->http, "Host:")`
  - find_hdr() must do strlen(3) on literal constant to find out how long it is.
- `find_hdr(sp->http, "Host:", 5)`
  - Better, but programmer has to get it right.
- `#define HTTP_HOST "\005Host:"`
  `find_hdr(sp->http, HTTP_HOST);`
  - Even better, programmer must only get it right once.

# Remember the caches

- Reuse most recently used resource.

  - It is more likely to be in cache.

- Process, Thread, Memory, File, Directory, Disk-block etc.

- Avoid round-robin scheduling.

  - Use "Last-In-First-Out" (ie: stack)
  - Not "First-In-First-Out" (ie: queue)

# Measure & Benchmark

- Understand how often things happen
  - Use tcov(1) or add counters.
- Check your assumptions
  - Use ktrace(8) to see the system-calls.
- Study system statistics
  - Netstat(8), Vmstat(8), sysctl(8)
- Time things if in doubt
  - use clock_gettime(3) or RDTSC()