# Bringing Innovative Load Balancing to NGINX

Safiruddin Khan[1*], Syed Ubaid Ul Haq[2], Amir Khan[3]

[1,2,3]*Department of Computer Science and Engineering, Galgotias University, Greater Noida, India*

*Abstract*: **Load equalization remains a very important space of study in engineering for the most part thanks to the increasing demand on knowledge centers and net servers. However, it's rare to envision enhancements in load equalization algorithms enforced outside of pricy specialized hardware. This scientific research is a shot to bring these innovative techniques to NGINX, the business leading open supply load balancer and net server. In addition to implementing a brand new, native NGINX module, I even have developed a straightforward work flow to benchmark and compare the performance of accessible load equalization algorithms in any given production surroundings. My benchmarks indicate that it's possible to require advantage of a lot of refined load distribution techniques while not paying a major performance price in further overhead.**

*Keywords*: **Balancing, Innovative, NGINX, Load.**

## 1. Background

Ultimately, load equalisation could be a balls into bins problem: one should decide however best to distribute m balls into n bins such every bin has roughly identical range of balls. though this might sound easy, load equalisation has remained a troublesome drawback in engineering. the foremost difficulties ar thanks to the complexities of distributing tasks with 2 major unknowns: load and time. Load could be a task's demand on the server, whereas time is expounded to each the length of a task and its arrival. In short, load equalisation is difficult as a result of the arrival of a task, however long it'll fancy complete, and also the process resources it needs, are ne'er inevitable and invariably freelance of every alternative.

These factors not solely contribute to the complexities of designing load balancers, however they conjointly create it troublesome to model SOME surroundings for testing them. what is more, not all load balancers are identical. The balls into bins drawback shows up in several areas of computing, all over from computer hardware task programing to telecommunication depends on a load balancer to induce work done as efficiently as attainable. Figure one shows the everyday design for load equalization in high performance net server environments.

My analysis is driven by rising the performance of load equalisation on net servers as a result of there has not been the maximum amount innovation compared to figure done on the TCP/IP network stack or software schedulers. However, one thing of these areas have in common is that the underlying applied math model of however tasks arrive that need distribution. This model is most typically understood as a distribution [1], that is why I exploit them in my simulation environments to assign every request a singular weight representing their point in time and cargo on the net servers. Figure a pair of provides a visible illustration what Poisson streams seem like relative to the arrival times of requests at a given interval.
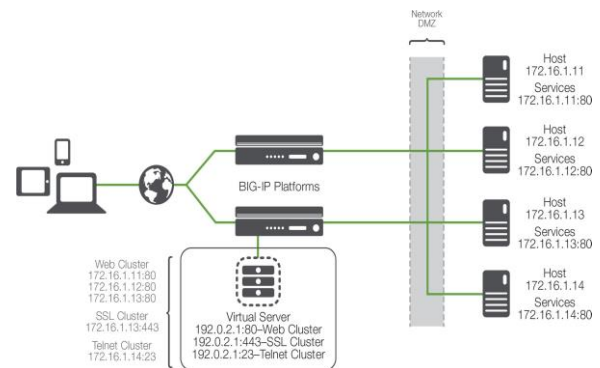

Fig. 1. High performance load balancing architecture

Web server load equalisation methods have hardly modified since their initial implementations. the 2 most well-liked algorithms are random and spherical robin (RR), the latter having a booming history in computer hardware programing, time-sharing systems, and DNS. These approaches work quite well underneath sure circumstances, however have vital drawbacks once considering however the net is employed these days. for instance, spherical robin works best only distributing requests of a homogenous length. once RR is employed as a computer hardware hardware, distinct time quanta are bonded, however this is often not the case for an online server, wherever requests have AN unknown length and cargo. Largely, these disadvantages are neglected as a result of random and RR appear to try and do a "good enough" job and a spotlight is primarily given to lower levels of networking and software style.

However, rising the power to distribute load as uniformly as attainable has many advantages that got to be thought-about. For one factor, an online application unfold across multiple servers victimization AN inefficient load balancer can end in one or 2 machines handling the bulk of the requests whereas others sit nearly idle. once this happens, it's common to feature another server into the surroundings as a result of it'll create it less doubtless for one machine to become full. this is often clearly not the simplest approach. By utilizing a much better

load equalisation algorithmic program, an online application will get the foremost out of every obtainable machine while not risking a premature upgrade. However, that's not all, reducing the overall range of further servers saves loads of cash, maintenance, and energy.

## 2. Project Description

There are a unit variety of load reconciliation algorithms that are shown to extend the performance of net servers once employed in place of random or RR, nonetheless few area units ever enforced in pre- vailing open supply comes. the most important advantage of mistreatment RR and random from a developer's purpose of read is that they area unit they're intuitive algorithms that are straightforward to implement and maintain. whereas dedicated hardware load balancers frequently cash in of recent innovation, the open supply community has been frequently left behind. My analysis is a trial to bring a number of the foremost recent and successful load reconciliation techniques into NGINX, one the leading open supply load balancer and net server.

Of these innovations, the formula especially that I would like to target originally comes from archangel Mitzenmacher's 2001 paper, the facility of 2 selections in randomised Load reconciliation. during this paper, Mitzenmacher outlines associate formula referred to as two-choices, that behaves exponentially superior to the standard methods like RR and random. Figure three illustrates the 2 selections formula in what Mitzenmacher presents because the "supermarket model", wherever a client desires to enter the smallest amount busy checkout queue. the thought behind 2-choices is that the economical shopper solely surveys two of the accessible queues and quickly enters the smallest amount huddled one. The less economical shopper fastidiously compares each queue before creating a choice. Mitzenmacher found that by choosing 2 random queues, it had been doable to avoid the ill-famed "thundering herd" drawback. If each client was seeking the smallest amount huddled queue, then at any given time, everybody is going to be sport towards one lane, mostly ignoring everything else. Once that queue fills up, another one is pursued down. With two-choices, multiple customers don't seem to be seemingly to be directed to a similar queue, however they're terribly seemingly to avoid the foremost huddled one.



```
import numpy.random as nr
import matplotlib.pyplot as plt

balls = 1000
bins = 8

s = nr.poisson(0.99, balls)
count, bins, ignored = plt.hist(s, bins, (0, bins))
plt.show()
```
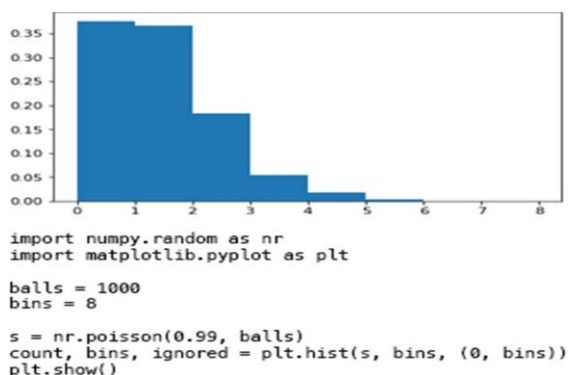
Fig. 2. Poisson distribution

The aim of my research is to review the behavior of those

breakthrough load reconciliation techniques in an exceedingly production environment. To accomplish this, I even have 2 goals: (1) Reproduce the work of Mitzenmacher et al. about the efficiency of varied load reconciliation methods. (2) Implement two-choices as associate NGINX module and take a look at it against the opposite accessible load balancers.

## 3. Experimental Setup

This project was ab initio impressed by a chat given by Tyler Mc- Mullen, titled Load reconciliation is not possible [5], wherever he outlines the challenges load balancers face once addressing the net as we all know it nowadays. I started my analysis by increasing the initial simulations given in his speak and shortly I used to be able to construct associate surroundings wherever I may reproduce the work given in analysis papers concerning the 2 selections formula.

I conducted my load reconciliation experimentation mistreatment associate Python notebook running within a python virtual surroundings as a result of it permits transportable and cross platform development. employing a Poisson stream with a mean of zero.99 as my request distribution model, I appointed a weight to every request to represent its arrival on the server. within the Python notebook I model the load reconciliation within the following way: there's a listing of length n representing the requests and a listing of length m representing the accessible servers. The re- quests area unit passed to a load reconciliation formula that increments a counter happiness to a specific server by that request's weight. in spite of everything requests area unit distributed, the quality deviation of requests among every server is compared between algorithms. an ideal load distribution would so have a customary deviation of zero.

The algorithms I enforced were random, round-robin, and two-choices: Random chooses a server for every request severally and uniformly arbitrarily, RR distributes the request to every server one by one, and 2-choices initial selects two servers severally and uniformly arbitrarily so chooses the server with the smallest amount load to method the request. Figure five provides smallest formula implementations employed in my initial testing surroundings and provides a far better sense of however my Python simulation was organized. The later stages of my analysis was done mistreatment special configuration files that enable my load reconciliation module to be dynamically joined to the system installation of NGINX. additionally, I utilised the Go artificial language to build an online server that compiles into a native binary for execution on multiple machines and ports.

All of the software system elements employed in my analysis area unit provided at intervals one organized dirty dog repository.

### A. NGINX Module Development

I developed 2 load equalisation modules for NGINX: random and two-choices. The underlying load balancer for NGINX is RR, however it additionally provides a module referred to as least_conn, which can distribute requests giving preference to the server with the smallest amount connections presently

established. The two choices module is enforced by incorporating the practicality provided by least_conn and my new random module. each modules are compiled and dynamically coupled into the system installation of NGINX as a result of it makes development abundant easier. However, each modules will be statically coupled if desired. though NGINX provides AN API for writing modules in perl, I selected to implement them directly in C to eliminate any potential overhead that will skew the results. I additionally take into account native NGINX module implementations a lot of helpful to the open supply community.
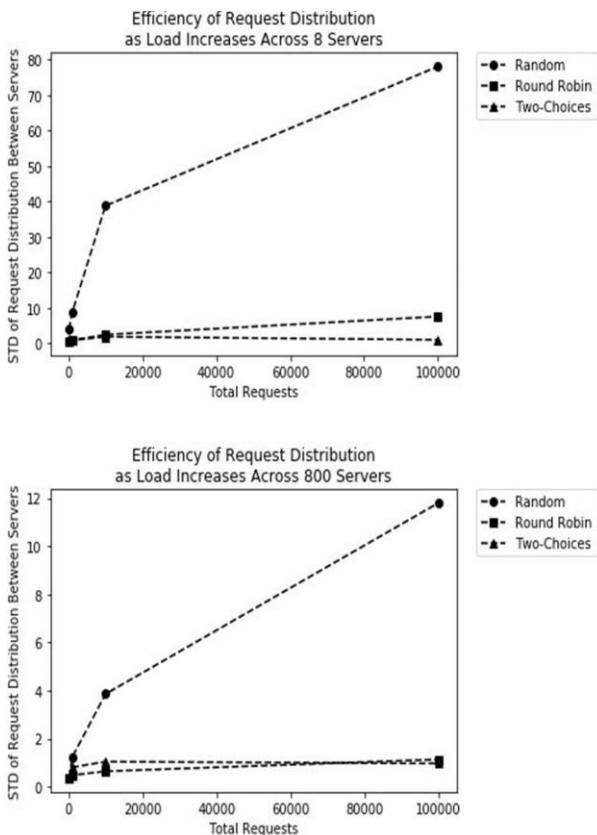




Fig. 3. IPython notebook simulation results for random, round robin, and two-choices

In order to check the effectiveness of the load equalisation algorithms, I created an easy webapp in Go which will simulate my production webserver setting. Go is a wonderful language to use for this task as a result of it's an intensive hypertext transfer protocol package within the customary library, compiles to native machine language, and doesn't would like any extra dependencies to host a webserver.

The Go webapp generates a Poisson random range for every incoming request. This range is then accustomed verify however long the webapp can sleep for before causation back a response. I try this to simulate the unpredictability of request length and cargo on the server. I selected to model my webserver setting with the Poisson method as a result of it's well understood and normally accustomed model the behavior of net track. Naturally, this may not offer AN correct model for all production net applications, however, I even have created a

workflow for benchmarking the performance of all NGINX load balancers, together with two-choices, on any given system. This workflow can permit anyone to look at the performance of every formula in their own production environments.

### B. Apache Bench Testing Strategy

The trade customary tool for benchmarking and measure net server performance could be a program line utility referred to as Apache Bench, 4 or ab. The interface is sort of straightforward, it permits you to specify several what percentages what number total requests to send to a web site and the way many ought to be created at the same time. when causation the requests, ab can offer some helpful info like the whole time to finish the requests, requests processed per second by the webserver, and therefore the average time spent per request. I exploit these metrics to measure the performance of the load balancers on NGINX additionally to graphing the latency of every request within the benchmark.

## 4. Results and Discussion
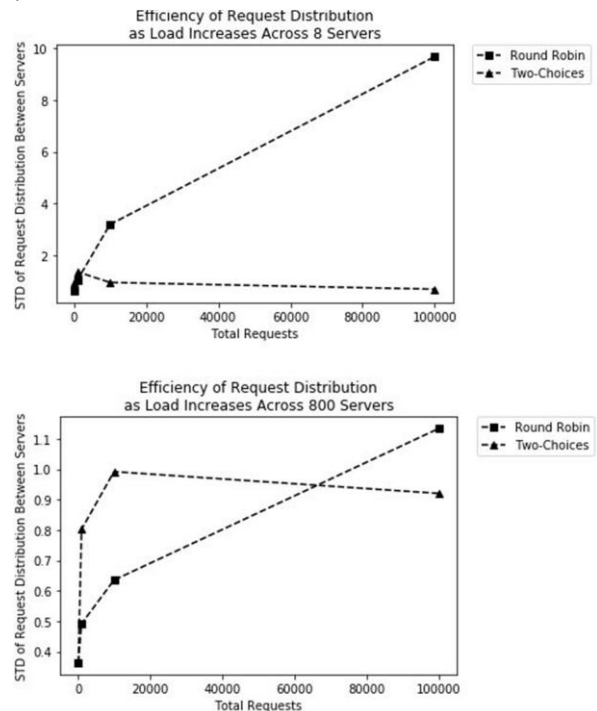
### A. Python Simulation Results





Fig. 4. A closer look at the load distribution capabilities of round robin and two-choices

My initial simulations rearmed the results bestowed by Mitzenmacher. once requests are weighted, the quality deviation of two-choices approaches zero because the quantity of requests being processed will increase. As Figures six indicates, RR will far better than random, however has AN increasing variance as requests in- crease. Figure eight highlights a very important observation: RR continually completes within the smallest amount of your time, whereas two-choices takes quite double as long to run. Additionally, value noting is that once the quantity of servers is accumulated, RR performs a lot of equally to two-choices, however, Figure

seven confirms that two-choices is clearly higher at maintaining a standardized distribution of requests across all out their servers. though my experiments rearm that two-choices is that the superior formula as so much as load distribution, the results raise a very important question. However, can the overhead of two-choices have an effect on the latency of a production net server?.

*B. NGINX Simulation Results*

My intensive benchmarking disclosed no obvious distinction be- tween load equalisation algorithms running in NGINX. despite the active module, performance remained regarding an equivalent. How- ever, there have been some general trends relating to coincidental and total requests that were anticipated, namely, after you fiood your net server with requests, it takes longer to reply.

What these results do indicate, is that the overhead of a load balancer could become negligible once taking under consideration the whole overhead related to finishing AN hypertext transfer protocol request. within the earlier simulations with python, I used to be involved that the accumulated latency of two-choices would build it AN inconvenient load balancer in a very production setting. However, my results show that we tend to could also be able to make the most of two-choice's uniform load distribution skills while not paying abundant performance penalty.

Yet, the shortage of a transparent distinction in algorithms could be a concern. it's an honest indication that my experimental setting isn't capable of simulating the conditions necessary to create high performance load equalisation noticeable. I'm not utterly afraid as a result of mistreatment ab to benchmark webserver performance is AN trade customary. Although, using a custom benchmarking technique for these experiments could have created a lot of obvious results. thereupon being aforementioned, I'm still assured within the viability of two-choices as a load balancer when running these experiments.

Additionally, the machine running all simulations will solely launch up to eight net servers, every process up to a hundred coincident connections from ab. five whereas it's doable that the load reconciliation modules ought to be tested with AN NGINX configuration containing many servers, it's going to be AN impossible expectation. once I at first contacted the NGINX list concerning my research, lead developer Maxim Dounin responded that algorithms like two-choices have not been thought-about for implementation as a result of it had been unlikely to result performance unless one was victimization NGINX during a very giant computing setting.

Most of my findings are summarized by Figure nine. once the quantity of coincident connections is unbroken comparatively low, every load reconciliation module behaves nearly identical. However, as we tend to increase the coincident connections, we tend to see that the overwhelming majority of requests are completed underneath five hundred ms, however close to five-hitter of requests take thousands of milliseconds longer to finish. This behavior could be a notable issue with victimization Apache Bench, however it conjointly addresses

the matter load reconciliation tries to resolve. That is, once an internet server becomes full, it's terribly exhausting for it to recover.

The stair-step pattern drawn in these graphs unsurprisingly correspond directly with my statistical distribution. every incoming request can pay either zero, 100, 200, or three hundred ms on the net server before obtaining a response. the very fact that we will visualize the Poisson stream nearly precisely is another indication that the overhead of load reconciliation is negligible underneath these testing conditions and NGINX.
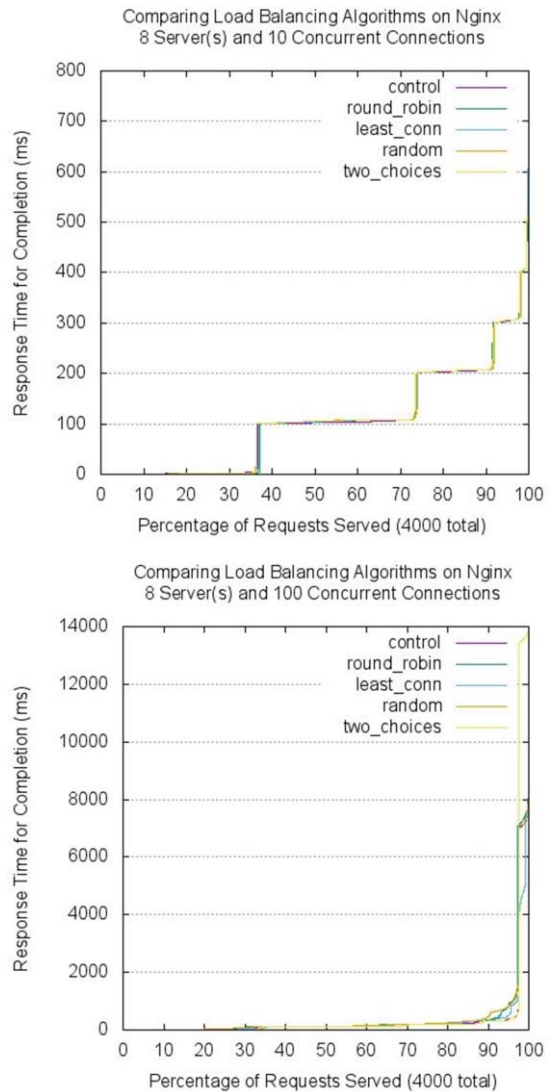


Fig. 5. Each load balancing algorithm has near identical performance in NGINX according to the ab results

In order to urge a much better sense of those apparently homogenous results, I created another mental image for examining the mini- mum, maximum, and average request latencies of every algorithmic rule. it's doable to look at some further trends victimization these new charts. Figure 6, rearms that underneath lower concurrency levels, performance is pretty uniform between algorithms. However, it remains unclear if any algorithmic rule is superior underneath high levels of

concurrency. whereas it seems two-choices could often have a plus, Figure eleven is a reminder however a number of latency outliers from Apache Bench will skew the graphs considerably.
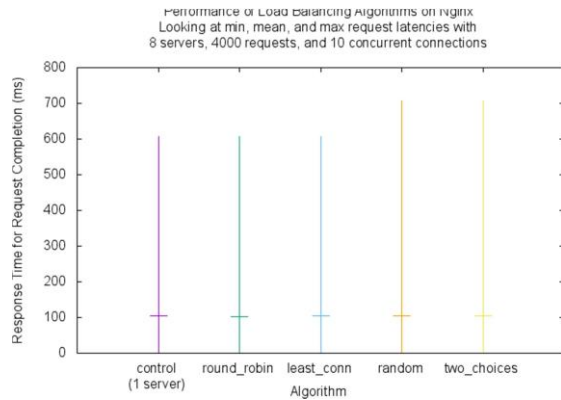


Fig. 6. Under low levels of concurrency, there are less outliers so it's possible to see the slight variations in performance

## 5. Related and Future Work

Overall, I'm excited by the outcomes of my capstone research. If I continue running experiments on a lot of subtle server environments, I hope to urge a lot of refined result set that may result in a much better understanding of NGINX load reconciliation performance. I attempt to contribute my 2 selections module upstream to the NGINX project further as answer any feedback I'll get from the opposite open supply developers. to boot, it'd be worthy to assemble a lot of information and analysis production net application server load a lot of completely. The statistical distribution could be a nice applied math model for a proof-of-concept, however my analysis would undoubtedly take pleasure in a richer applied math dataset. Load reconciliation for the foremost half is primarily a priority for big firms and information centers. For this reason, a lot of my background analysis concerned learning however the massive school firms are approaching this drawback. The prevailing ways to the load reconciliation drawback sometimes involves improvement deeper among the networking stack, wherever the matter may be a lot of discretely outlined and a lot of usually applied.

### A. Microsoft's JIQ

Join-Idle-Queue is that the latest and greatest load reconciliation algorithmic rule. it had been developed by Microsoft and achieves larger performance than two-choices and another competitive algorithmic rule referred to as join-shortest-queue. However, JIQ doesn't introduce communication overhead on the servers. this can be achieved by solely victimization native in- formation concerning server load. the concept behind JIQ is to "decouple discovery of gently loaded servers from job assignment" [3]. this can be achieved through utilizing idle CPUs to create the load reconciliation call. JIQ out-performs the competitive advanced load reconciliation algorithms and far like my results, Microsoft notes that these load reconciliation ways are most noticeable underneath very high server load.
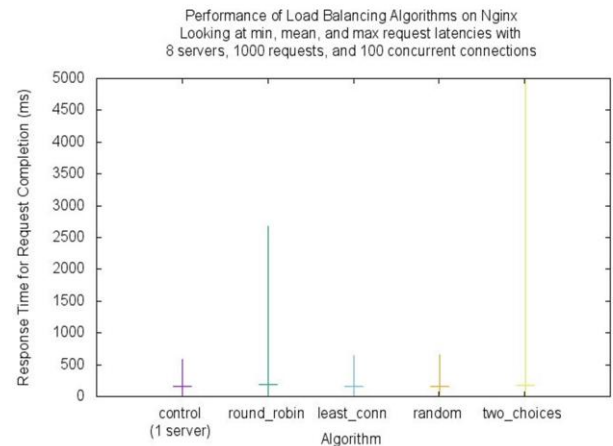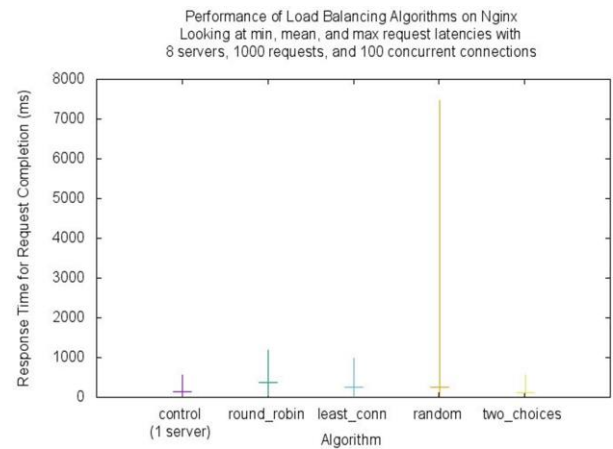




Fig. 7. Although ab is a great benchmarking tool, results are often inconsistent due to a few outliers

### B. Google's BBR

BBR stands for Bottleneck information measure and Round-trip propagation time. it's a brand new congestion management algorithmic rule developed and deployed by Google for increasing the output of communications protocol [2]. the aim of the algorithmic rule is to live the present bottleneck of the network and solely send enough information to "fill the pipe". The success of the algorithmic rule comes from activity network congestion in terms of its bottleneck rather than packet loss, that is however it's tradition- ally done. to boot, it had been found that most output is achieved once the loss rate was but the inverse sq. of the information measure delay product (BDP). BBR is already enforced within the Linux kernel for communications protocol.

### C. Facebook's Egress

Egress could be a method for evaluating network latency and congestion through "performance aware routing" on Facebook's network. The Egress paper explains some key components of running a network on an enormous scale that minimizes congestion. What Google did with communications protocol congestion, Facebook did with the border entrance protocol (BGP); they created it "capacity and performance aware". primarily, Facebook had to optimize its purpose of presence (PoP) servers to possess extremely efficient routing algorithms by establishing shorter ways, to deliver content to its billions of users. This paper illustrates a standard theme that

ancient implementations of networking protocols aren't any longer sufficient.

### *D. Linux Socket Balancing: Epoll-and-Accept*

An interesting downside concerning NGINX was mentioned by Marek Majkowski of CloudFlare, wherever he examines however UNIX schedules connections to sockets [4]. NGINX, like several applications might produce multiple employee processes to extend performance at scale. On Linux, these processes communicate over sockets. On NGINX, one socket "listens" for brand spanking new connections then distributes them to at least one of the out their employee processes. This behavior is strictly just like the load equalization mentioned during this paper, except that rather than process a call for participation on another webserver, at this level, NGINX distributes new connections among OS processes. it's additionally potential to own a model wherever their square measure multiple listening sockets and multiple employee processors. sadly for UNIX, once distributing connections between sockets victimisation epoll() to avoid obstruction on the accept() supervisor call instruction, the programming behavior becomes Last-In-First-Out (LIFO). That is, the busiest method is selected most frequently. a bit like the thundering herd downside, this ends up in Associate in Nursing unbalanced employee method load and a decrease in NGINX performance. However, by setting the SO_REUSEPORT socket possibility, every employee method can have a lot of uniform load at the price of upper latency.

### 6. Conclusion

Load balancers square measure a key part in trendy distributed systems. There square measure 2 general categories of load balancers: L4 and L7.

Both L4 and L7 load balancers square measure relevant in trendy architectures. L4 load balancers square measure moving towards horizontally ascendable distributed consistent hashing solutions. L7 load balancers square measure being heavily invested with in recently because of the proliferation of dynamic small service architectures.

Global load equalization and a split between the management plane and also the information plane is that the way forward for load equalization and wherever the bulk of future innovation and industrial opportunities are found.

The trade is sharply moving towards artifact OSS hardware and software system for networking solutions. I think ancient load equalization vendors like F5 are displaced 1st by OSS software system and cloud vendors. ancient router/switch vendors like Arista/Cumulus/etc. I feel have a bigger runway in on premise deployments however ultimately will be displaced by the general public cloud vendors and their native physical networks.

Overall, I feel this is often a desirable time in pc networking! The move towards OSS and software system for many systems is increasing the pace of iteration by orders of magnitude. what is more, as distributed systems continue their march to dynamism via "server-less" paradigms, the sophistication of the underlying network and cargo equalization systems can ought to be commensurately enlarged.

### References

[1]  Dimitri Aivaliotis, "Mastering NGINX," Second Edition.
[2]  Martin Fjordvald, Clement Nedelcu, "NGINX HTTP Server."
[3]  https://www.digitalocean.com/
[4]  https://docs.nginx.com/
[5]  https://upcloud.com/