# xh

Spencer Tipping

June 8, 2014

# Contents

# Part I

# xh runtime

# Chapter 1

# Self-replication

boot/xh-header

```perl
1   #!/usr/bin/env perl
2   BEGIN {eval(our $xh_bootstrap = q{
3   # xh | https://github.com/spencertipping/xh
4   # Copyright (C) 2014, Spencer Tipping
5   # Licensed under the terms of the MIT source code license
6
7   # For the benefit of HTML viewers (long story):
8   # <body style='display:none'>
9   # <script src='http://spencertipping.com/xh/page.js'></script>
10  use 5.014;
11  package xh;
12  our %modules;
13  our @module_ordering;
14  our %eval_numbers = (1 => '$xh_bootstrap');
15
16  sub with_eval_rewriting(&) {
17    my @result = eval {$_[0]->(@_[1..$#_])};
18    $@ =~ s/\(eval (\d+)\)/$eval_numbers{$1}/eg if $@;
19    die $@ if $@;
20    @result;
21  }
22
23  sub named_eval {
24    my ($name, $code) = @_;
25    $eval_numbers{$1 + 1} = $name if eval('__FILE__') =~ /\(eval (\d+)\)/;
26    with_eval_rewriting {eval $code};
27  }
28
29  our %compilers = (pl => sub {
```

```
30    my $package = $_[0] =~ s/\./::/gr;
31    named_eval $_[0], "{package ::$package;\n$_[1]\n}";
32    die "error compiling module $_[0]: $@" if $@;
33  });
34
35  sub defmodule {
36    my ($name, $code, @args) = @_;
37    chomp($modules{$name} = $code);
38    push @module_ordering, $name;
39    my ($base, $extension) = split /\.(\w+$)/, $name;
40    die "undefined module extension '$extension' for $name"
41      unless exists $compilers{$extension};
42    $compilers{$extension}->($base, $code, @args);
43  }
44
45  chomp($modules{bootstrap} = $::xh_bootstrap);
46  undef $::xh_bootstrap;
```

At this point we need a way to reproduce the image. Since the bootstrap code is already stored, we can just wrap it and each defined module into an appropriate BEGIN block.

Listing 1.2   boot/xh-header (continued)

```
1   sub image {
2     my @pieces = "#!/usr/bin/env perl";
3     push @pieces, "BEGIN {eval(our \$xh_bootstrap = <<'_')}",
4                   $modules{bootstrap},
5                   '_';
6     push @pieces, "BEGIN {xh::defmodule('$_', <<'_')}",
7                   $modules{$_},
8                   '_' for @module_ordering;
9     push @pieces, "xh::main::main;\n__DATA__";
10    join "\n", @pieces;
11  }
12  })}
```

# Chapter 2

# SSH routing fabric

xh does all of its distributed communication over SSH stdin/stdout tunnels (since remote hosts may have port forwarding disabled), which means that we need to implement a datagram format, routing logic, and a priority-aware traffic scheduler.

For simplicity, the only session type that's supported is RPC. The request and response each must fit into a single packet, which is size-limited to 64 kilobytes excluding the packet header. The fabric client will deal with larger requests and responses, but it will cause additional round-trips.

Listing 2.1  src/fabric.pl

```perl
1  package xh::fabric;
2  use Sys::Hostname;
3  use Time::HiRes qw/time/;
4  use Digest::SHA qw/sha256/;
5  # Mutable state space definition for the routing fabric. You should create
6  # one of these for every separate xh network you plan to interface with.
7
8  sub fabric_client {
9    my ($name, $bindings) = @_;
10   $name //= $ENV{USER} . '@' . hostname . '.local';
11   return {rpc_bindings     => $bindings,
12           instance_name    => $name,
13           instance_id      => 0,
14           edge_pipes       => {},
15           network_topology => {},
16           send_queue       => [],
17           blocked_rpcs     => {},
18           routing_cache    => {}};
19  }
20
21  sub fabric_rpc_bind {
```

```perl
22    my ($state, %bindings) = @_;
23    my $bindings = $state->{rpc_bindings};
24    $bindings->{$_} = $bindings{$_} for keys %bindings;
25    $state;
26  }
27
28  sub fabric_is_initialized {
29    my ($state) = @_;
30    !!$state->{instance_id};
31  }
32  use constant header_pack_format        => 'C32 N N d n C C n';
33  use constant signed_header_pack_format => 'C32 ' . header_pack_format;
34
35  use constant signed_header_length    => 32 + 32+4+4+8+2+1+1+2;
36  use constant header_signature_length => 32;
37
38  our $nonce_state = sha256(time . hostname);
39  sub packet_nonce {$nonce_state = sha256(time . $nonce_state)}
40
41  sub encode_packet {
42    my ($state, $destination_name, $message_type, $priority, $deadline) = @_;
43    die "data is too long: " . length($_[5]) . " (max is 65535 bytes)"
44      if length $_[5] >= 65536;
45
46    my $destination_id = $state->{routing_cache}{$destination_name};
47    return undef unless defined $destination_id;
48
49    my $header = pack header_pack_format, packet_nonce,
50                                          $state->{instance_id},
51                                          $destination_id->{endpoint_id},
52                                          time,
53                                          length $_[5],
54                                          $message_type,
55                                          $priority,
56                                          $deadline;
57    my $packet = $header . $_[5];
58    sha256($packet) . $packet;
59  }
60
61  sub decode_packet_header {
62    unpack signed_header_pack_format, $_[0];
63  }
64
65  sub signature_is_valid {
66    my ($sha) = decode_packet_header $_[0];
67    $sha eq sha256(substr $_[0], header_signature_length);
```

```perl
68  }
69  use constant {forgetful_rpc  => 0,
70                functional_rpc => 1,
71                rpc_reply      => 2,
72                rpc_error      => 3,
73                routing_error  => 4};
74  use constant {realtime_priority => 0,
75                high_priority     => 16,
76                normal_priority   => 256,
77                low_priority      => 32768};
78
79  use constant {realtime_deadline          => 0,
80                imperceptible_deadline     => 20,
81                short_interactive_deadline => 50,
82                long_interactive_deadline  => 100,
83                process_blocking_deadline  => 250,
84                background_deadline        => 2000,
85                far_deadline               => 32768};
86
87
88
89  \lstset{caption={src/fabric/dependencies.pl},name={src/fabric/dependencies.pl}}\begin{perlcode}
90  use Sys::Hostname;
91  use Time::HiRes qw/time/;
92  use Digest::SHA qw/sha256/;
```

**Listing 2.2**  src/fabric/state.pl

```perl
1   # Mutable state space definition for the routing fabric. You should create
2   # one of these for every separate xh network you plan to interface with.
3
4   sub fabric_client {
5     my ($name, $bindings) = @_;
6     $name //= $ENV{USER} . '@' . hostname . '.local';
7     return {rpc_bindings     => $bindings,
8             instance_name    => $name,
9             instance_id      => 0,
10            edge_pipes       => {},
11            network_topology => {},
12            send_queue       => [],
13            blocked_rpcs     => {},
14            routing_cache    => {}};
15  }
16
17  sub fabric_rpc_bind {
18    my ($state, %bindings) = @_;
```

```
19    my $bindings = $state->{rpc_bindings};
20    $bindings->{$_} = $bindings{$_} for keys %bindings;
21    $state;
22  }
23
24  sub fabric_is_initialized {
25    my ($state) = @_;
26    !!$state->{instance_id};
27  }
```

## 2.1  Packet format

Packets and headers are written in binary, and all multibyte numbers are big-endian. The structure of a packet is:

```
data+header SHA-256:      32 bytes
packet identity nonce:    32 bytes           \
source xh instance ID:    4 bytes            |
destination xh instance ID: 4 bytes          |
packet creation time:     8 bytes (double) |
data length:              2 bytes            | SHA applies to these bytes
message type:             1 byte             |
priority:                 1 byte             |
deadline:                 2 bytes            |
data:                     <= 65535 bytes   /
```

The only reason we represent packet creation time as a double rather than as a 64-bit integer is that 64-bit integer support is not guaranteed within Perl. As a result, we have a somewhat awkward situation where all absolute times are encoded as doubles and all deltas as integers.

Listing 2.3   src/fabric/packet-format.pl

```
1   use constant header_pack_format        => 'C32 N N d n C C n';
2   use constant signed_header_pack_format => 'C32 ' . header_pack_format;
3
4   use constant signed_header_length    => 32 + 32+4+4+8+2+1+1+2;
5   use constant header_signature_length => 32;
6
7   our $nonce_state = sha256(time . hostname);
8   sub packet_nonce {$nonce_state = sha256(time . $nonce_state)}
9
10  sub encode_packet {
11    my ($state, $destination_name, $message_type, $priority, $deadline) = @_;
12    die "data is too long: " . length($_[5]) . " (max is 65535 bytes)"
13      if length $_[5] >= 65536;
```

8

```
14
15    my $destination_id = $state->{routing_cache}{$destination_name};
16    return undef unless defined $destination_id;
17
18    my $header = pack header_pack_format, packet_nonce,
19                                          $state->{instance_id},
20                                          $destination_id->{endpoint_id},
21                                          time,
22                                          length $_[5],
23                                          $message_type,
24                                          $priority,
25                                          $deadline;
26    my $packet = $header . $_[5];
27    sha256($packet) . $packet;
28 }
29
30 sub decode_packet_header {
31    unpack signed_header_pack_format, $_[0];
32 }
33
34 sub signature_is_valid {
35    my ($sha) = decode_packet_header $_[0];
36    $sha eq sha256(substr $_[0], header_signature_length);
37 }
```

message type is one of the following values:

0  Forgetful RPC request. The receiver should execute the code, but the
   sender will not await a reply. This is used internally by xh to maintain
   routing graph information and clock offsets.

1  Functional RPC request. This indicates that the receiver should execute
   the given code, encoded as text, and send a reply. The code may contain
   references that require further RPCs to be issued.

2  RPC reply after a successful invocation. The return value of the function
   is encoded in quoted form, and may require further dereferencing via
   RPC.

3  Callee-side RPC error; the reply is a partially-evaluated quoted value,
   where any unevaluated pieces represent errors.

4  Routing error or timeout; the routing fabric generates this to indicate that
   it has given up on getting a successful reply. If this happens, the sender
   will automatically re-send the RPC unless the deadline has expired.

**Listing 2.4**  `src/fabric/message-types.pl`

```
1  use constant {forgetful_rpc  => 0,
2                 functional_rpc => 1,
3                 rpc_reply      => 2,
4                 rpc_error      => 3,
5                 routing_error  => 4};
```

`priority` and `deadline` are used for scheduling purposes. Zero is the highest priority, 65535 is the lowest. The deadline is used to indicate how time-sensitive the packet is; the queueing order function used by the scheduler is $\frac{2^c}{s}$, where:

$$c = \frac{\Delta t - d}{16 + p}$$

$\Delta t$ = ms since packet was originally sent

$d$ = the deadline

$p$ = the priority

$s$ = header + data size in bytes

$\Delta t$ is an estimated quantity, since hosts will not, in general, have synchronized clocks. However, xh uses a protocol similar to NTP to estimate clock offsets for each instance. These clock offsets are used to coordinate instances on different hosts. (See 2.3.)

Listing 2.5  `src/fabric/standard-priorities.pl`

```
1  use constant {realtime_priority => 0,
2                 high_priority      => 16,
3                 normal_priority    => 256,
4                 low_priority       => 32768};
5
6  use constant {realtime_deadline           => 0,
7                 imperceptible_deadline      => 20,
8                 short_interactive_deadline => 50,
9                 long_interactive_deadline  => 100,
10                process_blocking_deadline  => 250,
11                background_deadline         => 2000,
12                far_deadline                => 32768};
```

## 2.2 Routing logic

I assume the topology of xh instances will fit into memory. This won't be a problem for most installations; in practice, xh should be able to easily manage (and transfer data between) many hundreds of machines without slowing down. Each xh instance maintains a copy of the full routing graph, which includes information about edge timings.

The routing logic's job is to decide how to most effectively get packets from point A to point B, which, more formally, means minimizing the expected sum of delay costs. Doing this well involves a few factors:

1. An edge's average latency and throughput.

2. The variance in an edge's latency and throughput, absent xh traffic.

3. The impact of traffic on an edge's latency and throughput.

All of these are continuously measured and periodically propagated as network topology metadata.

Listing 2.6  `src/fabric/routing-rpcs.pl`

## 2.3   Clock offset estimation