

# lgmalloc: Predictive Memory Allocation

Luca M. Goddijn  
luca.goddijn@gmail.com  
Independent Researcher  
Amsterdam, Netherlands

## Abstract

Memory allocation represents a critical bottleneck in multi-threaded applications. Traditional allocators, like `tcmalloc`, suffer from central contention points that limit scalability, while thread-local implementations require substantial per-thread memory overhead. `lgmalloc` shifts the demand-reactive paradigm that standard allocation algorithms use, to a preemptive and predictive heuristic approach. The approach leverages the fundamental observation that call sites demonstrate predictable allocation patterns. By parsing the executable's ELF format before the program's entry point, `lgmalloc` performs static analysis of call sites to reveal allocation patterns and generate confidence scores for size classes. Based on the confidence scores, the system utilizes optimized code paths to the backend allocation API, reducing allocation latency. When confidence scores are low, the system falls back to reliable, standard allocation strategies. The thread-local design eliminates lock contention while substantially reducing per-thread memory overhead by optimizing the pre-allocated size classes each thread requires.

`lgmalloc` is available at: <https://github.com/Arty3/lgmalloc>

## Keywords

memory allocation, heap, malloc, allocator, static analysis, predictive optimization, thread-local, heuristics

## 1 Introduction

As performance requirements scale, memory allocators face a rising trade-off barrier between memory efficiency and thread scalability. Existing and standardized allocators demonstrate a clear division: optimizing for performance or minimizing memory overhead. Both strategies significantly compromise in the opposing dimension. `tcmalloc` [3] is Google's custom memory allocator, designed for high-performance infrastructure applications, while `mimalloc` [8] is Microsoft's general-purpose allocator which focuses on minimizing memory overhead. Both implementations clearly illustrate the fundamental trade-offs that limit current allocation strategies.

These trade-offs manifest in concrete design choices. Performance-focused allocators like `tcmalloc` employ fundamental strategies like thread-local caches and per-size-class free lists to minimize lock contention, which comes at the cost of substantial memory

overhead required to maintain these structures across threads. Conversely, memory-efficient implementations reduce per-thread allocation and utilize compact data structures. These sacrifice performance by bottlenecking critical contention points, which scale poorly at higher thread counts. The fundamental issue underlying these trade-offs is that the core principles allocators rely on are founded on a demand-reactive paradigm. This paradigm entails larger-than-needed allocations, organized in core data structures which are plagued by thread contention, or poor memory efficiency at scale.

`lgmalloc` is a general-purpose `malloc(3)` implementation that addresses this paradigm by leveraging the fundamental observation that allocation call sites exhibit significant predictability. This observation introduces three key innovations: Static analysis, confidence-based routing and a thread-local design whose effectiveness emerges through its integration with the other two approaches. `lgmalloc` statically analyzes the executable's structure before the program's main entry point, predicting allocation patterns and creating optimized code paths for allocation logic. This approach offers significant optimization potential for the key trade-off elements.

Consider an ideal scenario: a multi-threaded server where each thread allocates a compile-time constant size to process incoming data. The compile-time constant yields maximum confidence, constituting a perfectly optimized allocation system, where each thread only allocates exactly what it requires. The confidence-based routing system can now perform a direct jump to the respective backend allocation API function, which follows linear and predictable logic, offering minimal side effects. In contrast, existing allocators face bottlenecks at scale, for example `tcmalloc` experiences high contention on its global free list [7] and `jemalloc` [2] experiences significant memory overhead [5]

The key contributions of this work are:

- A novel static analysis system for predicting allocation patterns
- A confidence-based routing mechanism for optimized allocation paths
- A thread-local design that minimizes both contention and memory overhead
- Comprehensive evaluation showing [...]

The remainder of this paper is organized as follows: Section 2 surveys related work, Section 3 presents the `lgmalloc` architecture, Section 4 details key implementation aspects, Section 5 evaluates performance, Section 6 discusses results and implications, and Section 7 concludes.

Permission to make digital or hard copies of all or part of this work for personal or internal use, or for the internal or personal use of specific clients, is granted by ACM for users registered with ACM, provided that the fee of \$12.00 is paid directly to ACM. This fee code for users registered with ACM: 978-1-60959-123-4/18/06...\$12.00. Copyright 2018 ACM 978-1-60959-123-4/18/06...\$12.00. This work is published in the Proceedings of the ACM on Operating Systems and Security, Vol. 1, No. 1, 2018.

© 2025 Copyright held by the owner/author(s). Publication rights licensed to ACM.  
ACM ISBN 978-x-xxxx-xxxx-x/YYYY/MM  
<https://doi.org/10.1145/nnnnnnn.nnnnnnn>

2025-06-18 23:47. Page 1 of 1–2.

## 2 Background and Related Work

### 2.1 Evolution of Memory Allocation Strategies

### 2.2 Allocation Pattern Predictability

### 2.3 Static Analysis for Runtime Optimization

### 2.4 Compiler Techniques in Runtime Systems

### 2.5 Benchmarking Methodologies

### 2.6 Security Considerations

### 2.7 Constrained Systems

## 3 System Design

## 4 Implementation

### 4.1 Codebase Features

### 4.2 Configuration System

### 4.3 Build System and Compiler Integration

lmalloc focuses on maintaining a simple and comprehensive build system. The approach uses a single GNU Make-based makefile [10], prioritizing directness and uniformity over modular complexity. The makefile performs platform compliance verification to ensure Unix/POSIX compatibility and ELF support through runtime checks, since the implementation relies on those components being present for its innovative approach to memory allocation.

The release and debug targets generate optimized and development builds respectively, with release serving as the default target to prioritize production-ready binaries. The choice of compilers reflects distinct optimization and reliability philosophies: debug builds utilize Clang for its superior diagnostic capabilities, faster compilation times, and comprehensive error reporting [6, 9], making it ideal for development iterations where issue identification takes precedence over optimization. Conversely, release builds employ GCC for its more sophisticated optimization framework [1], leveraging `-flto=full` for link-time optimization. These compiler choices enable lmalloc's extensive use of compiler attributes and built-ins (Section 4.1) made possible by the GNU17 standard [4], while strict warning-as-error policies ensure high code quality standards.

## 5 Evaluation

### 5.1 Behavioral Unit Testing

## 6 Discussion

## 7 Conclusion

## References

- [1] Colfax Research. 2017. *Compiler Comparison: GCC vs. Intel vs. Clang*. Technical Report. Colfax Research. <https://colfaxresearch.com/compiler-comparison/> Performance comparison of GCC, Intel C++ Compiler, and Clang across various benchmarks.
- [2] Jason Evans. 2006. A scalable concurrent malloc(3) implementation for FreeBSD. In *Proceedings of the BSDCan Conference*. Ottawa, Canada.
- [3] Google. [n. d.]. TCMalloc: Thread-Caching Malloc. <https://google.github.io/tcmalloc/>
- [4] ISO/IEC. 2018. *Information technology – Programming languages – C*. Technical Report ISO/IEC 9899:2018. International Organization for Standardization. C17 standard with GNU extensions (GNU17).
- [5] Rivalino Matias Jr., Autran Macedo, Luciano C. Serafim, Rodrigo M. Rocha, and Felipe M. G. França. 2011. An Experimental Study on Memory Allocators in

Multicore and Multithreaded Applications. In *Proceedings of the 12th International Conference on Parallel and Distributed Computing, Applications and Technologies (PDCAT)*. IEEE, 92–98. doi:10.1109/PDCAT.2011.53 Comparative study showing jemalloc's thread-cache behavior and memory consumption patterns in multi-threaded scenarios.

- [6] Chris Lattner and Vikram Adve. 2004. LLVM: A Compilation Framework for Lifelong Program Analysis & Transformation. In *International Symposium on Code Generation and Optimization*. IEEE, 75–86. doi:10.1109/CGO.2004.1281665
- [7] Sangho Lee, Teresa L. Johnson, and Easwaran Raman. 2014. Feedback Directed Optimization of TCMalloc. In *Proceedings of the Workshop on Memory Systems Performance and Correctness*. ACM, New York, NY, USA. doi:10.1145/2618128.2618131 Study demonstrating performance bottlenecks in TCMalloc's central cache management and thread-local cache interactions.
- [8] Daan Leijen. 2019. mimalloc: Free List Sharding in Action. In *ISMM*.
- [9] LLVM Project. [n. d.]. Expressive Diagnostics. <https://clang.llvm.org/diagnostics.html> Clang's diagnostic system documentation.
- [10] Richard M. Stallman, Roland McGrath, and Paul D. Smith. 2020. *GNU Make: A Program for Directing Recompilation* (4.3 ed.). Free Software Foundation. <https://www.gnu.org/software/make/manual/> Version 4.3.