

More Bang for the Buck: Superlinear Scaling with Distributed Self-adjusting Systems

Paper #901, 12 + 7 pages

ABSTRACT

Conventional wisdom suggests that linear scaling of the worker pool in a distributed system can result in at most a linear performance improvement. In this paper we show that distributed systems can be *systematically* architected to achieve faster-than-linear (superlinear) scaling. Our insight is that dispatching jobs to parallel workers so that the locality of reference in the workers' input increases, and implementing the workers with a self-adjusting algorithm to take advantage of the higher locality, jointly yield superlinear scaling. We demonstrate the applicability of our methodology in extensive simulations: scaling textbook self-adjusting algorithms we obtain 100–3,300x speedup using only 48 CPU cores, up to 70x beyond linear scaling. Then, we present two operational case studies. Using our architectural blueprint to scale a Memcached+PostgreSQL storage system we attain 2.3x faster than linear scaling. Then we re-engineer the Linux packet classifier to self-adjust with load, obtaining 800x speedup on synthetic traces and 220x speedup on real firewall traces with 32 CPU cores, resulting 5–25x times raw performance improvement compared to the vanilla Linux kernel.

1 INTRODUCTION

With the end of Moore's law, computing power in modern systems increasingly comes in the form of parallel processing resources. A major obstacle faced by network engineers is how to harness this increasingly parallel computing power for scaling distributed systems [49, 72, 79, 86, 93].

In horizontally scaled applications a load balancer dispatches jobs across a fleet of workers that process the jobs in parallel [19]. In the context of *web applications*, HTTP load balancers [13, 23, 67] distribute requests across a swarm of backend web servers. Multicore *OS network stacks* [12, 50, 90] run multiple instances of the networking logic on different CPUs and leverage the NIC to dispatch packets to CPU cores. In *sharded key-value stores* [30] different servers handle different portions of the key-space. The system's overall goal is to achieve the greatest possible parallel speedup with a given number of workers, in order to minimize the execution time of a single task or maximize the number of completed tasks in a given time period.

Suppose a web app handles 100 requests per second using a single server. As we add another server we expect the throughput to increase to 200 requests per second, or slightly less if the

system is not perfectly parallel [4]. (Sub)linear scaling feels almost evident in this context: as the additional capacity can be consumed at 100 percent efficiency at best, we can “get at most equal bang for the capacity buck” [35]. Strikingly, faster-than-linear (or *superlinear*) growth seemingly defying this common wisdom has been observed experimentally in many high-performance computing applications and distributed systems [10, 18, 26, 39, 40, 42–44, 75, 80, 81, 87, 88]. Superlinear growth feels particularly alluring in this context, in that it assumes a system that somehow manages to produce more work than the computer capacity available to it [35].

In fact, faster-than-linear scaling is not supernatural at all. Suppose our sample web app serves a set of static assets (web pages, images, etc.) and assume each server is assigned a fixed subset of the assets, with a load balancer carefully routing client requests for each asset to the proper server. As the number of servers increases each server perceives requests to a progressively smaller subset of the assets, which may allow it to finish servicing requests faster, say, by *caching* the most popular assets in fast memory [26, 80]. The combined speedup resulting from the increase of web server capacity and the decrease of servers' “virtual job size” due to improved cache efficiency often yields superlinear scaling [18, 40, 40–42, 44, 75, 87, 88, 92]. This is, however, extremely sensitive to subtle design choices and a poor implementation can easily destroy the delicate superlinear scaling trend (see §2).

Despite that superlinearity is genuinely measured [10, 18, 39, 42, 43, 81] and thoroughly dissected [35, 40, 44, 75, 87, 88] in piecemeal applications, currently there is a lack of a general architectural blueprint that would help system architects to methodologically *engineer distributed systems towards superlinear scaling*. In this paper we aim to fill this gap. Our motivation is that networking applications are often embarrassingly parallel with little or no dependency between threads, promising massive (superlinear) parallel execution gains.

We observe that in order to achieve superlinearity one has to carefully combine an appropriate load balancing policy with a proper worker implementation. Indeed, load balancing in distributed systems is often non-arbitrary: web apps apply the “sticky sessions” rule to route all requests of a particular user to the same web server; networking code commonly uses IP 5-tuple hashing at the NIC to ensure that all packets of a flow are processed on the same CPU; and key-hashing in sharded key-value stores concentrates queries to a key at the same replica. Such policies tend to make the input streams

processed by the parallel workers more predictable, compared to the aggregate input processed by the system. Combining such a *locality boosting load balancer* with a *self-adjusting algorithm* so that workers can take advantage of the higher input predictability to adaptively improve their own performance will, as we show both theoretically and empirically, yield faster-than-linear speedup in a broad range of applications.

The power of the resultant *distributed self-adjusting system architecture* we advocate in this paper (§3) is *not* that it confirms the existence of superlinear scaling (this has been known for a while [87, 88]), neither that it would defy well-established scaling laws (it does not, see [25, 34, 35, 44]) nor that it produces the most efficient implementations possible (e.g., our packet classifier will not be as efficient as, say, a DPDK equivalent [78] just by the fact that it runs inside the Linux kernel [27]). Rather, our main contribution is that we precisely *identify the main architectural ingredients*, locality boosting and self-adjustment, *which in combination allow a distributed system to scale faster than linear*. Our system architecture will then provide a mental model that guides us in re-engineering several commonly used distributed systems with surprisingly little effort to attain real and tangible performance improvement, often in the range of several orders of magnitude as evidenced by our subsequent case studies.

First we confirm the viability of our system architecture in extensive simulations. Deploying well-known list and tree search algorithms from the literature we achieve 100–3300× speedup on 48 CPU cores, orders of magnitude surpassing plain, linear scaling. We support our empirical findings with a formal analysis and obtain a new generic scaling law for distributed self-adjusting systems (Appendix A). Then we present two fully operational case studies. As a major contribution we re-engineer the packet classifier built into the popular Linux kernel to reach superlinear scaling (§4). On synthetic and real-life firewall traces, our implementation exhibits up to 800× speedup with 32 CPU cores, 5–25× improvement beyond the default Linux firewall implementation which scales only (sub)linearly. We also apply our methodology to a combined Memcached+PostgreSQL storage system, yielding 2.3× faster than linear scaling (discussions moved to Appendix B for space reasons). We finally review related work (§5) and summarize the guidelines for designing and implementing superlinearly scaling distributed systems (§6). All code will be available as open source after publication. This work raises no ethical concerns.

2 BACKGROUND

There is an extensive background on scaling laws for characterizing the performance of a parallel system as the function

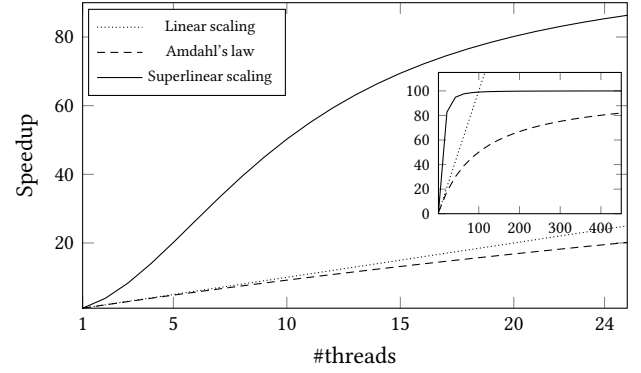


Figure 1: Linear scaling, Amdahl’s law and superlinear scaling ($s=0.01$). The inset shows the asymptotics.

of the computing/storage capacity available to it. In the following, we will use the terms “distributed” and “parallel” interchangeably to connote a networked system scaled to multiple independent compute threads (“workers”), e.g., scheduled to parallel CPUs of the same node, distributed to separate nodes, run in multiple datacenters, etc.

A cornerstone result in parallel computing, Amdahl’s law [4] establishes a firm limit on the performance gain one can obtain by distributing a computation task over multiple processors. Given a partially parallel program, denote the fraction of execution time spent in the sequential part of the code by s , and the parallel fraction by $(1-s)$. Here, some code is “sequential” if it cannot benefit from the improvement of parallel computing resources, like single-threaded code, critical sections guarded by exclusion locks, etc. Denote by $T(k)$ the runtime (in seconds) of the program when executed on k processors, and let $S(k) = \frac{T(1)}{T(k)}$ denote the performance improvement relative to a single-threaded execution (i.e., the *speedup*). Then, the following holds (see Fig. 1):

$$S(k) = \frac{T(1)}{T(k)} = \frac{1}{s + \frac{1-s}{k}}. \quad (1)$$

Here, $\frac{1-s}{k}$ establishes that the perfectly parallel part of the program executes k times faster on k processors than on a single core. By Amdahl’s law, (i) no code can scale faster than linear (i.e., $\frac{dS(k)}{dk} \leq 1$, with equality exactly when $s = 0$), (ii) throwing additional workers on a computation task yields diminishing returns ($\frac{dS(k)}{dk}$ is monotonically decreasing in k) and (iii) the asymptotics is limited by the sequential part only ($\lim_{k \rightarrow \infty} S(k) = \frac{1}{s}$). For different applications and extensions of Amdahl’s law, see [14, 20, 36, 46, 61, 69].

Curiously, there have been several reports from a broad range of applications indicating faster-than-linear scaling, e.g., database systems [26, 42], distributed storage systems

[18, 80, 88], SDN analytics [81], high-performance computing applications [39, 40, 75], multi-robot systems [43], information retrieval systems [87, 88], and large-scale network simulations [10] (see full taxonomies in [44, 75]). One way to reconcile these empirical observations and Amdahl’s law is the *scaled size model* [40]. Critical to Amdahl’s law is the assumption that the size of workers’ sub-problems remains constant as we scale the system [41]. Under this *fixed size* assumption [40], faster-than-linear scaling is impossible [25]. However, when this assumption fails, say, when the workers’ jobs get progressively smaller or execution gets gradually faster as we add more parallel workers (scaled size model), superlinear scaling often emerges [39, 42, 43, 81].

Sometimes faster-than-linear growth appears almost accidentally. Imagine a naive parallel dense matrix-multiplication algorithm that factors input matrices into multiple blocks, performs the multiplication of the blocks in parallel, and aggregates the results [75]. Easily, blocks will get smaller as we add more processors, so that after a certain point the entire input of workers will fit into CPU fast cache, yielding a disproportionately faster parallel execution. Conditions under which such superlinear (or “super-unitary” to be absolutely precise [44]) scaling emerges are widely discussed [40, 87, 88], analyzed [44, 75], and debated [25, 34, 35]. What is missing is a generic design methodology to *engineer* distributed systems for superlinear scaling. Such a model would also help identify the cases when superlinear scaling is possible, and when it is not. Our main contribution in this paper is a new system architecture to fill this gap.

3 DISTRIBUTED SELF-ADJUSTING SYSTEMS

We now present a general distributed systems architecture, which, as we show theoretically and empirically later, produces faster-than-linear scaling in several problem domains. Our main observation is that, whenever genuinely observed, superlinear scaling assumes two critical components: a policy to dispatch jobs to workers in a way to increase the locality of reference in workers’ input streams, plus an algorithm that can adaptively exploit the structure in the input to process it more efficiently. Our architecture is purely a software technique in that it does not require the addition of new cache space to a system. Nonetheless, it contains distributed caching as a special case and hence automatically takes advantage of additional fast memory, if available.

3.1 Locality-boosting load balancing

The first crucial component in our architecture is a locality-boosting load balancer. In this context, load balancing refers to the distribution of computational work or incoming network traffic across multiple parallel *workers* (servers, processors,

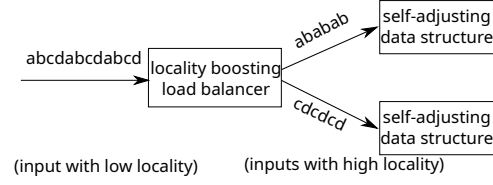


Figure 2: A locality boosting load balancer partitions the input sequence of a given locality into subsequences with higher locality. Self-adjusting data structures perform better on inputs with higher locality.

or nodes). A good load balancing strategy ensures optimal resource utilization, minimizes response time, avoids overloading any single resource, and, as we argue below, improves the locality in the input presented to the workers.

Locality of reference is the property of a sequence of inputs that subsequent items are statistically dependent on each other. A request set with minimal locality is uniformly distributed on the entire input domain and hence unpredictable, while one with maximal locality contains only a single item, i.e., maximally predictable. A *locality-boosting load balancing* policy is then a request dispatching strategy that can statistically or deterministically improve the locality of reference experienced by the worker threads, *turning an unpredictable system input into multiple streams of predictable input* to be processed by the workers (see Fig. 2).

We distinguish two types of locality in this context. *Spatial locality* means that the distribution of requests on the input domain is statistically biased towards a particular subset of the items. One way to ensure this in the load balancer is to *partition* the input domain into disjoint subsets, so that worker’s input distributions are concentrated on a smaller set of items. In contrast, a round robin or a uniform random load balancer will export its own spatial input locality unchanged to the workers. A related concept is *temporal locality*, which refers to the reuse of specific items in the input within a relatively small time duration. One way to boost temporal locality is to reorder items within a time window: e.g., Reframer applies controlled delays to order packet batches flow-wise, thereby enabling more efficient processing [29, 55].

3.2 Self-adjusting algorithms

The second critical enabler in our architecture is *self-adjusting algorithms*. Self-adjustment refers to the property of a dynamic data structure to *automatically reorganize itself based on the sequence of inputs it receives*, in order to optimize performance for future operations on frequently accessed items. Internal data reorganization always introduces extra complexity and overhead compared to a static data structure. Thus, self-adjustment can improve performance only if the input

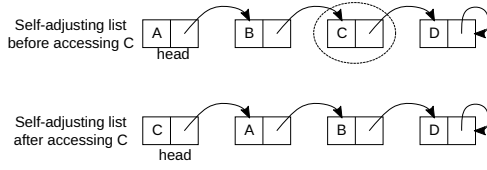


Figure 3: A self-adjusting list containing nodes A, B, C and D serves the request to C and moves C to the front of the list to speed up future accesses to C.

processed by the algorithm exhibits a certain amount of spatial or temporal locality. Next we review some prominent self-adjusting data structures from the literature.

Caches. The textbook example for self-adjustment is caching [76, 80, 95]. Caches can serve frequently accessed items fast by storing them in a software or a hardware fast memory. This is typically much faster than if we had to run the request through the full processing pipeline or the slow backing store. Caches have that almost magical capability of self-adaptation, without us having to engineer any prior knowledge of the input into the cache mechanism apart from a promise that it has nontrivial locality. When the promise is true, caches are an inexpensive way to improve throughput and response time. When there is no locality in the input, however, caches usually just add extra latency and overhead. Note that caches do not necessarily have to come in the form of hardware memory: a fast key-value store is a candidate cache for a slow database [26], a kernel fast-path flow cache is a useful way to speed up a slow user-space software switch [74], etc.

List lookup. One of the most widely used self-adjusting data structures is the *move-to-front list*. Suppose we wish to store a list of m items in a way so that reordering, insertion and deletion are fast, while lookup is also reasonably efficient. A straightforward choice is a static linked list. Here, the cost of accessing an item at position i is exactly i . Then, any linked list can be upgraded to a self-adjusting list using the move-to-front (MTF) heuristics: after accessing an item it is moved to the front, which improves lookup time for future requests of the same item at minimal cost (see Fig. 3). The MTF heuristics comes with appealing theoretical properties, namely that blindly moving the accessed item to the front of the list is close to the best reordering strategy one could design, even if one knew all future requests [84]. MTF lists handle both spatial and temporal locality. For uniformly distributed input MTF lists usually add nontrivial overhead compared to static lists due to the frequent and useless relinking of the list.

Classic applications of MTF lists are information retrieval systems, compression [16], etc. In general, any use case is a potential candidate application for MTF where the task is to match a request against a list of complex rules that do not lend themselves readily to be arranged into a fast lookup structure

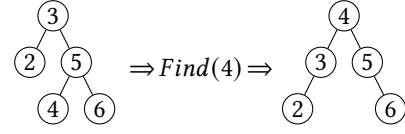


Figure 4: Splay-tree with elements 2, 3, 4, 5, 6. After accessing node 4 it is moved to the root making a subsequent lookup to the same node faster, while the tree is well-balanced.

(e.g., a search tree), like inference in explainable rule-based AI [22], rule matching in OpenFlow and P4 reference software switches [68], packet classification in networking (see later), etc. We note that caching is a subset of list lookup, in that every algorithm for list reorganization gives rise to a different cache management algorithm [84].

Search trees. A search tree is an efficient tree data structure for locating specific keys from within an ordered set. A *splay tree* is a self-adjusting version of a static search tree that can dynamically reorganize itself by moving popular items closer to the root and less frequently accessed elements to the bottom, while keeping the tree relatively well-balanced [9, 17, 85]. Since access time in a search tree is determined by the depth at which the requested item is to be found, splay trees can improve future access to the same or similar items when the input exhibits temporal or spatial locality (see Fig. 4). Note that a red-black tree, an AVL tree or any similar self-balancing tree is not self-adjusting, in that it can rearrange only with respect to the items *stored* in it but not with respect to the queries *posed* to it. Splay trees are widely used to adaptively speed up associative memory and data compression algorithms [48], as well as a building block for more complex self-adjusting algorithms.

3.3 Superlinear scaling

So how can locality-boosting load balancing and self-adjusting algorithms, when used together in a distributed system, produce superlinear scaling? First, we present a demonstration on a particular instantiation of the architecture, *distributed self-adjusting list lookup*, and then we provide a formal scaling characterization for general distributed self-adjusting systems.

Consider a partitioning load balancer (see Fig. 2) combined with a move-to-front list (see Fig. 3) implemented in the workers. Suppose that there are m items to be stored in the list and k workers, each maintaining an independent index into the list. To make things more complicated, we assume uniform request distribution on the entire input domain m at the system's input, which is, recall, the worst case for any self-adjusting algorithm by being totally *unpredictable*. Thus, for a single worker move-to-front reordering has no useful effect and the worst case access time is m , identical to that of a static linked list.

Now suppose we move from 1 worker to k parallel workers where $k \leq m$. This results, within our architecture, that the load balancer effectively partitions the uniformly distributed input on m items into k uniformly distributed input streams on only m/k different items (see Fig. 2). This means that the workers' input features a higher spatial locality than the system's input (which sports none). Had we used a random or a round robin load balancer the workers would still see all the m possible inputs, just with a sampled uniform distribution and no locality. After a while, each MTF list in the workers will have its specific subset of m/k items moved to the first m/k positions (in an arbitrary order), reducing the worst-case lookup time from m (1 worker) to m/k (k workers). This introduces $k \times$ speedup compared to the single-threaded case.

Then, superlinear speedup is merely a product of two simultaneous $k \times$ speedup factors: one $k \times$ factor comes from the self-adjusting list getting progressively faster as we add new workers, and another $k \times$ speedup as the total compute capacity available to the system grows k times. The effective speedup is then just the multiple of the two, yielding k^2 times speedup in total. Plugging into Amdahl's law we get the *scaling law for distributed MTF lists on uniform input* (see Fig. 1):

$$S_I(k) = \frac{T_I(1)}{T_I(k)} = \frac{1}{s + \frac{1-s}{k^2}} \quad k \leq m, \quad (2)$$

where s denotes the fraction of execution time spent in the sequential part of the code.

For small values of k , we obtain $O(k^2)$ scaling. As k grows sufficiently large, say, when $k = m$, the workers' input reduces to a singleton ($m/k = 1$). From this point the distributed MTF list reduces into a simple parallel hash table and superlinear speedup degrades into an "ordinary" Amdahl's scaling profile, until speedup eventually blocks on a serial bottleneck (e.g., the sequential load balancer). For anything between, the system adaptively finds the best combination of an MTF list and a hash-table, producing a quadratic scaling.

In general, superlinear speedup emerges as the superposition of two related speedup factors. First, by splitting the input into multiple input streams of improved locality, the locality-boosting load balancer reduces the "effective size" of the jobs workers will have to process (recall the "scaled size model", §2). Denote the "job size reduction" attainable with k workers by $\ell(k)$. Second, there is a "parallelizability" gain, denoted by $q(k)$, that is obtained by k self-adjusting workers processing the reduced workloads. In the Appendix we present a formal definition of these terms and define the below general *scaling law for distributed self-adjusting systems*:

$$S(k) = \frac{T(1)}{T(k)} = \frac{1}{s + \frac{1-s}{q(k) \cdot \ell(k)}}.$$

If $q(k) \cdot \ell(k) > k$ then we achieve superlinear scaling.

Our formal characterization describes the cases when superlinear scaling may be attainable (positive results) and when it can not (negative results). However, even if we attain superlinear scaling for a particular input sequence for certain values of k , this cannot be sustained infinitely: eventually, superlinear growth peters out as the system hits a bottleneck and growth falls back to (sub)linear or saturates.

It is important to stress that superlinear speedup is only possible with respect to a single-threaded/single-core baseline. Had we normalized with respect to a multi-threaded baseline constrained to a single core (see later) we would see only a linear speedup. Similarly, faster-than-linear growth appears only if both the virtual job size reduction and parallelization deliver actual performance improvement. Had we normalized running times with respect to the scaled job size we would obtain Amdahl's law. See Appendix A for a precise formal characterization.

3.4 Empirical evidence

Next, we present a series of simulation studies to confirm that locality-boosting load balancing combined with self-adjusting workers (but only this combination!) yields faster-than-linear scaling over a broad selection of load balancing policies, self-adjusting algorithms, and input distributions.

Our simulator was written in Go, using lightweight threads (goroutines) managed by the Go runtime to run a given number of workers in parallel, a home-grown implementation of static and MTF lists, and standard Go modules for LRU caches [33], static balanced trees [31] and splay trees [32]. In order to make the workload CPU-bounded, we used an "expensive" ordering operation underneath the search tree with every comparison costing a configurable w number of extra CPU cycles, and likewise LRU cache misses will cost ρ cycles. The simulator creates the specified combination of a load balancer, k worker threads running the selected lookup algorithm, and a random input sequence with a given request distribution, and then performs a configurable number of lookup operations and measures the total execution time with nanosecond precision. To obtain a full picture, the total execution time includes the transient time needed to warm up the self-adjusting algorithms as well as the overhead of request generation, goroutine scheduling, and memory management. For the specification of the evaluation platform, refer to §4.3.

Fig. 5 shows the results. First, the immediate observation is that *the right combination of a locality-boosting load balancer and a self-adjusting algorithm robustly delivers superlinear speedup*, irrespectively of the problem domain or the input distribution. Even for a worst-case uniform input we obtain $3,300 \times$ speedup for list access on 48 CPU cores, almost $70 \times$ of "ideal" linear speedup, $200 \times$ speedup on LRU caches and $65 \times$ speedup on tree search with 36 CPU cores. Second, *only*

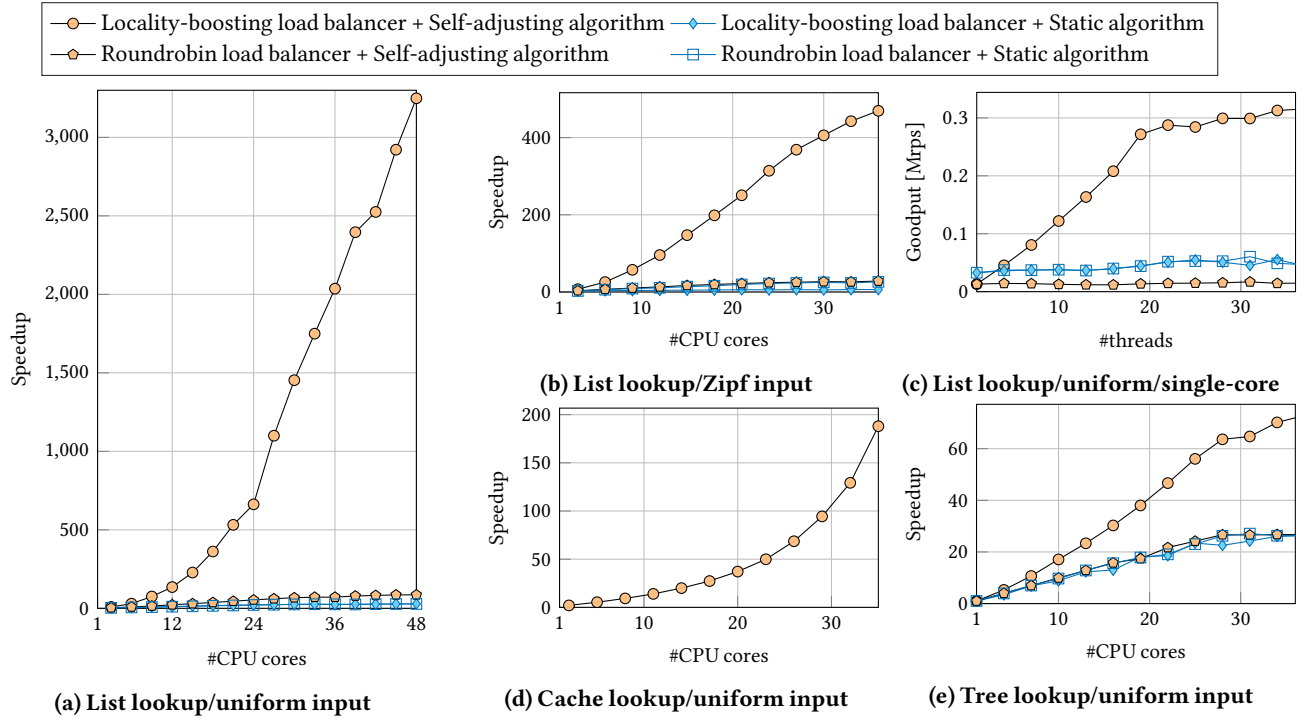


Figure 5: Static vs. self-adjusting distributed systems scaling laws with round-robin and hash-based load balancing: (a) static vs. MTF list access speedup on uniform input ($m = 100k$); (b) static vs. MTF list speedup on skewed input ($m = 100k$, Zipf power law with $\alpha = 1.01$); (c) static vs. MTF list access goodput with multiple threads running on a *single core* for uniform input ($m = 10k$); (d) cache access on uniform input ($m = 50k$, cache hit rate $\delta = 0.05$, $\rho = 100k$ cycles); and (e) static balanced vs. splay tree speedup ($m = 500$, $w = 100k$ cycles). Panels (a), (b), (d) and (e) show multicore speedup as the function of the number of CPU cores, each running a single worker, while (c) shows the single-core goodput (million requests per second).

the combination of locality-boosting load balancing and self-adjusting algorithms produces superlinear speedup, all other combinations (i.e., round robin with any algorithm or static algorithm with any load balancer) fall back to (sub)linear scaling. Third, *self-adjustment clearly has its overhead*. This can be observed in Fig. 5c that shows the absolute throughput instead of the relative speedup. Here, the single-threaded self-adjusting algorithm is slower than the static one due to processing an unpredictable input. Fourth, *the overhead of self-adjustment is irrelevant for more than one worker, or with skewed request distributions*. For instance, on a Zipf input distribution (Fig. 5b) even the single-threaded self-adjusting version is already 2–2.5× faster in an absolute term (not shown in the figure). However, *only combined with a locality-boosting load balancer it produces superlinear speedup*. Fifth, *our architecture yields visible parallel performance gain even if the CPU capacity is kept constant*. Fig. 5c shows an evaluation with an increasing number of parallel threads sharing a single CPU core, with a little surplus CPU for the load balancer. The results indicate that the

parallel self-adjusting system (but only this combination!) delivers linear speedup. Recall, in the multicore case superlinear speedup emerges due to the superposition of two independent $k \times$ speedup trends, one delivered by self-adjustment and another by the $k \times$ scaling of the total CPU power. When the total available CPU is limited only the first $k \times$ speedup factor is in effect, resulting in the observed linear scaling trend.

4 SUPERLINEAR SCALING IN THE LINUX KERNEL

Next we demonstrate our methodology by assembling *existing* techniques into a distributed self-adjusting scheme to understand when, and to what extent, superlinear scaling emerges. Here we present a case study for systematically applying the distributed self-adjusting systems architecture to a common networking problem: software packet classification [38]. Later in Appendix B we also apply our methodology to a popular distributed data storage setup using multiple Memcached servers as a fast cache to access a PostgreSQL

database. In both cases we consider the experiment successful if we can robustly reproduce faster-than-linear growth on some realistic workloads. Recall, it is a stated *nongoal* in this paper to conceive novel algorithms, let alone produce the fastest implementations possible with existing technology. Yet, the re-engineered self-adjusting packet classifier we present below will prove several times faster than the default Linux kernel implementation on a wide range of workloads.

To achieve superlinear scaling we need a self-adjusting algorithm in the first place (plus a locality-boosting load balancer). From the many potential use cases [9, 15, 16, 45, 45, 73, 85] we eventually chose packet classification for the following reasons. First, the default Linux firewall implementation, `nftables`, uses a static doubly linked list to evaluate classifier rules, which makes it an appealing candidate for applying the move-to-front (MTF) heuristics (but see ramifications related to handling rule-dependencies below). Second, underlying packet classification there is an infamously difficult theoretical problem [38, 52, 52, 70, 71, 83, 91], and achieving superlinear speedup on such a hard problem promises massive performance gain. Third, the Linux kernel network stack offers several flexible software and hardware based load balancers for dispatching packets to parallel classifier instances running on different CPU cores [77], which we will reuse to implement the locality-boosting load balancer component. And fourth, packet classifiers are very difficult to cache [21] (recall, caches are the “cheap” way to obtain superlinear scaling), which calls for a true self-adjusting packet classifier.

4.1 Self-adjusting packet classification

A network firewall is a means to control incoming and outgoing network traffic based on user-defined packet classifier rules (see Fig. 6). A classifier *rule* is a pair of a filter, a user-defined regular expression defined on specific fields of the packet header or metadata, and an action that decides what to do with the packets that match the filter (accept, drop, log, etc.). Rules are organized into linear chains ordered by rule priority. When a packet enters a chain, it is compared against the first rule. If there is a match, the corresponding action is executed and the lookup is over. Otherwise, subsequent rules are matched in priority order until the first match is found.

The `nftables` engine is a virtual machine that uses a Domain Specific Language for parsing and matching packet header fields [65]. This makes `nftables` agnostic to specific network protocols, in contrast to, e.g., `iptables`, which contains an embedded protocol parser. Currently, `nftables` is the default packet classifier in most Linux distributions.

One way to make `nftables` self-adjusting would be to replace the static linked list it uses internally for rule matching with a self-adjusting list. A naive application of MTF, however, would easily break the semantics of the firewall. This is

Prio	Proto	Src IP	Dst IP	Dst Port	Action
1	UDP	192.168.178.33	23.0.0.45/32	53	ACCEPT
2	TCP	10.10.10.0/24	23.0.0.45/32	443	DROP
3	UDP	192.168.178.0/24	23.0.0.45/32	53	DROP
4	TCP	10.10.10.10/32	23.0.0.45/32	ANY	ACCEPT
5	IP	192.168.0.0/16	23.0.0.0/8	ANY	ACCEPT

Figure 6: Sample firewall rule set. Source ports do not matter.

because rules in the chain may not be independent from each other, and hence may not be freely swapped [52].

Consider the example in Fig. 6 and suppose that, initially, rules are ordered priority-wise in the list: $\langle 1, 2, 3, 4, 5 \rangle$. Suppose that a packet with the IP 5-tuple $(192.168.0.1, 23.0.0.45, \text{UDP}, 1, 3478)$ enters the classifier, where the fields in the 5-tuple are IP source and destination address, protocol, and source and destination port, respectively. Rules are inspected in linear order until rule 5 is found as the first match, at which point the lookup terminates with the verdict ACCEPT. Now, a naive application of MTF would move rule 5 to the front of list, resulting in the order $\langle 5, 1, 2, 3, 4 \rangle$. Suppose a packet with the 5-tuple $(192.168.178.1, 23.0.0.45, \text{UDP}, 1, 53)$ is to be processed next: this will immediately match rule 5 at the front of the list yielding the verdict ACCEPT, despite that, if matched in priority order, rule 3 would be the correct match and the verdict should be DROP.

We say that rule u is *dependent* on another rule v if they have overlapping match criteria in all fields, v has a higher priority than u , and u and v define different actions. Such a dependency means that u is not allowed to be moved before v in the list, otherwise some packets may be erroneously classified. For instance, in the example of Fig. 6 rule 5 is dependent on rule 3, which is in turn dependent on rule 1, implying the dependency chain $5 \rightarrow 3 \rightarrow 1$. Similarly, rule 4 is dependent on rule 2.

A dependency-aware variant of the MTF heuristics, called the *Move-recursively-Forward* (MRF) algorithm, is defined in [1] (see Alg. 1). The idea is to push an accessed item forward in the list until the first dependency is reached. To prevent the item from blocking behind its direct dependency, the dependency is also moved forward until the first transitive dependency is hit. This process repeats until the head of the list is reached. Independent rules are however free to be moved without restrictions, to the point that if there are no dependencies then MRF simplifies into a plain MTF policy. Contrarily, if the entire rule set is a single dependency chain then no reordering is allowed and MRF degrades into a static list. In general, MRF moves frequently hit rules, with all their dependencies, to the first positions of the chain, which tends to improve lookup performance on high-locality input without jeopardizing the semantics of the classifier [1]. In addition, MRF is “almost” optimal in the same competitive sense as MTF,

Algorithm 1 Move Recursively Forward (MRF)

```

1: procedure MRF( $y$ )
2:   if  $y$  has no dependencies then
3:     Move  $y$  to the front of the list
4:   else
5:     Let  $z$  be the direct dependency of  $y$ 
6:     Move node  $y$  to position( $z$ )+1
7:     MRF( $z$ )

```

in that the best reordering one could obtain even if one knew the entire input sequence in advance would yield only a small constant factor improvement over MRF.

Going back to our earlier example, after rule 5 is hit in the list $\langle 1, 2, 3, 4, 5 \rangle$ MRF moves it immediately after the direct dependency 3 along the dependency chain $5 \rightarrow 3 \rightarrow 1$, 3 is moved to the position after 1, and the recursion ends resulting the order $\langle 1, 3, 2, 5, 4 \rangle$. If 5 was hit again, the lookup time would be only 4 instead of 5. Then, 5 would be moved forward again, yielding the order $\langle 1, 3, 5, 2, 4 \rangle$ and a lookup time of 3. Note that dependency chains can be moved by MRF independently from each other: e.g., if 4 was hit first then we would obtain $\langle 2, 1, 4, 3, 5 \rangle$ in the first iteration and eventually $\langle 2, 4, 1, 3, 5 \rangle$, with lookup time for 4 dropping from 4 to 2.

We created a comprehensive self-adjusting packet classifier implementation on top of `nftables` using the dependency-aware MRF algorithm [1]. Our implementation can run multiple MRF instances in parallel, each maintaining its own local rule order in a private per-CPU pointer array that indexes into a shared static rule list. Apart from lockless list reordering, this also enables lockless rule addition/deletion: every time the rule list is updated we simply allocate a new pointer array at each CPU and update the list head atomically.

The original MRF algorithm uses recursion (see Alg. 1), which may be expensive in the Linux kernel due to the overhead of maintaining the function call stack. To avoid this overhead, we defined an iterative version of the algorithm. When a rule is to be moved forward, we first check whether it can be swapped with the preceding rule. This is done by checking whether the two rules overlap using a range-based representation, which we extract from the rule's bytecode in the `nftables` virtual machine. If there is an overlap then the rule cannot be moved forward so we restart the process, this time trying to move the blocking dependency forward. Otherwise, the two rules are independent so they are immediately swapped and the iteration moves to the subsequent preceding rule. Reordering terminates when we reach the first position. A more efficient implementation would be to precompute dependencies on rule insertion/deletion and run the MRF algorithm using the cached dependencies; implementing this optimization is for further study.

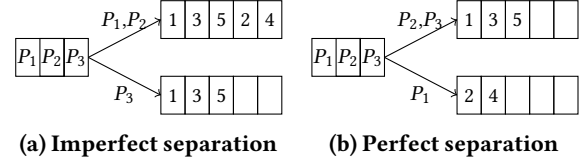


Figure 7: Locality-boosting load balancing over packet sequence P_1 (matching rule 4 of the sample classifier in Fig. 6) followed by P_2 and P_3 (both matching rule 5): (a) hash-based load balancing may assign P_2 and P_3 to different workers so that both will have to keep the dependency list $5 \rightarrow 3 \rightarrow 1$ in the active rule set, (b) perfect separation sends P_1 and P_2 to the same worker, yielding minimal active rule sets.

4.2 Locality-boosting load balancing

The other ingredient that we need to achieve faster-than-linear scaling is a locality-boosting load balancer. An ideal load balancer would partition the rule set into disjoint per-worker subsets. This would minimize the size of the *active rule set* at workers, which is defined as the set of rules for which a particular worker receives packets during a time window. The smaller the active rule set the fewer rules the classifier has to search through for each packet and the larger the contribution of self-adjustment to speedup. Contrarily, the larger the active rule set the more rules compete for the first positions in the list, which reduces the room for self-adjustment to reduce lookup time and erodes superlinear scaling.

There are several factors that may bloat workers' active rule sets. First, whenever a rule with nonzero dependencies is hit MRF adds its entire dependency chain to the active rule set. Second, packet classifiers often use wildcard rules, matching potentially a huge number of diverse traffic flows. If the load balancer dispatches two packets matching the same rule to two different workers, then both workers would have to include the same rule, with all its dependencies, in their active rule sets (see an example in Fig. 7). Note that the same rule duplication problem plagues many software packet classifier algorithms [37, 56, 83, 91].

Designing an ideal load balancer that minimizes workers' active rule sets, regardless of rule dependencies and flow diversity, seems difficult (but see a discussion in §5). Therefore, we adopt a simple hash-based load balancing scheme here that implements only "imperfect rule set partitioning". Our load balancer will however be fully implemented in hardware and run at line rate. This is crucial in order to minimize the overhead, which in our system entirely counts towards the sequential part of the workload and limits ultimate scaling. Later, we will show empirically that even this imperfect scheme is enough to reach superlinear speedup in many practical cases.

Our load balancer reuses the Receive Side Scaling (RSS, [12, 77]) function offered by most standard NICs. RSS evaluates a hash function over a selected set of header fields per each packet. The resultant hash value is then used to index into an indirection table to select a packet queue, and the corresponding CPU core, that will process the packet. The hash function can be configured to consider any combination of the IP 5-tuple header fields, which allows us to fine-tune locality-boosting in our load balancer.

4.3 Reproducing superlinear speedup

We conducted several experiments with the distributed self-adjusting packet classifier combined with the hash-based RSS load balancer. Our goal was to understand whether superlinear scaling can be robustly reproduced on a real network application using real packet I/O.

Testbed. The system-under-test (SUT) is a server equipped with a 32-core AMD EPYC 7502P@2.5 GHz CPU (64 cores with hyper-threading enabled), 128 GByte DDR4 main memory, 96 KB per-core L1 cache, 512 KB per-core L2 cache, and 128MB shared L3 cache. A server of similar configuration was used for traffic generation and measurement with DPDK/moongen [24], connected back-to-back to the SUT over Intel XL710 40GbE NICs. We used standard Ubuntu 22.04.4 LTS OS VMs with NIC-passthrough, running a patched v6.5 Linux kernel on the SUT replacing the `nftables` packet classifier with our own self-adjusting implementation. The benchmarks use the Topsy network testing automation and visualization tool [59]. Hyper-threading was disabled, unless otherwise noted.

The classifier rule sets come from two sources. A series of *realistic rule sets* was generated with ClassBench-ng [62, 89], which accurately model the characteristics of real access control lists and firewalls. ClassBench uses a seed file for describing the statistics of the generated 5-tuple rules, including address ranges, port distribution, and rule dependencies. For each rule set a matching input packet sequence was generated using the standard Classbench tools [62, 63]. We also used a series of *synthetic rule sets* and matching packet traces for conducting controlled microbenchmarks. For each synthetic rule set we generated a matching packet trace with uniform flow-size distribution, which, recall, represents the worst-case for self-adjustment. In all cases the rules and packets using unroutable IP addresses were manually removed (otherwise, Linux would drop some packets, distorting the results). Unless otherwise noted, the benchmarks run with an RSS-based hardware load balancer using an IP 5-tuple hash.

Macrobenchmarks. First, we asked whether superlinear scaling can be reproduced with real workloads. Fig. 8a, Fig. 8b and Fig. 8c give the speedup and the raw packet rate obtained with the default `nftables` packet classifier and our self-adjusting implementation on 3 ClassBench rule sets, each

containing 5000 rules, generated with the seeds `ac11`, `ipc1` and `fw1`, respectively. All rule and trace generation parameters were set to their default values.

Our observations are as follows. First, *superlinear scaling is indeed reproducible with our distributed self-adjusting packet classifier*, with maximum speedup on 32 cores ranging from 225× (about 7× faster than linear) for `ac11`, to 72× (about 2.2× of linear) with `ipc1` and 52× for `fw1` (1.6× faster than linear). In contrast, *the static `nftables` classifier scales almost linearly*. A closer analysis shows a slow sublinear trend representative of an Amdahl’s law profile for a very small sequential parameter ($s \sim 0.001$).

The speedup factor alone, however, does not reveal the full picture, as evidenced by Fig. 8c. The absolute packet rate of the self-adjusting classifier on the `fw1` seed is smaller than that of the static classifier, despite the superlinear speedup. In other words, a massive spurious speedup can be obtained by improving a slow baseline. Note, however, that this occurs only for the `fw1` seed (later we reveal why); for the rest of the benchmarks the self-adjusting version is robustly faster even in terms of raw performance (5.2× for `ac11` and 1.4× with `ipc1` on 32 cores). Nonetheless, with hyperthreading enabled we obtain $\sim 1.5\times$ absolute packet rate improvement on 64 cores even for the `fw1` seed (not shown in the figure), indicating that, with the sufficient amount of parallel resources, distributed self-adjustment eventually surpasses static algorithms even in terms of raw performance. In other words, when scaling is superlinear even a slow baseline becomes extremely fast ultimately.

The mean per-packet latency is shown in Fig. 9a. We observe that *superlinear speedup transforms into massive latency reduction*, resulting 52× smaller mean packet delay on 32 cores for the `ac11` seed using the self-adjusting algorithm. In contrast, the static `nftables` classifier produces a mostly flat latency profile, stabilizing at about 13ms per-packet delay.

Rule size. Next we turn to controlled microbenchmarks over synthetic input, which we fine-tune to highlight the effect of some specific characteristic of the classifier workload on scaling. The main factor affecting speedup is workers’ active rule set sizes, which determines the extent to which self-adjustment can arrange recently hit rules to the front of the rule list (see §4.2). We used the following template to generate synthetic rule-sets of configurable size:

Prio	Proto	Src IP	Dst IP	Dst Port	Action
1	UDP	A.B.C.D	E.F.G.H	1	ACCEPT
2	UDP	A.B.C.D	E.F.G.H	2	ACCEPT
...

The source and destination address are the same in each rule, and each action was set to accept. We obtained 3 rule sets this way, containing roughly 2k, 5k, and 10k rules, respectively (the real size is a close prime to minimize periodicity in the scaling profiles). Note that rules are independent and

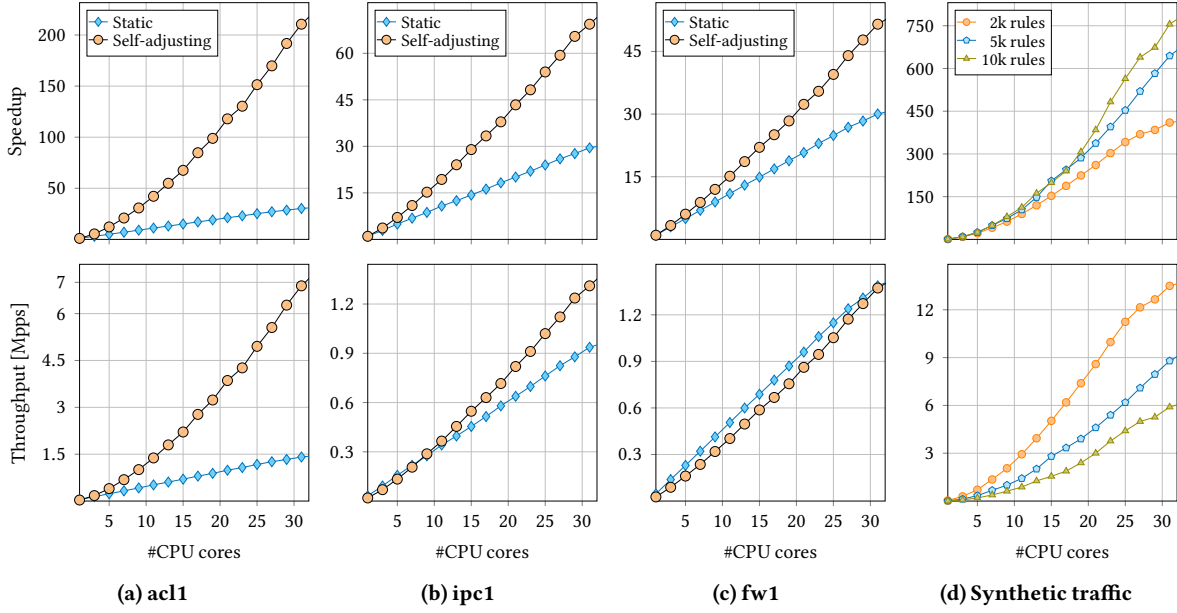


Figure 8: Macrobenchmarks: Scaling on 3 ClassBench rulesets generated from different seeds, containing 5000 rules each (panel (a), (b) and (c)), and synthetic rule set with uniform traffic and different rule sizes (panel (d)). Upper row shows relative speedup and the bottom row shows absolute throughput (packet rate in million packets per sec, mpps).

each rule matches exactly one flow, which represents the optimistic case for the self-adjusting classifier (see later for the pessimistic settings). We generated a matching packet trace containing one flow per rule.

Fig. 8d shows the results. The takeaway is that *superlinear speedup appears independently of the classifier size*, to the point that for 10k rules we see $> 800\times$ speedup on 32 cores. Again, the raw performance plot completes the picture: the larger the rule set the greater the superlinear speedup but the smaller the absolute packet rate. Nevertheless, superlinear scaling robustly appears in terms of the raw performance as well.

Rule dependencies. Rule-dependencies have a crucial role in self-adjustment, since for every rule with nonzero dependencies not just the rule but all its dependencies will also become active, bloating the active rule sets. To measure the effects of rule dependencies we created 3 synthetic rule sets with increasingly long dependency chains using the below template:

Prio	Proto	Src IP	Dst IP	Dst Port	Action
1	UDP	A.B.C.D/32	E.F.G.H	1	ACCEPT
2	UDP	A.B.C.D/31	E.F.G.H	1	DROP
...
...	UDP	A.B.C.D/0	E.F.G.H	1	ACCEPT
...	UDP	A.B.C.D/32	E.F.G.H	2	ACCEPT
...

For every rule in the synthetic rule set we add an extra d overlapping rules by varying the subnet prefix length in the source IP address filter between /32 (most specific, highest priority) and /0 (least specific, lowest priority). This creates

for every rule a chain of d increasingly more specific dependencies. Unfortunately, rule set size also increases d times, but this should not affect the basic superlinear speedup trends (recall Fig. 8d). We run the benchmarks with a 5k base rule set and add d dependencies per rule for $d=1$ (small-dependency), $d=2$, $d=4$ and $d=8$ (high-dependency). The packet trace contains a single flow per each “least specific” rule at the tail of the dependency chains.

Fig. 9b shows the absolute packet rate for the 4 synthetic rule sets. The most important observation is that, as expected, *the more dependencies the smaller the performance and the less visible the superlinear growth* (but note the simultaneous increase in the rule size). Manually checking the classifier statistics confirms that the MRF algorithm at each worker moves the active rules with all d dependencies to the front of the list, reducing the self-adjustment contribution to scaling for large settings of d . In terms of speedup, however, the trend is just the opposite (not shown here): the more dependencies the greater the speedup, again thanks to the slow baseline; e.g., we see $> 1,000\times$ speedup for $d=8$. We also found rule-dependencies to be the reason for the slow scaling on the fw1 ClassBench seed. We observed a similar slowdown when sending huge traffic to the final “catch-all” rule specifying the default action. As this rule depends on all other rules it cannot be moved forward, degrading the self-adjusting classifier into a static list.

Flow diversity. In this microbenchmark we vary the number of flows in the input packet trace per each rule. We used the

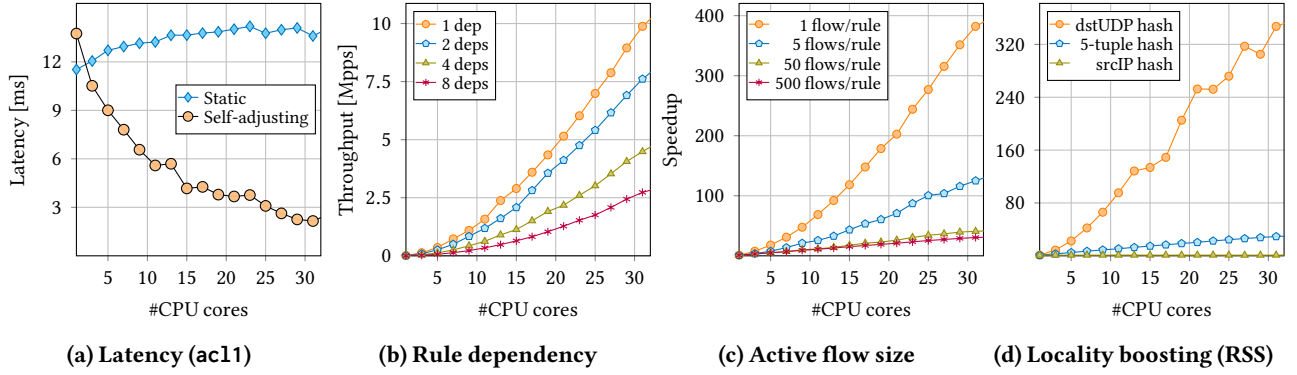


Figure 9: Microbenchmarks: (a) mean packet delay on the rule set generated from the ac11 Classbench speed (5k rules, uniform traffic); (b) raw packet rate for 4 synthetic rule sets with increasingly long dependency chains; (c) speedup for 4 packet traces with increasingly more active flows (independent rules); and (d) speedup with different RSS hash functions (same rules).

same synthetic rule set as previously, but we removed the dependencies. We generated 4 traces containing 1, 5, 50, and 500 uniformly distributed flows per rule, respectively. The results in Fig. 9c confirm that increasing flow diversity has negative impact on scaling: *the more flows per rule the less visible the superlinear speedup*. With a modest flow diversity (1–50 per rule) we observe 46–400 \times speedup on 32 cores. However, for 500 flows the superlinear trend disappears and scaling degrades to linear (32 \times speedup on 32 cores). We traced back the reason to the 5-tuple RSS load balancer. Recall, an optimal load balancing policy would dispatch all flows matching the same rule to the same worker, perfectly eliminating rule duplication at workers (see §4.2). However, the RSS-based 5-tuple hash only “imperfectly” partitions the rule set: manually verifying the classifier statistics reveals that for 500 flows per rule essentially every rule appears at every worker, completely removing the speedup contribution of self-adjustment.

Locality boosting. It seems that longer rule dependencies and growing flow diversity have negative impact on superlinear scaling. In this microbenchmark we show some clue that the negative impact can be removed using a proper locality-boosting load balancer. In particular, Fig. 9d shows the speedup for the previous high flow-diversity benchmark (5k independent rules, 500 uniform flows per rule) with different RSS-based hash functions. Our observations are as follows. An inadequate choice for the load balancing function removes scaling all together: e.g., the RSS hash matching on only the source IP address dispatches all input to the same worker (recall, the source IP is the same in all rules and flows), yielding no scaling at all. A better choice is a 5-tuple hash: this at least spreads the load but, as we checked above, causes massive rule duplication across workers, constraining scaling to linear. An optimal locality-boosting load balancer, however, would dispatch the packets matching the same rule to

the same worker, removing rule duplication. For our specific rule set, such “perfectly partitioning” policy is a hash function that uses only the UDP destination port. For this RSS hash, superlinear scaling is recovered in Fig. 9d, with roughly the same speedup as with no flow diversity in Fig. 9c. This confirms that faster-than-linear growth appears only if the load balancer is indeed “locality-boosting”.

5 RELATED WORK

Superlinear scaling. Amdahl’s famous scaling law [4], asserting sublinear speedup and diminishing returns for parallelization, is a cornerstone result in distributed computing [14, 36, 41, 41, 53]. During the almost 60 years since its first publication various use cases were reported that seemingly violate Amdahl’s scaling, triggering several useful extensions of the basic law [20, 36, 46, 61, 69]. One such phenomenon is faster-than-linear scaling, observed in a broad range of production workloads [10, 26, 39, 40, 42, 43, 75, 81, 87, 88]. For instance, [42] shows superlinear speedup for PostgreSQL and traces back the reason to a new “cache plan” for caching compiled SQL queries at each thread, [75] shows that dense matrix multiplication may exhibit faster-than-linear speedup when matrix rows/columns are optimized for CPU caches, etc. Superlinear growth is often found in Nature as well, e.g., describing the scaling of human communities to large cities [7].

There seem to be two common strategies to obtain superlinear scaling [44, 75]: either do disproportionately less work per worker as system is scaled [75], or add more resources per thread [44]. These techniques, however, are difficult to apply beyond specific use cases [40] or require adding more cache space [75]. Meanwhile, there have been heated debates on the controversies related to superlinear scaling: for instance, Gunther shows that an earlier report on faster-than-linear

scaling from a Hadoop MapReduce workload is attributable to a benchmarking error and, when measured the right way, reduces to sublinear scaling [25]. To the best of our knowledge, ours is the first general methodology that can systematically reproduce superlinear growth in a broad range of applications. **Locality-boosting load balancing.** In line with the recent trend to leverage NICs for intelligently moving data between the network, CPU, GPU and accelerators in computing systems [82], there have been several efforts to extend the static hash-based load balancing provided by RSS: Receive Flow Steering (RFS) is a mechanism to steer flows to the CPU on which the application that processes the flow is running [77] and RSS++ is a dynamic receive side scaling mechanism aiming to keep CPU load constant [12]. These mechanisms could be leveraged to implement more efficient locality boosting in the NIC: RFS can be used to direct all flows matching the same rule to the same CPU, RSS++ could be used to evenly spread load even for staggering workers that process the “difficult” high-dependency rules, etc. Furthermore, Reframer can be used to reorder packets for improving the temporal locality at workers’ input [29, 55], and SAX-PAC can be used to decompose a classifier rule set with many dependencies into multiple smaller but independent rule sets [52]. Hicuts [37], Hypercuts [83], Efficuts [91], and CutSplit [56] define “intelligent” packet header space cuts [51] to partition a rule set along a decision tree into smaller rule lists stored in the leaves of the tree. These schemes are complimentary to our approach: while [37, 56, 83, 91] use “smart” cuts with “dumb” lists in the leaves we rather use “dumb” cuts, implemented by hash-based load balancing, with “smart” rule lists in the workers to reach superlinear parallel scaling.

Self-adjusting data structures. Self-adjusting algorithms, the other ingredient for superlinear scaling, are widely applied in computer systems: caches are extensively used in predictive NFV state stores [54], database accelerators [26, 30, 66], distributed web caching and CDNs [95], and microservices [94]; move-to-front (MTF) lists are used for computing point maxima and convex hulls [15], program compilation and interpretation [45], detecting collisions in hash tables [45], and data compression [16]; further examples are splay trees [85], self-adjusting skip lists [17], push-down trees [9], or self-adjusting geometric data stores [73], etc. Another example for self-adjustment are runtime optimization frameworks which can just-in-time recompile code to specialize it to a particular structured input [29, 55, 58, 64]. All these are candidates to be used, along with a proper locality-boosting load balancer, to reach superlinear scaling in distributed applications. To what extent these algorithms *already* achieve superlinear scaling in production applications is perhaps one of the most intriguing open questions for future research.

6 CONCLUSIONS

In this paper, we present theoretical and empirical proof that locality-boosting load balancing combined with parallel self-adjusting algorithms together yield faster-than-linear speedup in many applications on a wide range of workloads. Our main contribution is that we identify the main architectural patterns commonly appearing in the use cases where superlinear scaling emerges and synthesize these into a comprehensive and universal methodology to *reproduce* it. We then show in extensive simulations that our methodology produces orders of magnitude faster scaling than previously observed. We also show superlinear scaling on applications used widely in production. We extend the default `nftables` Linux subsystem into a true self-adjusting packet classifier, which we use to identify the main workload characteristics (rule-dependency, flow diversity) that affect superlinear growth trends. We also reproduce faster-than-linear scaling on a Memcached+PostgreSQL distributed storage system. Future research will be needed to apply our methodology in a broader range of use cases: for instance, rule-based network intrusion detection systems like Snort or Suricata [47] or explainable AI inferencing seem like appealing application candidates.

Finally we summarize the experience we gathered in engineering systems towards superlinear scaling in a set of generic design guidelines. (1) *Optimize the load balancer for the worker implementation* (or the other way around) so that the load balancer boosts exactly the type of locality workers can exploit. A load balancer that boosts temporal locality will not improve the performance of a worker designed for spatial locality (and *vice versa*). (2) *Make workers’ internal data structures independent* so that each worker can autonomously rearrange itself with respect the locality of its own input. Shared data structures will not work. For instance, [30] maintains a single cache to offload popular Memcached queries that is shared across kernel threads. This blocks parallel self-adjustment by (re)mixing the locality in the threads’ input into a single unstructured workload. (3) *Workers must be CPU-bounded*, otherwise the system cannot benefit from parallelization. In many cases bounded memory is also needed for self-adjustment to count (e.g., in distributed caching). (4) *Avoid sequential bottlenecks* that may block speedup prematurely. For instance, we experimented with running the load-balancer for our packet classifier in the Linux kernel’s RPS function, which allows better locality boosting by letting us fine-tune the partitioning function, instead of the hardware RSS that supports only hash-based load balancing. Unfortunately, the single kernel thread quickly posed a firm sequential bottleneck, well before superlinear speedup could appear.

REFERENCES

- [1] V. Addanki, M. Pacut, A. Pourdamghani, G. Rétvári, S. Schmid, and J. Vanerio. Self-adjusting partially ordered lists. In *IEEE Conference*

- on *Computer Communications*, IEEE INFOCOM, pages 1–10, 2023.
- [2] S. Albers and S. Lauer. On list update with locality of reference. *J. Comput. Syst. Sci.*, 82(5):627–653, 2016.
 - [3] C. Alvarez, J. Corbal, and M. Valero. Fuzzy memoization for floating-point multimedia applications. *IEEE Trans. Comput.*, 54(7):922–927, 2005.
 - [4] G. M. Amdahl. Validity of the single processor approach to achieving large scale computing capabilities. In *Spring Joint Computer Conference, AFIPS '67 (Spring)*, page 483–485. Association for Computing Machinery, 1967.
 - [5] G. Andersson. An approximation algorithm for max p-section. In *STACS 99, 16th Annual Symposium on Theoretical Aspects of Computer Science*, pages 237–247, 1999.
 - [6] K. Andreev and H. Räcke. Balanced graph partitioning. *Theory Comput. Syst.*, 39(6):929–939, 2006.
 - [7] S. Arbesman, J. M. Kleinberg, and S. H. Strogatz. Superlinear scaling for innovation in cities. *Phys. Rev. E*, 79, 2009.
 - [8] C. Avin, M. Bienkowski, A. Loukas, M. Pacut, and S. Schmid. Dynamic balanced graph partitioning. *SIAM J. Discret. Math.*, 34(3):1791–1812, 2020.
 - [9] C. Avin, K. Mondal, and S. Schmid. Dynamically optimal self-adjusting single-source tree networks. In *LATIN 2020: Theoretical Informatics - Latin American Symposium*, volume 12118 of *Lecture Notes in Computer Science*, pages 143–154, 2020.
 - [10] S. Bai, H. Zheng, C. Tian, X. Wang, C. Liu, X. Jin, F. Xiao, Q. Xiang, W. Dou, and G. Chen. Unison: a parallel-efficient and user-transparent network simulation kernel. In *European Conference on Computer Systems, EuroSys*, page 115–131, 2024.
 - [11] A. Bar-Noy and M. Lampis. Online maximum directed cut. *J. Comb. Optim.*, 24(1):52–64, 2012.
 - [12] T. Barbette, G. P. Katsikas, G. Q. Maguire, and D. Kostić. RSS++: load and state-aware receive side scaling. In *International Conference on Emerging Networking Experiments And Technologies, ACM CoNEXT '19*, page 318–333, 2019.
 - [13] T. Barbette, E. Wu, D. Kostić, G. Q. Maguire, P. Papadimitratos, and M. Chiesa. Cheetah: a high-speed programmable load-balancer framework with guaranteed per-connection-consistency. *IEEE/ACM Transactions on Networking*, 30(1):354–367, 2022.
 - [14] G. Bell, J. Gray, and A. Szalay. Petascale computational systems. *Computer*, 39(1):110–112, 2006.
 - [15] J. L. Bentley, K. L. Clarkson, and D. B. Levine. Fast linear expected-time algorithms for computing maxima and convex hulls. *Algorithmica*, 9(2):168–183, 1993.
 - [16] J. L. Bentley, D. D. Sleator, R. E. Tarjan, and V. K. Wei. A locally adaptive data compression scheme. *Commun. ACM*, 29(4):320–330, 1986.
 - [17] P. Bose, K. Douieb, and S. Langerman. Dynamic optimality for skip lists and b-trees. In *ACM-SIAM Symposium on Discrete Algorithms, SODA*, pages 1106–1114, 2008.
 - [18] D. Brahneborg, W. Afzal, A. Čaušević, and M. Björkman. Superlinear and bandwidth friendly geo-replication for store-and-forward systems. In *Proceedings of the 15th International Conference on Software Technologies (ICSOF)*, pages 328–338. SciTePress, 2020.
 - [19] B. Burns. *Designing Distributed Systems: Patterns and Paradigms for Scalable, Reliable Services*. O'Reilly Media, Inc., 1st edition, 2018.
 - [20] K. W. Cameron and R. Ge. Generalizing Amdahl's Law for power and energy. *Computer*, 45(3):75–77, 2012.
 - [21] F. Chang, W. chang Feng, and K. Li. Approximate caches for packet classification. In *IEEE INFOCOM*, volume 4, pages 2196–2207, 2004.
 - [22] F. K. Došilović, M. Brčić, and N. Hlupić. Explainable artificial intelligence: A survey. In *International convention on information and communication technology, electronics and microelectronics (MIPRO)*, pages 210–215. IEEE, 2018.
 - [23] D. E. Eisenbud, C. Yi, C. Contavalli, C. Smith, R. Kononov, E. Mann-Hielscher, A. Cilingeroglu, B. Cheyney, W. Shang, and J. D. Hosein. Maglev: a fast and reliable software network load balancer. In *13th USENIX Symposium on Networked Systems Design and Implementation (NSDI 16)*, pages 523–535, Mar. 2016.
 - [24] P. Emmerich, S. Gallenmüller, D. Raumer, F. Wohlfart, and G. Carle. MoonGen: A Scriptable High-Speed Packet Generator. In *Internet Measurement Conference 2015 (IMC'15)*, Tokyo, Japan, Oct. 2015.
 - [25] V. Faber, O. M. Lubeck, and A. B. White. Superlinear speedup of an efficient sequential algorithm is not possible. *Parallel Comput.*, 3(3):259–260, 1986.
 - [26] B. Fitzpatrick. Distributed caching with memcached. *Linux J.*, 2004(124):5, 2004.
 - [27] J. Fried, G. I. Chaudhry, E. Saurez, E. Choukse, I. Goiri, S. Elnikety, R. Fonseca, and A. Belay. Making kernel bypass practical for the cloud with junction. In *21st USENIX Symposium on Networked Systems Design and Implementation (NSDI 24)*, pages 55–73. USENIX Association, 2024.
 - [28] A. Frieze and M. Jerrum. Improved approximation algorithms for MAX k-CUT and MAX BISECTION. In *Algorithmica* 18, pages 67–81, 1997.
 - [29] H. Ghasemirahni, T. Barbette, G. P. Katsikas, A. Farshin, A. Roozbeh, M. Gironi, M. Chiesa, G. Q. M. Jr., and D. Kostić. Packet order matters! improving application performance by deliberately delaying packets. In *USENIX Symposium on Networked Systems Design and Implementation (NSDI)*, pages 807–827, 2022.
 - [30] Y. Ghigoff, J. Sopena, K. Lazri, A. Blin, and G. Muller. BMC: accelerating memcached using safe in-kernel caching and pre-stack processing. In *18th USENIX Symposium on Networked Systems Design and Implementation (NSDI 21)*, pages 487–501, 2021.
 - [31] hashicorp/btree. <https://pkg.go.dev/github.com/google/btree>.
 - [32] golang-collections/collections. <https://pkg.go.dev/github.com/golang-collections/collections/collections#readme-splay-tree>.
 - [33] hashicorp/golang-lru. <https://pkg.go.dev/github.com/hashicorp/golang-lru/v2>.
 - [34] N. Gunther. Superlinear scalability. In *Symposium and Bootcamp on the Science of Security, HotSoS '13*, 2013.
 - [35] N. Gunther, P. Puglia, and K. Tomasette. Hadoop superlinear scalability: The perpetual motion of parallel performance. *Queue*, 13(5):20–42, 2015.
 - [36] N. J. Gunther. *Guerrilla Capacity Planning: A Tactical Approach to Planning for Highly Scalable Applications and Services*. Springer Publishing Company, Incorporated, 1st edition, 2010.
 - [37] P. Gupta and N. McKeown. Classifying packets with hierarchical intelligent cuttings. *IEEE Micro*, 20(1):34–41, 2000.
 - [38] P. Gupta and N. McKeown. Algorithms for packet classification. *IEEE Network*, 15(2):24–32, 2001.
 - [39] M. Gusev and S. Ristov. Superlinear speedup in Windows Azure cloud. In *IEEE International Conference on Cloud Networking (CLOUDNET)*, pages 173–175, 2012.
 - [40] J. Gustafson. Fixed time, tiered memory, and superlinear speedup. In *Distributed Memory Computing Conference*, volume 2, pages 1255–1260, 1990.
 - [41] J. L. Gustafson. Reevaluating Amdahl's Law. *Commun. ACM*, 31(5):532–533, may 1988.
 - [42] R. Haas. Scalability, in graphical form, analyzed. <http://rhaas.blogspot.com/2011/09/scalability-in-graphical-form-analyzed.html>, 2011.
 - [43] H. Hamann. Superlinear scalability in parallel computing and multi-robot systems: Shared resources, collaboration, and network topology. In *Architecture of Computing Systems (ARCS 2018)*, pages 31–42. Springer, 2018.
 - [44] D. Helmbold and C. McDowell. Modelling speedup (n) greater than n. *IEEE Transactions on Parallel and Distributed Systems*, 1(2):250–256, 1990.
 - [45] J. H. Hester and D. S. Hirschberg. Self-organizing linear search. *ACM Comput. Surv.*, 17(3):295–311, 1985.

- [46] M. D. Hill and M. R. Marty. Amdahl's Law in the multicore era. *Computer*, 41(7):33–38, 2008.
- [47] H. Jiang, G. Zhang, G. Xie, K. Salamatian, and L. Mathy. Scalable high-performance parallel design for network intrusion detection systems on many-core processors. In *Symposium on Architectures for Networking and Communications Systems*, ACM/IEEE ANCS, page 137–146, 2013.
- [48] D. W. Jones. Application of splay trees to data compression. *Communications of the ACM*, 31(8):996–1007, 1988.
- [49] M. Kablan, A. Alsudais, E. Keller, and F. Le. Stateless network functions: Breaking the tight coupling of state and processing. In *USENIX Conference on Networked Systems Design and Implementation*, NSDI'17, page 97–112, 2017.
- [50] G. P. Katsikas, T. Barbette, D. Kostić, R. Steinert, and G. Q. M. Jr. Metron: NFV service chains at the true speed of the underlying hardware. In *15th USENIX Symposium on Networked Systems Design and Implementation (NSDI 18)*, pages 171–186, Apr. 2018.
- [51] P. Kazemian, G. Varghese, and N. McKeown. Header space analysis: Static checking for networks. In *Symposium on Networked Systems Design and Implementation*, USENIX NSDI, pages 113–126, 2012.
- [52] K. Kogan, S. Nikolenko, O. Rottenstreich, W. Culhane, and P. Eugster. SAX-PAC (Scalable And EXpressive PAcet Classification). In *Conference of the ACM Special Interest Group on Data Communication*, SIGCOMM '14, page 15–26, 2014.
- [53] S. Krishnaprasad. Uses and abuses of Amdahl's Law. *J. Comput. Sci. Coll.*, 17(2):288–293, dec 2001.
- [54] J. Lei and V. Shrivastav. Seer: Enabling Future-Aware online caching in networked systems. In *USENIX Symposium on Networked Systems Design and Implementation (NSDI)*, pages 635–649, 2024.
- [55] T. Lévai, F. Németh, B. Raghavan, and G. Retvari. Batchy: batch-scheduling data flow graphs with service-level objectives. In *USENIX Symposium on Networked Systems Design and Implementation (NSDI)*, pages 633–649, 2020.
- [56] W. Li, X. Li, H. Li, and G. Xie. Cutsplit: A decision-tree combining cutting and splitting for scalable packet classification. In *IEEE INFOCOM*, pages 2645–2653, 2018.
- [57] H. Lim, D. Han, D. G. Andersen, and M. Kaminsky. MICA: A holistic approach to fast In-Memory Key-Value storage. In *11th USENIX Symposium on Networked Systems Design and Implementation (NSDI 14)*, pages 429–444, Apr. 2014.
- [58] L. Linguaglossa, S. Lange, S. Pontarelli, G. Rétvári, D. Rossi, T. Zinner, R. Bifulco, M. Jarschel, and G. Bianchi. Survey of performance acceleration techniques for Network Function Virtualization. *Proceedings of the IEEE*, 107(4):746–764, 2019.
- [59] T. Lévai, G. Pongrácz, P. Megyesi, P. Vörös, S. Laki, F. Németh, and G. Rétvári. The price for programmability in the software data plane: The vendor perspective. *IEEE Journal on Selected Areas in Communications*, 36(12):2621–2630, 2018.
- [60] S. Mahajan and J. Ramesh. Derandomizing semidefinite programming based approximation algorithms. In *Proc. 36th Ann. IEEE Symp. on Foundations of Comput. Sci. (FOCS)*, pages 162–169, 1995.
- [61] A. Marowka. Extending Amdahl's Law for heterogeneous computing. In *2012 IEEE 10th International Symposium on Parallel and Distributed Processing with Applications*, pages 309–316, 2012.
- [62] J. Matoušek, G. Antichi, A. Lučanský, A. W. Moore, and J. Kořenek. ClassBench-ng: recasting ClassBench after a decade of network evolution. In *Symposium on Architectures for Networking and Communications Systems*, IEEE/ACM ANCS, page 204–216, 2017.
- [63] S. Miano. sebymiano/pcap-utils. <https://github.com/sebymiano/pcap-utils>.
- [64] S. Miano, A. Sanaee, F. Risso, G. Rétvári, and G. Antichi. Domain specific run time optimization for software data planes. In *ACM International Conference on Architectural Support for Programming Languages and Operating Systems*, ASPLOS '22, page 1148–1164, 2022.
- [65] The nftables project. <https://wiki.nftables.org>.
- [66] R. Nishtala, H. Fugal, S. Grimm, M. Kwiatkowski, H. Lee, H. C. Li, R. McElroy, M. Paleczny, D. Peek, P. Saab, D. Stafford, T. Tung, and V. Venkataramani. Scaling memcache at Facebook. In *10th USENIX Symposium on Networked Systems Design and Implementation (NSDI 13)*, pages 385–398, Apr. 2013.
- [67] V. Olteanu, A. Agache, A. Voinescu, and C. Raiciu. Stateless datacenter load-balancing with Beamer. In *15th USENIX Symposium on Networked Systems Design and Implementation (NSDI 18)*, pages 125–139, 2018.
- [68] ONF. Openflow reference release. <https://github.com/mininet/openflow>, 2013.
- [69] I. Onyuksel and S. Hosseini. Amdahl's law: a generalization under processor failures. *IEEE Transactions on Reliability*, 44(3):455–462, 1995.
- [70] M. H. Overmars and F. A. van der Stappen. Range searching and point location among fat objects. *J. Algorithms*, 21(3):629–656, 1996.
- [71] M. Pacut, J. Vanerio, V. Addanki, A. Pourdamghani, G. Rétvári, and S. Schmid. Self-adjusting packet classification. *CoRR*, abs/2109.15090, 2021.
- [72] S. Palkar, C. Lan, S. Han, K. Jang, A. Panda, S. Ratnasamy, L. Rizzo, and S. Shenker. E2: a framework for NFV applications. In *Symposium on Operating Systems Principles*, SOSP '15, page 121–136, 2015.
- [73] E. Park and D. M. Mount. A self-adjusting data structure for multidimensional point sets. In *Algorithms - ESA Annual European Symposium*, volume 7501, pages 778–789, 2012.
- [74] B. Pfaff, J. Pettit, T. Koponen, E. Jackson, A. Zhou, J. Rajahalme, J. Gross, A. Wang, J. Stringer, P. Shelar, K. Amidon, and M. Casado. The design and implementation of Open vSwitch. In *12th USENIX Symposium on Networked Systems Design and Implementation (NSDI 15)*, pages 117–130, 2015.
- [75] S. Ristov, R. Prodan, M. Gusev, and K. Skala. Superlinear speedup in HPC systems: Why and when? In *Federated Conference on Computer Science and Information Systems (FedCSIS)*, pages 889–898, 2016.
- [76] O. Rottenstreich and J. Tapolcai. Optimal rule caching and lossy compression for longest prefix matching. *IEEE/ACM Transactions on Networking*, 25(2):864–878, 2016.
- [77] Scaling in the linux networking stack. <https://www.kernel.org/doc/Documentation/networking/scaling.txt>.
- [78] Packet classification and access control. https://doc.dpdk.org/guides/prog_guide/packet_classif_access_ctrl.html.
- [79] A. Sapio, M. Canini, C.-Y. Ho, J. Nelson, P. Kalnis, C. Kim, A. Krishnamurthy, M. Moshref, D. Ports, and P. Richtarik. Scaling distributed machine learning with In-Network aggregation. In *18th USENIX Symposium on Networked Systems Design and Implementation (NSDI 21)*, pages 785–808, 2021.
- [80] C. Sears. The elements of cache programming style. In *4th Annual Linux Showcase & Conference (ALS 2000)*. USENIX Association, 2000.
- [81] Sdn analytics and control: Superlinear. <https://blog.sflow.com/2010/09/superlinear.html>, 2010.
- [82] J. Sherry. The I/O driven server: From SmartNICs to data movement controllers. *Computer Communications Review (CCR)*, 53(3), 2023.
- [83] S. Singh, F. Baboescu, G. Varghese, and J. Wang. Packet classification using multidimensional cutting. In *Conference on Applications, Technologies, Architectures, and Protocols for Computer Communications*, page 213–224, 2003.
- [84] D. D. Sleator and R. E. Tarjan. Amortized efficiency of list update and paging rules. *Commun. ACM*, 28(2):202–208, Feb. 1985.
- [85] D. D. Sleator and R. E. Tarjan. Self-adjusting binary search trees. *J. ACM*, 32(3):652–686, 1985.
- [86] C. Sun, J. Bi, Z. Zheng, H. Yu, and H. Hu. Nfp: Enabling network function parallelism in nf. In *Conference of the ACM Special Interest Group on Data Communication*, SIGCOMM '17, page 43–56, 2017.

- [87] H. Sutter. Going superlinear. Dr. Dobb's J., 2008. <https://www.drdoobs.com/cpp/going-superlinear/206100542>.
- [88] H. Sutter. Super linearity and the bigger machine. Dr. Dobb's J., 2008. <https://www.drdoobs.com/parallel/super-linearity-and-the-bigger-machine/206903306>.
- [89] D. E. Taylor and J. S. Turner. ClassBench: a packet classification benchmark. *IEEE/ACM Transactions on Networking*, 15(3):499–511, 2007.
- [90] W. Tu, Y.-H. Wei, G. Antichi, and B. Pfaff. Revisiting the Open VSwitch Dataplane ten years later. In *Conference of the ACM Special Interest Group on Data Communication*, SIGCOMM '21, page 245–257, 2021.
- [91] B. Vamanan, G. Voskuilen, and T. N. Vijaykumar. EffiCuts: optimizing packet classification for memory and throughput. In *Conference of the ACM Special Interest Group on Data Communication*, SIGCOMM '10, page 207–218, 2010.
- [92] Wikipedia. Speedup: Super-linear speedup. https://en.wikipedia.org/wiki/Speedup#Super-linear_speedup.
- [93] S. Woo, J. Sherry, S. Han, S. Moon, S. Ratnasamy, and S. Shenker. Elastic scaling of stateful network functions. In *USENIX Conference on Networked Systems Design and Implementation*, NSDI'18, page 299–312, 2018.
- [94] H. Zhang, K. Kallas, S. Pavlatos, R. Alur, S. Angel, and V. Liu. MuCache: A general framework for caching in microservice graphs. In *USENIX Symposium on Networked Systems Design and Implementation (NSDI)*, pages 221–238, 2024.
- [95] Y. Zhang, J. Yang, Y. Yue, Y. Vigfusson, and K. Rashmi. SIEVE is simpler than LRU: an efficient Turn-Key eviction algorithm for web caches. In *USENIX Symposium on Networked Systems Design and Implementation (NSDI)*, pages 1229–1246, 2024.

APPENDIX A ANALYSIS

Scaling in distributed systems refers to the improvement in performance achieved by adding more machines to a system. It is defined as the ratio of the completion time of the baseline system with a single machine to the completion time of the system with k machines, $S(k) = T(1)/T(k)$.

We claim that our system scales along two dimensions: load balancer efficiency $\ell(k)$ and parallelizability $q(k)$.

$$S(k) = \frac{T(1)}{T(k)} = \frac{1}{s + \frac{1-s}{q(k) \cdot \ell(k)}} ,$$

for some values of $q(k)$ and $\ell(k)$ that depend on the input σ and the load balancer. The value $\ell(k)$ captures the reduction in the total work of the system by using a load balancer, and $q(k)$ captures how well the reduced workload can be parallelized on k machines. If $q(k) \cdot \ell(k) > k$, we say that our system *scales superlinearly*.

We draw the following conclusions from our analysis.

- (1) The system can scale superlinearly for certain input streams combined with the right load balancer. The uniform input example is just one such example.
- (2) On the contrary, some other input streams cannot even achieve linear scaling with any load balancer.
- (3) The parallelization factor $q(k)$ can be at most k (reaching $q(k) = k$ for uniform input).
- (4) The workload reduction factor $\ell(k)$ depends on the input sequence σ and the load balancer.

- (5) The workload reduction factor $\ell(k)$ cannot be reduced indefinitely with growth of k . Hence, the system can scale superlinearly only for small values of k .

Our analysis holds for a wide range of self-adjusting data structures (including LRU caches).

A.1 The model

Architecture. Consider k identical parallel machines $M_1, M_2, M_3, \dots, M_k$, each having its own isolated memory and running an instance of a self-adjusting list D . The stream of requests σ arriving at a load balancer is partitioned into k streams $\sigma(M_1), \sigma(M_2), \dots, \sigma(M_k)$ and dispatched to the machines. The load balancer dispatches the requests to machines based solely on the request itself, ignoring the state of the system. The load balancer is a function $f_k: \mathcal{U} \rightarrow \{1, 2, \dots, k\}$ from the universe of all items \mathcal{U} to the machines, and this function partitions the universe \mathcal{U} into k subsets $\mathcal{U}_1, \mathcal{U}_2, \dots, \mathcal{U}_k$ (often referred to as affinity domains).

Cost model. The time to process a request σ_t at time t includes both the computational overhead of the load balancer (denoted $T(f_k(\sigma_t))$) and the processing time by D at machine $M_{f_k(\sigma_t)}$ (denoted $T(g_D(\sigma_t))$). Notably, the processing time g_D varies over time due to the self-adjusting nature of the data structure.

Self-adjusting data structures often achieve some variant of *working set property* that links the input history to the cost of processing a request. In our work, the working set for an item σ_t request at time t is defined as the set of distinct items requested since the last request to the item σ_t . With each data structure D , there exists an associated cost function g_D

$$\text{cost}(x, t, \sigma) \leq g_D(|W_t(x)|) ,$$

where $W_t(x)$ is the number of distinct requests to items other than x since the last request to x . With these assumptions, we capture e.g. LRU caches, Move-to-Front lists, splay trees and more.

Objective. Our goal is to minimize the completion time of the schedule of jobs induced by the stream of requests σ executed on parallel machines. The schedule finishes when all requests from σ are processed. The load may be uneven, and some machines may be idle throughout execution, but the system is not allowed to reassign the requests to other machines.

Benchmark. Our benchmark $T(1)$ is a single self-adjusting data structure D that runs on a single machine M_1 and processes the entire stream σ , with a trivial load balancer $f_1(\cdot) = M_1$. A single self-adjusting data structure is the most natural baseline choice for self-adjusting data structures.

A.2 Superlinear Scaling of Self-adjusting Distributed Systems (a positive result)

The load balancer f_k partitions the input stream σ into k more local streams $\sigma(M_1), \sigma(M_2), \dots, \sigma(M_k)$, and reduces the sum of machine's workloads by a factor of $\ell(k)$. The workload is then executed on k machines, with the objective to reduce the schedule completion time, which brings speedup of the factor of $q(k)$, $q(k) \leq k$.

A.3 Study of $\ell(k)$: how workload is reduced

The value of $\ell(k)$ depends on the input σ and the load balancer f_k . Hence, $\ell(k)$ indirectly relies on k through f_k (these are tied together in our architecture). Furthermore, for technical reasons, $\ell(k)$ depends on the serial portion of the workload s . To estimate $\ell(k)$ for a given stream σ , we need to relate σ with the load balancer f_k for the given data structure D .

Self-adjusting data structures and working sets. Our law captures various self-adjusting data structures, such as lists, caches and their generalizations. In these data structures, the cost of accessing an item depends on the internal structure and changes over time depending on the history of requests. Self-adjusting data structures have a property that the cost of accessing an item x at time t depends on the number of distinct items requested since the last access of x . This is often referred to as the *working set property*, and it can hold in an amortized sense. A working set for an item σ_t request at time t is defined as the set of distinct items requested since the last request to the item σ_t .

With each data structure D , there exists an associated cost function g_D

$$\text{cost}(x, t, \sigma) \leq g_D(|W_t(x)|) ,$$

where $W_t(x)$ is the number of distinct requests to items other than x since the last request to x . With these assumptions, we still capture parallel extensions of LRU caches, Move-to-Front lists, splay trees.

We illustrate g_D for Move-to-Front and LRU. In Move-to-Front the cost of accessing an item is linear: $g_{\text{MTF}}(x, t, \sigma) = |W_t(x)| + 1$. LRU is the algorithm Move-to-Front casted into the cost model of caching with a generalized cost function g_{LRU} that is non-linear: for a cache of size B , the cost is 1 if the working set size is $B + 1$, and 0 otherwise. We note that this generalized setting introduced by Sleator and Tarjan [84] captures more general data structures than just lists and caching under a common characterization of g_D .

Load balancer isolates working sets. Recall that the load balancer partitions the stream into k streams $\sigma(M_1), \sigma(M_2), \dots, \sigma(M_k)$ and dispatches them to the machines. These streams

may have reduced working set sizes at the machines in comparison to the working set sizes of the original stream σ . Precisely, a load balancer $f_k : \mathcal{U} \rightarrow \{1, 2, \dots, k\}$ partitions the universe \mathcal{U} into k subsets $\mathcal{U}_1, \mathcal{U}_2, \dots, \mathcal{U}_k$, and the working set for machine i at the time $x = \sigma_t$ is requested is $W_t(x, t, \sigma(M_i)) = W_t(x, t, \sigma) \cap \mathcal{U}_i$.

The cost for the parallel workload of the baseline is

$$T(1) = \sum_t g(|W_t(x)|) .$$

The cost of the parallel workload for the distributed system is as follows. In total, we observe cost savings from load balancing for the stream σ , data structure D characterized by a function g , and a load balancer f_k .

$$T(k) = \sum_t g(|W_t(x) \cap \mathcal{U}_{f_k(\sigma_t)}|) .$$

Note the total work decreases by r for the stream σ .

The speedup ratio in the dimension of $\ell(k)$ is then

$$1/\ell(k) = \frac{\sum_t g(|W_t(x) \cap \mathcal{U}_{f_k(\sigma_t)}|)}{\sum_t g(|W_t(x)|)} .$$

A.4 The study of $q(k)$: how parallelizable the reduced workload is

Recall that our objective is not to minimize the total work, but to minimize the completion time of the schedule. Fix a reduced parallel workload from the previous section of total size $(1 - s)/\ell(k)$, and let's look at its components. We are interested in minimizing the completion time of the last machine. The speedup ratio in the dimension of $q(k)$ is given by

$$1/q(k) = \frac{\sum_t g(|W_t(x) \cap \mathcal{U}_{f_k(\sigma_t)}|)}{\max_i \sum_{t: f_k(\sigma_t)=i} g(|W_t(x) \cap \mathcal{U}_{f_k(\sigma_t)}|)} .$$

Some request streams and load balancer pairs are better at achieving $q(k)$ close to its theoretical limit k (given by Amdahl's law). Our uniform input example achieves perfect parallelization of $q = k$, since all jobs are the same size and all machines process the same number of jobs. On the other hand, some unparallelizable streams such as a repeated request to a single item have $q = 1$, because they use only one machine in our architecture¹. Heterogeneous workload can also be parallelizable, for example a single machine M_1 can process a majority of the stream σ if the stream σ_{M_1} is local, while the rest of the machines process longer jobs to finish at the same time as M_1 . It should however be obvious that the streams that achieve good parallelization $q(k)$ need to have roughly equal workloads for each machine, and the only streams that can achieve that consist of requests to multiple affinity domains.

¹To remedy that, we could look into load balancers that distribute jobs to multiple machines (generalizations of functions f_k , to functions from \mathcal{U} to sets of machines). This goes beyond the scope of this paper. See RSS+ paper for reference [12].

To maximize $q(k)$, the load balancer needs to distribute the workload evenly among the machines.

A.5 Characterizing the speedup

Definition A.1. Fix any stream σ and load balancer f_k . The value $\ell(k)$ is defined as the ratio of the total parallel work of the distributed system to the total parallel work of the baseline.

Definition A.2. Fix any stream σ and load balancer f_k . The value $q(k)$ is defined the ratio of the total reduced parallel work $(1-s) \cdot \ell(k)$ to the completion time of the last machine.

Then, the speedup of the distributed system is given by the following theorem.

THEOREM A.3. Consider a data structure D with a cost that is upper-bounded by a non-decreasing function g_D of the working set size.

Consider a load balancer that dispatches inputs σ taken from a universe \mathcal{U} to k identical parallel workers W_1, W_2, \dots, W_k , each running an instance of a self-adjusting algorithm D , using a deterministic function $f_k : \mathcal{U} \rightarrow \{1, 2, \dots, k\}$ that partitions the input universe \mathcal{U} into k disjoint subsets $\mathcal{U}_1, \mathcal{U}_2, \dots, \mathcal{U}_k$. For an input σ , self-adjusting distributed systems scaling with k is characterized in two dimensions: load balancer efficiency $\ell(k)$ and parallelizability $q(k)$. This gives us speedup

$$\frac{T(1)}{T(k)} \leq \frac{1}{s + \frac{1-s}{q(k) \cdot \ell(k)}} ,$$

for $q(k)$ and $\ell(k)$ defined in Definition A.2 and A.1, which depend on the input σ , load balancer f_k and the cost function g_D .

The proof of the above theorem is a direct consequence of executing a reduced workload (reduced by a factor of $\ell(k)$) on k machines with the parallelization $q(k)$. Note that s is the serial fraction of the workload, which is common to $T(1)$ and $T(k)$.

Our theorem applies to e.g. Move-to-Front lists and LRU caches, since their cost functions g_{MTF} and g_{LRU} are monotonically increasing in the working set size.

We dedicate the next two subsections to partitioning the speedup into these two dimensions and independently analyzing them.

A.6 Impossibility of scaling for self-adjusting distributed systems (a negative result)

We conclude by outlining a scaling impossibility law, analogous to Amdahl's law for self-adjusting distributed systems. We conclude that for any stream σ , scaling is limited to the initial phase and cannot continue indefinitely with k . Therefore, the superlinear scaling observed in practice is only merely a transient phenomenon.

THEOREM A.4. Assume that the cost of a data structure D is lower-bounded in terms of the working set size as a monotonically non-decreasing function g_D . Then, our distributed self-adjusting system cannot scale better than

$$S(k) = \frac{T(1)}{T(k)} = \frac{1}{s + \frac{1-s}{k \cdot \ell(k)}} ,$$

where $\ell(k)$ depends on g_D , σ and f_k .

The proof of this theorem is a direct consequence of the Amdahl's law to the reduced workload: the reduced parallel workload $(1-s)/\ell(k)$ can be executed at most k times faster. We leave the proof of to the full version of the paper.

Two key consequences arise from the above theorem.

(1) Superlinearity is an initial-only phenomenon

For any input sequence σ , the maximum achievable multiplicative $\ell(k)$ is fixed. This occurs when k is large enough so the load balancer f_k isolates all working sets. However, for many inputs σ , the improvements can dry up even for smaller k (the more local the sequences are, the more effective the load balancer can be). As a corollary, the system can scale with the parameter k due to additional resources, but the savings from load balancing do not grow indefinitely with k .

OBSERVATION 1. For each input sequence σ , there exists a constant K such that for all $k > K$, for each load balancer f_k , the savings $\ell(k)$ are fixed do not increase with k . Combined with the fact that $q \leq k$, this implies that superlinear scaling cannot continue indefinitely with k .

(2) Tight analysis for LRU caches and MTF lists

In case of Move-to-Front lists and LRU caches, the cost function $\ell(k)$ is both upper- and lower-bounded as a function of the working set size. Hence, the analysis of scaling is tight for these algorithms.

In particular, the LRU algorithm for caching cannot scale superlinearly indefinitely. Therefore, the scaling observed in the literature is only an initial effect caused by the combined influence of the increased number of machines and the load balancer.

A.7 How to find good load balancers?

There are many constraints for the load balancer, e.g. it should be efficiently computable and should parallelize the workload well (measured by the parameter q). Now, we focus solely on the locality-boosting aspect of the load balancer, how well the workload is reduced by cutting affinity domains with f_k . We can visualize the saving from load balancing with help of a weighted complete graph, where each edge weight represents the saved cost by isolating the affinity domains of the two machines. For each input stream, the costs of a self-adjusting

data structure can be decomposed into the sum of costs accounted to pairs of nodes, see the work of Albers and Lauer for details [2]. The optimal load balancer uses the heaviest cut in such a graph.

First, we consider load balancers that remain fixed over time, and we additionally assume that the input stream is known in advance. The optimal k -cut is known as *maximum k -cut*, and is known to be NP-hard problem [28, 60]. If the cuts are balanced (a natural choice), then the problem is known as *maximum k -section* [5] (from the perspective of maximizing the cut) and *minimum graph k -balanced partitioning* [6] (from the perspective of minimizing the non-cut edges). In the former model, we have a polynomial time algorithm that achieves a constant-factor approximation [5], and in the latter model, we have a polylogarithmic approximation [6].

It is often unrealistic to know the entire input sequence in advance, and online variants of the problem are studied, where additionally the load balancer can change the assignment of the requests to the machines over time. We refer to online variants of the above problems: online k -cut [11] and online graph partitioning [8].

APPENDIX B SUPERLINEAR SCALING IN DISTRIBUTED CACHING

Superlinear scaling often emerges in systems where a “fast” *distributed cache* is deployed in front of “slow” processing system or storage engine [42, 81, 88]. Examples include multi-processor CPUs with unshared Level-1 fast cache memory that make access to program arguments more efficient [75], runtimes that selectively “memoize” the results of costly computations [3], FIB caches in OS network stacks that maintain the most recent IP routes in fast memory to sidestep longest prefix matching [76], hierarchical (mega)flow caches that serve as a fast-path in programmable software switches [74], etc. All these workloads may benefit from caches becoming more efficient as the system is scaled and, potentially, show superlinear speedup on certain workloads. Below we reproduce this finding using Memcached as a fast cache for the PostgreSQL database management system [26, 30, 57, 66].

It is instructive to quantify superlinear speedup in this context using a simple model. Suppose a source emits uniformly distributed random requests for m items and requests are distributed among k workers, each using a separate cache of size c , by hashing on the request id. Initially, the cache hit rate for a single worker that processes all m possible requests is $\delta := c/m$. Adding k workers effectively partitions the requests into k random buckets so that each worker will perceive uniformly distributed requests for only m/k items, which improves the cache hit rate at each worker to $\frac{c}{m/k} = k\delta$ ($k\delta \leq 1$). This puts

the lookup time of the system of k parallel caches to

$$T_c(k) = \begin{cases} s + \frac{1-s}{k} (k\delta + (1-k\delta)\rho) & \text{if } k\delta \leq 1 \\ s + \frac{(1-s)}{k} & \text{otherwise} \end{cases}, \quad (3)$$

where δ is the single-threaded cache hit rate, ρ is the penalty for a cache miss, and s denotes the fraction of execution time spent in the sequential part of the code.

The speedup $S_c(k) = \frac{T_c(1)}{T_c(k)}$ for the parameters $s=0.1$, $\delta=0.1$ and $\rho=10$ is depicted in Fig. 10. The lower envelope of the scaling profile is given by Amdahl’s law for the system with random or round robin load-balancing. As k grows the scaling profile progresses over a superlinear curve to an elevated Amdahl’s law profile, representative of a system serving *all* requests from fast memory. Note that this occurs *only* if request dispatching is chosen carefully to partition the item space. Modulo hashing assigns the same item to the same worker deterministically, so that workers process only a subset of the items that may have a greater chance to fit into the cache. In contrast, a random or a round robin load balancer may assign any item to any worker, which defeats the purpose of improving workers’ cache hit rate.

Note that for a system to match this scaling profile the fast cache and the slow backend must be scaled jointly. In the below we will use this setup, by increasing the number of Memcached instances and PostgreSQL client threads *simultaneously*. The version where *only* the cache is scaled while the backend runs with constant CPU resources would match a different scaling profile characterized by $T_c(k) = s + (1-s)(\delta + (1-k\delta)\rho)$, $k\delta \leq 1$. This also produces superlinear speedup, just with a slower ramp-up.

In our case study Memcached is gradually scaled from a single replica to 15 replicas, playing the role of a distributed cache for a “slow” PostgreSQL v14 database. PostgreSQL is scaled proportionally to the number of Memcached replicas; in particular we run 4 PostgreSQL client threads per cache replica. Note that each Memcached replica runs on a different port letting clients to address each one separately, which is critical to implement key-hashing. We wrote a custom multi-threaded client that performs a configurable number of cache-aside read iterations: first it tries to read from Memcached and, on a cache miss, makes a read from PostgreSQL and writes the result back into the cache. We used two different load balancing schemes to route key requests: random load balancing reads from a random Memcached replica, while key-hashing always reads/writes the same key from/to the same Memcached replica.

Fig. 11 shows the results for PostgreSQL pre-filled with 1,000,000 key-value pairs of 16 byte keys and 48 byte values, with a configurable number of PostgreSQL threads and Memcached replicas with 4 MB of cache each. As expected, superlinear scaling emerges with key-hashing, yielding 35× speedup with 15 Memcached replicas, 2.3× higher than linear

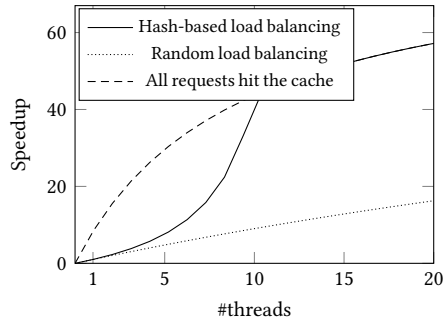


Figure 10: Scaling laws for distributed caching.

scaling. In contrast, random request routing exhibits only linear scaling. The reason is the improving cache hit rate as Memcached is scaled: with 14 replicas we reach close to 100% cache hit rate and speedup falls back into the linear range, as predicted by the analysis.

We note that superlinear scaling in this context is extremely sensitive to certain benchmark parameters, like the number of Memcached replicas, PostgreSQL threads, and client

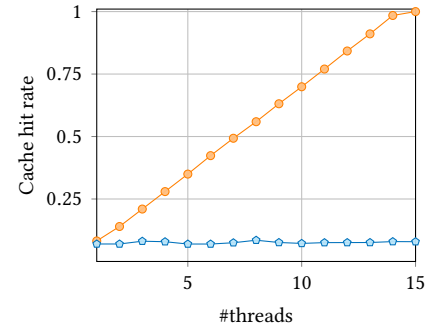
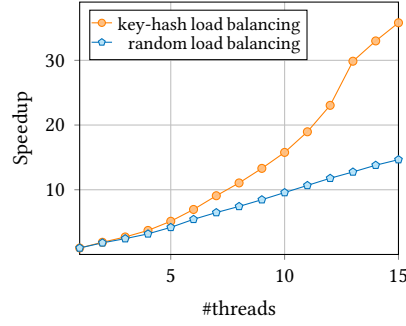


Figure 11: Results for a joint scaling of Memcached+PostgreSQL with and without key-hashing: speedup and cache-hit rate.

threads. This is because for faster-than-linear scaling to appear Memcached replicas must be both CPU-bounded (so that adding more replicas will improve throughput) *and* memory-bounded (so that improving cache hit rate will cause speedup) at the same time. Earlier reports indicate that this occurs surprisingly commonly in practice [40, 42, 44, 75, 87, 88].