# Detecting and dealing with Denial of Service in all its forms

By Andreas Aaris-Larsen

W/TH secure

# Introduction and definitions

Denial of Service (DoS) can be split into many categories, but for the purposes of this document the following basic and conceptual categories will be used:

• Traffic congestion
• Resource depletion
• Logical

The most common one, and the one most often associated with Distributed Denial of Servce (DDoS) attacks, are the traffic congestion ones, often also referred to as "volumetric" DoS. The premise is fairly simple:
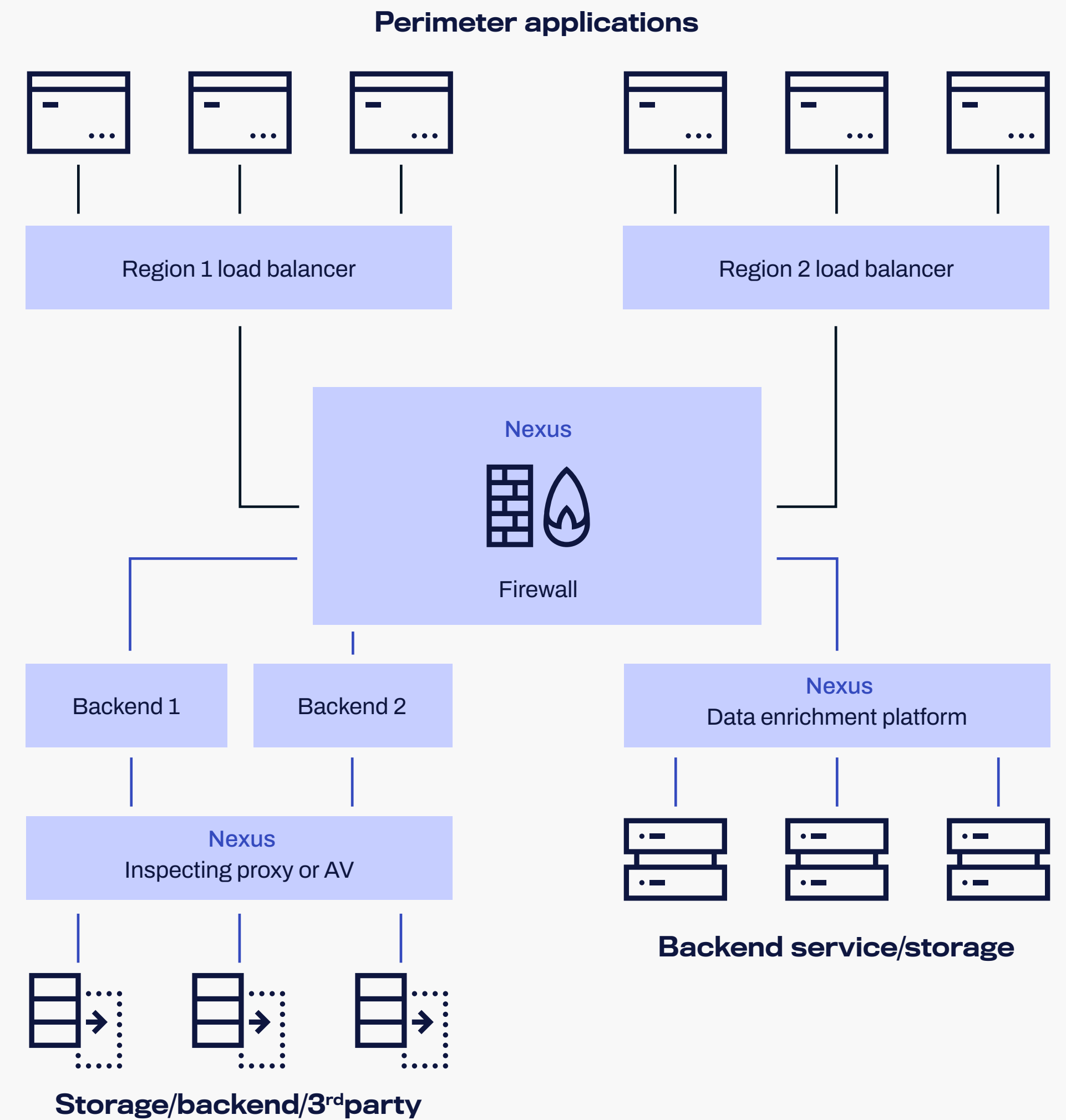
• The threat actors' Internet connection has more bandwidth than yours (or the actor controls an army of bots whose combined bandwidth exceeds yours)
• The threat actor can produce more traffic than your Internet connection can handle, for instances through amplification attacks with a spoofed origin address
• You device becomes unable to respond to other users' request, because all the bandwidth is used on responding to threat actor traffic

While the idea of traffic volume exceeding available bandwidth is a relatively simple concept, for the depletion one, the concept is that every bit of traffic sent to your server incurs a cost in terms of system resources (CPU, RAM, storage, etc.) needed to process and respond to the incoming request. This can be anything from the web server needing to parse a basic HTTP GET request and formulate a response with some content, to complex data modelling and processing involving multiple internal systems and $3^{rd}$ parties. Each request will require a small fraction of the servers available CPU, RAM and storage, and in fact might require resources on several systems depending on the type of application. Multiply that by the number of expected concurrent users, and you have the formula for how developers and infrastructure operators calculate the specifications for the systems that will be hosting their solution. Submit more simultaneous requests than the developers and operators were expecting, and the system is likely to run out of resources and react in unforeseen ways or simply stop working altogether due to resource depletion[1].

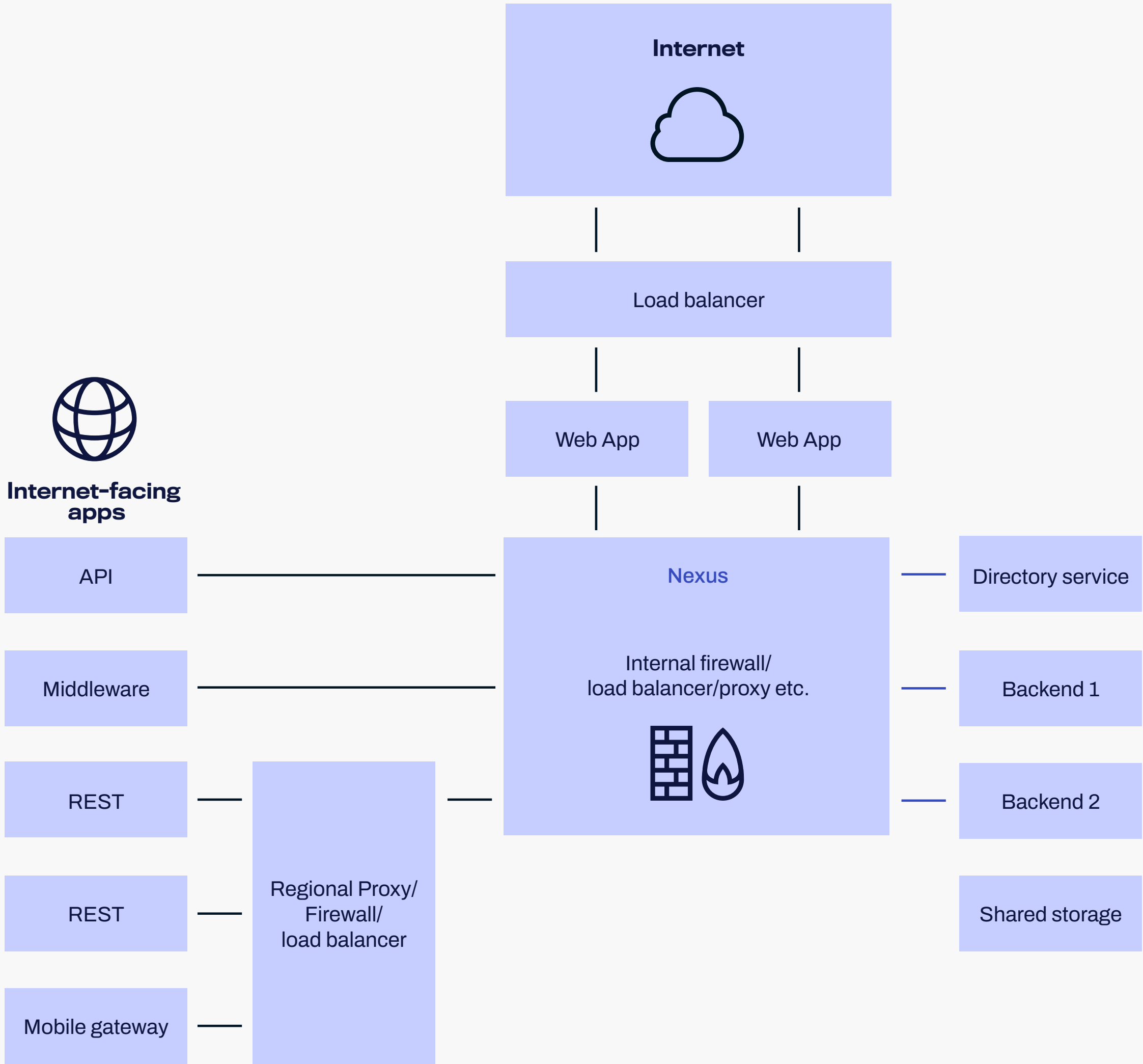[1]DDoS on Dyn Impacts Twitter, Spotify, Reddit – Krebs on Security

The third and final category, is a bit of a beast; logical Denial of Service conditions. While traditionally discussed as application programming mistakes that can lead to the scenarios of traffic congestion and resource depletion (discussed at length in this OWASP resource, this category also includes design and architectural conditions resulting in DoS. Consider for instance a large multinational organisation, with physical locations around the globe and an extensive external perimeter. While DoS conditions may be sought to be prevented by placing load balancers, web application firewalls and content-inspecting proxies along the edge of the estate, it is quite normal to see traffic from the perimeter make its way through the infrastructure and arrive at one or more nexuses or intersections, along with traffic from other edges.

**Perimeter applications**

Region 1 load balancer

Region 2 load balancer

Nexus

Firewall

Backend 1

Backend 2

Nexus
Data enrichment platform

Nexus
Inspecting proxy or AV

**Backend service/storage**

**Storage/backend/3ʳᵈparty**

This could, for instance, be as part of a second set of firewalls, a message queuing system or some other data enrichment solution. While the intended use may be to have traffic traverse the edge systems and go through the second layer, before returning the relevant response, when this happens from multiple edges at the same time it often results in those second layer components receiving more traffic than they were ever designed for.

This can in turn lead to the two other forms of DoS discussed earlier but falls into the category of logical DoS as the logical topology of the environment and the logical placement of these nexus' are what induces the DoS condition. And as such a nexus is often a critical or expensive component, it is not abnormal for there to only be one such appliance or solution, which at the same time makes it a single-point-of-failure, compounding the impact of the DoS.

# Detection

The first step in managing the risk and potential impact of a DoS situation, is to be alerted to its imminent presence. Preferably ahead of time. To do this, detection capabilities must be present, and these can take a number of forms. Using Content Delivery Network (CDN) providers specialising in anti-DDoS services (such as Akamai) is currently part of everyone's' recommendations, and very much works as advertised, although they usually only trigger alerts when the DoS attack is in full effect. While this part of the paper seeks to identify additional controls aside from CDNs and traffic sink holing, the following outlines the key requirements that should be posed to prospective anti-DDoS vendors, and the insights a high-quality vendor should provide as part of the service:

## Prevention

- Automated rate controls that block traffic based on vendor and client-aligned thresholds
- Standard and custom Web Application Firewall (WAF) rules based on public and proprietary signature feeds, as well as custom attack signatures
- Tools to establish a baseline of customer-related traffic, and rapidly compare current data to said baseline to identify indicators of attacks. This should include real-time monitoring of items such as headers (especially user agents) to identify emerging deviations from the baseline. Such deviations could trigger pager-style alerts when high or peak traffic volumes are observed
- Ability to redirect or sinkhole traffic, or otherwise perform traffic reduction

## Analysis and insights

The following are the key insights that any anti-DDoS vendor should be expected to provide as part of a high-quality offering within the anti-DDoS space, following an observed DDoS event:

- What was the source of the traffic?
  - by IPs participating in the attack
  - by ISPs participating in the attack
  - by geographical location
- What was the traffic type
  - by protocol and service
  - by attack type (amplification, fragmentation, flooding, botnet, etc)
- What was the traffic load?
- What was the duration of the maximum load?
- What was the ebb and flow of the DDoS? (Indicating a ramp-up in load, or that traffic-generating nodes were rapidly detected and dismantled)
- Who were the primary contributors to the attack?
  - Were the contributors part of known botnets?
  - Does any known signatures or Tactics, Techniques and Procedures (TTPs) match the traffic generated? Which ones?
  - If part of known clusters, what is the maximum expected throughput from the threat actor?
    - Does the vendor expect to be able to handle this?
    - What preliminary steps should be taken to prepare for an event that utilises the maximum throughput?
  - Do any of the participating nodes correlate with other DDoS attacks observed by the vendor within the past three months?
- Is the vendor aware of any clients in similar verticals currently experiencing similar attacks?
- Does the attack appear to be coordinated and executed in synchronisation with similar clients protected by the vendor?
- Is the vendor able to provide any attribution or suggestions for mitigating controls to reduce the impact of future events?
- Considering the geographical origin of the attack and the geographical location of the target, does the event coincide with any regional holidays or special events? (Christmas, New Years, Black Friday, etc.)

# Early detection

## DNS log monitoring

As most larger organisations will be in control of the authoritative DNS server for their domain, live collection of DNS logs from said server should be in place. If a botnet was to be tasked with targeting any domain or subdomain belonging to the organisation, each individual node would first need to resolve the domain to an IP address to attack. To reduce latency and the amount of needed DNS lookups, subsequent attack activity will traditionally be targeting that IP. Due to the traditionally dispersed nature of botnets, nodes would be geographically spread out across the globe and by extension be subject to any number of local or regional DNS caches and servers. As such, when a botnet is tasked with targeting a specific domain, the initial domain name resolution from all nodes would result in a sudden and markedly rise in resolutions on the authoritative DNS server (as it is statistically unlikely that all local or regional DNS servers will

have the company domain/subdomain cached or available), which could be observed by the victim company. As botnets consist of many different devices, of varying performance, synchronicity is not necessarily part of the botnet design, and a spike in DNS resolutions might be experienced suddenly or as rising over time. If the spike is sudden, it could be indicative of an imminent event, whereas a rising spike would indicate

that the botnet is still in the process of gathering strength and preparatory actions can be started. As the DNS logs will not only include the source IP addresses of the regions (regional DNS servers or nodes in a botnet) trying to resolve a domain, they will also reveal the potential target of the attack in the form of the domain or subdomain being resolved.
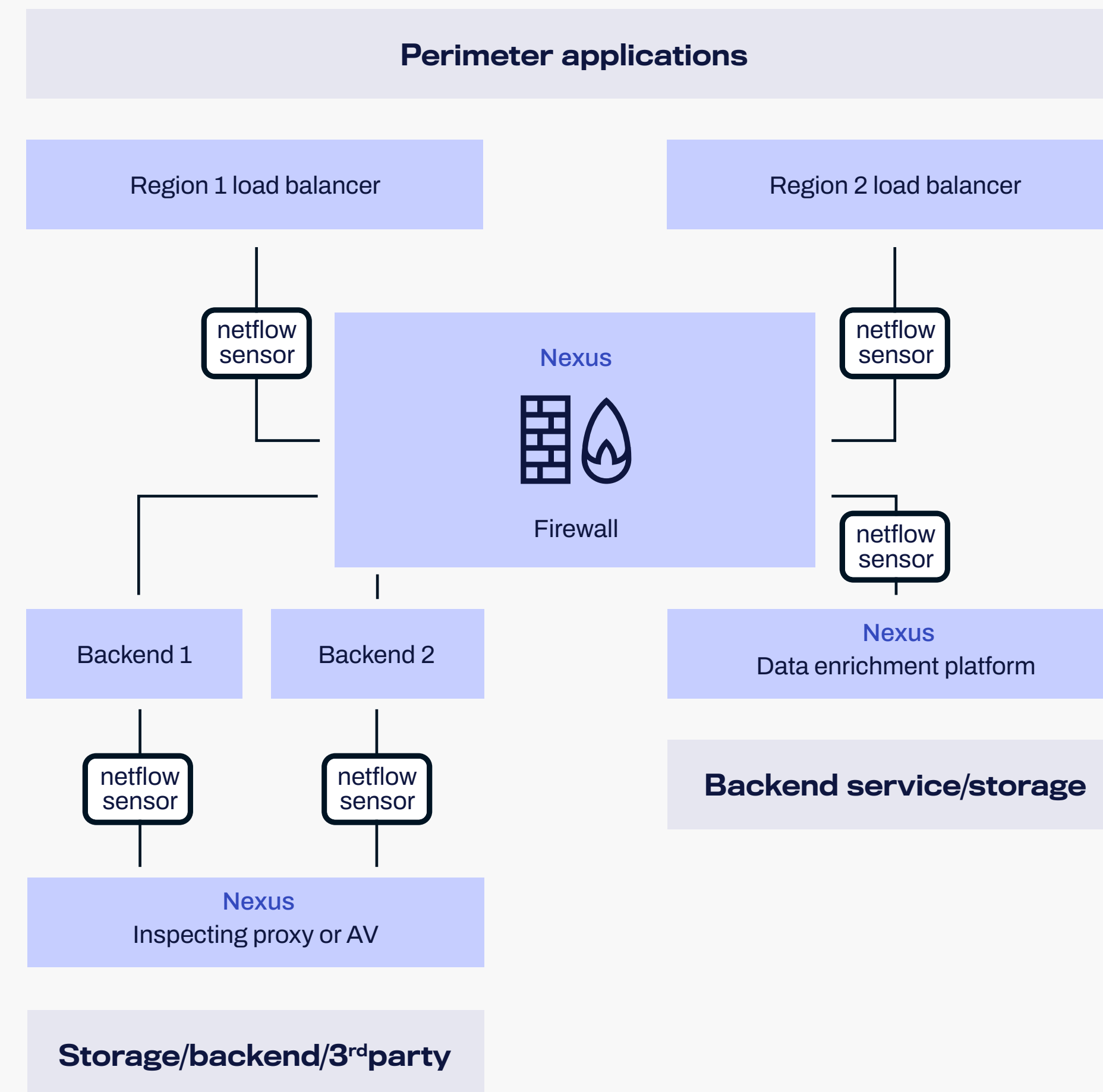
# NetFlow monitoring

While it is impractical to log and parse the content of every single connection that reaches the edge of the organisation's network perimeter, the collection of NetFlow logs can provide valuable information with a much smaller footprint in terms of resources and cost. Collecting and monitoring NetFlow from perimeter devices will allow the company to identify sudden spikes or increases in incoming traffic from specific regions (or just overall traffic volume in general). As NetFlow are unidirectional, the footprint of the collection can also be reduced by only collecting incoming sessions, which make them ideal for long-term collection and allow for pattern identification in old data. Based on timestamps, any spikes in traffic flows can be correlated with spikes in DNS resolutions as previously described, which in combination can reveal what asset is being targeted and if the incident is originating from a specific region.

NetFlow records can likewise be collected at each step between the edge perimeter and infrastructure components. This will let you identify what components are at risk of being oversaturated when traffic flows from several edge devices and converge on a single device on the internal network. Sudden spikes in traffic on such internal devices can then be correlated based on timestamps with the flows of all edge devices, which can again in turn be correlated with spikes in DNS monitoring as previously described. Such information can also inform decisions about potentially restricting access from a geographical region to reduce impact of the rest of the network, although it should be noted that such a response would be vulnerable to false-flag operations and could be abused to significant effect by threat actors. For this reason, the decision to restrict access from a geographical region should not be automated, and the origin and type of traffic should be close reviewed before making such a decision, to ensure that this is not a threat actor abusing this process to deny services to their chosen region.

By subscribing to threat intelligence sources from regional and/or industry-specific SIRTs/CERTs or partnering CSOCs, this type of collection also allows the organisation to proactively and retroactively review any interactions between the estate and known-bad actors, and commensurate steps can be taken to block access ahead of time.

**Perimeter applications**

Region 1 load balancer

Region 2 load balancer

netflow sensor

netflow sensor

**Nexus**

Firewall

netflow sensor

Backend 1

Backend 2

**Nexus**
Data enrichment platform

netflow sensor

netflow sensor

**Backend service/storage**

**Nexus**
Inspecting proxy or AV

**Storage/backend/3rdparty**

# Resource monitoring

A common reason for logical DoS is when lower-tier network devices, or core appliances inside the perimeter, crash. Examples of these are distribution-layer switches and routers, content-inspecting proxies, stateful firewalls, or hypervisors and similar devices hosting virtualized components. For all such infrastructure components, basic health-monitoring should be in place to measure things such as CPU and RAM utilization, power consumption, bandwidth saturation and available storage, and baseline values should be established for all of these. This should apply to both virtualized and physical devices in the network path between the Internet and the responding system. Alerts should trigger when significant deviations from the baseline are observed, or an upward trend in utilization is experienced. On such occasions, these observations should be correlated with NetFlow logs and analysis, to determine if the event is triggered by external stimulants such as an attack, or if this is unrelated.

# Managing Risk

This paper assumes that all relevant hardening has already been done on edge devices. This should include things such as ensuring the edge devices (firewall / load balancer etc.)

are correctly resourced with respect to the expected number of users and volume of traffic (while keeping in mind the best practice of always designing networking solutions to be 20% bigger than what was actually requested), include the devices in your corporate patch management policies and procedures, disabling un-needed features and blocking administrative access across all interfaces, as well as employing centralised and monitored AAA solutions. It is also assumed relevant mitigation configuration changes are in place such as;

• timeout values have been set for such things as TCP sessions (to avoid half-open or SYN-flood attacks)
• lowering of ICMP and UDP flooding thresholds
• rate limiting traffic to specific endpoints
• dropping spoofed or malformed packets
• only relevant traffic is allowed to traverse the devices

Do not rely on vendor default settings on these point, as they are in the majority of cases designed for performance and usability, not security. As such, in dealing with an active DoS or DDoS scenario, this paper focuses on what can be done as the attack is in its early stages  and when it is fully in progress. The focus is on responding to it, and weathering the storm because no matter what advance mitigations are in place, it is impossible to achieve effective and long-term prevention when dealing with motivated and well-resourced threat actors.

# Your friendly neighbourhood ISP

While this will obviously vary depending on your region, Internet Service Providers (ISP) have no more love for DDoS attackers than you do, and will often be in a position to help you (at least part of the way) or this could and should be negotiated as part of entering into a contract with the ISP. Using the baselines discussed in the NetFlow monitoring and resource monitoring sections, visibility into when the limit of your Internet connection(s) is reached will be available to you. When nearing that limit, reach out to your ISP and request a temporary increase in bandwidth (if possible), and/or activate a secondary subscription as part of a high-availability or redundant setup. While such actions will buy you additional time, it is ultimately a matter of how much you are able to spend on procuring more bandwidth compared to the threat actors ability to do the same, and with DDoS you are likely to lose that race. Choosing a response action needs to be determined based on what level of attack (or what level of motivated and financed opponent) you wish to be able to withstand over time and is largely a contractual issue that needs to be negotiated with the ISP ahead of time and what additional cost such actions may incur.

When spikes in traffic are observed as described in this paper, selected traffic captures should be collected, and attempts should be made to fingerprint the traffic (user agents, IP addresses, headers, patterns etc.) so that a potential signature/rule for the offending traffic can be created. The objective being that the signature/rule can be applied to the edge devices themselves to help mitigate the impact. It is also worth reaching out to the ISP and share the signatures and rules, requesting their assistance in blocking the traffic further up stream. The ISP may also be able to fingerprint the traffic in more detail, and provider better high-volume filtering than what can be achieved with the company's own equipment.

# Downgrading the user experience

When dealing with DoS, the issue is not always the data flowing into an organisation, sometimes it is the data flowing out. Modern systems and applications are becoming extraordinarily complex beasts, and the relationship between requests and responses are no longer as symmetrical as they once were. Many applications now respond with a sizable JSON response with encoded graphic, file or other media content. In fact, if you observe any highly interactive content that is returned from servers following just a single small request, you quickly see how the outbound connection can be bogged down just as easily as the incoming one. This is why caching proxies were recruited into the mix, and any setup that hope to resist any form of DoS attack should have this configured and monitored religiously. This is not only for the transfer of content that requires resources but also for any dynamically generated content. The generation itself also requires resources, and this is highly asymmetrical compared to the effort required to trigger such a creation. This is why a common DoS vector is to identify requests that triggers exactly such resource-expensive activities and submit consecutive or concurrent requests of that type.

While monitoring the general health and resource consumption of application servers across the estate (as described in the previous sections), an effective albeit temporary solution to spikes is to trigger a switch between the normal, highly dynamic and highly interactive application, and a much simpler and potentially more static version of the application. Think of it as downgrading the user experience from the visual standards of the 2020's, to the text-based experience of the early 1990's. It might not look as good, and not all the features will be available, but you'll still be able to provide a basic service set for a lot longer, due to the instant reduction in required outbound bandwidth and local system resources. This can also be combined with advanced caching of data per user sessions, allowing you to present cached data that is unlikely to be need updating (such as profile data or terms and conditions) very rapidly, without needing to constantly query backend systems that might be under very high load. The net result is that it becomes possible to handle more requests, without needing to handle more actual traffic volume or incur system expenditure.

# Hot failover load balancing

Whether deployed at the edge of the estate, or as part of a second layer or internal infrastructure components, firewalls, proxies, load balancers etc. are critical components and should almost always be deployed in a high-availability setup with either a cold, warm or hot standby device ready to go. While the standby device is obviously intended to be a backup, it also represents un-utilized resources and un-earned Return-On-Investment (ROI), which should be considered for activation in an emergency or in a crisis such as DoS instance. Depending on the configuration and the appliance in question, the backup device may be configured to operate in load-balancing configuration between the original device and the standby-device, effectively doubling the capacity of the infrastructure node.

An alternative to this, which can make this type of ad-hoc load balancing simpler to enact, is to configure all infrastructure nodes to run as part of a clustered setup, in which new nodes can be added and an appropriate store of extra nodes can be kept in lieu of standby devices. This can however be cost-prohibitive for all but the most critical network nexuses, and will also require that the cluster itself be considered in terms of cold, warm, or hot standbys for continuity assurance.

# DNS whackamole

If, from the DNS and NetFlow monitoring, you observe a specific domain or specific IP being targeted by an observed climbing spike, a decision needs to be made;

*Do we fight, or do we run?*

If the choice is made to stand our ground, the options of contacting the ISP for additional bandwidth, signature creation and blocking, recruiting un-utilized standby network devices and downgrading the user experience, this needs to happen now. But, knowing that this is likely not a fight that can be sustained, the glory might be found in running. When a rising spike in traffic for a specific domain is observed (or to an IP with multiple domains), a proposed strategy is to start lowering the DNS Time To Live (TTL) for the domain(s). The DNS TTL controls how long a domain-to-IP mapping is cached in a resolver, which in turn determines how much of a delay exists between a domain changing its IP address and this information being available to end-users. This exists to reduce the load on DNS resolvers and DNS servers. As the spike in DNS and NetFlow traffic continues to climb, the TTL should continuously be lowered, so that when the attack reaches critical mass the domain can be switched to another

IP address, leaving the DDoS attack to fire at an IP with nothing there (effectively sink holing it). With the TTL being as small as possible at this point, the number of legitimate users experiencing downtime as a result of the switch, should be minimal. With the domain now residing on a new IP, the attacking nodes would need to be re-tasked, which in turn would result in a new series of DNS resolutions, which you would be able to monitor for. In this fashion, as a defender, it becomes possible to continuously move the targeted victim domain around with minimal impact on end-users, while frustrating the threat actor and creating overhead until they decide to end the attack effort.

# Messaging

All of the detection and managing of a DoS attack described in this paper, are designed for one purpose: Buying time.

"Death. Taxes. Denial of Service" are the only three certainties in life, and there is no good way to beat any of them.

But given time and warning, the victim company can control the narrative, and keep the reputational (if not operational) damages of the attack to a minimum. A generic message for each of the following groups should be drafted ahead of time (with the advice of legal counsel as needed):

• employees
• stakeholders
• end-users/clients/the public
• Partners and supporters

Standardised messages are timesavers, but they read exactly as what they are; generic and non-specific to the case at hand. So, they should only serve as a starting point, and the early warnings provided by the controls described in this paper, should buy legal counsel and process owners time to polish and refine the message (and get started on potential take-down notices to ISPs). This ensures that you have a statement ready by the time the full impact of the attack is felt, that the wording is relevant and as accurate as possible, employees, stakeholders and other constituents can be put

at ease by knowing what is happening and what is being done to address the issue. This also permits you to inform users that the company is under sustained attack, which is why there is a reduced user experience. The document can also advise on what they can expect in terms of service returning to normal, and of the expected impact. Transparency should be the only policy and is more likely to generate support from all concerned.

Denial of Service attacks in all their forms are employed by threat actors of all stripes, from individuals with zero technical knowledge to state-sponsored Advanced Persistent Threat (APT) groups with what might appear to be unlimited resources. These efforts should help drive up the cost of attack for all levels of threat actors. The hallmark of an APT is not that their techniques are advanced, but that they are as advanced as they need to be. They are also persistent, and they will keep trying, increasing the level of complexity of their attacks. But as they do so, they expend time, they expend resources to a point where cost vs reward is unsustainable even for the most egregious group, and they might just realise that their time and effort is best spent targeting someone other than you. And that is the objective of dealing with DoS.

The details in this document will help deal with all the threat group types by driving up the cost of an attack and it is based on our collective experience in advising clients over decades across industry verticals.

# Who We Are

WithSecure™, formerly F–Secure Business, is cyber security's reliable partner. IT service providers, MSSPs and businesses – along with the largest financial institutions, manufacturers, and thousands of the world's most advanced communications and technology providers – trust us for outcome-based cyber security that protects and enables their operations. Our AI-driven protection secures endpoints and cloud collaboration, and our intelligent detection and response are powered by experts who identify business risks by proactively hunting for threats and confronting live attacks. Our consultants partner with enterprises and tech challengers to build resilience through evidence-based security advice. With more than 30 years of experience in building technology that meets business objectives, we've built our portfolio to grow with our partners through flexible commercial models.

WithSecure™ Corporation was founded in 1988, and is listed on NASDAQ OMX Helsinki Ltd.

W / T H®
secure