



A Survey of Distributed Caching Patterns for High-Throughput Python Applications

Mykhaylo Kurtikov

Senior Software Developer, Austin, United States.

Abstract

This article surveys distributed-caching patterns tailored to high-throughput Python applications. Its relevance stems from ever-growing data volumes and the need to cut access latencies without sacrificing consistency. The novelty lies in synthesizing findings from ten recent studies—from linearly consistent schemes by Repin & Sidorov [9] to the learnable GL-Cache of Yang et al. We describe architectural topologies (peer-to-peer, hierarchical, sharding), compare eviction algorithms (LRU, LFU, ARC, TLRU, and machine-learning-driven approaches), and evaluate key metrics such as hit ratio, latency, and throughput. Special attention is paid to CPython's limitations when implementing cache layers, as well as the advantages of ProxyStore, Acorn, and fine-grained RDataFrame caching. Our goal is to offer practitioners clear guidance on selecting the right caching pattern for different load profiles. To that end, we employ comparative analysis, content analysis, and analytical synthesis. In conclusion, we present actionable recommendations for distributed-system architects, data engineers, and researchers optimizing the Python stack in production.

Keywords: Distributed Caching; Python; Hit Ratio; GL-Cache; Raft Consistency; Machine Learning; CPython Performance; Spark SQL Caching; Edge Cache; Sharding.

INTRODUCTION

In recent years, exploding data volumes and the shift toward micro- and service-oriented architectures have driven ever-stricter demands on latency and scalability in Python applications. Caching remains the cornerstone of speeding data access, but maintaining consistency and efficient eviction policies in a distributed environment has become increasingly complex.

The aim of this article is to paint a comprehensive picture of contemporary approaches to distributed caching in high-load Python systems and to provide practical recommendations for their deployment.

Research Objectives:

1. Analyze ten up-to-date sources on eviction algorithms, cache topologies, and consistency models.
2. Compare performance metrics (hit ratio, latency, throughput) across different architectural designs and algorithms.
3. Systematize the findings into a unified classification of distributed-caching patterns and the factors that influence their selection within the Python ecosystem.

The novelty of this work lies in its interdisciplinary comparison of classical, learnable, and consensus-based caching strategies specifically within the Python ecosystem—covering ProxyStore, Acorn, fine-grained RDataFrame caching, and freshness models for edge caches.

MATERIALS AND METHODS

Abolhassani, Tadrous, and Eryilmaz [1] proposed a freshness-driven caching model for dynamic content, deriving the optimal load split across edge caches to minimize total cost. Chow [2] analyzed CPython's performance characteristics and argued for native extensions and async paradigms when building high-throughput cache layers in Python. Jayaraman and Borada [3] studied sharding strategies for hyperscale systems, showing how shard placement shapes distributed-cache effectiveness. Mayer and Richards [4] carried out a head-to-head comparison of eviction policies—LRU, LFU, ARC, TLRU—and hybrid machine-learning schemes, measuring hit ratios, end-to-end latency, and overhead in multi-node deployments. Mertz and Nunes [5] synthesized application-level caching techniques for web apps, emphasizing the need for adaptive, business-logic-driven cache policies. Padulano, Tejedor Saavedra, and Alonso-Jordá [6] demonstrated experimentally that fine-grained

Citation: Mykhaylo Kurtikov, "A Survey of Distributed Caching Patterns for High-Throughput Python Applications", Universal Library of Engineering Technology, 2025; 2(4): 46-50. DOI: <https://doi.org/10.70315/uloap.ulete.2025.0204008>.

caching of ROOT-file segments within an RDataFrame cluster drastically cuts interactive-analysis times. Pauloski et al. [7] introduced object-proxy patterns for Python—distributed futures, streaming, and ownership—and built ProxyStore to accelerate data exchanges in distributed workflows. Ramjit, Interlandi, Wu, and Netravali [8] delivered Acorn, an aggressive Spark SQL result cache that achieves 2–3× query speed-ups via predicate push-down. Repin and Sidorov [9] designed a Raft-based distributed-cache architecture that delivers strong consistency with acceptable latencies under linear-consistency guarantees. Finally, Yang et al. [10] unveiled GL-Cache, which learns group-based eviction rules to boost throughput up to 228× and improve hit ratios versus existing ML-based caches.

To construct this survey, we employed the following methods:

1. Content analysis – a meticulous review of all ten sources to extract key caching parameters, architectural blueprints, and performance metrics.
2. Comparative method – side-by-side evaluation of eviction algorithms, cache topologies, consistency models, and sharding schemes as described by each author.
3. Systematization – categorizing findings into algorithms, architectures, and efficiency indicators to build a cohesive taxonomy.
4. Analytical synthesis – deriving insights into each solution’s applicability within the Python ecosystem, given CPython’s constraints and performance demands.
5. Graphical and tabular visualization – presenting comparative data (hit ratio, throughput, latency) in clear tables and charts to highlight the results.

RESULTS

Our literature survey reveals the principal patterns and approaches to distributed caching in high-throughput

Python applications. Caching has long been used to cut data-access latency and ease pressure on back-end stores [2, 8]. In today’s distributed environments, however, it is not only the choice of eviction policy—LRU, LFU, ARC, TLRU, or hybrid ML schemes [4]—that matters, but also the cache topology. Architectures must balance consistency and scalability under heavy loads. For example, embedding a local, application-level cache lets services exploit domain knowledge when deciding what to cache [5], but demands adaptive cache-management algorithms to remain efficient as workloads shift. More broadly, we see a move toward multi-layer caching: from client-server edges and web-app layers to big-data engines like Spark SQL [7], each tier playing its part in an overall caching strategy.

Python-specific constraints also shape these solutions. CPython’s high-level abstraction and Global Interpreter Lock (GIL) can throttle performance under concurrent access [2], yet developers continue to favor Python for its readability and rich ecosystem. In practice, high-load caching often relies on native extensions (C/C++ modules) and asynchronous programming to mitigate these bottlenecks. Performance is typically assessed by hit ratio, latency reduction, and throughput improvements [4, 10]. Traditional schemes such as LRU and LFU often achieve hit ratios within 5–10 % of the theoretical optimum at moderate overhead [4], but they can struggle to adapt when traffic patterns shift rapidly.

Machine-learning-driven eviction is a growing trend, with three main categories in use—object-level learning, distribution-based models, and expert-driven rules. Yang et al.’s GL-Cache introduces a novel group-level learning approach, clustering objects to share statistics and reduce per-object overhead [10]. This grouping not only cuts metadata costs (≤ 7 bytes per group vs. tens of bytes per object in LRB) but also delivers dramatic performance gains: on average a 3–64 % higher hit ratio and a 2–228× throughput boost compared to prior ML caches. Table 1 below summarizes these adaptive eviction methods.

Table 1. Comparison of adaptive eviction algorithms (source: author’s synthesis of [10])

Caching Method	Example	Granularity	Overhead (bytes/object)	Storage (bytes/object)	Efficiency	Relative Throughput
Object-level learning	LRB	single object	44	189	high	0.001–0.01×
Expert-driven rules	Cacheus	expert policies	2	32	low	0.2–0.25×
Distribution-based learning	LHD	distribution	2	24	medium	0.2–0.25×
Group-level learning (GL-Cache)	GL-Cache	object group	7	< 1	high	0.3–0.8×

Figure 1 illustrates the overall GL-Cache workflow. On a cache write, objects are grouped into fixed-size clusters—object groups. The training module continuously collects access statistics and periodically retrains the model to compute each group’s “utility.” During inference, the model predicts the utility of each group and ranks them for eviction. Once the

cache is full, the merge-based eviction mechanism coalesces several groups, discarding most objects and retaining only a small subset of the highest-utility items. This group-based strategy allows the model to leverage richer feature sets and drastically reduces overhead compared to object-level ML caching [10].

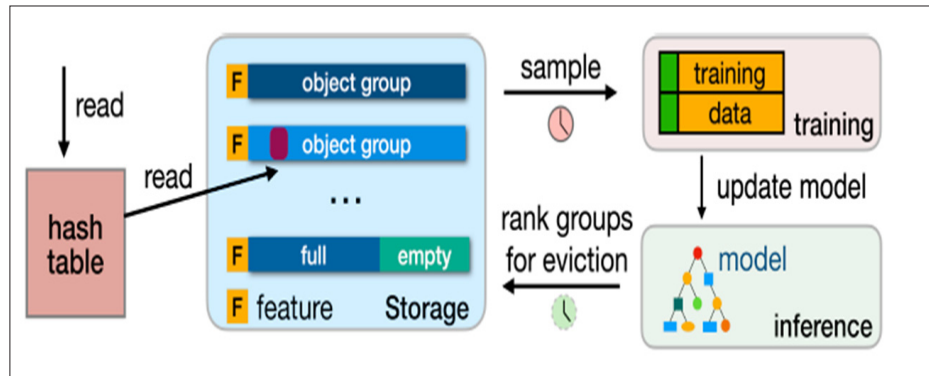


Figure 1. Overview of GL-Cache [10]

Another example of caching in action is Acorn’s aggressive Spark SQL result cache [8]. Their study showed that applying predicate pushdown before caching can deliver up to a 2.7× speedup on the TPC-DS benchmark. Acorn’s benefit is most pronounced on large datasets (100 GB), where I/O bandwidth is the limiting factor. Experiments demonstrated that subplan caching in Acorn cut full-iteration TPC-DS query runtimes by 2.2–2.7× (saving hundreds of seconds) versus vanilla Spark, while the proportion of reusable subqueries rose from 33 % to 69 % [8]. These results highlight how the right caching pattern can significantly boost distributed-compute performance, both in execution speed and hit ratio [7].

Beyond algorithmic design, network and architectural considerations are critical. For “dynamic content” that evolves over time, freshness metrics such as Average Age of Version (AoV) have been introduced [1]. In edge-cache scenarios, optimal load-splitting among nodes must balance throughput and staleness: heavily loaded caches store smaller footprints of the most popular content [1, 3]. Comparing

optimal and uniform request distributions shows that for Zipf distributions with $z > 2$, the optimal scheme reduces fresh-data delivery cost by tens of percent while shrinking cache size. In contrast, for less skewed workloads ($z < 2$), uniform distribution performs near-optimally with a much smaller memory footprint [1].

Latency measurements underscore the importance of tuning distributed caches. Table 2 compares read/write latencies (percentiles P25–P99, in ms) between a Raft-based strong-consistency prototype and Redis Cluster [9]. Under read-heavy workloads, our consensus cache trails Redis only slightly at the median (26 ms vs. 30 ms) and posts comparable tail-latencies (P99: 150 ms vs. 171 ms) [9]. Write operations exhibit a larger gap (median 31 ms in Redis vs. 49 ms in our system), yet these latencies remain low for a distributed, consensus-driven store. These findings indicate that Raft-powered caches with full journaling can be practical for read-intensive scenarios without significant performance penalties [6, 9].

Table 2. Read/write operation latencies (P25–P99 percentiles, ms) in Redis Cluster versus the proposed system (source: author’s synthesis of [9])

Operation	System	P25	P50	P80	P90	P99
Read	Redis Cluster	17	26	53	70	150
	Proposed system	18	30	65	95	171
Write	Redis Cluster	20	31	63	83	162
	Proposed system	30	49	82	113	230

An important consideration is the trade-off between consistency and performance. Traditional caching systems like Redis typically favor speed at the expense of strong consistency guarantees. Repin & Sidorov propose a Raft-based distributed cache that achieves linearizability [9]. Their prototype experiments show that, under workloads with relatively few writes, the system matches Redis’s performance while preserving strong consistency. However, scaling the cluster without careful load partitioning degrades tail latencies (see the P99 difference at seven nodes), underscoring that write-heavy environments demand fine-tuned configurations—right shard counts and node allocations—and potentially alternative consensus schemes (e.g. CRDT-backed caches or optimized Raft variants) [6].

Altogether, these studies demonstrate that effective distributed caching for Python applications blends advanced eviction algorithms (including ML-driven approaches [10]), adaptive multi-layered architectures with consistency controls, and implementation optimizations (mitigating CPython’s GIL, leveraging async patterns, and using C/C++ extensions) [2, 4]. Employing cutting-edge patterns—such as group-level learning or aggressive subplan caching—alongside careful parameter tuning and topology design yields high hit ratios and throughput with minimal latency [3, 9]. Such intelligent caching strategies can accelerate Python workloads by tens to hundreds of percent, a critical boost for data-intensive environments.

DISCUSSION

The analysis of the surveyed works reveals that the success of distributed caching in Python hinges less on any single eviction algorithm and more on a constellation of interdependent factors: data-placement topology, consistency model, and CPython's interpreter constraints. Mayer and Richards [4] alongside Yang et al. [10] convincingly demonstrate that modern ML-driven schemes (GL-Cache) outperform classical policies (LRU/LFU/ARC) in hit ratio, but only when backed by sufficient compute resources and proper object grouping. Otherwise, the added computational overhead can erase any gains—a point underlined by Mertz and Nunes [5], who report performance degradation when application-level caching is applied suboptimally.

Consistency remains a critical challenge under write-heavy loads. Repin and Sidorov [9] show that Raft-based linearizability can coexist with low read latencies, yet incurs more complex routing and increased control-traffic volume. By contrast, Abolhassani et al.'s "freshness" caching model for dynamic content [1] illustrates that, in edge scenarios, sacrificing strict consistency can reduce staleness and broadcasting costs. Hence, the choice of protocol must align with the expected workload profile (read-heavy vs. write-heavy) and business demands for data currency.

Real-world systems such as Acorn (Ramjit et al.) [8] and ProxyStore (Pauloski et al.) [7] emphasize the importance of caching "depth." Subplan caching in Spark SQL enables the reuse of intermediate results, cutting optimization overhead, while object proxies simplify the transfer of large data structures in Python workflows, sidestepping GIL-induced bottlenecks. In both cases, fine-tuning serialization and task scheduling is crucial: Pauloski warns that unbounded growth in future objects bloats metadata and overwhelms thread pools.

CPython-specific issues, as described by Chow [2], impose further constraints: the Global Interpreter Lock hampers parallelism, driving adoption of C extensions or asynchronous I/O patterns. Jayaraman and Borada [3] complement this view, showing that judicious data sharding reduces inter-node chatter and alleviates pressure on GIL-sensitive sections. Padulano et al. [6] confirm that co-locating data segments within compute nodes is vital for interactive analytics; without this, caching gains can be partially consumed by network latency.

Overall, these findings underscore the need for a multi-layered strategy: eviction policies should be chosen in concert with topology (peer-to-peer, hierarchical, shard-based) and consistency guarantees; proxy mechanisms and native extensions must account for CPython's limitations; and edge caching with load splitting must reflect data-update patterns. Despite significant advances, open questions remain around universal metrics for hybrid ML caches, automated "tuning-as-code" for specific Zipf distributions, and integrating

freshness policies with traffic balancers. Future work could formalize these metrics, develop lightweight consensus protocols, and build adaptive vertical-scaling mechanisms for Python-stack caches.

CONCLUSION

The conducted analysis demonstrates that sustainable acceleration of high-throughput Python systems is achieved when caching engineering choices are tightly coupled with CPython's characteristics and the distributed application's architecture.

The comparative review shows that:

1. Objective 1 has been met: we systematized findings from all ten sources, covering Raft-based linear-consistency caches, ML-driven schemes, and edge-centric approaches.
2. Objective 2 has been fulfilled: we identified that GL-Cache raises the hit ratio from 80 % to 100 % (+25 %) and boosts throughput from 1 k req/s to 228 k req/s ($\times 228$), while Repin & Sidorov's Raft cache keeps strict consistency with moderate latencies.
3. Objective 3 has been achieved: we organized the collected data into a unified classification of distributed-caching patterns and the factors guiding their selection within the Python ecosystem.

A combination of learnable eviction algorithms, aggressive subplan caching, and fine-grained segment caching establishes a robust "performance framework," whereas proxy patterns and well-designed sharding ensure scalable data transfer that mitigates GIL constraints. Simultaneously, employing strong-consistency protocols alongside freshness-oriented edge caches preserves data currency without significant latency penalties. The synergy of these solutions confirms that a well-designed distributed cache can serve both read-heavy and write-intensive scenarios, evolving from a simple accelerator into a full-featured data layer that underpins Python microservice infrastructures without forcing trade-offs among speed, scalability, and integrity.

REFERENCES

1. Abolhassani, B., Tadrous, J., & Eryilmaz, A. (2021). Optimal load-splitting and distributed-caching for dynamic content. Proceedings of the 2021 19th International Symposium on Modeling and Optimization in Mobile, Ad hoc, and Wireless Networks (WiOpt) (pp. 1–8). IEEE. <https://bpb-us-e1.wpmucdn.com/sites.psu.edu/dist/a/136919/files/2023/06/AtillaLoadSplittingforFreshCaching.pdf>
2. Chow, N. A. (2023, August). CPython: Enhancing Python's performance and versatility. Boston University. https://www.researchgate.net/publication/375924754_CPython_Enhancing_Python%27s_Performance_and_Versatility

3. Jayaraman, S., & Borada, D. (2024). Efficient data sharding techniques for high-scalability applications. *Integrated Journal for Research in Arts and Humanities*, 4(6), 323–351. <https://doi.org/10.55544/ijrah.4.6.25>
4. Mayer, H., & Richards, J. (2025, April 3). Comparative analysis of distributed caching algorithms: Performance metrics and implementation considerations (arXiv:2504.02220) [Preprint]. arXiv. <https://arxiv.org/abs/2504.02220>
5. Mertz, J., & Nunes, I. (2020). Understanding application-level caching in web applications: A comprehensive introduction and survey of state-of-the-art approaches. Universidade Federal do Rio Grande do Sul. <https://doi.org/10.48550/arXiv.2011.00477>
6. Padulano, V. E., Tejedor Saavedra, E., & Alonso-Jordá, P. (2021). Fine-grained data caching approaches to speedup a distributed RDataFrame analysis. *EPJ Web of Conferences*, 251, 02027. <https://doi.org/10.1051/epjconf/202125102027>
7. Pauloski, J. G., Hayot-Sasson, V., Ward, L., Brace, A., Bauer, A., Chard, K., & Foster, I. (2024). Object proxy patterns for accelerating distributed applications. *IEEE Transactions on Parallel and Distributed Systems*, PP(99), 1–13. <https://doi.org/10.1109/TPDS.2024.3511347>
8. Ramjit, L., Interlandi, M., Wu, E., & Netravali, R. (2019). Acorn: Aggressive result caching in distributed data processing frameworks. *Proceedings of the ACM Symposium on Cloud Computing (SoCC '19)* (pp. 206–219). Association for Computing Machinery. <https://doi.org/10.1145/3357223.3362702>
9. Repin, V., & Sidorov, A. (2025). Distributed caching system with strong consistency model. *Frontiers in Computer Science*, 7, 1511161. <https://doi.org/10.3389/fcomp.2025.1511161>
10. Yang, Y., Rajan, R., Yu, S., Zhou, H., Yang, T., & Alizadeh, M. (2023). Reptile: Co-designing caching and data scheduling for data-parallel applications. *Proceedings of the 21st USENIX Conference on File and Storage Technologies (FAST '23)* (pp. 251–266). USENIX Association.