



IE2050

Operating Systems

2nd Year, 2nd Semester

Assignment

In-Depth Analysis of Operating System Concepts through Research Paper Exploration

Submitted to

Sri Lanka Institute of Information Technology

Student ID	Name
IT23169180	G.K.T.M.Abesinghe
IT23227590	N.D.Kariyawasam
IT23363052	K.H.S.Damsiluni

In partial fulfillment of the requirements for the
Bachelor of Science Special Honors Degree in Information Technology

25/04/2025

Table of Contents

Introduction.....	3
Problem Being Addressed.....	3
Proposed Solution	4
1. Lifetime Prediction Using Machine Learning	4
2. LLAMA (Learned Lifetime-Aware Memory Allocator)	4
Results and Implications	4
Critical Evaluation	5
Strengths	5
Weaknesses	5
Areas for Future Research	5
Relating the Paper to Broader OS Concepts	6
Relevance to Modern Operating Systems	6
Applications and Future Extensions	6
References.....	7
Maas, M., Andersen, D. G., Isard, M., Javanmard, M. M., McKinley, K. S., & Raffel, C. (2024). <i>Combining Machine Learning and Lifetime-Based Resource Management for Memory Allocation and Beyond</i> . Communications of the ACM. Combining Machine Learning and Lifetime-Based Resource Management for Memory Allocation and Beyond – Communications of the ACM	7

Introduction

Memory management is a fundamental operating system design principle whose effects lie directly upon the performance, efficiency, and dependability of general-purpose and specialty computing ecosystems. Memory allocator issues have emerged due to the increased workload sophistication of modern workloads, especially in cloud environments, data center installations, and other evolving computing scenarios. The best optimization in those environments is often the use of large memory pages - i.e., 2MB "huge pages." Huge pages reduce address translation overhead by improving TLB coverage, consequently decreasing CPU stalls and improving execution throughput.

However, as beneficial as it is, massive page use also reveals massive memory management problems, primarily memory fragmentation. It happens when long-lived objects are mixed with huge pages in a manner that memory cannot be freed even when most objects in a page are short-lived and their lifespan has ended. In their paper "Combining Machine Learning and Lifetime-Based Resource Management for Memory Allocation and Beyond" the authors solve this issue by proposing a memory allocation strategy that employs machine learning to predict object lifetimes at allocation time, thus enabling improved page utilization and memory reclamation.

Problem Being Addressed

Traditional memory allocators, especially in C++ systems like TCMalloc, typically operate with memory by splitting allocations into size classes. But they don't look at how long a specific allocated object will be around. This creates a significant inefficiency in using huge pages: a page cannot be reclaimed until all objects of the page have been deallocated. Even a single long-lived object on an otherwise empty page can prevent its reclamation, resulting in wasted memory and decreased performance.

This inefficiency is compounded for server workloads that typically consist of many small-lived allocations punctuated with some long-lived objects like session state or cached data. The inability of the allocator to group objects of similar lifetimes together results in short-lived and long-lived objects being grouped together, and thus memory compaction or reclamation becomes almost impossible without object movement—an operation not possible in C++.

The main problem, then, is a lack of information regarding object lifetimes at allocation time, which makes it difficult to organize memory layout in a way that maximizes the use of large pages. Current profiling methods either incur too much overhead (for online profiling) or cannot provide full coverage (for offline profiling), especially across application versions and configurations.

Proposed Solution

To clear up this, the authors propose a two-stage solution of predictive modeling and a redesigned allocator.

1. Lifetime Prediction Using Machine Learning:

The authors employ an LSTM model that is trained on stack traces observed during sampling runs. These stack traces—essentially the call stacks leading to memory allocations—are encoded as token sequences, equivalent to natural language. The model learns to predict for every allocation context one of several lifetime buckets ($\leq 10\text{ms}$, 100ms , 1s , etc.). The allocator employs this classification to infer the lifetime of new objects at runtime.

2. LLAMA (Learned Lifetime-Aware Memory Allocator):

LLAMA is a memory allocator that organizes memory not by size, but by estimated object lifetimes. Huge pages are assigned to some lifetime classes, and objects with similar expected lengths are grouped together. The organization enables more pages to be fully released and returned to the OS when all of their objects die, reducing fragmentation considerably.

The mechanism also includes techniques for detecting mispredictions. If an object lives longer than its predicted class, LLAMA can reclassify the page and raise its expected lifetime. The opposite can be performed earlier when lifetimes are shorter than expected, with pages being reclassified, which is efficient. A hashing-based cache makes most predictions at very low overhead, invoking the full ML model only on cache misses.

Results and Implications

The authors evaluate LLAMA on four workloads that are representative of common use cases: an image-processing server, TensorFlow inference using InceptionV3, a large data processing pipeline, and the Redis key-value store. These workloads exercise a wide range of memory allocation patterns and lifetimes, and stress memory systems with large thread counts, allocation sizes, and allocation frequencies.

The key findings are:

- Decrease in fragmentation: LLAMA reduces memory fragmentation by 19% to 78% for diverse workloads. Redis and the data pipeline, in particular, see the most drastic improvements.
- Memory efficiency: Aggregate memory usage, for example, drops in the image server workload from 664MB (with TCMalloc) to 446MB with LLAMA.
- Prediction accuracy: Classification accuracy of lifetime ranges from 94% to 100% when the weighting is based on allocation frequency.

- Low performance overhead: Despite ML predictions, LLAMA introduces little latency due to its cache approach. Cache-hitting average allocations are only marginally slower than with traditional allocators.

These results demonstrate that ML lifetime prediction is not only feasible at runtime but also beneficial for improving memory utilization and reclaimability without paying a high performance penalty.

Critical Evaluation

The approach of the paper is innovative but pragmatic, with robust interaction between systems engineering and machine learning. Specifically, it is a paradigm shift in memory allocator design from reactive heuristics to proactive prediction-based approaches.

Strengths:

- Introduces a new paradigm for memory allocators that synchronizes allocation with object lifetime, reducing fragmentation by orders of magnitude.
- Uses machine learning in a lean, practical fashion, demonstrating that inference can be fast enough for system-level pieces.
- Provides a highly generalizable solution that might be extended from memory allocation to storage, scheduling, and power management.

Weaknesses:

- Relying on training data from past workloads can limit generalizability to completely new applications or radically different versions of existing ones.
- The effectiveness of the ML model is contingent on the quality and representativeness of sampled stack traces.
- Misprediction handling, while performing well, might still leave some fragmentation in edge cases.

Areas for Future Research:

- Exploring more sophisticated models (e.g., Graph Neural Networks or Transformers) that could lead to better generalization across yet-to-be-known stack traces.
- Extending the idea of lifetime-aware allocation to other OS resources such as file systems, process scheduling, or network buffer management.
- Investigating the integration of LLAMA with dynamic observability systems and system telemetry to enable real-time adaptation without retraining.

In summary, this work is a promising step toward intelligent operating systems that learn to adapt to application behavior using machine learning. LLAMA is an operational proof of concept of how predictive analytics can be embedded in the fabric of system software to yield measurable performance and efficiency gains.

Relating the Paper to Broader OS Concepts

Contribution to Operating Systems Understanding

This paper provides a comprehensive analysis of how intelligent memory management can evolve with the help of machine learning. Rather than using advanced heuristics to enhance allocators, it proposes the use of intelligence based on prediction. This is a paradigm shift in memory allocator design—moving from reactive management to predictive, proactive allocation.

It also points to the importance of relating software behavior (stack traces, code paths) with physical memory layout—a realm that's commonly segmented. The allocator has knowledge of application semantics with allocation contexts in LLAMA, bringing user-space behavior together with kernel-level memory management.

Relevance to Modern Operating Systems

Modern operating systems must serve more complex workloads—AI/ML inference, web-scale micro services, and multi-tenant configurations. LLAMA's ML-driven design demonstrates how OSes in the future can leverage intelligent, data-driven mechanisms for improved efficiency and performance.

Thanks to the availability of observability tools and tracing systems, it is highly plausible to integrate ML for optimization - beyond memory - into file systems, I/O scheduling, and process management.

Applications and Future Extensions

LLAMA's algorithm can be used in any computationally orientated space-time bin packing problem. Future potential areas include

- **Storage Systems:** Estimating file lifetimes to avoid wasting SSD wear.
- **Process Management:** Grouping short-lived and long-lived processes on different cores for better cache usage.
- **Cloud Resource Allocation:** Assignments of containers or VMs to hosts based on resource duration forecast.
- **Power Management:** Power gating idle memory banks based on expected access patterns.

References

Maas, M., Andersen, D. G., Isard, M., Javanmard, M. M., McKinley, K. S., & Raffel, C. (2024). *Combining Machine Learning and Lifetime-Based Resource Management for Memory Allocation and Beyond*. Communications of the ACM. [Combining Machine Learning and Lifetime-Based Resource Management for Memory Allocation and Beyond – Communications of the ACM](#)