```
                    ------------------------
                    HAProxy Starter Guide
                    ------------------------

                          version 1.6
```

This document is an introduction to HAProxy for all those who don't know it, as
well as for those who want to re-discover it when they know older versions. Its
primary focus is to provide users with all the elements to decide if HAProxy is
the product they're looking for or not. Advanced users may find here some parts
of solutions to some ideas they had just because they were not aware of a given
new feature. Some sizing information are also provided, the product's lifecycle
is explained, and comparisons with partially overlapping products are provided.

This document doesn't provide any configuration help nor hint, but it explains
where to find the relevant documents. The summary below is meant to help you
search sections by name and navigate through the document.

Note to documentation contributors :
    This document is formatted with 80 columns per line, with even number of
    spaces for indentation and without tabs. Please follow these rules strictly
    so that it remains easily printable everywhere. If you add sections, please
    update the summary below for easier searching.

1. Available documentation
--------------------------

The complete HAProxy documentation is contained in the following documents.
Please ensure to consult the relevant documentation to save time and to get the
most accurate response to your needs. Also please refrain from sending questions
to the mailing list whose responses are present in these documents.

  - intro.txt (this document) : it presents the basics of load balancing,
    HAProxy as a product, what it does, what it doesn't do, some known traps to
    avoid, some OS-specific limitations, how to get it, how it evolves, how to
    ensure you're running with all known fixes how to update it, complements and
    alternatives.

  - management.txt : it explains how to start haproxy, how to manage it at
    runtime, how to manage it on multiple nodes, how to proceed with seamless
    upgrades.

  - configuration.txt : the reference manual details all configuration keywords
    and their options. It is used when a configuration change is needed.

  - architecture.txt : the architecture manual explains how to best architect a
    load-balanced infrastructure and how to interact with third party products.

  - coding-style.txt : this is for developers who want to propose some code to
    the project. It explains the style to adopt for the code. It's not very
    strict and not all the code base completely respects it but contributions
    which diverge too much from it will be rejected.

  - proxy-protocol.txt : this is the de-facto specification of the PROXY
    protocol which is implemented by HAProxy and a number of third party
    products.

  - README : how to build haproxy from sources

2. Quick introduction to load balancing and load balancers
----------------------------------------------------------

Load balancing consists in aggregating multiple components in order to achieve
a total processing capacity above each component's individual capacity, without
any intervention from the end user and in a scalable way. This results in more
operations being performed simultaneously by the time it takes a component to
perform only one. A single operation however will still be performed on a single
component at a time and will not get faster than without load balancing. It
always requires at least as many operations as available components and an
efficient load balancing mechanism to make use of all components and to fully
benefit from the load balancing. A good example of this is the number of lanes
on a highway which allows as many cars to pass during the same time frame
without increasing their individual speed.

Examples of load balancing :

  - Process scheduling in multi-processor systems
  - Link load balancing (eg: EtherChannel, Bonding)
  - IP address load balancing (eg: ECMP, DNS roundrobin)
  - Server load balancing (via load balancers)

The mechanism or component which performs the load balancing operation is
called a load balancer. In web environments these components are called a
"network load balancer", and more commonly a "load balancer" given that this
activity is by far the best known case of load balancing.

A load balancer may act :

131    - at the link level : this is called link load balancing, and it consists in
132      chosing what network link to send a packet to;
133
134    - at the network level : this is called network load balancing, and it
135      consists in chosing what route a series of packets will follow;
136
137
138    - at the server level : this is called server load balancing and it consists
139      in deciding what server will process a connection or request.
140
141    Two distinct technologies exist and address different needs, though with some
142    overlapping. In each case it is important to keep in mind that load balancing
143    consists in diverting the traffic from its natural flow and that doing so always
144    requires a minimum of care to maintain the required level of consistency between
145    all routing decisions.
146
147    The first one acts at the packet level and processes packets more or less
148    individually. There is a 1-to-1 relation between input and output packets, so
149    it is possible to follow the traffic on both sides of the load balancer using a
150    regular network sniffer. This technology can be very cheap and extremely fast.
151    It is usually implemented in hardware (ASICs) allowing to reach line rate, such
152    as switches doing ECMP. Usually stateless, it can also be stateful (consider
153    the session a packet belongs to and called layer4-LB or L4), may support DSR
154    (direct server return, without passing through the LB again) if the packets
155    were not modified, but provides almost no content awareness. This technology is
156    very well suited to network-level load balancing, though it is sometimes used
157    for very basic server load balancing at high speed.
158
159    The second one acts on session contents. It requires that the input streams is
160    reassembled and processed as a whole. The contents may be modified, and the
161    output stream is segmented into new packets. For this reason it is generally
162    performed by proxies and they're often called layer 7 load balancers or L7.
163    This implies that there are two distinct connections on each side, and that
164    there is no relation between input and output packets sizes nor counts. Clients
165    and servers are not required to use the same protocol (for example IPv4 vs
166    IPv6, clear vs SSL). The operations are always stateful, and the return traffic
167    must pass through the load balancer. The extra processing comes with a cost so
168    it's not always possible to achieve line rate, especially with small packets.
169    On the other hand, it offers wide possibilities and is generally achieved by
170    pure software, even if embedded into hardware appliances. This technology is
171    very well suited for server load balancing.
172
173    Packet-based load balancers are generally deployed in cut-through mode, so they
174    are installed on the normal path of the traffic and divert it according to the
175    configuration. The return traffic doesn't necessarily pass through the load
176    balancer. Some modifications may be applied to the network destination address
177    in order to direct the traffic to the proper destination. In this case, it is
178    mandatory that the return traffic passes through the load balancer. If the
179    routes doesn't make this possible, the load balancer may also replace the
180    packets' source address with its own in order to force the return traffic to
181    pass through it.
182
183    Proxy-based load balancers are deployed as a server with their own IP address
184    and ports, without architecture changes. Sometimes this requires to perform some
185    adaptations to the applications so that clients are properly directed to the
186    load balancer's IP address and not directly to the server's. Some load balancers
187    may have to adjust some servers' responses to make this possible (eg: the HTTP
188    Location header field used in HTTP redirects). Some proxy-based load balancers
189    may intercept traffic for an address they don't own, and spoof the client's
190    address when connecting to the server. This allows them to be deployed as if
191    they were a regular router or firewall, in a cut-through mode very similar to
192    the packet based load balancers. This is particularly appreciated for products
193    which combine both packet mode and proxy mode. In this case DSR is obviously
194    still not possible and the return traffic still has to be routed back to the
195    load balancer.

196    A very scalable layered approach would consist in having a front router which
197    receives traffic from multiple load balanced links, and uses ECMP to distribute
198    this traffic to a first layer of multiple stateful packet-based load balancers
199    (L4). These L4 load balancers in turn pass the traffic to an even larger number
200    of proxy-based load balancers (L7), which have to parse the contents to decide
201    what server will ultimately receive the traffic.
202
203
204    The number of components and possible paths for the traffic increases the risk
205    of failure; in very large environments, it is even normal to permanently have
206    a few faulty components being fixed or replaced. Load balancing done without
207    awareness of the whole stack's health significantly degrades availability. For
208    this reason, any sane load balancer will verify that the components it intends
209    to deliver the traffic to are still alive and reachable, and it will stop
210    delivering traffic to faulty ones. This can be achieved using various methods.
211
212    The most common one consists in periodically sending probes to ensure the
213    component is still operational. These probes are called "health checks". They
214    must be representative of the type of failure to address. For example a ping-
215    based check will not detect that a web server has crashed and doesn't listen to
216    a port anymore, while a connection to the port will verify this, and a more
217    advanced request may even validate that the server still works and that the
218    database it relies on is still accessible. Health checks often involve a few
219    retries to cover for occasional measuring errors. The period between checks
220    must be small enough to ensure the faulty component is not used for too long
221    after an error occurs.
222
223    Other methods consist in sampling the production traffic sent to a destination
224    to observe if it is processed correctly or not, and to evince the components
225    which return inappropriate responses. However this requires to sacrify a part
226    of the production traffic and this is not always acceptable. A combination of
227    these two mechanisms provides the best of both worlds, with both of them being
228    used to detect a fault, and only health checks to detect the end of the fault.
229    A last method involves centralized reporting : a central monitoring agent
230    periodically updates all load balancers about all components' state. This gives
231    a global view of the infrastructure to all components, though sometimes with
232    less accuracy or responsiveness. It's best suited for environments with many
233    load balancers and many servers.
234
235    Layer 7 load balancers also face another challenge known as stickiness or
236    persistence. The principle is that they generally have to direct multiple
237    subsequent requests or connections from a same origin (such as an end user) to
238    the same target. The best known example is the shopping cart on an online
239    store. If each click leads to a new connection, the user must always be sent
240    to the server which holds his shopping cart. Content-awareness makes it easier
241    to spot some elements in the request to identify the server to deliver it to,
242    but that's not always enough. For example if the source address is used as a
243    key to pick a server, it can be decided that a hash-based algorithm will be
244    used and that a given IP address will always be sent to the same server based
245    on a divide of the address by the number of available servers. But if one
246    server fails, the result changes and all users are suddenly sent to a different
247    server and lose their shopping cart. The solution against this issue consists
248    in memorizing the chosen target so that each time the same visitor is seen,
249    he's directed to the same server regardless of the number of available servers.
250    The information may be stored in the load balancer's memory, in which case it
251    may have to be replicated to other load balancers if it's not alone, or it may
252    be stored in the client's memory using various methods provided that the client
253    is able to present this information back with every request (cookie insertion,
254    redirection to a sub-domain, etc). This mechanism provides the extra benefit of
255    not having to rely on unstable or unevenly distributed information (such as the
256    source IP address). This is in fact the strongest reason to adopt a layer 7
257    load balancer instead of a layer 4 one.
258
259    In order to extract information such as a cookie, a host header field, a URL
260    or whatever, a load balancer may need to decrypt SSL/TLS traffic and even

possibly to reencrypt it when passing it to the server. This expensive task
explains why in some high-traffic infrastructures, sometimes there may be a
lot of load balancers.

Since a layer 7 load balancer may perform a number of complex operations on the
traffic (decrypt, parse, modify, match cookies, decide what server to send to,
etc), it can definitely cause some trouble and will very commonly be accused of
being responsible for a lot of trouble that it only revealed. Often it will be
discovered that servers are unstable and periodically go up and down, or for
web servers, that they deliver pages with some hard-coded links forcing the
clients to connect directly to one specific server without passing via the load
balancer, or that they take ages to respond under high load causing timeouts.
That's why logging is an extremely important aspect of layer 7 load balancing.
Once a trouble is reported, it is important to figure if the load balancer took
a wrong decision and if so why so that it doesn't happen anymore.

3. Introduction to HAProxy
--------------------------

HAProxy is written "HAProxy" to designate the product, "haproxy" to designate
the executable program, software package or a process, though both are commonly
used for both purposes, and is pronounced H-A-Proxy. Very early it used to stand
for "high availability proxy" and the name was written in two separate words,
though by now it means nothing else than "HAProxy".

3.1. What HAProxy is and is not
-------------------------------

HAProxy is :

  - a TCP proxy : it can accept a TCP connection from a listening socket,
    connect to a server and attach these sockets together allowing traffic to
    flow in both directions;

  - an HTTP reverse-proxy (called a "gateway" in HTTP terminology) : it presents
    itself as a server, receives HTTP requests over connections accepted on a
    listening TCP socket, and passes the requests from these connections to
    servers using different connections.

  - an SSL terminator / initiator / offloader : SSL/TLS may be used on the
    connection coming from the client, on the connection going to the server,
    or even on both connections.

  - a TCP normalizer : since connections are locally terminated by the operating
    system, there is no relation between both sides, so abnormal traffic such as
    invalid packets, flag combinations, window advertisements, sequence numbers,
    incomplete connections (SYN floods), or so will not be passed to the other
    side. This protects fragile TCP stacks from protocol attacks, and also
    allows to optimize the connection parameters with the client without having
    to modify the servers' TCP stack settings.

  - an HTTP normalizer : when configured to process HTTP traffic, only valid
    complete requests are passed. This protects against a lot of protocol-based
    attacks. Additionally, protocol deviations for which there is a tolerance
    in the specification are fixed so that they don't cause problem on the
    servers (eg: multiple-line headers).

  - an HTTP fixing tool : it can modify / fix / add / remove / rewrite the URL
    or any request or response header. This helps fixing interoperability issues
    in complex environments.

  - a content-based switch : it can consider any element from the request to
    decide what server to pass the request or connection to. Thus it is possible

to handle multiple protocols over a same port (eg: http, https, ssh).

  - a server load balancer : it can load balance TCP connections and HTTP
    requests. In TCP mode, load balancing decisions are taken for the whole
    connection. In HTTP mode, decisions are taken per request.

  - a traffic regulator : it can apply some rate limiting at various points,
    protect the servers against overloading, adjust traffic priorities based on
    the contents, and even pass such information to lower layers and outer
    network components by marking packets.

  - a protection against DDoS and service abuse : it can maintain a wide number
    of statistics per IP address, URL, cookie, etc and detect when an abuse is
    happening, then take action (slow down the offenders, block them, send them
    to outdated contents, etc).

  - an observation point for network troubleshooting : due to the precision of
    the information reported in logs, it is often used to narrow down some
    network-related issues.

  - an HTTP compression offloader : it can compress responses which were not
    compressed by the server, thus reducing the page load time for clients with
    poor connectivity or using high-latency, mobile networks.

HAProxy is not :

  - an explicit HTTP proxy, ie, the proxy that browsers use to reach the
    internet. There are excellent open-source software dedicated for this task,
    such as Squid. However HAProxy can be installed in front of such a proxy to
    provide load balancing and high availability.

  - a caching proxy : it will return as-is the contents its received from the
    server and will not interfere with any caching policy. There are excellent
    open-source software for this task such as Varnish. HAProxy can be installed
    in front of such a cache to provide SSL offloading, and scalability through
    smart load balancing.

  - a data scrubber : it will not modify the body of requests nor responses.

  - a web server : during startup, it isolates itself inside a chroot jail and
    drops its privileges, so that it will not perform any single file-system
    access once started. As such it cannot be turned into a web server. There
    are excellent open-source software for this such as Apache or Nginx, and
    HAProxy can be installed in front of them to provide load balancing and
    high availability.

  - a packet-based load balancer : it will not see IP packets nor UDP datagrams,
    will not perform NAT or even less DSR. These are tasks for lower layers.
    Some kernel-based components such as IPVS (Linux Virtual Server) already do
    this pretty well and complement perfectly with HAProxy.

3.2. How HAProxy works
----------------------

HAProxy is a single-threaded, event-driven, non-blocking engine combining a very
fast I/O layer with a priority-based scheduler. As it is designed with a data
forwarding goal in mind, its architecture is optimized to move data as fast as
possible with the least possible operations. As such it implements a layered
model offering bypass mechanisms at each level ensuring data don't reach higher
levels when not needed. Most of the processing is performed in the kernel, and
HAProxy does its best to help the kernel do the work as fast as possible by
giving some hints or by avoiding certain operation when it guesses they could
be grouped later. As a result, typical figures show 15% of the processing time
spent in HAProxy versus 85% in the kernel in TCP or HTTP close mode, and about

```
391  30% for HAProxy versus 70% for the kernel in HTTP keep-alive mode.
392
393  A single process can run many proxy instances; configurations as large as
394  300000 distinct proxies in a single process were reported to run fine. Thus
395  there is usually no need to start more than one process for all instances.
396
397  It is possible to make HAProxy run over multiple processes, but it comes with
398  a few limitations. In general it doesn't make sense in HTTP close or TCP modes
399  because the kernel-side doesn't scale very well with some operations such as
400  connect(). It scales pretty well for HTTP keep-alive mode but the performance
401  that can be achieved out of a single process generally outperforms common needs
402  by an order of magnitude. It does however make sense when used as an SSL
403  offloader, and this feature is well supported in multi-process mode.
404
405  HAProxy only requires the haproxy executable and a configuration file to run.
406  For logging it is highly recommended to have a properly configured syslog daemon
407  and log rotations in place. The configuration files are parsed before starting,
408  then HAProxy tries to bind all listening sockets, and refuses to start if
409  anything fails. Past this point it cannot fail anymore. This means that there
410  are no runtime failures and that if it accepts to start, it will work until it
411  is stopped.
412
413  Once HAProxy is started, it does exactly 3 things :
414
415    - process incoming connections;
416
417    - periodically check the servers' status (known as health checks);
418
419    - exchange information with other haproxy nodes.
420
421  Processing incoming connections is by far the most complex task as it depends
422  on a lot of configuration possibilities, but it can be summarized as the 9 steps
423  below :
424
425    - accept incoming connections from listening sockets that belong to a
426      configuration entity known as a "frontend", which references one or multiple
427      listening addresses;
428
429    - apply the frontend-specific processing rules to these connections that may
430      result in blocking them, modifying some headers, or intercepting them to
431      execute some internal applets such as the statistics page or the CLI;
432
433    - pass these incoming connections to another configuration entity representing
434      a server farm known as a "backend", which contains the list of servers and
435      the load balancing strategy for this server farm;
436
437    - apply the backend-specific processing rules to these connections;
438
439    - decide which server to forward the connection to according to the load
440      balancing strategy;
441
442    - apply the backend-specific processing rules to the response data;
443
444    - apply the frontend-specific processing rules to the response data;
445
446    - emit a log to report what happened in fine details;
447
448    - in HTTP, loop back to the second step to wait for a new request, otherwise
449      close the connection.
450
451  Frontends and backends are sometimes considered as half-proxies, since they only
452  look at one side of an end-to-end connection; the frontend only cares about the
453  clients while the backend only cares about the servers. HAProxy also supports
454  full proxies which are exactly the union of a frontend and a backend. When HTTP
455  processing is desired, the configuration will generally be split into frontends
```

```
456  and backends as they open a lot of possibilities since any frontend may pass a
457  connection to any backend. With TCP-only proxies, using frontends and backends
458  rarely provides a benefit and the configuration can be more readable with full
459  proxies.
460
461
462  3.3. Basic features
463  ---------------
464
465  This section will enumerate a number of features that HAProxy implements, some
466  of which are generally expected from any modern load balancer, and some of
467  which are a direct benefit of HAProxy's architecture. More advanced features
468  will be detailed in the next section.
469
470
471  3.3.1. Basic features : Proxying
472  --------------------------
473
474  Proxying is the action of transferring data between a client and a server over
475  two independant connections. The following basic features are supported by
476  HAProxy regarding proxying and connection management :
477
478    - Provide the server with a clean connection to protect them against any
479      client-side defect or attack;
480
481    - Listen to multiple IP address and/or ports, even port ranges;
482
483    - Transparent accept : intercept traffic targetting any arbitrary IP address
484      that doesn't even belong to the local system;
485
486    - Server port doesn't need to be related to listening port, and may even be
487      translated by a fixed offset (useful with ranges);
488
489    - Transparent connect : spoof the client's (or any) IP address if needed
490      when connecting to the server;
491
492    - Provide a reliable return IP address to the servers in multi-site LBs;
493
494    - Offload the server thanks to buffers and possibly short-lived connections
495      to reduce their concurrent connection count and their memory footprint;
496
497    - Optimize TCP stacks (eg: SACK), congestion control, and reduce RTT impacts;
498
499    - Support different protocol families on both sides (eg: IPv4/IPv6/Unix);
500
501    - Timeout enforcement : HAProxy supports multiple levels of timeouts depending
502      on the stage the connection is, so that a dead client or server, or an
503      attacker cannot be granted resources for too long;
504
505    - Protocol validation: HTTP, SSL, or payload are inspected and invalid
506      protocol elements are rejected, unless instructed to accept them anyway;
507
508    - Policy enforcement : ensure that only what is allowed may be forwarded;
509
510    - Both incoming and outgoing connections may be limited to certain network
511      namespaces (Linux only), making it easy to build a cross-container,
512      multi-tenant load balancer;
513
514    - PROXY protocol presents the client's IP address to the server even for
515      non-HTTP traffic. This is an HAProxy extension that was adopted by a number
516      of third-party products by now, at least these ones at the time of writing :
517      - client : haproxy, stud, stunnel, exaproxy, ELB, squid
518      - server : haproxy, stud, postfix, exim, nginx, squid, node.js, varnish
519
520
```

```
521
522 3.3.2. Basic features : SSL
523 ---------------------------
524 HAProxy's SSL stack is recognized as one of the most featureful according to
525 Google's engineers (http://istlsfastyet.com/). The most commonly used features
526 making it quite complete are :
527
528 - SNI-based multi-hosting with no limit on sites count and focus on
529   performance. At least one deployment is known for running 50000 domains
530   with their respective certificates;
531
532 - support for wildcard certificates reduces the need for many certificates ;
533
534 - certificate-based client authentication with configurable policies on
535   failure to present a valid certificate. This allows to present a different
536   server farm to regenerate the client certificate for example;
537
538 - authentication of the backend server ensures the backend server is the real
539   one and not a man in the middle;
540
541 - authentication with the backend server lets the backend server it's really
542   the expected haproxy node that is connecting to it;
543
544 - TLS NPN and ALPN extensions make it possible to reliably offload SPDY/HTTP2
545   connections and pass them in clear text to backend servers;
546
547 - OCSP stapling further reduces first page load time by delivering inline an
548   OCSP response when the client requests a Certificate Status Request;
549
550 - Dynamic record sizing provides both high performance and low latency, and
551   significantly reduces page load time by letting the browser start to fetch
552   new objects while packets are still in flight;
553
554 - permanent access to all relevant SSL/TLS layer information for logging,
555   access control, reporting etc... These elements can be embedded into HTTP
556   header or even as a PROXY protocol extension so that the offloaded server
557   gets all the information it would have had if it performed the SSL
558   termination itself.
559
560 - Detect, log and block certain known attacks even on vulnerable SSL libs,
561   such as the Heartbleed attack affecting certain versions of OpenSSL.
562
563 - support for stateless session resumption (RFC 5077 TLS Ticket extension).
564   TLS tickets can be updated from CLI which provides them means to implement
565   Perfect Forward Secrecy by frequently rotating the tickets.
566
567
568 3.3.3. Basic features : Monitoring
569 ----------------------------------
570 HAProxy focuses a lot on availability. As such it cares about servers state,
571 and about reporting its own state to other network components :
572
573
574 - Servers state is continuously monitored using per-server parameters. This
575   ensures the path to the server is operational for regular traffic;
576
577 - Health checks support two hysteresis for up and down transitions in order
578   to protect against state flapping;
579
580 - Checks can be sent to a different address/port/protocol : this makes it
581   easy to check a single service that is considered representative of multiple
582   ones, for example the HTTPS port for an HTTP+HTTPS server.
583
584 - Servers can track other servers and go down simultaneously : this ensures
585   that servers hosting multiple services can fail atomically and that noone
```

```
586   will be sent to a partially failed server;
587
588 - Agents may be deployed on the server to monitor load and health : a server
589   may be interested in reporting its load, operational status, administrative
590   status independantly from what health checks can see. By running a simple
591   agent on the server, it's possible to consider the server's view of its own
592   health in addition to the health checks validating the whole path;
593
594 - Various check methods are available : TCP connect, HTTP request, SMTP hello,
595   SSL hello, LDAP, SQL, Redis, send/expect scripts, all with/without SSL;
596
597 - State change is notified in the logs and stats page with the failure reason
598   (eg: the HTTP response received at the moment the failure was detected). An
599   e-mail can also be sent to a configurable address upon such a change ;
600
601 - Server state is also reported on the stats interface and can be used to take
602   routing decisions so that traffic may be sent to different farms depending
603   on their sizes and/or health (eg: loss of an inter-DC link);
604
605 - HAProxy can use health check requests to pass information to the servers,
606   such as their names, weight, the number of other servers in the farm etc...
607   so that servers can adjust their response and decisions based on this
608   knowledge (eg: postpone backups to keep more CPU available);
609
610 - Servers can use health checks to report more detailed state than just on/off
611   (eg: I would like to stop, please stop sending new visitors);
612
613 - HAProxy itself can report its state to external components such as routers
614   or other load balancers, allowing to build very complete multi-path and
615   multi-layer infrastructures.
616
617
618 3.3.4. Basic features : High availability
619 -----------------------------------------
620
621 Just like any serious load balancer, HAProxy cares a lot about availability to
622 ensure the best global service continuity :
623
624 - Only valid servers are used ; the other ones are automatically evinced from
625   load balancing farms ; under certain conditions it is still possible to
626   force to use them though;
627
628 - Support for a graceful shutdown so that it is possible to take servers out
629   of a farm without affecting any connection;
630
631 - Backup servers are automatically used when active servers are down and
632   replace them so that sessions are not lost when possible. This also allows
633   to build multiple paths to reach the same server (eg: multiple interfaces);
634
635 - Ability to return a global failed status for a farm when too many servers
636   are down. This, combined with the monitoring capabilities makes it possible
637   for an upstream component to choose a different LB node for a given service;
638
639 - Stateless design makes it easy to build clusters : by design, HAProxy does
640   its best to ensure the highest service continuity without having to store
641   information that could be lost in the event of a failure. This ensures that
642   a takeover is the most seamless possible;
643
644 - Integrates well with standard VRRP daemon keepalived : HAProxy easily tells
645   keepalived about its state and copes very well with floating virtual IP
646   addresses. Note: only use IP redundancy protocols (VRRP/CARP) over cluster-
647   based solutions (Heartbeat, ...) as they're the ones offering the fastest,
648   most seamless, and most reliable switchover.
649
650
```

### 3.3.5. Basic features : Load balancing
----------------------------------

HAProxy offers a fairly complete set of load balancing features, most of which
are unfortunately not available in a number of other load balancing products :

- no less than 9 load balancing algorithms are supported, some of which apply
  to input data to offer an infinite list of possibilities. The most common
  ones are round-robin (for short connections, pick each server in turn),
  leastconn (for long connections, pick the least recently used of the servers
  with the lowest connection count), source (for SSL farms or terminal server
  farms, the server directly depends on the client's source address), uri (for
  HTTP caches, the server directly depends on the HTTP URI), hdr (the server
  directly depends on the contents of a specific HTTP header field), first
  (for short-lived virtual machines, all connections are packed on the
  smallest possible subset of servers so that unused ones can be powered
  down);

- all algorithms above support per-server weights so that it is possible to
  accommodate from different server generations in a farm, or direct a small
  fraction of the traffic to specific servers (debug mode, running the next
  version of the software, etc);

- dynamic weights are supported for round-robin, leastconn and consistent
  hashing ; this allows server weights to be modified on the fly from the CLI
  or even by an agent running on the server;

- slow-start is supported whenever a dynamic weight is supported; this allows
  a server to progressively take the traffic. This is an important feature
  for fragile application servers which require to compile classes at runtime
  as well as cold caches which need to fill up before being run at full
  throttle;

- hashing can apply to various elements such as client's source address, URL
  components, query string element, header field values, POST parameter, RDP
  cookie;

- consistent hashing protects server farms against massive redistribution when
  adding or removing servers in a farm. That's very important in large cache
  farms and it allows slow-start to be used to refill cold caches;

- a number of internal metrics such as the number of connections per server,
  per backend, the amount of available connection slots in a backend etc makes
  it possible to build very advanced load balancing strategies.


### 3.3.6. Basic features : Stickiness
----------------------------------

Application load balancing would be useless without stickiness. HAProxy provides
a fairly comprehensive set of possibilities to maintain a visitor on the same
server even across various events such as server addition/removal, down/up
cycles, and some methods are designed to be resistant to the distance between
multiple load balancing nodes in that they don't require any replication :

- stickiness information can be individually matched and learned from
  different places if desired. For example a JSESSIONID cookie may be matched
  both in a cookie and in the URL. Up to 8 parallel sources can be learned at
  the same time and each of them may point to a different stick-table;

- stickiness information can come from anything that can be seen within a
  request or response, including source address, TCP payload offset and
  length, HTTP query string elements, header field values, cookies, and so
  on...

- stick-tables are replicated between all nodes in a multi-master fashion ;

- commonly used elements such as SSL-ID or RDP cookies (for TSE farms) are
  directly accessible to ease manipulation;

- all sticking rules may be dynamically conditioned by ACLs;

- it is possible to decide not to stick to certain servers, such as backup
  servers, so that when the nominal server comes back, it automatically takes
  the load back. This is often used in multi-path environments;

- in HTTP it is often prefered not to learn anything and instead manipulate
  a cookie dedicated to stickiness. For this, it's possible to detect,
  rewrite, insert or prefix such a cookie to let the client remember what
  server was assigned;

- the server may decide to change or clean the stickiness cookie on logout,
  so that leaving visitors are automatically unbound from the server;

- using ACL-based rules it is also possible to selectively ignore or enforce
  stickiness regardless of the server's state; combined with advanced health
  checks, that helps admins verify that the server they're installing is up
  and running before presenting it to the whole world;

- an innovative mechanism to set a maximum idle time and duration on cookies
  ensures that stickiness can be smoothly stopped on devices which are never
  closed (smartphones, TVs, home appliances) without having to store them on
  persistent storage;

- multiple server entries may share the same stickiness keys so that
  stickiness is not lost in multi-path environments when one path goes down;

- soft-stop ensures that only users with stickiness information will continue
  to reach the server they've been assigned to but no new users will go there.


### 3.3.7. Basic features : Sampling and converting information
----------------------------------------------------------

HAProxy supports information sampling using a wide set of "sample fetch
functions". The principle is to extract pieces of information known as samples,
for immediate use. This is used for stickiness, to build conditions, to produce
information in logs or to enrich HTTP headers.

Samples can be fetched from various sources :

- constants : integers, strings, IP addresses, binary blocks;

- the process : date, environment variables, server/frontend/backend/process
  state, byte/connection counts/rates, queue length, random generator, ...

- variables : per-session, per-request, per-response variables;

- the client connection : source and destination addresses and ports, and all
  related statistics counters;

- the SSL client session : protocol, version, algorithm, cipher, key size,
  session ID, all client and server certificate fields, certificate serial,
  SNI, ALPN, NPN, client support for certain extensions;

- request and response buffers contents : arbitrary payload at offset/length,
  data length, RDP cookie, decoding of SSL hello type, decoding of TLS SNI;

- HTTP (request and response) : method, URI, path, query string arguments,
  status code, headers values, positionnal header value, cookies, captures,

```
781                   authentication, body elements;
782 A sample may then pass through a number of operators known as "converters" to
783 experience some transformation. A converter consumes a sample and produces a
784 new one, possibly of a completely different type. For example, a converter may
785 be used to return only the integer length of the input string, or could turn a
786 string to upper case. Any arbitrary number of converters may be applied in
787 series to a sample before final use. Among all available sample converters, the
788 following ones are the most commonly used :
789
790  - arithmetic and logic operators : they make it possible to perform advanced
791    computation on input data, such as computing ratios, percentages or simply
792    converting from one unit to another one;
793
794  - IP address masks are useful when some addresses need to be grouped by larger
795    networks;
796
797  - data representation : url-decode, base64, hex, JSON strings, hashing;
798
799  - string conversion : extract substrings at fixed positions, fixed length,
800    extract specific fields around certain delimiters, extract certain words,
801    change case, apply regex-based substitution ;
802
803  - date conversion : convert to http date format, convert local to UTC and
804    conversely, add or remove offset;
805
806  - lookup an entry in a stick table to find statistics or assigned server;
807
808  - map-based key-to-value conversion from a file (mostly used for geolocation).
809
810
811 3.3.8. Basic features : Maps
812 ----------------------------
813
814 Maps are a powerful type of converter consisting in loading a two-columns file
815 into memory at boot time, then looking up each input sample from the first
816 column and either returning the corresponding pattern on the second column if
817 the entry was found, or returning a default value. The output information also
818 being a sample, it can in turn experience other transformations including other
819 map lookups. Maps are most commonly used to translate the client's IP address
820 to an AS number or country code since they support a longest match for network
821 addresses but they can be used for various other purposes.
822
823 Part of their strength comes from being updatable on the fly either from the CLI
824 or from certain actions using other samples, making them capable of storing and
825 retrieving information between subsequent accesses. Another strength comes from
826 the binary tree based indexation which makes them extremely fast event when they
827 contain hundreds of thousands of entries, making geolocation very cheap and easy
828 to set up.
829
830
831 3.3.9. Basic features : ACLs and conditions
832 -------------------------------------------
833
834 Most operations in HAProxy can be made conditional. Conditions are built by
835 combining multiple ACLs using logic operators (AND, OR, NOT). Each ACL is a
836 series of tests based on the following elements :
837
838  - a sample fetch method to retrieve the element to test ;
839
840  - an optional series of converters to transform the element ;
841
842  - a list of patterns to match against ;
843
844  - a matching method to indicate how to compare the patterns with the sample
845
```

```
846 For example, the sample may be taken from the HTTP "Host" header, it could then
847 be converted to lower case, then matched against a number of regex patterns
848 using the regex matching method.
849
850 Technically, ACLs are built on the same core as the maps, they share the exact
851 same internal structure, pattern matching methods and performance. The only real
852 difference is that instead of returning a sample, they only return "found" or
853 "not found". In terms of usage, ACL patterns may be declared inline in the
854 configuration file and do not require their own file. ACLs may be named for ease
855 of use or to make configurations understandable. A named ACL may be declared
856 multiple times and it will evaluate all definitions in turn until one matches.
857
858 About 13 different pattern matching methods are provided, among which IP address
859 mask, integer ranges, substrings, regex. They work like functions, and just like
860 with any programming language, only what is needed is evaluated, so when a
861 condition involving an OR is already true, next ones are not evaluated, and
862 similarly when a condition involving an AND is already false, the rest of the
863 condition is not evaluated.
864
865 There is no practical limit to the number of declared ACLs, and a handful of
866 commonly used ones are provided. However experience has shown that setups using
867 a lot of named ACLs are quite hard to troubleshoot and that sometimes using
868 anonymous ACLs inline is easier as it requires less references out of the scope
869 being analysed.
870
871
872 3.3.10. Basic features : Content switching
873 ------------------------------------------
874
875 HAProxy implements a mechanism known as content-based switching. The principle
876 is that a connection or request arrives on a frontend, then the information
877 carried with this request or connection are processed, and at this point it is
878 possible to write ACLs-based conditions making use of these information to
879 decide what backend will process the request. Thus the traffic is directed to
880 one backend or another based on the request's contents. The most common example
881 consists in using the Host header and/or elements from the path (sub-directories
882 or file-name extensions) to decide whether an HTTP request targets a static
883 object or the application, and to route static objects traffic to a backend made
884 of fast and light servers, and all the remaining traffic to a more complex
885 application server, thus constituting a fine-grained virtual hosting solution.
886 This is quite convenient to make multiple technologies coexist as a more global
887 solution.
888
889 Another use case of content-switching consists in using different load balancing
890 algorithms depending on various criteria. A cache may use a URI hash while an
891 application would use round robin.
892
893 Last but not least, it allows multiple customers to use a small share of a
894 common resource by enforcing per-backend (thus per-customer connection limits).
895
896 Content switching rules scale very well, though their performance may depend on
897 the number and complexity of the ACLs in use. But it is also possible to write
898 dynamic content switching rules where a sample value directly turns into a
899 backend name and without making use of ACLs at all. Such configurations have
900 been reported to work fine at least with 300000 backends in production.
901
902
903 3.3.11. Basic features : Stick-tables
904 -------------------------------------
905
906 Stick-tables are commonly used to store stickiness information, that is, to keep
907 a reference to the server a certain visitor was directed to. The key is then the
908 identifier associated with the visitor (its source address, the SSL ID of the
909 connection, an HTTP or RDP cookie, the customer number extracted from the URL or
910
```

from the payload, ...) and the stored value is then the server's identifier.

Stick tables may use 3 different types of samples for their keys : integers, strings and addresses. Only one stick-table may be referenced in a proxy, and it is designated everywhere with the proxy name. Up to 8 key may be tracked in parallel. The server identifier is committed during request or response processing once both the key and the server are known.

Stick-table contents may be replicated in active-active mode with other HAProxy nodes known as "peers" as well as with the new process during a reload operation so that all load balancing nodes share the same information and take the same routing decision if a client's requests are spread over multiple nodes.

Since stick-tables are indexed on what allows to recognize a client, they are often also used to store extra information such as per-client statistics. The extra statistics take some extra space and need to be explicitly declared. The type of statistics that may be stored includes the input and output bandwidth, the number of concurrent connections, the connection rate and count over a period, the amount and frequency of errors, some specific tags and counters, etc...In order to support keeping such information without being forced to stick to a given server, a special "tracking" feature is implemented and allows to track up to 3 simultaneous keys from different tables at the same time regardless of stickiness rules. Each stored statistics may be searched, dumped and cleared from the CLI and adds to the live troubleshooting capabilities.

While this mechanism can be used to surclass a returning visitor or to adjust the delivered quality of service depending on good or bad behaviour, it is mostly used to fight against service abuse and more generally DDoS as it allows to build complex models to detect certain bad behaviours at a high processing speed.


3.3.12. Basic features : Formated strings
--------------------------------

There are many places where HAProxy needs to manipulate character strings, such as logs, redirects, header additions, and so on. In order to provide the greatest flexibility, the notion of formated strings was introduced, initially for logging purposes, which explains why it's still called "log-format". These strings contain escape characters allowing to introduce various dynamic data including variables and sample fetch expressions into strings, and even to adjust the encoding while the result is being turned into a string (for example, adding quotes). This provides a powerful way to build header contents or to customize log lines. Additionally, in order to remain simple to build most common strings, about 50 special tags are provided as shortcuts for information commonly used in logs.


3.3.13. Basic features : HTTP rewriting and redirection
--------------------------------

Installing a load balancer in front of an application that was never designed for this can be a challenging task without the proper tools. One of the most commonly requested operation in this case is to adjust requests and response headers to make the load balancer appear as the origin server and to fix hard coded information. This comes with changing the path in requests (which is strongly advised against), modifying Host header field, modifying the Location response header field for redirects, modifying the path and domain attribute for cookies, and so on. It also happens that a number of servers are somewhat verbose and tend to leak too much information in the response, making them more vulnerable to targetted attacks. While it's theorically not the role of a load balancer to clean this up, in practice it's located at the best place in the infrastructure to guarantee that everything is cleaned up.

Similarly, sometimes the load balancer will have to intercept some requests and

respond with a redirect to a new target URL. While some people tend to confuse redirects and rewriting, these are two completely different concepts, since the rewriting makes the client and the server see different things (and disagree on the location of the page being visited) while redirects ask the client to visit the new URL so that it sees the same location as the server.

In order to do this, HAProxy supports various possibilities for rewriting and redirect, among which :

- regex-based URL and header rewriting in requests and responses. Regex are the most commonly used tool to modify header values since they're easy to manipulate and well understood;

- headers may also be appended, deleted or replaced based on formated strings so that it is possible to pass information there (eg: client side TLS algorithm and cipher);

- HTTP redirects can use any 3xx code to a relative, absolute, or completely dynamic (formated string) URI;

- HTTP redirects also support some extra options such as setting or clearing a specific cookie, dropping the query string, appending a slash if missing, and so on;

- all operations support ACL-based conditions;


3.3.14. Basic features : Server protection
--------------------------------

HAProxy does a lot to maximize service availability, and for this it deploys large efforts to protect servers against overloading and attacks. The first and most important point is that only complete and valid requests are forwarded to the servers. The initial reason is that HAProxy needs to find the protocol elements it needs to stay synchronized with the byte stream, and the second reason is that until the request is complete, there is no way to know if some elements will change its semantics. The direct benefit from this is that servers are not exposed to invalid or incomplete requests. This is a very effective protection against slowloris attacks, which have almost no impact on HAProxy.

Another important point is that HAProxy contains buffers to store requests and responses, and that by only sending a request to a server when it's complete and by reading the whole response very quickly from the local network, the server side connection is used for a very short time and this preserves server resources as much as possible.

A direct extension to this is that HAProxy can artificially limit the number of concurrent connections or outstanding requests to a server, which guarantees that the server will never be overloaded even if it continuously runs at 100% of its capacity during traffic spikes. All excess requests will simply be queued to be processed when one slot is released. In the end, this huge resource savings most often ensures so much better server response times that it ends up actually being faster than by overloading the server. Queued requests may be redispatched to other servers, or even aborted in queue when the client aborts, which also protects the servers against the "reload effect", where each click on "reload" by a visitor on a slow-loading page usually induces a new request and maintains the server in an overloaded state.

The slow-start mechanism also protects restarting servers against high traffic levels while they're still finalizing their startup or compiling some classes.

Regarding the protocol-level protection, it is possible to relax the HTTP parser to accept non standard-compliant but harmless requests or responses and even to fix them. This allows bogus applications to be accessible while a fix is being developed. In parallel, offending messages are completely captured with a

detailed report that help developers spot the issue in the application. The most
dangerous protocol violations are properly detected and dealt with and fixed.
For example malformed requests or responses with two Content-length headers are
either fixed if the values are exactly the same, or rejected if they differ,
since it becomes a security problem. Protocol inspection is not limited to HTTP,
it is also available for other protocols like TLS or RDP.

When a protocol violation or attack is detected, there are various options to
respond to the user, such as returning the common "HTTP 400 bad request",
closing the connection with a TCP reset, faking an error after a long delay
("tarpit") to confuse the attacker. All of these contribute to protecting the
servers by discouraging the offending client from pursuing an attack that
becomes very expensive to maintain.

HAProxy also proposes some more advanced options to protect against accidental
data leaks and session crossing. Not only it can log suspicious server responses
but it will also log and optionally block a response which might affect a given
visitors' confidentiality. One such example is a cacheable cookie appearing in a
cacheable response and which may result in an intermediary cache to deliver it
to another visitor, causing an accidental session sharing.


### 3.3.15. Basic features : Logging

Logging is an extremely important feature for a load balancer, first because a
load balancer is often accused of the trouble it reveals, and second because it
is placed at a critical point in an infrastructure where all normal and abnormal
activity needs to be analysed and correlated with other components.

HAProxy provides very detailed logs, with millisecond accuracy and the exact
connection accept time that can be searched in firewalls logs (eg: for NAT
correlation). By default, TCP and HTTP logs are quite detailed an contain
everything needed for troubleshooting, such as source IP address and port,
frontend, backend, server, timers (request receipt duration, queue duration,
connection setup time, response headers time, data transfer time), global
process state, connection counts, queue status, retries count, detailed
stickiness actions and disconnect reasons, header captures with a safe output
encoding. It is then possible to extend or replace this format to include any
sampled data, variables, captures, resulting in very detailed information. For
example it is possible to log the number cumulated requests for this client or
the number of different URLs for the client.

The log level may be adjusted per request using standard ACLs, so it is possible
to automatically silent some logs considered as pollution and instead raise
warnings when some abnormal behaviour happen for a small part of the traffic
(eg: too many URLs or HTTP errors for a source address). Administrative logs are
also emitted with their own levels to inform about the loss or recovery of a
server for example.

Each frontend and backend may use multiple independant log outputs, which eases
multi-tenancy. Logs are preferably sent over UDP, maybe JSON-encoded, and are
truncated after a configurable line length in order to guarantee delivery.


### 3.3.16. Basic features : Statistics

HAProxy provides a web-based statistics reporting interface with authentication,
security levels and scopes. It is thus possible to provide each hosted customer
with his own page showing only his own instances. This page can be located in a
hidden URL part of the regular web site so that no new port needs to be opened.
This page may also report the availability of other HAProxy nodes so that it is
easy to spot if everything works as expected at a glance. The view is synthetic
with a lot of details accessible (such as error causes, last access and last

change duration, etc), which are also accessible as a CSV table that other tools
may import to draw graphs. The page may self-refresh to be used as a monitoring
page on a large display. In administration mode, the page also allows to change
server state to ease maintenance operations.


## 3.4. Advanced features

### 3.4.1. Advanced features : Management

HAProxy is designed to remain extremely stable and safe to manage in a regular
production environment. It is provided as a single executable file which doesn't
require any installation process. Multiple versions can easily coexist, meaning
that it's possible (and recommended) to upgrade instances progressively by
order of criticity instead of migrating all of them at once. Configuration files
are easily versionned. Configuration checking is done off-line so it doesn't
require to restart a service that will possibly fail. During configuration
checks, a number of advanced mistakes may be detected (eg: for example, a rule
hiding another one, or stickiness that will not work) and detailed warnings and
configuration hints are proposed to fix them. Backwards configuration file
compatibility goes very far away in time, with version 1.5 still fully
supporting configurations for versions 1.1 written 13 years before, and 1.6
only dropping support for almost unused, obsolete keywords that can be done
differently. The configuration and software upgrade mechanism is smooth and non
disruptive in that it allows old and new processes to coexist on the system,
each handling its own connections. System status, build options and library
compatibility are reported on startup.

Some advanced features allow an application administrator to smoothly stop a
server, detect when there's no activity on it anymore, then take it off-line,
stop it, upgrade it and ensure it doesn't take any traffic while being upgraded,
then test it again through the normal path without opening it to the public, and
all of this without touching HAProxy at all. This ensures that even complicated
production operations may be done during opening hours with all technical
resources available.

The process tries to save resources as much as possible, uses memory pools to
save on allocation time and limit memory fragmentation, releases payload buffers
as soon as their contents are sent, and supports enforcing strong memory limits
above which connections have to wait for a buffer to become available instead of
allocating more memory. This system helps guarantee memory usage in certain
strict environments.

A command line interface (CLI) is available as a UNIX or TCP socket, to perform
a number of operations and to retrieve troubleshooting information. Everything
done on this socket doesn't require a configuration change, so it is mostly used
for temporary changes. Using this interface it is possible to change a server's
address, weight and status, to consult statistics and clear counters, dump and
clear stickiness tables, possibly selectively by key criteria, dump and kill
client-side and server-side connections, dump captured errors with a detailed
analysis of the exact cause and location of the error, dump, add and remove
entries from ACLs and maps, update TLS shared secrets, apply connection limits
and rate limits on the fly to arbitrary frontends (useful in shared hosting
environments), and disable a specific frontend to release a listening port
(useful when daytime operations are forbidden and a fix is needed nonetheless).

For environments where SNMP is mandatory, at least two agents exist, one is
provided with the HAProxy sources and relies on the Net-SNMP perl module.
Another one is provided with the commercial packages and doesn't require Perl.
Both are roughly equivalent in terms of coverage.

It is often recommended to install 4 utilities on the machine where HAProxy is
deployed :

```
1171  - socat (in order to connect to the CLI, though certain forks of netcat can
1172    also do it to some extents);
1173
1174  - halog from the latest HAProxy version : this is the log analysis tool, it
1175    parses native TCP and HTTP logs extremely fast (1 to 2 GB per second) and
1176    extracts useful information and statistics such as requests per URL, per
1177    source address, URLs sorted by response time or error rate, termination
1178    codes etc... It was designed to be deployed on the production servers to
1179    help troubleshoot live issues so it has to be there ready to be used;
1180
1181  - tcpdump : this is highly recommended to take the network traces needed to
1182    troubleshoot an issue that was made visible in the logs. There is a moment
1183    where application and haproxy's analysis will diverge and the network traces
1184    are the only way to say who's right and who's wrong. It's also fairly common
1185    to detect bugs in network stacks and hypervisors thanks to tcpdump;
1186
1187
1188  - strace : it is tcpdump's companion. It will report what HAProxy really sees
1189    and will help sort out the issues the operating system is responsible for
1190    from the ones HAProxy is responsible for. Strace is often requested when a
1191    bug in HAProxy is suspected;
1192
1193
1194  3.4.2. Advanced features : System-specific capabilities
1195  ------------------------------------------------------
1196
1197  Depending on the operating system HAProxy is deployed on, certain extra features
1198  may be available or needed. While it is supported on a number of platforms,
1199  HAProxy is primarily developed on Linux, which explains why some features are
1200  only available on this platform.
1201
1202  The transparent bind and connect features, the support for binding connections
1203  to a specific network interface, as well as the ability to bind multiple
1204  processes to the same IP address and ports are only available on Linux and BSD
1205  systems, though only Linux performs a kernel-side load balancing of the incoming
1206  requests between the available processes.
1207
1208  On Linux, there are also a number of extra features and optimizations including
1209  support for network namespaces (also known as "containers") allowing HAProxy to
1210  be a gateway between all containers, the ability to set the MSS, Netfilter marks
1211  and IP TOS field on the client side connection, support for TCP FastOpen on the
1212  listening side, TCP user timeouts to let the kernel quickly kill connections
1213  when it detects the client has disappeared before the configured timeouts, TCP
1214  splicing to let the kernel forward data between the two sides of a connections
1215  thus avoiding multiple memory copies, the ability to enable the "defer-accept"
1216  bind option to only get notified of an incoming connection once data become
1217  available in the kernel buffers, and the ability to send the request with the
1218  ACK confirming a connect (sometimes called "biggy-back") which is enabled with
1219  the "tcp-smart-connect" option. On Linux, HAProxy also takes great care of
1220  manipulating the TCP delayed ACKs to save as many packets as possible on the
1221  network.
1222
1223  Some systems have an unreliable clock which jumps back and forth in the past
1224  and in the future. This used to happen with some NUMA systems where multiple
1225  processors didn't see the exact same time of day, and recently it became more
1226  common in virtualized environments where the virtual clock has no relation with
1227  the real clock, resulting in huge time jumps (sometimes up to 30 seconds have
1228  been observed). This causes a lot of trouble with respect to timeout enforcement
1229  in general. Due to this flaw of these systems, HAProxy maintains its own
1230  monotonic clock which is based on the system's clock but where drift is measured
1231  and compensated for. This ensures that even with a very bad system clock, timers
1232  remain reasonably accurate and timeouts continue to work. Note that this problem
1233  affects all the software running on such systems and is not specific to HAProxy.
1234  The common effects are spurious timeouts or application freezes. Thus if this
1235  behaviour is detected on a system, it must be fixed, regardless of the fact that
```

```
1236  HAProxy protects itself against it.
1237
1238  3.4.3. Advanced features : Scripting
1239  ------------------------------------
1240
1241  HAProxy can be built with support for the Lua embedded language, which opens a
1242  wide area of new possibilities related to complex manipulation of requests or
1243  responses, routing decisions, statistics processing and so on. Using Lua it is
1244  even possible to establish parallel connections to other servers to exchange
1245  information. This way it becomes possible (though complex) to develop an
1246  authentication system for example. Please refer to the documentation in the file
1247  "doc/lua-api/index.rst" for more information on how to use Lua.
1248
1249
1250  3.5. Sizing
1251  -----------
1252
1253  Typical CPU usage figures show 15% of the processing time spent in HAProxy
1254  versus 85% in the kernel in TCP or HTTP close mode, and about 30% for HAProxy
1255  versus 70% for the kernel in HTTP keep-alive mode. This means that the operating
1256  system and its tuning have a strong impact on the global performance.
1257
1258
1259  Usages vary a lot between users, some focus on bandwidth, other ones on request
1260  rate, others on connection concurrency, others on SSL performance. this section
1261  aims at providing a few elements to help in this task.
1262
1263  It is important to keep in mind that every operation comes with a cost, so each
1264  individual operation adds its overhead on top of the other ones, which may be
1265  negligible in certain circumstances, and which may dominate in other cases.
1266
1267  When processing the requests from a connection, we can say that :
1268
1269  - forwarding data costs less than parsing request or response headers;
1270
1271  - parsing request or response headers cost less than establishing then closing
1272    a connection to a server;
1273
1274  - establishing an closing a connection costs less than a TLS resume operation;
1275
1276  - a TLS resume operation costs less than a full TLS handshake with a key
1277    computation;
1278
1279  - an idle connection costs less CPU than a connection whose buffers hold data;
1280
1281  - a TLS context costs even more memory than a connection with data;
1282
1283  So in practice, it is cheaper to process payload bytes than header bytes, thus
1284  it is easier to achieve high network bandwidth with large objects (few requests
1285  per volume unit) than with small objects (many requests per volume unit). This
1286  explains why maximum bandwidth is always measured with large objects, while
1287  request rate or connection rates are measured with small objects.
1288
1289  Some operations scale well on multiple process spread over multiple processors,
1290  and others don't scale as well. Network bandwidth doesn't scale very far because
1291  the CPU is rarely the bottleneck for large objects, it's mostly the network
1292  bandwidth and data busses to reach the network interfaces. The connection rate
1293  doesn't scale well over multiple processors due to a few locks in the system
1294  when dealing with the local ports table. The request rate over persistent
1295  connections scales very well as it doesn't involve much memory nor network
1296  bandwidth and doesn't require to access locked structures. TLS key computation
1297  scales very well as it's totally CPU-bound. TLS resume scales moderately well,
1298  but reaches its limits around 4 processes where the overhead of accessing the
1299  shared table offsets the small gains expected from more power.
1300
```

1301   The performance numbers one can expect from a very well tuned system are in the
1302   following range. It is important to take them as orders of magnitude and to
1303   expect significant variations in any direction based on the processor, IRQ
1304   setting, memory type, network interface type, operating system tuning and so on.
1305
1306   The following numbers were found on a Core i7 running at 3.7 GHz equipped with
1307   a dual-port 10 Gbps NICs running Linux kernel 3.10, HAProxy 1.6 and OpenSSL
1308   1.0.2. HAProxy was running as a single process on a single dedicated CPU core,
1309   and two extra cores were dedicated to network interrupts :
1310
1311   - 20 Gbps of maximum network bandwidth in clear text for objects 256 kB or
1312     higher, 10 Gbps for 41kB or higher;
1313
1314   - 4.6 Gbps of TLS traffic using AES256-GCM cipher with large objects;
1315
1316   - 83000 TCP connections per second from client to server;
1317
1318   - 82000 HTTP connections per second from client to server;
1319
1320   - 97000 HTTP requests per second in server-close mode (keep-alive with the
1321     client, close with the server);
1322
1323   - 243000 HTTP requests per second in end-to-end keep-alive mode;
1324
1325   - 300000 filtered TCP connections per second (anti-DDoS)
1326
1327   - 160000 HTTPS requests per second in keep-alive mode over persistent TLS
1328     connections;
1329
1330   - 13100 HTTPS requests per second using TLS resumed connections;
1331
1332   - 1300 HTTPS connections per second using TLS connections renegociated with
1333     RSA2048;
1334
1335   - 20000 concurrent saturated connections per GB of RAM, including the memory
1336     required for system buffers; it is possible to do better with careful tuning
1337     but this setting it easy to achieve.
1338
1339   - about 8000 concurrent TLS connections (client-side only) per GB of RAM,
1340     including the memory required for system buffers;
1341
1342   - about 5000 concurrent end-to-end TLS connections (both sides) per GB of
1343     RAM including the memory required for system buffers;
1344
1345   Thus a good rule of thumb to keep in mind is that the request rate is divided
1346   by 10 between TLS keep-alive and TLS resume, and between TLS resume and TLS
1347   renegociation, while it's only divided by 3 between HTTP keep-alive and HTTP
1348   close. Another good rule of thumb is to remember that a high frequency core
1349   with AES instructions can do around 5 Gbps of AES-GCM per core.
1350
1351   Having more core rarely helps (except for TLS) and is even counter-productive
1352   due to the lower frequency. In general a small number of high frequency cores
1353   is better.
1354
1355   Another good rule of thumb is to consider that on the same server, HAProxy will
1356   be able to saturate :
1357
1358   - about 5-10 static file servers or caching proxies;
1359
1360   - about 100 anti-virus proxies;
1361
1362   - and about 100-1000 application servers depending on the technology in use.
1363
1364
1365   3.6. How to get HAProxy

---

  ---------------------
1366
1367   HAProxy is an opensource project covered by the GPLv2 license, meaning that
1368   everyone is allowed to redistribute it provided that access to the sources is
1369   also provided upon request, especially if any modifications were made.
1370
1371   HAProxy evolves as a main development branch called "master" or "mainline", from
1372   which new branches are derived once the code is considered stable. A lot of web
1373   sites run some development branches in production on a voluntarily basis, either
1374   to participate to the project or because they need a bleeding edge feature, and
1375   their feedback is highly valuable to fix bugs and judge the overall quality and
1376   stability of the version being developed.
1377
1378
1379   The new branches that are created when the code is stable enough constitute a
1380   stable version and are generally maintained for several years, so that there is
1381   no emergency to migrate to a newer branch even when you're not on the latest.
1382   Once a stable branch is issued, it may only receive bug fixes, and very rarely
1383   minor feature updates when that makes users' life easier. All fixes that go into
1384   a stable branch necessarily come from the master branch. This guarantees that no
1385   fix will be lost after an upgrade. For this reason, if you fix a bug, please
1386   make the patch against the master branch, not the stable branch. You may even
1387   discover it was already fixed. This process also ensures that regressions in a
1388   stable branch are extremely rare, so there is never any excuse for not upgrading
1389   to the latest version in your current branch.
1390
1391   Branches are numbered with two digits delimited with a dot, such as "1.6". A
1392   complete version includes one or two sub-version numbers indicating the level of
1393   fix. For example, version 1.5.14 is the 14th fix release in branch 1.5 after
1394   version 1.5.0 was issued. It contains 126 fixes for individual bugs, 24 updates
1395   on the documentation, and 75 other backported patches, most of which were needed
1396   to fix the aforementionned 126 bugs. An existing feature may never be modified
1397   nor removed in a stable branch, in order to guarantee that upgrades within the
1398   same branch will always be harmless.
1399
1400   HAProxy is available from multiple sources, at different release rhythms :
1401
1402   - The official community web site : http://www.haproxy.org/ ; this site
1403     provides the sources of the latest development release, all stable releases,
1404     as well as nightly snapshots for each branch. The release cycle is not fast,
1405     several months between stable releases, or between development snapshots.
1406     Very old versions are still supported there. Everything is provided as
1407     sources only, so whatever comes from there needs to be rebuilt and/or
1408     repackaged;
1409
1410   - A number of operating systems such as Linux distributions and BSD ports.
1411     These systems generally provide long-term maintained versions which do not
1412     always contain all the fixes from the official ones, but which at least
1413     contain the critical fixes. It often is a good option for most users who do
1414     not seek advanced configurations and just want to keep updates easy;
1415
1416   - Commercial versions from http://www.haproxy.com/ : these are supported
1417     professional packages built for various operating systems or provided as
1418     appliances, based on the latest stable versions and including a number of
1419     features backported from the next release for which there is a strong
1420     demand. It is the best option for users seeking the latest features with
1421     the reliability of a stable branch, the fastest response time to fix bugs,
1422     or simply support contracts on top of an opensource product;
1423
1424
1425   In order to ensure that the version you're using is the latest one in your
1426   branch, you need to proceed this way :
1427
1428   - verify which HAProxy executable you're running : some systems ship it by
1429     default and administrators install their versions somewhere else on the
1430     system, so it is important to verify in the startup scripts which one is

```
1431 used;
1432
1433 - determine which source your HAProxy version comes from. For this, it's
1434   generally sufficient to type "haproxy -v". A development version will
1435   appear like this, with the "dev" word after the branch number :
1436
1437       HA-Proxy version 1.6-dev3-385ecc-68 2015/08/18
1438
1439 A stable version will appear like this, as well as unmodified stable
1440 versions provided by operating system vendors :
1441
1442       HA-Proxy version 1.5.14 2015/07/02
1443
1444 And a nightly snapshot of a stable version will appear like this with an
1445 hexadecimal sequence after the version, and with the date of the snapshot
1446 instead of the date of the release :
1447
1448       HA-Proxy version 1.5.14-e4766ba 2015/07/29
1449
1450 Any other format may indicate a system-specific package with its own
1451 patch set. For example HAProxy Enterprise versions will appear with the
1452 following format (<branch>-<latest commit>-<revision>) :
1453
1454       HA-Proxy version 1.5.0-994126-357 2015/07/02
1455
1456 - for system-specific packages, you have to check with your vendor's package
1457   repository or update system to ensure that your system is still supported,
1458   and that fixes are still provided for your branch. For community versions
1459   coming from haproxy.org, just visit the site, verify the status of your
1460   branch and compare the latest version with yours to see if you're on the
1461   latest one. If not you can upgrade. If your branch is not maintained
1462   anymore, you're definitely very late and will have to consider an upgrade
1463   to a more recent branch (carefully read the README when doing so).
1464
1465 HAProxy will have to be updated according to the source it came from. Usually it
1466 follows the system vendor's way of upgrading a package. If it was taken from
1467 sources, please read the README file in the sources directory after extracting
1468 the sources and follow the instructions for your operating system.
1469
1470
1471 4. Companion products and alternatives
1472 --------------------------------------
1473
1474 HAProxy integrates fairly well with certain products listed below, which is why
1475 they are mentioned here even if not directly related to HAProxy.
1476
1477
1478 4.1. Apache HTTP server
1479 -----------------------
1480
1481 Apache is the de-facto standard HTTP server. It's a very complete and modular
1482 project supporting both file serving and dynamic contents. It can serve as a
1483 frontend for some application servers. In can even proxy requests and cache
1484 responses. In all of these use cases, a front load balancer is commonly needed.
1485 Apache can work in various modes, certain being heavier than other ones. Certain
1486 modules still require the heavier pre-forked model and will prevent Apache from
1487 scaling well with a high number of connections. In this case HAProxy can provide
1488 a tremendous help by enforcing the per-server connection limits to a safe value
1489 and will significantly speed up the server and preserve its resources that will
1490 be better used by the application.
1491
1492 Apache can extract the client's address from the X-Forwarded-For header by using
1493 the "mod_rpaf" extension. HAProxy will automatically feed this header when
1494 "option forwardfor" is specified in its configuration. HAProxy may also offer a
1495 nice protection to Apache when exposed to the internet, where it will better
```

```
1496 resist to a wide number of types of DoS.
1497
1498 4.2. NGINX
1499 ----------
1500
1501
1502 NGINX is the second de-facto standard HTTP server. Just like Apache, it covers a
1503 wide range of features. NGINX is built on a similar model as HAProxy so it has
1504 no problem dealing with tens of thousands of concurrent connections. When used
1505 as a gateway to some applications (eg: using the included PHP FPM), it can often
1506 be beneficial to set up some frontend connection limiting to reduce the load
1507 on the PHP application. HAProxy will clearly be useful there both as a regular
1508 load balancer and as the traffic regulator to speed up PHP by decongestionning
1509 it. Also since both products use very little CPU thanks to their event-driven
1510 architecture, it's often easy to install both of them on the same system. NGINX
1511 implements HAProxy's PROXY protocol, thus it is easy for HAProxy to pass the
1512 client's connection information to NGINX so that the application gets all the
1513 relevant information. Some benchmarks have also shown that for large static
1514 file serving, implementing consistent hash on HAProxy in front of NGINX can be
1515 beneficial by optimizing the OS' cache hit ratio, which is basically multiplied
1516 by the number of server nodes.
1517
1518
1519 4.3. Varnish
1520 ------------
1521
1522 Varnish is a smart caching reverse-proxy, probably best described as a web
1523 application accelerator. Varnish doesn't implement SSL/TLS and wants to dedicate
1524 all of its CPU cycles to what it does best. Varnish also implements HAProxy's
1525 PROXY protocol so that HAProxy can very easily be deployed in front of Varnish
1526 as an SSL offloader as well as a load balancer and pass it all relevant client
1527 information. Also, Varnish naturally supports decompression from the cache when
1528 a server has provided a compressed object, but doesn't compress however. HAProxy
1529 can then be used to compress outgoing data when backend servers do not implement
1530 compression, though it's rarely a good idea to compress on the load balancer
1531 unless the traffic is low.
1532
1533 When building large caching farms across multiple nodes, HAProxy can make use of
1534 consistent URL hashing to intelligently distribute the load to the caching nodes
1535 and avoid cache duplication, resulting in a total cache size which is the sum of
1536 all caching nodes.
1537
1538
1539 4.4. Alternatives
1540 -----------------
1541
1542 Linux Virtual Server (LVS or IPVS) is the layer 4 load balancer included within
1543 the Linux kernel. It works at the packet level and handles TCP and UDP. In most
1544 cases it's more a complement than an alternative since it doesn't have layer 7
1545 knowledge at all.
1546
1547 Pound is another well-known load balancer. It's much simpler and has much less
1548 features than HAProxy but for many very basic setups both can be used. Its
1549 author has always focused on code auditability first and wants to maintain the
1550 set of features low. Its thread-based architecture scales well with high
1551 connection counts, but it's a good product.
1552
1553 Pen is a quite light load balancer. It supports SSL, maintains persistence using
1554 a fixed-size table of its clients' IP addresses. It supports a packet-oriented
1555 mode allowing it to support direct server return and UDP to some extents. It is
1556 meant for small loads (the persistence table only has 2048 entries).
1557
1558 NGINX can do some load balancing to some extents, though it's clearly not its
1559 primary function. Production traffic is used to detect server failures, the
1560 load balancing algorithms are more limited, and the stickiness is very limited.
```

```
1561    But it can make sense in some simple deployment scenarios where it is already
1562    present. The good thing is that since it integrates very well with HAProxy,
1563    there's nothing wrong with adding HAProxy later when its limits have been faced.
1564
1565    Varnish also does some load balancing of its backend servers and does support
1566    real health checks. It doesn't implement stickiness however, so just like with
1567    NGINX, as long as stickiness is not needed that can be enough to start with.
1568    And similarly, since HAProxy and Varnish integrate so well together, it's easy
1569    to add it later into the mix to complement the feature set.
1570
```