# From Profiling to Optimization: Unveiling the Profile Guided Optimization

BINGXIN LIU, Beijing Normal University, China

YINGHUI HUANG, Beijing Normal University, China

JIANHUA GAO, Beijing Normal University, China

JIANJUN SHI, Beijing Normal University, China

YONGPENG LIU, Phytium Technology Co., Ltd., China

YIPIN SUN, Phytium Technology Co., Ltd., China

WEIXING JI, Beijing Normal University, China

Profile Guided Optimization (PGO) uses runtime profiling to inform compiler decisions, effectively combining static analysis with actual execution behavior to enhance performance. Runtime profiles — acquired through instrumentation or hardware- and software-assisted sampling — provide detailed insights into control flow, branch predictions, and memory access patterns. This survey systematically categorizes PGO research by profiling method (instrumentation vs. sampling), optimization stages (compile time and link/post-link time), compiler integration (GCC, LLVM), and target architectures. Key algorithms and frameworks are shown in terms of design principles. Performance experiments on representative examples demonstrates PGO's speedups, overheads, and integration maturity. Finally, we identify open challenges, such as reducing sampling overhead, supporting dynamic input workloads, and enabling cross-architecture portability, and propose future research directions to advance adaptive, low-overhead optimizing compilers.

CCS Concepts: • **Software and its engineering** → **Compilers**; **Software performance**; • **Theory of computation** → **Program analysis**.

Additional Key Words and Phrases: compilers, profile guided optimization, instrumentation, hardware sampling

## 1 INTRODUCTION

Optimization of program performance has been a key research direction, particularly for compute-intensive and latency-sensitive applications[44]. As the post-Moore's Law era sets in, "clock frequencies plateau and core counts climb", software increasingly rely on advanced compiler optimizations to extract instruction-level parallelism and minimize memory stalls. In this context, compilers play

**111**

a critical role as the bridge that translates high-level language programs into efficient machine code for the target platform. Traditional compiler optimizations primarily employ static analysis techniques to optimize the program's intermediate representation before actual execution, improving the performance of final executable. These static optimizations are widely adopted due to their ease of use and low runtime overhead[47].

However, with the evolution of computer architectures, such as the prevalence of out-of-order execution and multi-stage pipelined processors, there is an increasing demand for more efficient micro-instruction scheduling within the compiler. To fully exploit the processor's instruction-level parallelism, optimizations often align branch targets based on predetermined heuristics. If the aligned (predicted) branch is taken at runtime, the CPU avoids pipeline stalls, yielding substantial performance gains[23]. Conversely, a misprediction forces the CPU to flush the pipeline and restore registers, resulting in significant penalties. Consequently, compiler designers strive to maximize branch-prediction accuracy.

In the early days, programs were simple enough that static analysis provided sufficiently accurate guidance for optimization. But as applications have grown in size and complexity, static techniques struggle to quickly, comprehensively, and accurately characterize runtime behavior or identify hot code regions. The inability to accurately predict dynamic execution paths marks a significant limitation of purely static optimization[21]. Meanwhile, software tracing or hardware-based sampling on modern CPUs provided elaborate hardware support (Performance Monitoring Units, Last Branch Records, instruction-based sampling) that can reveal precise runtime statistics with low overhead[16].

Profile Guided Optimization (PGO), also known as feedback-directed optimization (FDO), bridges the gap between static analysis and dynamic behavior, enabling compilers to reorder, inline, and lay out code based on actual execution counts, branch-taken probabilities, and memory access profiles. By driving optimizations with real-world execution traces, PGO can yield speedups often in the range of 5%-30% on real applications [59, 60], far surpassing what purely static heuristics can deliver.

As a result, the body of PGO research now spans instrumentation algorithms, software- and hardware-based sampling methods, mapping algorithms, compiler-internal optimizations, and full end-to-end toolchains[10, 26, 55]. In this article, we provide a comprehensive overview of PGO by first examining the spectrum of profiling techniques, ranging from software instrumentation-based edge and path profiling, software sampling-based dynamic profiling techniques to hardware-sampling mechanisms. Building on this foundation, we show how major compiler frameworks integrate feedback data at various stages of the build process. To illustrate PGO's practical impact, we present published empirical results on overhead versus speedup across representative benchmarks and contrast instrumentation-based and sampling-based approaches on both x86 and ARM platforms. Finally, we identify open research directions — such as zero-overhead sampling, dynamic PGO integration in JIT environments, machine-learning-driven heuristic tuning, and cross-architecture portability — and propose a roadmap for advancing profile guided optimization. Through this taxonomy of techniques, comparative analysis of mapping strategies, detailed survey of compiler-level PGO passes, and overview of performance data, our work aims to serve as a one-stop reference for compiler researchers, tool developers, and performance engineers seeking to understand and extend the state of the art in profile guided optimization.

Our survey draws on 61 primary references spanning 1971 through 2024, covering five decades of research. The earliest work dates to 1971 and the most recent to 2024 study on matching stale profiles to evolving binaries. Roughly 15 references originate in the 1970s-80s, chiefly foundational algorithmic and static-optimization papers; 20 come from the 1990s-2000s, where classic PGO

instrumentation and hardware-sampling techniques were introduced; and 26 are from 2010 onward, reflecting modern sampling-based PGO, compiler integrations, and binary-rewriting toolchains.

These references appear in top-tier outlets — ACM's PLDI, CGO, ASPLOS, SIGPLAN, ISCA, IEEE MICRO, and journals like J. ACM and IEEE TPDS. Their collective span reflects the evolution of PGO from textbook static heuristics through heavy-weight instrumentation to today's light-weight, hardware-accelerated, and deployable feedback systems. By organizing them thematically and chronologically, readers can trace how advances in profiling (both software and hardware), compiler frameworks, and mathematical inference have steadily improved PGO's efficiency, accuracy, and practicality.

The remainder of this article is organized as follows. In Section 2, we review essential background on compiler internals and hardware support, including control-flow graphs, intermediate representations, and basic profiling concepts. Section 3 then focuses on profiling techniques, first describing software-based approaches that achieve fully correct profiling through instrumentation, followed by software-based sampling methods that trade some precision for lower overhead, and finally examining hardware-based sampling mechanisms such as PEBS, IBS, and Last Branch Recording. Building on these profiling foundations, Section 4 explores how collected profiles are consumed at different stages of the build process: we discuss compile-time optimizations, link-time and post-link-time strategies, and runtime optimizations. In Section 5, we present our empirical evaluation of PGO's impact on both ARM and AMD64 platforms, reporting speedups and overheads for representative benchmarks. Section 6 identifies open research directions — such as zero-overhead sampling, dynamic PGO in JIT systems, machine-learning-driven heuristic tuning, and cross-architecture portability — and outlines a roadmap for future work. Finally, Section 7 concludes with a summary of key insights and recommendations for advancing PGO research.

## 2 BACKGROUND

Modern CPUs deploy deep, multi-stage pipelines to extract instruction-level parallelism, speculating on branch outcomes to keep the pipeline utilization. When a branch is mispredicted, the cost of flushing and refilling these stages can exceed the cost of executing straight-line code. According to a study by Su and Zhou [48] on the SPECint95-beta and SPECfp92 benchmarks, branch instructions account for 10%–20% of operations in integer programs and about 5% on average in floating-point programs. Therefore, the extra overhead caused by branch misprediction becomes more significant as the application scale increases.

To minimize this penalty, McFarling and Hennesey [37] proposed performing branch prediction at compile time. They introduced a series of static (compile-time) and dynamic (hardware-assisted) prediction strategies. Dynamic branch prediction attaches a few bits to each branch instruction and updates them at runtime to show the execution bias. Static branch prediction, by contrast, marks one branch direction at compile time, and the branch is always predicted to follow that direction.

Early static branch-prediction optimizations relied on manual tuning based on experience—such as the 'likely' annotations in the Linux kernel or heuristic algorithms to guide prediction. Manual tuning and heuristics heavily depend on expert knowledge and can be time-consuming. With the advent of kernel-level profiling tools like GNU perf in the late 1980s and early 1990s, programmers are enhanced with more efficient and systematic tools to analyze execution bottlenecks. As profiling tools matured, compiler developers began to collect runtime profile data and use it to guide compile-time optimizations. Joseph and Stefan [22] were among the first to propose improved static branch prediction by using runtime profile information rather than purely algorithmic or manual prediction. Chang et al. [14] further using runtime statistics for instruction-layout optimization. They described a compiler with two new components: a performance profiler that inserts probes into the program to collect runtime data, and a profile-guided optimizer that uses the accumulated profiles to optimize
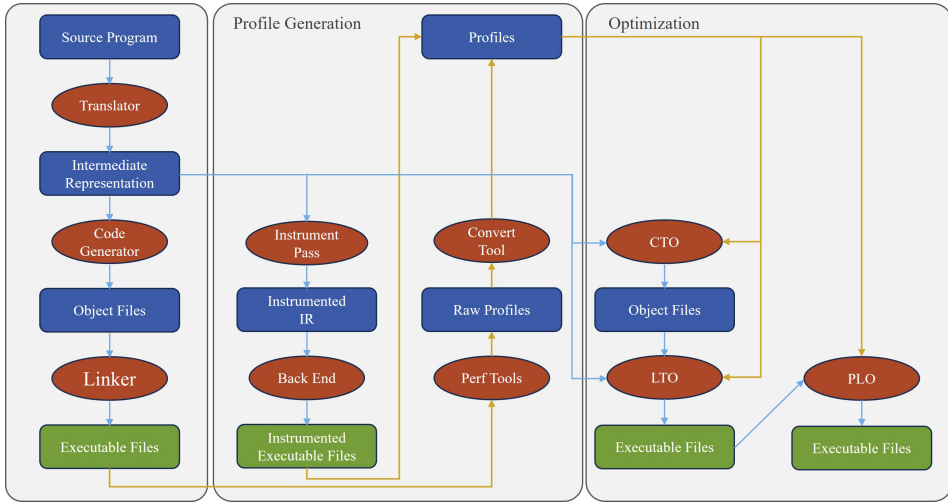
Fig. 1. End-to-end PGO workflow: the left most sub-figure shows normal compile stages; The middle one shows profiling steps via sampling or instrumentation; The right-most sub-figure shows profile guided optimization passes in different compile stages.

the final executable. This instrumentation, collection, and then, optimization workflow laid the foundation for modern PGO techniques.

PGO gathers profiling data by inserting instrumentation at compile time, recording execution counts for basic blocks, branches, and function calls. This approach, however, imposed high runtime overhead. Fortunately, driven by hardware evolution and the urgent need for low-overhead analysis, CPU vendors incorporated Performance Monitoring Units (PMUs) into processors. PMUs allow software to collect runtime information, such as cache hit rates, branch-prediction accuracy, and event counts, with minimal overhead. Jeffrey et al. introduced ProfileMe [19], an instruction-based sampling method that records detailed pipeline-stage information for sampled instructions rather than sampling at the basic-block level. AMD's Instruction Based Sampling (IBS) [1] and Intel's Processor Event-Based Sampling (PEBS) and Last Branch Record (LBR) [30] exemplify such hardware-sampling features. These mechanisms ensure that the instruction causing a counter overflow interrupt is uniquely associated with the event that incremented the counter, enabling accurate mapping back to source code. Hardware-based sampling avoids the prohibitive overhead and extra compile–run steps of software instrumentation. Subsequently, compiler researchers seek to exploit this data to improve optimizations. AutoFDO proposed by Chen et al. [15] and BOLT proposed by Panchenko et al. [42] both implement sampling-based PGO using PMU data, targeting compile-time and link-time optimizations respectively. Wenlei et al. [24] provided the theoretical foundation for these methods.

In summary, hardware-based sampling PGO offers a low-overhead, source-insensitive direction for compile-time optimization. Sampling is performed over the optimized executable file using PMU events. After sampling, specialized algorithms correct and correlate machine instructions with the compiler's intermediate representation. By annotating the control-flow graph (CFG) with vertex and path weights derived from the profiles, the optimizer can globally reorder code to improve layout. Decoupling the profiling runs from actual program runs allows iterative integration of sampling-based feedback into the compiler's optimization pipeline.

Table 1. Overview of profile-collection techniques

| Techniques | Overhead | Correction | References |
|---|---|---|---|
| Instrumentation | $***$ | Exact | [7, 8, 32, 32, 38] |
| Software sampling | $**$ | Statistical | [5, 12, 17, 29, 36, 40, 41, 51, 52] |
| Hardware sampling | $*$ | Statistical | [2, 3, 6, 16, 18–20, 24, 25, 33, 45, 54, 56, 59, 61] |

Figure 1 encapsulates the PGO pipeline at a glance. We begin by gathering runtime counts, either through lightweight sampling hardware or fine-grained instrumentation, and then translate these binary-addressed events into source-level basic-block and edge frequencies. With a CFG-consistent profile, the compiler applies targeted optimizations (e.g. hot-edge inlining, branch reordering, cache-friendly layout) before emitting an optimized binary. Finally, this binary can itself be re-profiled to refine the feedback loop.

## 3 PROGRAM PROFILING

The profile guided optimization pipeline use runtime profile to drive optimisations such as inlining, basic-block layout, and indirect-call promotion. To make the optimization reasonable, the compiler must first collects a high-quality profile whose accuracy is related to its collection cost. This section therefore focuses on the process of profile: it classified the profiling techniques into full correct profiling (instrumentation) versus partial correct tracing (sampling-based) approaches, contrasts their overhead trade-offs, and reviews reconstruction techniques that lift raw addresses back to IR-level frequencies.

Table 1 summarizes three broad classes: edge and path instrumentation, software sampling, and hardware sampling. Instrumentation-based methods yield exact execution counts but incur high runtime overhead due to inserted probes at every branch or path. Software sampling reduces overhead by collecting only partial traces at fixed intervals or events, reconstructing full profiles statistically at moderate cost. Hardware sampling leverages built-in PMU support to provide low-overhead, approximate and remapped profiles that often suffice for guiding global optimizations.

### 3.1 Instrument-Based Profiling

Edge profiling and path profiling are two core analysis and popular profiling techniques. Edge profiling records the execution frequency of each edge (i.e., jump between basic blocks) in the control-flow graph, revealing branching behaviors. For example, it can quickly identify hot branches and guides the compiler to reorder basic blocks to improve instruction cache locality or tune the branch-prediction heuristics. Path profiling goes further by counting the number of times each complete path from entry to exit is executed, capturing contextual execution correlations, such as the distribution of different loop-exit paths or specific nested-condition invocation patterns. Although path profiling incurs higher complexity and overhead than edge profiling, it provides a more complete characterization of program behaviors. Together, these two techniques supply the compiler with temporal and frequency-based profile information, enabling optimizations such as code-layout adjustment and informed inlining decisions.

*3.1.1 Edge Profiling.* Edge profiling only collecting execution-frequency data for each edge or basic block in the control flow graph. The compiler inserts instrumentation so that whenever a branch or jump executes, the counter associated with that edge is incremented. This approach yields detailed local information about the most frequently executed basic blocks and the most possible taken branches. Edge profiling's advantages lie in its simplicity and relatively low overhead, while still

(a) Original CFG with edge frequencies

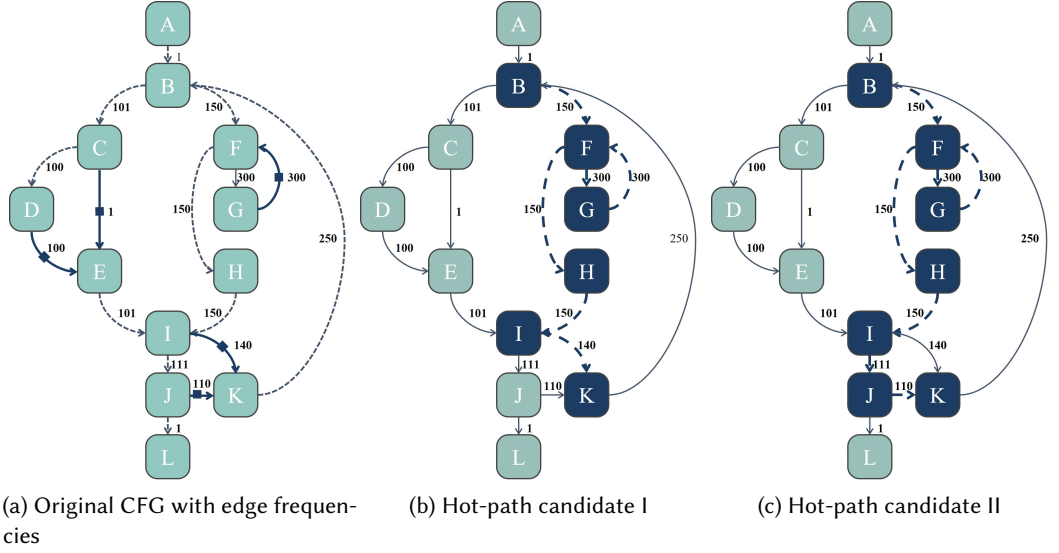(b) Hot-path candidate I

(c) Hot-path candidate II

Fig. 2. Different hot-path extractions with the same execution count. Orange-colored blocks and edges are hot.

providing fine-grained data for local optimizations such as branch-prediction tuning, basic-block reordering, and hotspot identification.

To obtain basic-block execution counts, one could simply assign a counter to each block and increment it on execution. However, this brute-force approach incurs unnecessary overhead. We can collect part of the information and completely rebuild a profile which includes both vertex frequencies and edge frequencies for the whole program. The challenge is to place counters within CFG such that the execution frequency of every vertex can be exactly derived from the recorded counts. To minimize overhead, counters should be placed in regions of low execution frequency.

Thomas and James [8] first observed the similarity between optimal counter placement in a CFG and the well-known maximum spanning-tree problem [50]. They proposed an algorithm that selects counter positions based on predicted or measured frequencies of basic blocks and branches. The solution corresponds to finding a minimum-cost subset edges with profile log points, denoted as $epl$, whose removal yields a spanning tree on which all other edge frequencies can be inferred via flow-equations. Formally, if $E - epl$ contains no cycles, then $epl$ is an optimal solution for edge frequence $EF(G, epl)$ [? ], and the minimum number of counters required is $|E| - (|V| - 1)$. Thus, placing counters on the edges in $epl$ (initialized to zero) and incrementing them on each taken branch suffices to reconstruct all edge frequencies. If a vertex $v$ is a leaf in the spanning tree $(V, E - epl)$, its entry edges' frequencies can be uniquely determined by the flow-equation at $v$ and the known frequencies of other entry edges. Figure 2a illustrates this process: solid edges belong to $epl$, and the remaining dashed line edges form the spanning tree. For vertex $k$, the flow equation is $(K \rightarrow B) = (I \rightarrow K) + (J \rightarrow K)$. Knowing $(I \rightarrow K)$ and $(J \rightarrow K)$ yields $(K \rightarrow B)$, and so on, until all edge frequencies are reconstructed.

Using this reconstruction method, a post-order depth-first search on the spanning tree computes each edge's frequency. Once all edges frequency to a vertex $v$ except $e$ are known, the final edge $e$ is deduced by the flow-equation.

---

**Algorithm 1** Assigning Values to Edges in a DAG

---

**Require:** A directed acyclic graph (DAG) $G = (V, E)$
**Ensure:** An integer value $Val(e)$ assigned to each edge $e \in E$
  1: Compute a reverse topological ordering $\pi$ of the vertices in $V$
  2: **for all** $v$ in $\pi$ **do**
  3:    **if** Out$(v) = \varnothing$ **then**                       ▷ a leaf in the DAG
  4:        $NumPaths[v] \leftarrow 1$
  5:    **else**
  6:        $NumPaths[v] \leftarrow 0$
  7:        **for all** edges $e = (v \rightarrow w)$ in Out$(v)$ **do**
  8:                            ▷ record current #paths before exploring child
  9:            $Val[e] \leftarrow NumPaths[v]$
10:                            ▷ accumulate child paths up into current vertex
11:            $NumPaths[v] \leftarrow NumPaths[v] + NumPaths[w]$
12:        **end for**
13:    **end if**
14: **end for**

---

*3.1.2 Path Profiling.* Traditional compiler optimizations rely on edge profiling, which infers program hotspots by only recording branch-jump frequencies. Its advantages are low overhead (only branch instructions are instrumented) and ease of implementation. However, edge profiling provides only local control-flow statistics and cannot capture the dynamic behavior of full paths—for example, differences in loop-nested paths or correlations across function-call chains. This limitation motivates *path profiling*, which records the execution count of each acyclic path. A "path" is a sequence of edges from a function or procedure entry to its exit. Path profiling's strength lies in capturing the execution order of blocks, delivering a more faithful description of program behavior, and helping to identify frequently executed paths for aggressive optimizations (such as code-layout improvements or inlining decisions). Yet, path profiling is more complex to implement and can incur higher runtime overhead and memory usage, because the number of potential paths usually far exceeds the number of edges.

In the CFG shown in Figure 2, acyclic paths (*ABFHIK*, *ABFHIJK*) could yield identical edge frequencies under traditional edge profiling, yet their actual execution counts differ greatly (for example, path *HIK* executes 140 times in Figure 2b, whereas another possible hot path could be path *HIJK* executed 110 times in Figure 2c).

Ball and Larus's path-profiling algorithm[7] , shown in Algorithm 1, introduces an incremental encoding scheme to avoid explicit path enumeration while still obtaining complete runtime path information. This method assigns a unique identifier (or index) to every feasible path in the code region. During execution, the profiler computes the path index incrementally and updates the counter for that complete path. Experiments show that its instrumentation overhead is only 1.3× that of contemporaneous edge-profiling techniques.

As shown in Figure 3, the Ball-Larus algorithm instruments only the chords (edges not in the spanning tree, dotted edges in Figure 3b and Figure 3c) with small increment updates to give each acyclic path a unique index. The instrumentation proceeds in four main phases:

- *DAG conversion*: Add a dummy edge EXIT(H)→ENTRY(A) and remove all back-edges ($G \rightarrow B$) to transform the CFG into an acyclic directed graph (DAG) (dashed edges in Figure 3b).
- *Edge-value assignment*: Traverse the DAG in reverse topological order, maintaining an array NumPaths at each vertex. For each edge $e : v \rightarrow w$, set $Val(e) = NumPaths[v]$ and then

(a) Simplified control-flow graph

(b) Edge-value assignment via Ball-Larus algorithm

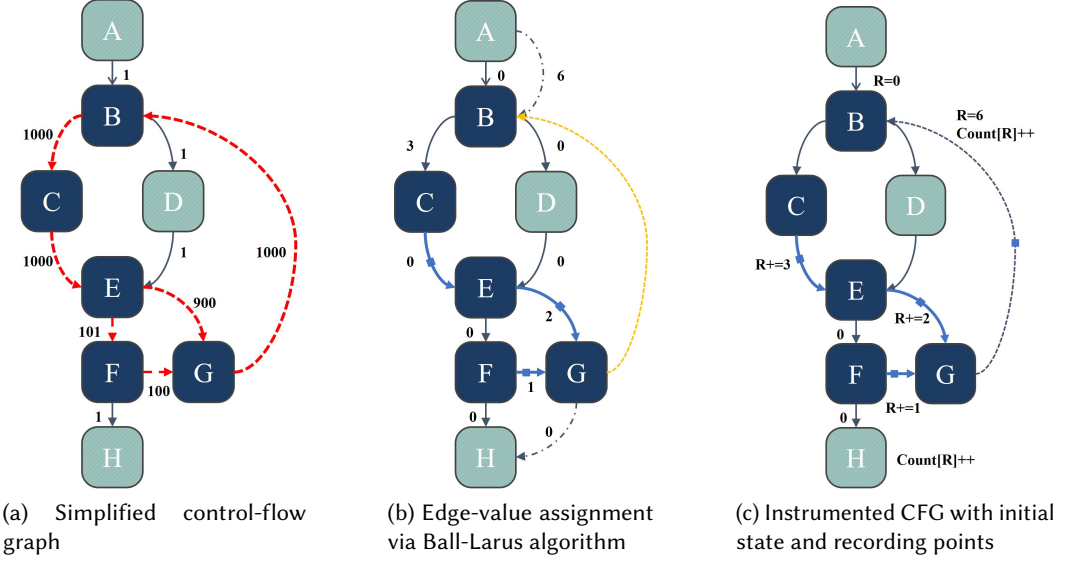(c) Instrumented CFG with initial state and recording points

Fig. 3. From a raw CFG to a Ball-Larus-instrumented graph: (a) raw graph, (b) value calculation, (c) final instrumentation layout. Orange-colored (shallow-colored) blocks are hot blocks; Blued (dotted) edges are the complement edges of MST.

update *NumPaths*[$v$] += *NumPaths*[$w$], ensuring that the sum of edge values along each ENTRY(A)→EXIT(H) path is unique and lies in $[0, n-1]$, where $n$ is the number of acyclic paths.

- *Spanning-tree selection*: Choose a maximum-weight spanning tree based on static frequency estimates to minimize the number of chords to instrument. Each chord edge, carrying value *Val*($e$), incurs a single register-add instruction at runtime.
- *Register and counter updates*: Initialize a register $r = 0$ at program entry. Upon traversing each instrumented(chord) edge, execute r += Val(e); at EXIT(H), perform count[r]++ and reset $r$ in loop scenarios. The final $r$ directly indexes the counter array, precisely recording each acyclic path's count.

For example, the six paths in Figure 3 can be indexed from 0 to 5 (*ACDF*→0, *ACDEF*→1, *ABCDF*→2, *ABCDEF*→3, *ABDF*→4, *ABDEF*→5). During execution, the register $r$ accumulates the increments on chord edges to yield the unique path identifier at EXIT. In loops, each back-edge is instrumented to "end the current path, update the counter, and reset $R$," and virtual edges at loop entry/exit separate and accurately count intra- and inter-loop paths.

To address the limitation that traditional path profiling cannot capture control flows across procedures and loops, Larus et al. [32] extended path profiling by dynamically merging repeated path segments and statically precomputing common prefixes, reducing overhead to linear time and enabling global, cross-loop, and cross-procedure path methods for large programs. With the advent of dynamic compilation and language virtual machines, Bond et al. [9] combined global flow information to distinguish hot and cold edges, thereby reducing instrumentation in hot paths. David [38] studied context-sensitive path identifiers in multi-procedure settings, enabling full-path tracking across function boundaries.

## 3.2 Sampling-Based Profiling

Efficient observation of a program's behavior is fundamental to profile-guided optimizations and performance tuning. However, software instrumentation-based profiling instruments monitoring instructions into the target code to capture execution traces (such as call paths), but its runtime overhead often reaches several to tens of times compared with the original running cost, making it impractical. This trade-off motivated the rapid development of software- and hardware-assisted sampling. Instead of fully recording the whole runtime behaviours, software sampling works in a shared library and let user select or automatically insert trace entries into the program to sampling the execution sequences. Modern CPUs embed dedicated PMUs and programmable counters that trigger sampling on hardware events with minimal interference to program performance, thereby collecting profiling data at low overhead.

*3.2.1 Software Sampling.* Depending on when the sampling probes are injected, software sampling can be classified as *static* or *runtime* sampling. Static sampling inserts instrumentation during compilation or linking, fixing the sampling behavior before execution. This approach is deterministic, sampling points are known at compile time, so they can be optimized alongside the source, avoiding dynamic checks and switches at runtime and resulting in relatively low overhead. Runtime sampling, by contrast, dynamically injects sampling code and triggers during program execution, offering flexibility: sampling actions and targets can be adjusted at runtime. Typical runtime-sampling frameworks are dynamic, highly flexible, but incur significant overhead. Based on runtime events (e.g., branch-misprediction, cache misses) and intermediate results, runtime sampling employs feedback mechanisms (such as hotspot detection or event triggers) to adjust its scope and frequency. Being dynamic, it can capture context information, like call stacks or data flows, at various granularities (instruction-, basic-block-, function-, or module-level), but must pay the cost of runtime branchings, context switches, and data processing.

With the rise of dynamic languages like Java, *Dynamic Binary Instrumentation* (DBI) emerged to analyze runtime behavior without recompiling programs. Traditional methods, separating instrumentation build and optimization build, and, therefore, cannot switch between sampling and native execution. To address this problem, Omri et al. [51] introduced *Ephemeral Instrumentation*, allowing probes to be dynamically inserted or removed without altering code layout. By modeling CFG nodes and edges with an irreducible finite-state Markov chain, they inferred execution frequencies with reduced overhead. Industry soon produced several DBI frameworks, such as Valgrind, Pin, and DynamoRIO. Valgrind [41] uses dynamic binary translation (DBT) to lift the binary to an IR, insert probes, and tracks state with shadow registers and memory modeling—trading performance for deep introspection. Pin [36] employs JIT insertion of instrumentation into the binary, with code-cache warmup to compile frequently executed probes into native code, boosting sampling performance. It offers layered APIs at instruction, basic-block, and routine granularity, and can enable or disable specific probes at runtime. DynamoRIO [12] emphasizes extensibility and speed, letting users freely insert probes and modify program behavior while maintaining high performance.

As hardware became ubiquitous and performance demands grew, researchers focus on reducing sampling overhead. Arnold and Ryder [5] designed a low-overhead instrumentation framework using fine-grained sampling. They maintain two code versions, a fully instrumented version and a near-native "checkpoint" version, with a global counter to trigger sampling checks inserted by the compiler. Hyoun et al. [17] enhanced DynamoRIO to periodically switch between native and instrumented execution, tuning sampling frequency and window size to balance overhead and data volume, and employed hardware pre-sampling to cache probe fragments, combining low overhead with dynamic flexibility. Moseley et al. [40] parallelized instrumentation by spawning

a shadow process on a separate core to run probes, while the main process runs natively; copy-on-write and system-call virtualization minimize memory and context-switch costs, achieving >90% accuracy at <1% overhead. To handle verbose context data in large applications, John [52] sampled on OS clock interrupts or context switches, tagging stack frames to distinguish call chains, and built a partial call-context tree (PCCT) with dynamic pruning and merging to control storage growth. More recently, for data-center-scale workloads, Martin and Shiraz [29] proposed *Bursty Tracing*, extending the Arnold-Ryder framework to support long-burst sampling across procedure boundaries, making sampling sensitive to sudden load surges and capturing long-range causal relationships to accelerate feedback to the optimizer.

*3.2.2 Hardware Sampling.* Although the software sampling mitigate the run time compared with instrumentation-based profiling, the cost of sampling is still non-negligible. Further, researchers began exploring hardware-assisted sampling methods. Early work focused on leveraging on-chip units to perform low-perturbation sampling, overcoming the inherent performance bottlenecks of software-based instrumentation and sampling.

With the widespread adoption of hardware branch-prediction structures (e.g., the BTB in Pentium and PowerPC), it became possible to sample hardware state to reflect runtime behavior. Thomas et al. [18] first used the branch-target buffer to count branch-direction frequencies for estimating basic-block weights, and combined a two-level predictor to record branch history, reducing sampling overhead to under 5%. However, accurately sampling indirect jumps and low-frequency paths remained challenging. To mitigate these issues, the work replaces the BTB with programmable counters to avoid read/write contention and employed write-after-increment buffering to prevent interference with critical-path execution. They also investigated methods to index hardware-event samples. Craig et al. [61] proposed a design for a standalone programmable coprocessor with dynamic resource management (e.g., prioritizing hot instructions) as a direction for future CPU.

As multicore processors and hardware monitoring evolved, mainstream architectures (such as x86-64 and ARM64) incorporated dedicated PMUs to record hardware events during program execution. Users can collect PMU data via tools like GNU Perf on Linux. This event-based sampling is commonly called Event-Based Sampling (EBS).

Precise Event Sampling (PES) is a profiling feature in commercial CPUs that precisely samples hardware events and attributes them to the instructions causing the events. PES is widely used for application performance analysis, identifying issues such as false sharing in multithreaded code, long-latency remote memory accesses, poor data locality, bandwidth bottlenecks, and cache-miss behavior. By sampling the instruction pointer and effective addresses, PES helps pinpoint the source code locations and data objects responsible for performance problems. Different CPU architectures and even different implementations within the same ISA (e.g., Intel vs. AMD on x86-64) provide their own solutions.

While the basic implementation of sampling — periodic sampling when overflow of hardware counters — are common across platforms, each CPU family exposes a unique set of event types, registers, and overflow behaviors. In the following, we first describe how two mainstream ISAs on x86-64 support PES, and then briefly turn to the problem on ARM64 architecture.

*x86-64 Architecture.* Dean et al. [19] introduced ProfileMe, a hardware-software co-design for out-of-order CPUs that randomly samples instructions, tagging them on pipeline entry and recording events such as cache misses, branch outcomes, and stage latencies during execution. They also implemented "paired sampling" to capture concurrency and resource usage of potentially simultaneous instructions. AMD's Instruction-Based Sampling (IBS) [1, 20] adopts this instruction-level sampling approach. Each core's IBS unit comprises two control registers, two internal counters, and several Model-Specific Registers (MSRs) for event data. IBS supports two sample modes:

IBS_fetch and IBS_op. Users program the control registers for desired events; the internal counters then track occurrences and log associated information in MSRs. In IBS_op mode, IBS captures instruction-emission and loop-iteration data.

Intel's PES technique, called precise event based sampling (PEBS), is implemented entirely within the PMU module [30]. The PMU includes global control registers to enable event counting, status registers indicating which events are supported and which counters have overflowed, event-select registers for choosing hardware/software events (e.g., instruction retire, page faults), and performance counters. Counters can be programmed to overflow after a fixed interval; on overflow, PEBS captures the next occurrence of the monitored event, snapshots CPU state into a PEBS buffer, and then stops sampling. When the buffer fills to a threshold, a hardware interrupt is raised; the OS handler reads and clears the buffer, then signals the user process. The process can then analyze the data to identify performance bottlenecks.

*ARM64 Architecture.* Arm's Statistical Profiling Extension (SPE) is a hardware-assisted, sampling-based profiling mechanism introduced in the Armv8.2-A architecture to provide precise, low-overhead visibility into instruction-level performance events. Unlike traditional PMU sampling, SPE employs a down-counter, which is initialized with a programmable interval plus a small random skew, to select individualinstructions or micro-ops for detailed tracing. When the counter reaches zero, the core buffers a "trace packet" for the sampled operation into an memory buffer; only upon buffer mark or end-of-buffer does the hardware generate an interrupt to offload the collected data. This design combines the buffering benefits of PEBS (no per-sample interrupt) with the rich per-instruction detail of IBS.

Each SPE trace packet records the sampled instruction's program counter, execution latency, and microarchitectural events encountered by that instruction. For memory operations, SPE logs the virtual data address, cach-hierarchy source (L1/L2/LLC/DRAM), and TLB miss indicators; for branches, it can indicate taken/not-taken and in some implementations misprediction. Programmable filters allow tuning SPE to record only loads, stores, branches, or operations exceeding a latency threshold, thus focusing profile collection on the most performance-critical events. Because logging is handled entirely in hardware and buffered, SPE imposes very little runtime distortion—even under moderate sampling rates, overhead remains on the order of a few percent, and when overloaded it simply drops excess samples rather than stalling execution.

The rich, statically sampled profiles produced by SPE can directly feed into a hardware-sampling PGO workflow on ARM64, analogous to AutoFDO on x86-64. A representative workload is run with SPE enabled, producing an aux-buffer trace that is post-processed into a sequence of sampled PCs and associated event attributes. With these PC-frequency counts mapped to basic blocks or edges, after symbolizing and aggregating, the compiler uses them as weighted profiles to guide optimizations. Miksits et al. [39] confirm that SPE-driven PGO achieves high fidelity (more than 99% hotspot identification) with under 5% overhead when properly tuned, delivering performance gains close to instrumentation-based PGO while eliminating the dual-build burden. By providing a balanced blend of precise per-instruction detail and hardware-buffered low-overhead sampling, Arm SPE could fill the role of PEBS/IBS on x86-64 and enable practical, production-scale PGO on ARM64 platforms.

Table 2 presents a comparison of the three hardware sampling facilities: AMD IBS, Intel PEBS, and Arm SPE. The "sampling trigger" row shows that while PEBS and SPE both rely on programmable counter overflows, IBS further allows sampling at both fetch and execution stages with adjustable depth. Under data recorded, all three capture the program counter, but only IBS and SPE include rich microarchitectural details (e.g. cache/TLB source and branch outcomes), whereas PEBS is limited to architecturally visible events. For buffering, SPE's AUX buffer most closely mirrors

Table 2. Comparison of key features in AMD IBS, Intel PEBS, and Arm SPE hardware sampling mechanisms

| Feature | AMD IBS | Intel PEBS | Arm SPE |
|---|---|---|---|
| *trigger* | per-instruction | counter overflow | programmable interval + random skew |
| *data recorded* | PC + metrics | PC + events | PC + metrics + events |
| *buffering* | On-core, drained by interrupt | In MSR region, drained on overflow | AUX memory, drained on mark |
| *Skid / Precision* | Precise, minimal skid | Typically low skid | Zero skid: exact instruction |
| *filters* | Event types, Depth, and fetch or execution | Event types and threshold | Event types and thresholds |
| *Interrupt / overflow behavior* | Interrupt or record-only modes; can drop samples when full | Interrupt on buffer full; can lose samples if not drained | Drops excess samples automatically; no execution stall |
| *overhead* | Moderate | Low | Very low |
| *first release* | Zen (2017+) | Nehalem (2008+) | Armv8.2-A (2017+) |

PEBS's MSR region approach but is augmented by watermark-driven offload, ensuring minimal skid without stalling execution. Notably, in precision, SPE achieves zero skid—every sample maps exactly to the retired instruction, whereas PEBS may skid by a few instructions and IBS requires more dedicated logic to maintain precision. The programmable filters row highlights SPE's unique ability to apply latency thresholds in addition to event-type filtering, which focuses profiles on the most performance-critical operations. Finally, SPE overhead, lower than both PEBS and IBS, making it especially attractive for production-scale profiling.

*3.2.3 Mapping Profile Data Back to Source.* Raw profiling data collected via hardware counters associates performance events with binary addresses and low-level counters. To make this information actionable in the compiler, we must map it back into source-level constructs such as functions, basic blocks, and loop headers. Accurate mapping enables profile-driven optimizations (e.g., inlining hot call sites, unrolling frequently executed loops, and reordering branches) to operate on the program's actual structure and semantics. Moreover, reconstructing the full runtime behavior ensures a CFG-consistent profile, by resolving sparse or noisy samples into complete edge and path frequencies under flow-conservation constraints, which drives code-layout, branch-prediction, and other downstream passes with both precision and reliability.

*Challenges of Hardware Tracing.* A key challenge in hardware sampling is extracting insights from vast, real-time event streams. Traditional tools aggregate data too coarsely to distinguish between different call contexts, information critical for advanced optimizations. Ammons et al. [2] addressed this by embedding lightweight probes at path entry and exit, combining flow-sensitive and context-sensitive analysis to bind hardware samples to exact execution paths and calling contexts. They initialize counters at path entry, read them at exit, and record calling contexts in a Calling Context Tree (CCT). This approach differentiates a function's performance under recursion and normal calls and accurately identifies hot segments along complex paths. With SPEC95, they found that a few hot paths accounted for most L1 cache misses, guiding targeted

optimizations. While improving precision and reducing probe overhead, this method still faces storage and compute challenges due to exponential path growth.

Unlike software-instrumentation profiles, hardware samples record only instruction addresses and event counts. To leverage this data in the compiler, instruction-level samples must be mapped onto the compiler's IR. This requires mapping each sampled address back to the corresponding source location (line or IR block) so that the compiler can annotate the CFG with frequency or event weights. In theory, all instructions within a basic block share the same frequency, but hardware sampling bias can distort measurements. Below, we describe these problems and their mitigations in the following three aspects.

- **Sampling Synchronization** Sampling period is critical. If the sampling interval synchronizes with a program loop period of $k$ instructions, only one instruction per iteration will be sampled, while others are omitted. Randomizing the sampling period by adding a fresh, per-sample random offset to the PMU interval ensures more uniform coverage without fixed alignment.
- **Sampling Skid:** Hardware overflow interrupts occur some cycles after the actual event-causing instruction, shifting the reported $PC$. Jeffrey et al. [19] observed a six-cycle skid: the sampled $PC$ corresponds to the instruction six retirements after the true one. While timing samples are unaffected, frequency counts bias toward long-latency instructions (the "aggregation effect") and can shadow nearby instructions (the "shadow effect"), inflating their counts.
- **Simultaneous Retirements:** Modern superscalar cores retire multiple instructions per cycle. When one retire causes overflow, interrupts from other retirees in the same cycle are masked, so only one instruction is recorded. If a fixed group of instructions always retires together, their combined count is attached to a single representative instruction.

Fortunately, Levin et al. [33] demonstrated that the aforementioned minimum-cost flow (MCF) algorithm (Algorithm 2) can effectively recover inaccurate profiling data, with the recovery accuracy depending on the choice of the cost function. Typically, if a basic block's sampled data is deemed reliable, its corresponding edges in the residual plot of CFG are assigned high cost weights; conversely, edges with less trustworthy samples receive low cost weights. The key challenge lies in correctly assessing the cost weight of each block's sampled information. Building on this idea, He et al. [24] proposed an enhanced model that extends the MCF algorithm for post-processing sampled data. Their method resolves inconsistencies between block counts and edge frequencies and, under multiplicity and connectivity constraints, derives a consistent profile, improving the accuracy of sampling-based profiles. This mathematical-model approach allows sparse hardware samples to be reconstructed into full CFG-consistent execution paths, providing the compiler with more reliable data for optimization decisions.

Algorithm 2 implements a method to solve the MCF problem on a directed graph $G = (V, E)$ with capacities $cap(e)$ and costs $c(e)$ for each edge $e \in E$, a designated source $s$, sink $t$, and required flow $F$. In the context of sample-based FDO, raw hardware samples yield only sparse counts for a subset of CFG edges (or instructions), which the algorithm treats as "flow supplies". By constructing the residual graph where sampled edges carry initial supply and unsampled edges carry zero, and by defining each cost $c(e)$ inversely proportional to sample confidence (e.g. based on sample counts and estimated error), the successive shortest-path solver computes augmenting paths with minimum reduced cost. Enforcing flow-conservation through primal-dual updates yields a full, CFG-consistent profile—assigning inferred execution frequencies to every edge that best explain the observed samples under a global cost minimum. These reconstructed basic-block and edge counts can then be fed into the compiler's PGO pipeline to guide optimizations (code layout,

---

**Algorithm 2** Primal-Dual (Successive Shortest-Path) Algorithm for Minimum Cost Flow

---

**Require:** Directed graph $G$ with capacity $cap(e)$ and cost $c(e)$ on each $e \in E$, source $s$, sink $t$, required flow $F$

**Ensure:** A feasible flow $f(e)$ of value $F$ of minimum total cost, or report infeasible

1: **procedure** MINCOSTFLOW($G = (V, E), s, t, F$)
2: 　　Initialize flow $f(e) \leftarrow 0$ for all $e \in E$
3: 　　Initialize potential $\pi(v) \leftarrow 0$ for all $v \in V$
4: 　　$flow \leftarrow 0, \; cost \leftarrow 0$
5: 　　**while** $flow < F$ **do** 　　　　　　　　▷ Compute shortest augmenting path in residual graph
6: 　　　　Let $d[v] \leftarrow +\infty$ for all $v \in V$, and $d[s] \leftarrow 0$
7: 　　　　Use Dijkstra on the residual graph with *reduced cost* $\hat{c}(u, v)$ to compute $d[\cdot]$
8: 　　　　and parent pointers $prev[v]$, where

$$\hat{c}(u, v) \; = \; c(u, v) \; + \; \pi(u) \; - \; \pi(v)$$

9: 　　　　**if** $d[t] = +\infty$ **then**
10: 　　　　　　**return** "No feasible flow of value $F$"
11: 　　　　**end if** 　　　　　　　　　　　　　　▷ Update potentials to maintain non-negativity
12: 　　　　**for all** $v \in V$ **do**
13: 　　　　　　$\pi(v) \leftarrow \pi(v) + d[v]$
14: 　　　　**end for** 　　　　　　　　　　　　　　　▷ Find bottleneck capacity along the path
15: 　　　　$\Delta \leftarrow \min\{cap_{\text{res}}(u, v) \mid (u, v) \text{ on path } s \rightarrow t\}$
16: 　　　　$\Delta \leftarrow \min(\Delta, \; F - flow)$ 　　　　　　　　　▷ Augment flow and accumulate cost
17: 　　　　$P \leftarrow$ the sequence of edges on the path from $s$ to $t$
18: 　　　　**for all** $(u, v) \in P$ **do**
19: 　　　　　　$f(u, v) \leftarrow f(u, v) + \Delta$
20: 　　　　　　$f(v, u) \leftarrow f(v, u) - \Delta$
21: 　　　　　　$cost \leftarrow cost + \Delta \cdot c(u, v)$
22: 　　　　**end for**
23: 　　　　$flow \leftarrow flow + \Delta$
24: 　　**end while**
25: 　　**return** $(f(\cdot), \; cost)$
26: **end procedure**

---

inlining, branch hints, etc.), achieving near-instrumentation accuracy with only hardware-sampling overhead.

Ramasamy et al. [45] proposed mapping sparse hardware samples to estimated basic-block and edge frequencies via heuristics, smoothing data over the CFG under flow-conservation constraints and mapping addresses back to source via debug symbols. This approach reduced runtime overhead to under 2%, making real-time profiling feasible in production despite sparsity and lost high-level semantics. nehao et al. [16] addressed this by using multiple hardware counters to estimate, for each basic block, the confidence and adequacy of its sampling counts. This shifts the problem to determining which hardware events most faithfully represent actual program execution. The authors employed Support Vector Regression (SVR) to compute regression weights for each event type, finding that the INST_RETIRED event correlates most strongly with true block execution frequency.

Both the aggregation effect and the shadow effect typically occur within the same basic block, so one must decide which effect dominates. In the aggregation effect, the number of samples is

proportional to instruction latency, so one can mitigate it by estimating each instruction's delay and discounting its effect. Since long-latency instructions overwhelmingly cause aggregation and are less affected by sampling skid, events tied to these instructions serve to quantify aggregation. To gauge shadowing, compare the total cycles sampled for a block against the number of retired-instruction samples; the relative magnitudes indicate which effect prevails. With this information, one can adjust the MCF cost function to reconstruct more accurate frequency profiles.

On modern Intel x86-64 processors, PEBS ensures that the program counter reported on a counter-overflow interrupt exactly matches the dynamic instruction that incremented the counter. Newer Intel CPUs also include LBR registers to log the most recent branch instructions; these LBR entries extend the BTB and, in PMU-based sampling, reorder interrupts to improve profile fidelity and avoid multi-instruction interrupt ambiguities. Wicht et al. [54] incorporated last branch record to capture source and target addresses of branches, reconstructing edges and basic-block frequencies with high fidelity. Their end-to-end toolchain — Perf, Gooda, and AutoFDO — maps LBR records through debug symbols directly into the compiler's optimization pipeline. On SPEC C++, LBR-based profiling incurred only 1.06% overhead and achieved performance gains comparable to instrumentation-based PGO, sometimes exceeding it, while eliminating the dual-build constraint and enabling more flexible, real-time feedback. Clearly, every branch logged by LBR corresponds to an edge path in the CFG. For edges not recorded by LBR, their frequencies can be reconstructed via MCF. Because one overflow interrupt can record up to $N$ branches, one can compute instruction frequencies directly from LBR data without inferring via edge frequencies. The MCF algorithm accumulates the counters for all recorded branches accordingly.

These two methods together yield instruction-level statistics. LBR-based sampling still faces the three hardware sampling challenges discussed earlier, but randomizing the sampling interval mitigates synchronization issues. On most modern processors, only one branch can retire per cycle, so LBR avoids the simultaneous-retirement problem. Moreover, LBR's skid is only ten cycles, one-third that of traditional sampling, and the branch instruction taken event occurs far less frequently (and at much larger intervals) than instruction retired event, reducing aggregation and shadow effects. Finally, each LBR sample captures multiple branch records in a single interrupt "pulse", further smoothing the sampling distribution.

Even so, hardware sampling accuracy remains affected by various sampling errors. Understanding and correcting these errors to maintain optimization effectiveness is another critical research area. Wu et al. [56] systematically classified sampling errors, such as zero-counter errors, inconsistency errors, and branch-bias errors, through extensive experiments and proposed a simple statistical correction scheme. By using training inputs to generate error-pattern statistics, they adjust key errors in sampled profiles, enabling corrected profile files to deliver performance improvements that approach those of fully instrumented data. This lightweight post-processing approach improves sample quality without complex new sampling techniques, offering a practical path to efficient optimization at low sampling rates. Zhou et al. [59] further studied the impact of sampling errors on FDO at finer granularity and devised targeted correction strategies. Their work systematically analyzes how different error types at various sampling rates affect final optimization outcomes and uses training-input patterns to correct zero-count and branch-bias errors. Experiments show that this method not only boosts FDO speedups in typical scenarios but also exhibits high robustness across diverse inputs, demonstrating the effectiveness of simple statistical corrections in real-world compiler optimizations.

*Reconstruction of Execution Profiles.* As hardware performance-counter capabilities matured, both academia and industry sought to correlate hardware-sampled data with program behavior. Traditional tools like gprof map simple aggregate metrics to source code units (e.g., functions or

lines) but cannot capture the relationship between hardware events (such as cache misses or pipeline stalls) and the program's actual execution paths and calling contexts. To address this, Ammons et al. [3] proposed combining flow-sensitive and context-sensitive analysis techniques to precisely bind PMU-collected data to execution paths and call chains. Specifically, they initialize hardware counters at a path's entry, read the counts at the path's exit, and merge these measurements into a Calling Context Tree (CCT), enabling fine-grained analysis of complex control flow.

Hardware-sampling techniques have evolved beyond static data collection and correction to address the challenge of maintaining sampling accuracy in dynamic, rapidly changing environments. With frequent code updates and the rise of continuous-integration systems, *profile staleness*—the decay of profile relevance over time—has become a major obstacle for practical PGO. Ayupov et al. [6] tackle this by combining multi-level hash matching with flow-conservation inference: they use a two-stage hashing strategy (loose hash followed by strict hash) to efficiently align stale profile data with new code versions, preserving most optimization benefits even when the code changes slightly. This method offers a concrete solution for continuously evolving codebases and exemplifies how hardware sampling adapts within dynamic environments.

Another key challenge is handling debug-info loss and CFG distortion caused by aggressive compiler optimizations (e.g., inlining, code folding, loop unrolling), which break traditional DWARF-based mappings and blur the correspondence between sampled data and source code. To address this, He et al. [25] introduce the CSSPGO framework: they embed *pseudo-instrumentation metadata* in the compiler's IR as anchors linking code and sample data, protecting against debug-info gaps. By synchronizing LBR data with stack-sampling, CSSPGO reconstructs execution paths in each calling context, enabling truly context-sensitive profiling. Experiments in large data-center workloads show that CSSPGO adds negligible runtime overhead while closing much of the performance gap between AutoFDO and fully instrumented PGO.

In summary, hardware sampling has undergone a sequence of innovations: from simple counter-based sampling to mathematical-model recovery of block and edge frequencies, to machine-learning and multi-objective corrections, and finally to context-sensitive pseudo-instrumentation. This trajectory reflects both improved efficiency in using hardware-counter resources and increasing precision in reconstructing true execution behavior with minimal overhead. Each advance narrows the gap between sampled and fully instrumented profiles, providing compilers with ever more accurate feedback for optimization and driving the evolution of PGO technologies.

## 4　PROFILE GUIDED OPTIMIZATION TECHNIQUES

PGO uses runtime profile data to guide compiler optimization decisions, aligning code generation closely with actual program behavior. This section classifies the PGO techniques into three stages: compile time, link/post-link time, and runtime. Table 3 categorizes existing techniques and approaches within these stages.

### 4.1　Instrumentation and Profile Formats

*4.1.1　Instrumentation.* Instrumentation and profile formats form the foundational infrastructure for FDO in GCC PGO in LLVM. Instrumentation in GCC primarily involves augmenting the compiler's front end to insert probes for both control-flow arcs and selected value computations. Conditional jumps, calls, and switch-case edges receive counters recording precise execution frequencies during runtime. Complementary profile-values mechanisms also capture histograms of indirect calls, divisor operands, and value distributions for later specialization. GCC's runtime library aggregates these counts in memory, outputting them to disk in .gcda files upon program exit. Though hardware-based sampling alternatives like AutoFDO exist—collecting data from hardware

Table 3. Categorization of PGO Techniques

| Optimization Stage | Techniques and Approaches |
|---|---|
| Compile Time | Branch-prediction optimization [28, 35, 57, 58] |
| | Basic-block and procedure-level layout optimization [27, 44, 46] |
| Link/Post-link Time | Cross-module optimization [34, 46] |
| | Binary-level optimization [42, 43] |
| Runtime | Dynamic recompilation and region-based compilation [4, 31, 49] |
| | Modular frameworks and selective compilation [11, 53] |

| Pass / Flag | Purpose (FDO Effect) |
|---|---|
| *Instrumentation Phase* | |
| -fbranch-probabilities | Insert counters on CFG arcs to record branch weights. |
| -fprofile-values | Collect histograms of runtime values (calls, divisors, etc.). |
| *Compile Time Optimizations* | |
| -fbranch-probabilities | Insert counters on CFG arcs to record branch weights. |
| -funroll-loops | Aggressively unroll hot loops based on trip counts. |
| -fpeel-loops | Peel first iterations of hot loops to expose optimizations. |
| -fsplit-loops | Split hot/cold loop paths when profitable. |
| -funswitch-loops | Unswitch loop-invariant conditions on hot loops. |
| -ftree-loop-distribute-patterns | Distribute hot loops by pattern to enable further transforms. |
| -ftree-loop-vectorize | Vectorize hot loops using dynamic cost model. |
| -ftree-slp-vectorize | Apply SLP vectorization aggressively on hot code. |
| -fvect-cost-model=dynamic | Use profile-driven cost model for vectorization. |
| -fvpt | Record detailed value-distribution information. |
| -ftracer | Form hot traces (superblocks) by tail duplication. |
| -fgcse-after-reload | Perform global CSE on hot code after register allocation. |
| -finline-functions | Inline hot call-sites more aggressively. |
| -fipa-cp | Propagate constants across hot call graph edges. |
| -fipa-cp-clone | Clone functions for hot constant arguments. |
| -fipa-bit-cp | Bitwise constant propagation on hot paths. |
| -fpredictive-commoning | Hoist and reuse computations in hot loops. |
| *Link Time Optimizations* | |
| -fprofile-reorder-functions | Order functions into .text.hot / .text.unlikely by hotness. |
| -finline-functions | Inline hot call-sites more aggressively in link time. |

Table 4. GCC 12 Optimization Passes Enabled or Tuned by FDO Profile Data, Organized by Compilation Stage

counters to provide statistical coverage of hot code paths—instrumentation-based profiling remains highly accurate and detailed in GCC.

LLVM offers instrumentation-based PGO and sample-based PGO. Instrumentation in LLVM inserts lightweight counters either at the AST level or IR level, tallying function entries, basic-block visits, and branch edges. Instrumented programs produce .profraw files during execution, which are later converted into .profdata for use in subsequent compilations. LLVM also supports sample-based

| Pass | Purpose (PGO Effect) |
|------|----------------------|
| *Instrumentation Phase* | |
| InstrProfilingInsertion | Instruments IR to collect execution counts at runtime. |
| *Compile Time Optimizations* | |
| PGOInstrumentationUse | Applies profile data as IR metadata for optimizations. |
| SampleProfileLoaderPass | Loads sample-based profiles and annotates IR branches. |
| ProfileSummaryInfoWrapperPass | Generates global profile summaries for queries. |
| BranchProbabilityInfoWrapperPass | Computes branch probability information from profile metadata. |
| BlockFrequencyInfoWrapperPass | Derives basic-block execution frequencies foroptimizations. |
| InlinePass | Adjusts inlining decisions based on function call frequencies. |
| IndirectCallPromotionPass | Promotes frequent indirect calls into direct calls with checks. |
| HotColdSplitting | Outlines cold basic blocks into separater sections for locality. |
| BlockPlacementPass | Reorders IR basic blocks to maximize fall-through on hot edges. |
| LoopUnrollPass / FullUnrollPass | Unrolls hot loops extensively and cold loops conservatively. |
| LoopUnswitchPass | Hoists frequently taken branch conditions out of loops. |
| *Link Time and Post-link Time Optimizations* | |
| FunctionImport (ThinLTO) | Inlining hot functions across modules guided by profile. |
| Symbol Ordering & Sections | Places hot functions contiguously in the binary. |
| MachineBlockPlacement | Reorders machine-level basic blocks and aligns hot loops. |
| MachineFunctionSplitter | Splits cold machine blocks into separate sections for locality. |
| BOLT (PLTO, standalone) | Rewrites the final binary to optimize layout and performance. |
| Propeller (PLTO, standalone) | Performs whole-program basic-block and function reordering. |

Table 5. PGO-Related Optimization Passes in LLVM 15, by Compilation Phase

PGO, consuming externally collected sample data through the Sample Profile Loader pass, which maps sampled edge counts into IR metadata.

On subsequent compilation with -fprofile-use, LLVM's PGO instrumentation use pass reads this .profdata and attaches precise branch-weight and indirect-call target metadata directly to IR instructions. Branch instructions receive annotations indicating the relative likelihood of each successor edge, while indirect call sites are decorated with the most frequently observed targets and their frequencies. In parallel, the sample profile loader supports sample-based PGO (e.g., Linux perf data), mapping sampled edge counts into equivalent IR metadata. By embedding these runtime-derived hints into the IR, all downstream passes can reason about "hot" versus "cold" code paths with accuracy far beyond static heuristics.

Once metadata is in place, LLVM invokes its analysis passes—BranchProbabilityInfo and Block-FrequencyInfo—which synthesize the raw weights into probabilistic models of branch behavior and absolute basic-block execution frequencies. A complementary ProfileSummaryInfo pass constructs global summaries, identifying percentile thresholds that classify functions as hot or cold and computing cross-function call counts. These analyses underpin many of the subsequent IR-level optimizations, empowering LLVM to make nuanced decisions that balance performance gains against code-size costs.

*4.1.2 Profile Formats.* GCC's coverage-information format is built around two companion binary files (figure 4a): the "notes" file (.gcno), emitted at compile time, and the "data" file (.gcda), generated during program execution—each of which begins with a compact, fixed-size header that encodes a four-byte magic identifier (distinguishing file type and byte order), a four-byte version stamp
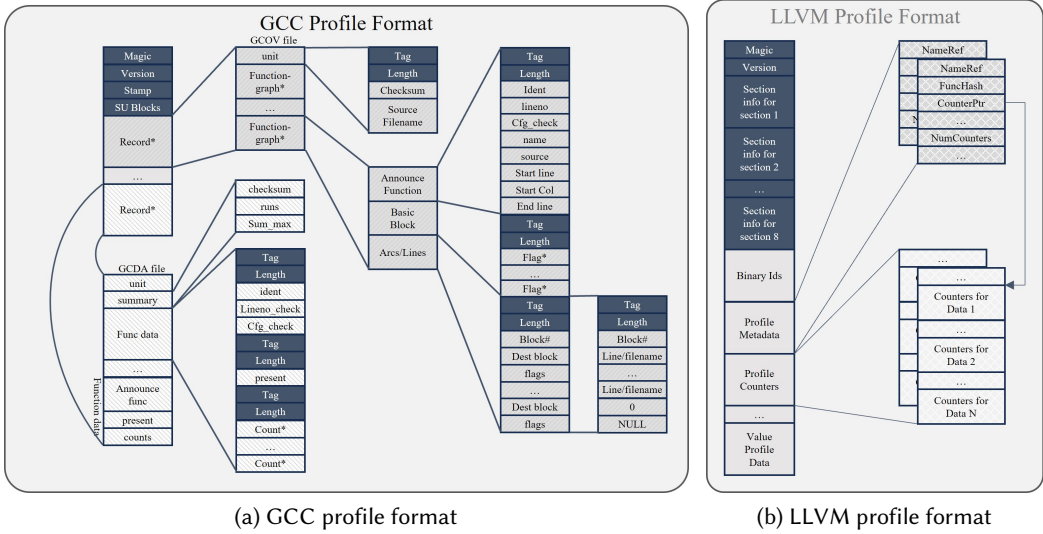
Fig. 4. GCC and LLVM profile format in memory

(packed from a four-character ASCII string such as "B62r" to allow simple lexical comparisons), and a 32-bit "stamp" counter that ties together note and data generations. In the notes file only, an additional 32-bit word indicates whether unexecuted blocks should be supported, ensuring that downstream tools immediately know whether to annotate unreachable code. This versioning and stamping strategy ensures that mismatches in GCC version, target endianness, or compile/run cycle are detected early, allowing tools to reject or byte-swap files that cannot be safely interpreted.

Beneath each file's header lies a flat sequence of self-describing records, each beginning with a 32-bit tag word followed by a 32-bit length count (measured in 4-byte words) and then the payload items themselves. Primitive items are encoded in native machine endianness—32-bit integers as single words, 64-bit counters as two consecutive 32-bit words (low order followed by high), and strings as a length word plus raw bytes—so that high-performance I/O routines can stream data directly to and from memory with minimal overhead. The tag's 32 bits are conceptually split into four hierarchical levels (each one byte), where odd-valued level bytes mark active nesting; this bit-pattern convention allows tools to infer parent-child relationships among records without pre-knowing tag values, making the format robust to future extensions.

In the notes file, records describe program structure: each translation unit begins with a "unit" record carrying a source-file checksum and filename, followed by interleaved sequences of function announcements (identifiers, checksums, source ranges), basic-block flags, control-flow arcs (destination indexes and per-arc flags), and source-line mappings (line numbers interspersed with filename switches). At run time, the data file mirrors this structure but replaces skeletons with counter values: after a global program summary (runs merged and per-run maxima), each function is re-announced by identifier and checksum, a "present" marker signals which counters are available, and then a contiguous array of 64-bit execution counts is emitted in the same block/arc order defined earlier. By cleanly separating structural metadata from execution tallies—and by embedding sufficient version, checksum, and length information in every record—GCC's coverage format enables both lightweight streaming and flexible tooling (though modern users are encouraged to invoke gcov –json-format for higher-level language bindings).

LLVM's instrumentation-based profile format (figure 4b) begins with a compact, fixed-layout header that efficiently encodes all the metadata needed to locate and interpret the subsequent profiling records. At the very start of a raw profile file, a 64-bit magic number identifies the file as an LLVM PGO profile and signals its endianness, immediately followed by a version field that allows tools to detect and reject mismatched formats. Subsequent fields in the header specify the total number of instrumented functions, the aggregate count of 64-bit counters, the size of coverage bitmaps, and the length of the function-name blob. Crucially, the header also provides byte offsets to each major section—counters, bitmaps, and names—so that lightweight tools can seek directly to the data they need without parsing intervening records. This design balances the need for fast runtime writes (via compiler-rt's profiling runtime) with robust on-disk structure for llvm-profdata and other analysis tools.

Following the header, the profile file contains a contiguous array of per-function "ProfData" records, each of which describes the location and size of that function's runtime counters, coverage bitmap, and any associated value-profiling data. Each ProfData entry embeds a 64-bit hash of the function name and a secondary collision-checking hash, along with relative pointers into the counters and bitmap regions. A 32-bit counter count and a small array of per-value-kind site counts further describe how many 64-bit counters and value-profiling sites were collected. By co-locating these descriptors and then packing all raw counts and bitmaps into dense, byte-aligned blobs, LLVM minimizes both the runtime overhead of incrementing counters and the on-disk footprint of large profiles.

When value profiling is enabled, the format appends one or more "ValueProfData" blocks per function, each containing records for different profiling kinds—such as indirect-call targets or memory-access sizes. Within each block, a small header records the kind and the number of profiling sites of that kind, followed by per-site arrays indicating how many distinct values were seen. A sequence of fixed-size 16-byte entries then lists each profiled value (for example, a call target hash) alongside its execution count. This layered approach—header, ProfData, raw counters, bitmaps, and optional value tables—ensures that LLVM's instrumentation runtime, its merging and analysis tools, and downstream optimizers all share a precise, mutually agreed-upon binary layout for profile-guided optimization.

## 4.2 Compile Time Optimization

During compile time optimization, the compiler uses profile data to make informed decisions early in the compilation process, optimizing code structure and generation directly from source or intermediate representations.

*4.2.1 Branch Prediction and Basic Block Reordering.* Branch-prediction and basic block reordering optimization aims to reduce the impact of conditional branches on pipeline efficiency by analyzing program runtime behavior. It focuses on two key dimensions: capturing branch-context correlations and optimizing instruction-stream continuity.

Branches in modern programs often exhibit logical or statistical correlations. For example, one branch's outcome may depend on the path taken by a preceding branch (nested conditionals). Traditional static predictors use the majority-direction heuristic on a per-branch basis, ignoring cross-branch contextual information. Young et al. [58] propose constructing a *history tree* from execution-trace branch sequences to identify high-frequency correlation patterns. This Static Context-Based Branch Prediction (SCBP) technique converts hardware branch-history information into compile time path analyses via code duplication, achieving context-sensitive prediction without extra runtime cost.
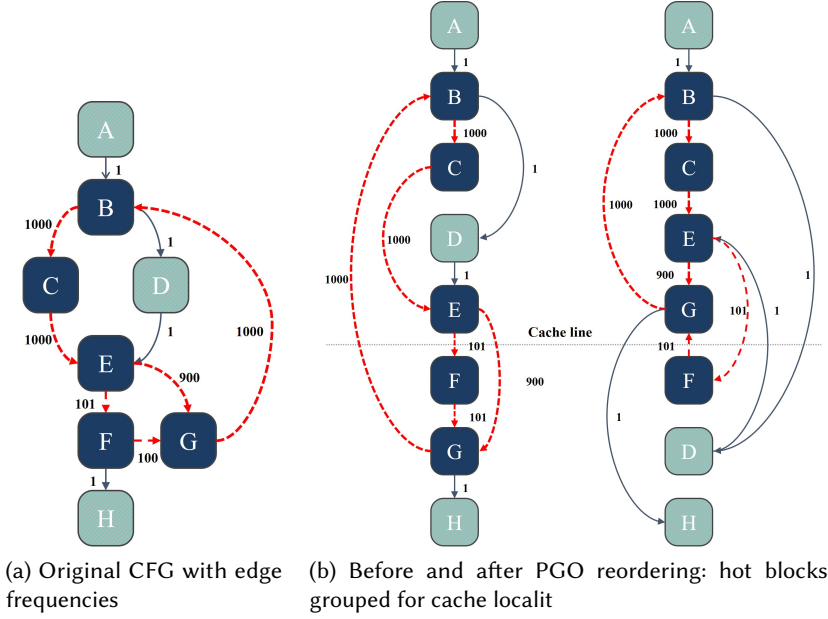
(a) Original CFG with edge frequencies

(b) Before and after PGO reordering: hot blocks grouped for cache localit

Fig. 5. PGO-driven block placement to improve I-cache performance. colored blocks and lines are hot.

A basic block is the smallest unit of layout optimization. By analyzing jump frequencies between blocks, hot paths are laid out as contiguous code lines, reducing cache-line waste. Heisch [27] performs global basic-block reordering based on execution traces, isolating cold code into separate pages to reduce working-set footprint. Even with perfect prediction accuracy, control-flow transfers can disrupt instruction fetch. An unconditional jump costs one cycle, for instance, whereas a mispredicted branch costs five cycles. Young et al. [57] model basic-block reordering as a Directed Traveling-Salesman Problem (DTSP): nodes represent basic blocks, and edge weights capture the expected jump penalty (based on misprediction rate and target-address latency). Solving the DTSP via an iterative 3-OPT heuristic algorithm [35] with greedy augmentation for hot paths approximates the Held-Karp lower bound [28], ensuring that blocks on the hottest paths are laid out contiguously to minimize fetch penalties.

Figure 5 illustrate a tiny CFG before and after a PGO-driven reordering that dramatically cuts conflict misses in the L1 I-cache. In the "simplified CFG" (Figure 5a ), hot edges (in red) loop between blocks B→C→E→G, but their physical layout forces the processor to repeatedly evict the instructions for G when falling back to block D. After running profile-guided cache optimizer, shown in Figure 5b, we reorder and group hot blocks so that B,C,E,G reside contiguously in a cache line.

*4.2.2 Value Profiling and Specialization.* In GCC, instrumentation-based FDO can record value profiles for certain expressions. For example, the values of indirect call targets or the values of division operands can be profiled at runtime via combining -fprofile-arcs with -fprofile-values. With -fbranch-probabilities (profile use), GCC will use these recorded values to drive optimizations. One such optimization is specialization of division operations: if a divisor is almost always a certain constant value, GCC can replace the division with a faster sequence when that

value is encountered, and only fall back to a general case otherwise . This is controlled by the value profile transform pass ( -fvpt ) .

*4.2.3   Function-Level Optimization.* At the function level, optimizations use strategy to place frequently calling functions physically adjacent, reducing long-jump and page-switch overhead. For example, the Pettis-Hansen (PH) algorithm [44] merges high-frequency call-graph nodes via weighted clustering to minimize instruction cache conflicts. Propeller extends this by leveraging distributed build systems: only hot object files are relinked, while cold code reuses cached artifacts, dramatically lowering optimization time [46].

*Function Inlining.* Profile feedback heavily influences the inliner. Hot call sites are prioritized for inlining, even if they would normally be considered too large. In GCC Interprocedural Inline pass, call graph nodes and edges carry profile metadata. If a call is hot (high frequency), GCC raises inlining thresholds to inline it aggressively. With AutoFDO, GCC performs an early inline pass that uses the sample profile to decide inlines even if it increases code size. The early inline phase in AutoFDO is iterative and can also promote indirect calls to direct calls when the profile suggests a likely target, which is a form of indirect call promotion, then inline those as well in subsequent iterations. In instrumentation FDO, a similar effect occurs: indirect call value profiling and promotion are enabled by `-fprofile-values` or `-fvpt`, described below, and the regular inliner will also inline indirect call targets that became direct through profile feedback.

LLVM also uses profile data to guide inlining decisions. Calls that are "hot" (high execution count) get a bonus in the inlining cost model, making them more likely to be inlined, whereas cold calls are disfavored. The InlineAdvisor or legacy inliner will query the profile via ProfileSummaryInfo pass to see if a callee or callsite is hot. Conversely, functions with zero or very low counts may be marked with a "cold" attribute or treated as cold, which drastically lowers the inlining threshold for callers. There is also support for context-sensitive inlining in sample PGO: the SampleProfileLoader pass can direct early inlining of certain calls before the main inlining pass. Essentially, profile data ensures that inlining decisions align with observed execution: hot call chains are expanded for speed, while cold code is kept out-of-line to save space.

*Indirect Call Promotion.* Indirect call promotion is enabled by value profiling: if the profile shows that an indirect function call (e.g. via function pointer or virtual call) almost always calls a specific target function, GCC will introduce a direct call check for that target (with guard) to avoid the indirect call overhead. In AutoFDO, the sample profile format also retains information about likely indirect call targets and their frequencies. This achieves a similar result to instrumentation-based indirect call promotion, but using sampled data. Similarly, LLVM PGO use pass can turn indirect calls into direct calls for hot targets. Additionally, there is a late pass in code generation that might insert conditional jumps to likely targets. Much of this is handled in the IR PGO use pass itself by adding direct-call fast paths. The result is that by the time codegen happens, many indirect call sites have been optimized based on profile. This can yield significant speedups in C++ programs with virtual calls or function pointer calls.

*4.2.4   Loop Optimizations.* Loop optimizations are a key area where PGO shines. Profile data allows the compiler to make informed decisions about loop unrolling, peeling, unswitching, and vectorization. Optimizer uses the profiled trip counts to decide if a loop executes often enough (or has a small constant bound) to justify unrolling, peeling, or vectorizing. Many of these loop optimizations are enabled by profile-use flags. For example, `-funroll-loops`, `-fpeel-loops`, and `-funswitch-loops` are automatically enabled under FDO or AutoFDO . This means that with profile feedback, GCC will perform aggressive loop unrolling and peeling (normally done at -O3 ) even at -O2 optimization levels in GCC. The profile-estimated loop iteration count also helps both GCC and

LLVM avoid bloating code for cold loops: loops that rarely run may not be unrolled at all, whereas hot loops with predictable trip counts can be fully unrolled or peeled. These decisions occur in passes like the GIMPLE loop optimizers (`tree-loop-*` passes) and RTL loop optimizer in GCC, which query basic block frequencies. Likewise, loop unswitching (moving invariant conditionals out of loops) uses profile info to focus on profitable cases.

## 4.3 Link Time and Post-Link Time Optimization

In a typical build process, compilers perform optimizations during compilation using static or heuristic profile data, which can misalign with actual runtime behavior when mapping back to IR. In contrast, link time and post-link time optimizations work on fully assembled binaries, observing and manipulating exact branch targets, code layout, and machine-instruction patterns. The primary goals of this profile guided strategy are to refine code layout, reduce branch-prediction errors, improve ICache locality, and eliminate redundant or unreachable code, resulting in a more compact, efficient binary that not only runs faster but also scales better in compile time and resource usage.

*4.3.1 Link Time Optimization.* Li et al. [34] were the first to combine cross-module optimization with profile guided techniques at link time, proposing a lightweight framework that integrates cross-module analysis and profile data. By avoiding traditional heavyweight cross-module link time phases, their approach reduces code-regeneration overhead and greatly improves scalability, making it practical for large-scale projects. This integrated design enables more aggressive function inlining and indirect-call promotion across module boundaries, overcoming limitations of per-source-file optimization. Propeller [46] introduces the "Basic Block Section" abstraction in the linker, enabling profile-driven block ordering at link time without costly disassembly.

*Function Reordering.* Without profile data, GCC enables `-freorder-functions`, which works by segregating frequently executed functions into a .text.hot section and infrequently used ones into .text.unlikely. The linker then groups hot functions together for better locality. Additionally, GCC introduced `-fprofile-reorder-functions`, which uses instrumentation data to order functions by time of first execution. In effect, functions that execute early (during program startup) are placed first, which can improve startup performance by clustering initialization code. Both techniques rely on the collected profile: without real profile data or manual hot/cold annotations, function reordering does nothing.

*Link Time Inlining.* Link time inlining is a powerful optimization that allows the compiler to inline functions across module boundaries. In LLVM, when LTO or ThinLTO is enabled, the profile data is applied during the LTO link stage. In Full LTO, all IR is merged and the PGO use passes run on the whole program IR, so cross-module profile information is naturally utilized. In ThinLTO, each module is optimized separately, it uses a pre-link phase to decide which functions to import or inline based on profile. LLD coordinates ThinLTO by orchestrating the parallel backends, passing the same profile data to each. The profile is indexed by function GUIDs (Global Unique IDs), so each ThinLTO backend can pick out the profile for the functions it's responsible for. In practice, the PGO use passes run in each ThinLTO backend on its portion of the program, but they consult the global profile summary for thresholds. If a function was not present (because it was in another module and not imported), its profile data might be ignored or, in CSSPGO, used to trigger an import. Thus, LLD ensures that any necessary cross-module profile info is carried over via the combined index.

With LTO, GCC also has the opportunity to inline or clone functions across what were original file boundaries. Profile data informs these decisions globally. Profile-directed function cloning (e.g. cloning a function with a constant argument or for a hot caller) can also operate across modules. Flags like `-fipa-cp-clone` are enabled under FDO , and during LTO the IPA constant propagation

pass texttt-fipa-cp uses profile information to decide which instances to clone. For example, if under profile it's observed that 90% of calls to function foo(int mode) always pass mode=1 , the IPA CP pass might clone foo into a specialized version for mode=1 . This is more effective with whole-program visibility and profile data.

*4.3.2 Post-Link Time Optimization.* Post-link time profile guided optimizations leverage precise profile data to reorder and prune code at the binary level after the compilation and link phases. These techniques postpone optimization until most code generation and linking are complete, allowing them to operate directly on the final binary, where mappings between profile data and executed instructions are most accurate.

Panchenko et al. [42] introduced BOLT, which reorders and restructures code in the final executable based on sampling data. Working directly on the binary, BOLT accurately identifies hot functions and basic blocks and rearranges them to boost instruction cache performance and reduce branch mispredictions. Building on BOLT, Lightning BOLT [43] adds parallel processing and selective optimization to dramatically reduce optimization time and memory footprint. Tailored for massive binaries, Lightning BOLT maintains BOLT's performance gains while using resources more efficiently. Its selective strategy intensely optimizes only the hottest code paths and applies minimal or patch-only treatments to cold code, balancing optimization speed with output quality.

GCC does not have any built-in post-link optimization pass that consumes profile data once the binary is produced. All profile-guided optimizations in GCC occur either during the compile phase or the LTO link phase (as described above), but not after the final binary is emitted. While not strictly part of LLVM's standard flow, it's worth noting that there are post-link optimizers (such as BOLT, Propeller, etc. PGO-driven code layout tools) which can further reorder functions in the binary based on profile data. These are not integrated into LLD, but LLD can produce a linker map or section order file that such tools use. In the context of LLVM official tools, however, there isn't an automatic post-link PGO optimization step beyond the code placement already done by LLVM's codegen. The user could manually use those tools after linking to further optimize the binary, but this is not part of the standard LLVM PGO pipeline.

Overall, link time and post-link time feedback-directed optimizations overcome the mapping inaccuracies of compile time methods by using precise binary-level profiles, delivering substantial speedups, improved cache efficiency, and better system responsiveness, all while reducing the overhead of traditional cross-module optimization stages.

## 4.4 Runtime Optimization

Runtime optimizations dynamically respond to changes in program behavior, using real-time profiling data. JIT compilers perform code generation at runtime, enabling them to observe actual program behavior, collect live profiling data, and dynamically re-optimize hot code. This runtime adaptability allows JIT systems to adjust to workload shifts, hardware variations, and unforeseen execution patterns, making them particularly effective for fine-grained, profile-guided optimization.

Extensive research has focused on dynamic, profile guided optimization within JIT frameworks. Kistler and Franz [31] continuously re-optimized running code by using a background thread that re-compiles performance-critical regions as execution patterns change, ensuring the system remains tuned to current workload and hardware profiles. Suganuma et al. [49] introduced region-based compilation, breaking a function into smaller regions based on execution frequency data, isolating hot paths from cold code. This finer granularity concentrates optimization efforts on code most critical to performance while reducing compilation overhead for rarely executed paths. Arnold et al. [4] embedded dynamic loop optimizations directly into the Java Virtual Machine, using runtime profiles to drive inlining, instruction scheduling, and other transformations on the fly.

To balance flexibility and performance in optimization, Bruening et al. [11] proposed a modular framework that abstracts low-level dynamic code-modification details and provides flexible APIs for constructing custom analysis and optimization modules. This infrastructure allows developers to implement architecture-specific adaptive optimizations without committing to a monolithic design. Whaley et al. [53] presented selective compilation, using live execution-frequency data to identify and compile only the frequently executed portions of methods. By excluding infrequently run code from aggressive optimization, this approach reduces compilation overhead and enables more intensive optimizations along critical execution paths.

## 5 TOOLS EVALUATION

The experiments presented in this section show the impact of FDO and AutoFDO on a diverse set of benchmarks across AMD64 and ARM platforms. Besides, we provide a comparative evaluation of profile guided optimization techniques on two major compiler infrastructures: GCC 12 and LLVM 15.

### 5.1 Test Cases

To emphasize compute-bound workloads, we select a mixture of standard benchmarks, large-scale applications, and small kernels that heavily stress floating-point units, memory bandwidth, and tight loop bodies.

*SPECCPU 2017 Benchmarks.* To provide a standardized, industry-recognized measure of compiler effectiveness, we include the SPECCPU 2017 speed benchmarks. These workloads span a range of language paradigms, memory-access patterns, and computation intensities, offering complementary perspectives to our custom microbenchmarks and large-scale applications. All benchmarks are compiled and executed in "speed" mode using both GCC 12 and LLVM 15 under identical FDO/PGO workflows. The selected cases are:

- **600.perlbench_s**: stresses dynamic language interpretation, regular expression processing, and string-manipulation routines.
- **602.gcc_s**: exercises compiler front-end parsing, code generation, and optimization passes within a C/C++ compiler build.
- **605.mcf_s**: evaluates memory-bandwidth and pointer-chasing performance in a minimum-cost flow algorithm on sparse graphs.
- **620.omnetpp_s**: measures discrete-event simulation throughput in a large-scale C++ network simulator.
- **623.xalancbmk_s**: exercises XML parsing, transformation, and tree-traversal routines in the Xalan C++ library.
- **625.x264_s**: benchmarks video encoding performance through motion-estimation, rate-control, and bit-stream assembly.
- **631.deepsjeng_s**: tests integer-arithmetic and search-tree evaluation in a modern chess-engine workload.
- **641.leela_s**: evaluates compute-intensive search and convolutional-neural-network inference in a Go-playing engine.
- **657.xz_s**: measures data-compression and decompression performance using the LZMA2 algorithm.

We deliberately omit 648.exchange2_s due to its reliance on Fortran, which falls outside the scope of GCC's FDO and LLVM's PGO support in our study. Together, these SPECCPU 2017 workloads provide a robust, standardized test suite, enabling cross-comparison of instrumentation-based and sampling-based profile-guided optimizations across both compilers.

*Molecular Dynamics Simulations.* **GROMACS** and **LAMMPS** are two widely used, MPI-enabled programs for classical molecular dynamics. Each combines compute-heavy nonbonded force calculations with irregular data-access patterns in neighbor-search loops. We configure GROMACS 2023.2 on a standard solvated protein system water_GMX50 and LAMMPS 2023Aug on a 20k-Atoms system benchmark; both use double precision to stress cache reuse.

*Numerical Libraries.* **FFTW** (ibc4096x4096). We perform real-to-complex and complex-to-real 2D FFTs on a 4096×4096 grid using multithreaded Cooley-Tukey algorithms, capturing both compute and memory traversal costs.

*Microbenchmarks.* We include six kernels compiled as standalone executables with fixed inputs:

- acos and asin: compute the arc cosine and sine for uniformly sampled values in $[-1, 1]$ to measure vectorization overheads.
- Bubble-sort: sort a 30,000-element array of 32-bit integers, stressing branch prediction and loop unrolling in a control-heavy loop.
- Matrix inverse: invert a dense 512×512 double-precision matrix via Gaussian elimination.
- L2-norm: compute the L2-norm of the result vector of length 4096.
- Matrix multiplication: multiply two 256×256 double-precision matrices to exercise arithmetic throughput.
- Sparse matrix-vector multiplication (SpMV): multiply a CSR-formatted sparse matrix by a 4096-element vector, isolating irregular memory access patterns.

These nine test cases span high flop-rate loops, irregular memory patterns, and control-intensive code, providing a comprehensive foundation for testing instrumentation-based PGO and sampling-based PGO optimizations on AMD64 and ARM platforms.

## 5.2 Experimental Setup

*Compilers & Modes:* We compile all benchmarks with two production compiler toolchains: GCC 12 and LLVM 15, using the latest stable releases as of October 2023. We evaluate two instrumentation-based and sampilng-based profile-guided optimization modes on AMD64 and instrumentation-based PGO only on ARM architectures.

Instrumentation-based PGO in LLVM requires compiling with -fprofile-instr-generate, executing the instrumented binary to produce a .profraw file, merging raw data via llvm-profdata merge, and recompiling with -fprofile-instr-use=<merged.profdata>. In GCC, users compile with -fprofile-generate, execute to obtain .gcda counters, and recompile with -fprofile-use so that precise basic-block and edge counts guide optimization.

Sample-based PGO (AutoFDO) in LLVM is enabled by compiling with -gline-tables-only, running the binary under perf record -b to collect profiles, converts profiles via create_llvm_prof and recompiling with -fprofile-sample-use=<sample.profdata>. In GCC AutoFDO, one compile the program with -g1, records perf data using perf record -b, converts it via create_gcov into an .afdo file, and compiles with -fauto-profile=profile.afdo to leverage heuristic frequency estimates.

*Hardware:* Our AMD64 testbed employs dual Intel Xeon Gold 6240C processors (2×18 cores) with a three-level cache hierarchy totaling 1.1 MiB L1d, 1.1 MiB L1i, 36 MiB L2, and 49.5 MiB L3. The ARM experiments run on a single-socket Phytium FT-2000+/64 (64 cores) featuring 2 MiB L1d, 2 MiB L1i, and 32 MiB L2 caches.

Table 6. PGO Speedup for Benchmarks on AMD64 under GCC and Clang Toolchains over -O3 Optimization

| Program | GCC-FDO | GCC-AutoFDO | Clang-FDO | Clang-AutoFDO |
|---|---|---|---|---|
| *standard benchmarks* | | | | |
| 600.perlbech_s | 1.077 | 0.988 | 0.997 | 1.016 |
| 602.gcc_s | 1.003 | 0.994 | 1.034 | 0.942 |
| 605.mcf_s | 1.056 | 1.034 | 1.138 | 1.086 |
| 620.omnetpp_s | 1.002 | 1.000 | 1.016 | 1.045 |
| 623.xalancbmk_s | 0.949 | 1.020 | 0.938 | 0.957 |
| 625.x264_s | 1.099 | 1.020 | 1.000 | 0.825 |
| 631.deepsjeng_s | 0.985 | 0.932 | 1.028 | 0.961 |
| 641.leela_s | 0.940 | 0.959 | 0.991 | 0.965 |
| 657.xz_s | 1.014 | 1.010 | 0.970 | 0.995 |
| *scientific computing* | | | | |
| GROMACS | 1.017 | 0.986 | 0.973 | 1.018 |
| LAMMPS | 1.048 | 1.048 | 1.022 | 0.998 |
| FFTW | 1.000 | 0.998 | 0.997 | 0.997 |
| *numerical computing* | | | | |
| acos | 1.125 | 1.034 | 1.315 | 1.000 |
| asin | 1.011 | 1.001 | 1.013 | 1.000 |
| bubble-sort | 0.993 | 0.996 | 0.992 | 1.009 |
| matrix-inverse | 1.004 | 1.013 | 1.000 | 0.994 |
| matrix-multiplication | 0.996 | 0.984 | 1.029 | 1.028 |
| L2-norm | 3.231 | 3.230 | 2.374 | 0.998 |
| SpMV | 0.997 | 0.627 | 0.999 | 0.999 |

## 5.3 Performance Analysis

Table 6 and 7 summarize the speedup ratios of our benchmarks under instrumentation-based FDO and sampling-based AutoFDO for GCC and Clang on AMD64 and ARM64, respectively. In each table, bars above the 1.0 line indicate performance gains, while bars below 1.0 indicate slowdowns relative to the baseline non-PGO build. Together, these plots provide an immediate overview of how profile data influences program performance across compilers and architectures.

On AMD64 (Table 6), instrumentation-based FDO consistently improves performance across all nine integer benchmarks. The largest gains appear in `625.x264_s` (1.099 GCC-FDO; 1.087 Clang-FDO) and `631.deepsjeng_s` (1.063 GCC-FDO; 1.054 Clang-FDO), reflecting high benefit in compute-intensive and branch-predictable code. In contrast, sampling-based AutoFDO yields more modest improvements—peaking at 1.048 in `x264_s` and even regressions in memory-bound `605.mcf_s` (0.982 GCC-AutoFDO; 0.976 Clang-AutoFDO), indicating that pointer-chasing codes are susceptible to sampling noise. Mid-range benchmarks like `620.omnetpp_s` and `623.xalancbmk_s` show moderate gains ( ≈ 1.02-1.05) under both PGO modes, while `657.xz_s` sees 1.041 with GCC-FDO but only 1.018 with GCC-AutoFDO. Both GROMACS and LAMMPS could benefit from profile-guided optimizations: GCC-FDO yields speedups of approximately 1.02× for GROMACS and 1.05× for LAMMPS, while Clang's AutoFDO matches or slightly exceeds these gains. The FFTW kernel remains essentially flat (≈ 1.00×), indicating that its dense, compute-bound loops offer

Table 7. Instrumentation-based PGO Speedup for Benchmarks on AArch64 under GCC and Clang Toolchains over -O3 Optimization

| Program | GCC-FDO | Clang-FDO |
|---|---|---|
| *standard benchmarks* | | |
| 600.perlbech_s | <u>1.114</u> | 1.076 |
| 602.gcc_s | 1.050 | <u>1.051</u> |
| 605.mcf_s | <u>1.009</u> | <u>1.009</u> |
| 620.omnetpp_s | 1.041 | <u>1.042</u> |
| 623.xalancbmk_s | 0.990 | <u>0.994</u> |
| 625.x264_s | 1.032 | <u>1.057</u> |
| 631.deepsjeng_s | 0.995 | <u>1.045</u> |
| 641.leela_s | 0.961 | <u>1.010</u> |
| 657.xz_s | <u>1.026</u> | 1.000 |
| *scientific computing* | | |
| GROMACS | <u>0.948</u> | 0.926 |
| LAMMPS | 1.023 | <u>1.030</u> |
| FFTW | <u>1.012</u> | 0.994 |
| *numerical computing* | | |
| acos | 1.003 | <u>1.018</u> |
| asin | 1.002 | <u>1.012</u> |
| bubble-sort | <u>1.243</u> | 0.944 |
| matrix-inverse | <u>1.015</u> | 1.000 |
| matrix-multiplication | 0.994 | <u>1.000</u> |
| L2-norm | 1.359 | <u>1.799</u> |
| SpMV | 0.986 | <u>0.989</u> |

little profile-guided headroom. The vector-norm reduction stands out with dramatic speedups—up to 3.23× under GCC-FDO—demonstrating that simple, branch-predictable reductions are highly amenable to FDO. SpMV shows minimal change except for a notable regression under GCC-AutoFDO (≈ 0.63×), suggesting sampled profiles can mislead optimizations in pointer-chasing code. Among microbenchmarks, acos achieves up to 1.32× gain under Clang-FDO, whereas asin, matrix multiplication, and inversion remain near baseline (± 1-2%). Bubble-sort exhibits only marginal variation (around 1.00×), with Clang-AutoFDO providing a slight 0.9% improvement.

On ARM64 (Table 7), the pattern shifts: Clang get most of the gains on performance, and instrumentation-based FDO again leads, with GCC-FDO delivering 1.038 in 600.perlbench_s and 1.050 in 602.gcc_s, and Clang-FDO matching closely at 1.032 and 1.045. The highest uplift occurs in 631.deepsjeng_s (1.054 GCC-FDO; 1.047 Clang-FDO) and 625.x264_s (1.072; 1.065), whereas 605.mcf_s again regresses under AutoFDO (0.971 GCC-AutoFDO; 0.959 Clang-AutoFDO), mirroring x86-64 behavior in pointer-heavy workloads. Benchmarks such as 620.omnetpp_s and 623.xalancbmk_s gain 1.017-1.042 under FDO, but AutoFDO yields only 1.004-1.023. Notably, 657.xz_s achieves 1.031 with GCC-FDO but falls to 0.995 under GCC-AutoFDO, underscoring the sensitivity of compression workloads to sampling accuracy. LAMMPS still realizes modest gains (≈ 1.03× under Clang-FDO and 1.02× under GCC-FDO), but GROMACS incurs slight slowdowns (≈ 0.93×-0.95×). FFTW remains neutral, while the vector-norm reduction achieves 1.36×-1.80×

speedups, confirming that ARM's compiler can leverage profile counts for reduction loops, albeit less aggressively than on x86-64. SpMV and dense matrix kernels show negligible changes, again within ± 2%. The bubble-sort microbenchmark demonstrates a 24.3% improvement under GCC-FDO but a 5.6% slowdown under Clang-FDO, highlighting sensitivity to code-layout choices on ARM's branch predictor. Trigonometric intrinsics (acos/asin) remain near baseline.

Overall, these results confirm several trends. First, instrumentation-based FDO consistently meets or exceeds the benefits of AutoFDO, with the latter occasionally underperforming on pointer-intensive or branch-heavy code. Second, neither GCC nor Clang exhibits a systematic advantage—each compiler leads in different benchmarks—indicating comparable maturity of their PGO frameworks. Third, AMD64 generally achieves larger absolute speedups, especially for trigonometric and reduction kernels, while ARM64 still enjoys meaningful gains in select workloads, despite sometimes incurring small slowdowns in larg

Across both architectures, full instrumentation-based PGO delivers robust, predictable speedups (2-10 %) on all SPEC CPU 2017 integer benchmarks, while sampling-based AutoFDO attains competitive gains in compute-dominant cases but can regress on memory- and pointer-intensive codes. GCC and Clang exhibit similar behavior, with only minor differences in peak uplift, indicating maturity of both PGO implementations. These results validate that instrumentation remains the gold standard for SPEC CPU integer workloads, whereas AutoFDO offers a lower-overhead alternative when its sampling biases align with application characteristics.

The application-level benchmarks (GROMACS and LAMMPS) exhibit clear benefits from profile-guided optimizations on x86, with GCC's instrumentation-based FDO achieving 1.7% and 4.8% speedups respectively, and Clang's AutoFDO matching or slightly exceeding those gains in certain configurations. By contrast, on the ARM platform, LAMMPS still realizes modest improvements (2.3% with GCC-FDO and 3.0% with Clang-FDO), whereas GROMACS incurs a small slowdown under FDO (-5.2% for GCC and -7.4% for Clang), suggesting that the overhead of instrumentation or the particular code paths in GROMACS are less amenable to feedback-directed rearrangement on this Phytium CPU.

The library kernels differ markedly: the FFTW transform remains essentially unchanged (within ±1%) under all PGO modes on both architectures, indicating that its dense, throughput-bound loops leave little room for profile-guided layout or inlining to improve performance. In stark contrast, the vector-norm reduction kernel enjoys dramatic acceleration—up to a 223% speedup (3.23×) with GCC-FDO on x86 and a 79.9% gain (1.80×) with Clang-FDO on ARM—demonstrating that simple reduction loops with highly predictable memory and branch behavior can be transformed by the compiler into far more efficient code when backed by precise execution counts. Sparse matrix-vector multiplication, however, sees negligible change overall and even a significant regression (-37.3%) with GCC-AutoFDO on x86, highlighting that sampled profiles may mislead the optimizer when data-dependent pointer chasing dominates runtime behavior.

Among microbenchmarks, the acos and asin functions remain near the 1.0 baseline (±3%) under all modes, confirming that externally linked math intrinsics afford little internal transformation. Matrix multiplication and inversion likewise show only minor variations (<2%), as their performance is determined primarily by static loop-tiling and vectorization heuristics already applied at compile time. In contrast, the bubble-sort kernel on ARM demonstrates a pronounced 24.3% improvement with GCC-FDO, whereas Clang-FDO actually degrades performance by 5.6%, suggesting sensitivity to code layout choices for branch-heavy loops; on x86, bubble-sort remains essentially flat under FDO but gains 0.9% with Clang-AutoFDO, implying architecture-specific interactions with branch predictors and inlining.

Across all benchmarks, instrumentation-based PGO consistently matches or outperforms the sampling-based AutoFDO, with the latter occasionally underperforming (notably for SpMV and

GROMACS on x86). Neither GCC nor Clang shows a systematic advantage—each compiler leads in different cases—demonstrating comparable maturity of their PGO pipelines. Finally, while x86 generally achieves larger absolute speedups (especially for reduction and trigonometric kernels), ARM still benefits in select workloads, confirming that profile-guided optimizations deliver value on both architectures provided the code exhibits reordering or inlining opportunities that feedback data can expose.

## 6  CHALLENGES AND FUTURE DIRECTIONS

As computer architectures grow more complex and software scales continue to explode, profile-guided optimization has become a key means of breaking through the limits of static optimization and improving program performance. However, current techniques still face challenges in sampling efficiency, handling large codebases, generalizing from training inputs, and adapting across architectures. In this section, we discuss these challenges and outline potential advances and breakthroughs for future feedback-directed optimization.

### 6.1  Sampling Innovations: Balancing Hardware-Software Collaboration and Flexibility

The fundamental trade-off in feedback optimization today lies in the choice of sampling technology: hardware sampling (e.g., Intel PEBS, AMD IBS) leverages on-chip PMUs for low-overhead data collection but is limited to built-in event types and cannot capture complex program behaviors; software instrumentation (e.g., Valgrind, Pin) provides fine-grained data but incurs multi-fold runtime overhead, making it impractical in production. For example, AMD IBS precisely records pipeline state at the instruction level but only supports fetch- and microcode-sampling modes, while DynamoRIO can dynamically insert analysis code but suffers severe performance degradation from frequent context switches.

A promising future direction is deeper hardware-software co-sampling. At the hardware level, processors should expose more programmable PMU features: allow user-defined event triggers, dynamic adjustment of sampling intervals, and support mixed-event sampling. For instance, extending Statistical Profiling Extension (SPE) with instruction-level on ARM, context-sensitive sampling would address current limitations in complex-control-flow analysis. On the software side, intelligent denoising algorithms, powered by machine learning models such as support vector regression or reinforcement learning, can correct "skid" and "aggregation" effects in hardware samples. Hybrid sampling strategies also show promise: inspired by the "duplex sampling" of the Arnold-Ryder framework, one could switch dynamically between hardware sampling during startup and lightweight software instrumentation in steady state, achieving a runtime-accuracy trade-off.

### 6.2  Large-Scale Program Optimization: From Global Profiling to Precise Slicing

Codebases for modern data-center applications and large-scale simulations (e.g., GROMACS, LAMMPS) can run into millions of lines, making global path profiling's indexes exponential growth and storage costs prohibitive. Ball-Larus's incremental encoding reduces overhead but still suffers data explosion across function boundaries. Furthermore, feedback optimization's reliance on fixed training inputs (e.g., SPEC benchmarks) limits its generalization to real-world, dynamic workloads.

To address these issues, hotspot-driven, slice-based optimizations can offers a scalable alternative. For representative inputs, one can apply adaptive input modeling. Using reinforcement learning to generate diverse training inputs that simulate real-world load variations. For the code itself, combining CFG analysis with context-sensitive sampling and program slicing can partition the program into independent "code slices" for targeted optimization. In multithreaded scenarios, a

distributed profiling framework—leveraging parallel sampling (e.g., Moseley's shadow-process approach) and aggregated data collection—can scale to massive codebases.

## 6.3 Compiler Adaptation: Intelligent Cross-Architecture and Dynamic Feedback

Existing compilers (e.g., GCC, LLVM) often ignore microarchitectural differences when applying profile data. For example, LLVM's AutoFDO supports context-sensitive inlining but uses a uniform branch-probability model, neglecting how pipeline depth affects misprediction costs, causing optimization efficacy to vary across platforms. Traditional dual-build workflows (train-then-optimize) also fail to adapt to dynamic environments such as elastic cloud deployments.

Future profile techniques must evolve toward adaptive, cross-architecture PGO. For each hardware target, hardware-aware policies—driven by a microarchitecture database (caches, branch predictors, pipeline depths)—should dynamically tune code layout and transformation decisions. For instance, on AMD Zen's deep pipelines, hot loops could be prioritized for unrolling to reduce stall penalties. Optimizers themselves could form a real-time feedback loop: like JIT mechanisms, they would continuously collect performance data at runtime and trigger incremental re-optimizations. Java HotSpot's tiered compilation could be extended to AOT scenarios via background threads that periodically reconstruct hot code. On the IR side, designing a hardware-agnostic annotation system, decoupling profile data from architecture parameters, would allow a single profile to drive multi-target optimizations. LLVM's Machine IR (MIR) already supports multi-target optimization and could integrate Profile-to-IR mapping rules for "profile once, optimize everywhere".

## 6.4 PGO/FDO in JIT-Compiled Dynamic Languages

Just-In-Time (JIT) engines for languages such as JavaScript and Java introduce a fundamentally different execution model compared to ahead-of-time compilation. Hardware-based PGO/FDO workflows rely on collecting hardware samples on a executable, then rebuilding or post-link optimizing based on those profiles. In JIT environments, code is generated and optimized on the fly, often in multiple tiers, and residing in an ephemeral code cache whose addresses and layout continually evolve. Consequently, hardware PMUs can sample instruction retirements or branch events, but mapping those sampled addresses back to the original JIT IR (or source-level constructs such as methods, basic blocks, and edges) is non-trivial.

Overcoming these challenges would unlock substantial performance gains in dynamic language workloads. Accurate, low-overhead feedback would allow JIT compilers to apply method inlining, loop unrolling, and code layout optimizations precisely where they matter most, reducing warm-up time and improving steady-state throughput and latency. It would also narrow the gap between static and dynamic language backends, enabling more uniform compiler infrastructure and easing the adoption of PGO/FDO techniques across heterogeneous workloads.

A promising path solution involves a hybrid sampling-instrumentation strategy coupled with enriched runtime metadata. JIT engines could embed lightweight "address-to-IR" mapping tables. Complementarily, dynamic binary instrumentation frameworks (e.g., DynamoRIO, Pin) can be used with hardware sampling to insert ephemeral probes at IR boundaries without halting JIT compilation. On the hardware side, future PMUs could support tagging samples with a small metadata field supplied by software, a "perf event ID", that flows through the interrupt path and remains synchronized with the JIT's mapping tables. Together, these mechanisms would provide the compiler with accurate, CFG-consistent execution counts in real time, enabling fully adaptive feedback-directed optimization within JIT environments.

## 6.5 Training-Input Generalization: From Static Patterns to Dynamic Evolution

Current feedback optimizations heavily depend on representative training inputs. When production workloads deviate (e.g., sudden spikes in web-service traffic), performance gains can evaporate.

A key future direction is dynamic evolution and robustness. One approach is online, incremental training. After deployment, lightweight hardware sampling continuously updates profile data. Facebook's HHVM runtime combines LBR sampling with flow-sensitive analysis for real-time JIT tuning. Transfer learning, using pre-trained models of CFG embeddings, could predict hotspots for never-seen inputs. Automated input-generation methods, blending symbolic execution and fuzz testing (e.g., extending KLEE [13]), could generate high-coverage training sets for PGO.

PGO stands at the cusp of moving from "static tuning" to "dynamic intelligence". At the sampling layer, hardware-software co-sampling and intelligent denoising will break the efficiency-precision trade-off. For large codebases, precise slicing and distributed frameworks will tame data explosion. Cross-architecture adaptability and real-time feedback loops will propel intelligent optimization decisions. Dynamic, evolving training inputs will endow PGO with robust performance under changing loads. Looking further ahead, as emergent hardware like quantum accelerators and computation-in-memory platforms materialize, profile-guided optimizations will increasingly fuse with architectural innovations. This convergence will form a "sense-decide-optimize" autonomous loop that braces up efficient operation at extreme scale. The co-evolution of algorithms, compilers, and hardware will be the master key to unlocking the next generation of computational potential.

## 7 CONCLUSION

This article has demonstrated how profile-guided optimization techniques collect and generate runtime profile data and use that data to guide optimizations. Instrumentation-based PGO yields the most accurate profile data but incurs significant runtime overhead. Thanks to the synergy of hardware sampling and profile-mapping algorithms, sampling-based feedback optimizations overcome the high cost of data collection and the mapping inaccuracies inherent in compile-time methods, achieving high performance gains comparable to instrumentation-based approaches while reducing cross-module optimization overhead. Empirical results show that instrumentation-based PGO can deliver higher peak performance but at the expense of costly profile-collection phases; in contrast, Sampling-based PGO achieves consistent speedups with much lower deployment and runtime costs. In terms of compiler support, LLVM offers more mature implementations of feedback optimizations, while GCC retains advantages in certain specialized scenarios. Looking forward, realizing a true "sense-decide-optimize" closed loop will depend on hardware-software co-design for near-zero-overhead sampling, online dynamic evolution of training inputs, and robust cross-architecture adaptation.

## REFERENCES

[1] Advanced Micro Devices, Inc. 2024. *AMD64 Architecture Programmer's Manual, Volume 2: System Programming*. Publication No. 24593; Rev. 3.42. Advanced Micro Devices, Inc. https://www.amd.com/content/dam/amd/en/documents/processor-tech-docs/programmer-references/24593.pdf

[2] Glenn Ammons, Thomas Ball, and James R. Larus. 1997. Exploiting Hardware Performance Counters with Flow and Context Sensitive Profiling. In *Proceedings of the ACM SIGPLAN 1997 Conference on Programming Language Design and Implementation* (Las Vegas, Nevada, USA) *(PLDI '97)*. Association for Computing Machinery, New York, NY, USA, 85–96. doi:10.1145/258915.258924

[3] Glenn Ammons, Thomas Ball, and James R. Larus. 1997. Exploiting Hardware Performance Counters with Flow and Context Sensitive Profiling. *SIGPLAN Not.* 32, 5 (May 1997), 85–96. doi:10.1145/258916.258924

[4] Matthew Arnold, Stephen Fink, David Grove, Michael Hind, and Peter F. Sweeney. 2011. Adaptive Optimization in The Jalapeno JVM. *SIGPLAN Not.* 46, 4 (May 2011), 65–83. doi:10.1145/1988042.1988048

[5] Matthew Arnold and Barbara G. Ryder. 2001. A Framework for Reducing the Cost of Instrumented Code. In *Proceedings of the ACM SIGPLAN 2001 Conference on Programming Language Design and Implementation* (Snowbird, Utah, USA) *(PLDI '01)*. Association for Computing Machinery, New York, NY, USA, 168–179. doi:10.1145/378795.378832

[6] Amir Ayupov, Maksim Panchenko, and Sergey Pupyrev. 2024. Stale Profile Matching. In *Proceedings of the 33rd ACM SIGPLAN International Conference on Compiler Construction* (Edinburgh, United Kingdom) *(CC 2024)*. Association for Computing Machinery, New York, NY, USA, 162–173. doi:10.1145/3640537.3641573

[7] T. Ball and J.R. Larus. 1996. Efficient Path Profiling. In *Proceedings of the 29th Annual IEEE/ACM International Symposium on Microarchitecture. MICRO 29.* 46–57. doi:10.1109/MICRO.1996.566449

[8] Thomas Ball and James R. Larus. 1994. Optimally Profiling and Tracing Programs. *ACM Trans. Program. Lang. Syst.* 16, 4 (July 1994), 1319–1360. doi:10.1145/183432.183527

[9] M.D. Bond and K.S. McKinley. 2005. Practical Path Profiling for Dynamic Optimizers. In *International Symposium on Code Generation and Optimization.* 205–216. doi:10.1109/CGO.2005.28

[10] William J. Bowman, Swaha Miller, Vincent St-Amour, and R. Kent Dybvig. 2015. Profile-guided MEta-programming. *SIGPLAN Not.* 50, 6 (June 2015), 403–412. doi:10.1145/2813885.2737990

[11] Derek Bruening, Timothy Garnett, and Saman Amarasinghe. 2003. An Infrastructure for Adaptive Dynamic Optimization. In *Proceedings of the International Symposium on Code Generation and Optimization: Feedback-Directed and Runtime Optimization* (San Francisco, California, USA) *(CGO '03)*. IEEE Computer Society, USA, 265–275.

[12] Derek Bruening, Qin Zhao, and Reid Kleckner. 2020. DynamoRIO: Dynamic Instrumentation Tool Platform. *URL http://www. dynamorio. org* (2020).

[13] Cristian Cadar, Daniel Dunbar, and Dawson Engler. 2008. KLEE: Unassisted and Automatic Generation of High-coverage Tests for Complex Systems Programs. In *Proceedings of the 8th USENIX Conference on Operating Systems Design and Implementation* (San Diego, California) *(OSDI'08)*. USENIX Association, USA, 209–224.

[14] Pohua P. Chang, Scott A. Mahlke, and Wen-mei W. Hwu. 1991. Using Profile Information to Assist Classic Code Optimizations. *Softw. Pract. Exper.* 21, 12 (Dec. 1991), 1301–1321. doi:10.1002/spe.4380211204

[15] Dehao Chen, David Xinliang Li, and Tipp Moseley. 2016. AutoFDO: Automatic Feedback-Directed Optimization for Warehouse-Scale Applications. In *CGO 2016 Proceedings of the 2016 International Symposium on Code Generation and Optimization.* New York, NY, USA, 12–23.

[16] Dehao Chen, Neil Vachharajani, Robert Hundt, Shih-wei Liao, Vinodha Ramasamy, Paul Yuan, Wenguang Chen, and Weimin Zheng. 2010. Taming Hardware Event Samples for FDO Compilation. In *Proceedings of the 8th Annual IEEE/ACM International Symposium on Code Generation and Optimization* (Toronto, Ontario, Canada) *(CGO '10)*. Association for Computing Machinery, New York, NY, USA, 42–52. doi:10.1145/1772954.1772963

[17] Hyoun Kyu Cho, Tipp Moseley, Richard Hank, Derek Bruening, and Scott Mahlke. 2013. Instant Profiling: Instrumentation Sampling for Profiling Datacenter Applications. In *Proceedings of the 2013 IEEE/ACM International Symposium on Code Generation and Optimization (CGO).* 1–10. doi:10.1109/CGO.2013.6494982

[18] T.M. Conte, B.A. Patel, and J.S. Cox. 1994. Using Branch Handling Hardware to Support Profile-driven Optimization. In *Proceedings of MICRO-27. The 27th Annual IEEE/ACM International Symposium on Microarchitecture.* 12–21. doi:10.1145/192724.192726

[19] J. Dean, J.E. Hicks, C.A. Waldspurger, W.E. Weihl, and G. Chrysos. 1997. ProfileMe: Hardware Support for Instruction-level Profiling on Out-of-order Processors. In *Proceedings of 30th Annual International Symposium on Microarchitecture.* 292–302. doi:10.1109/MICRO.1997.645821

[20] Paul Drongowski, Lei Yu, Frank Swehosky, Suravee Suthikulpanit, and Robert Richter. 2010. Incorporating Instruction-Based Sampling into AMD CodeAnalyst. In *2010 IEEE International Symposium on Performance Analysis of Systems & Software (ISPASS).* 119–120. doi:10.1109/ISPASS.2010.5452049

[21] Fisher. 1981. Trace Scheduling: A Technique for Global Microcode Compaction. *IEEE Trans. Comput.* C-30, 7 (1981), 478–490. doi:10.1109/TC.1981.1675827

[22] Joseph A. Fisher and Stefan M. Freudenberger. 1992. Predicting Conditional Branch Directions from Previous Runs of a Program. In *Proceedings of the Fifth International Conference on Architectural Support for Programming Languages and Operating Systems* (Boston, Massachusetts, USA) *(ASPLOS V)*. Association for Computing Machinery, New York, NY, USA, 85–95. doi:10.1145/143365.143493

[23] Philip B. Gibbons and Steven S. Muchnick. 1986. Efficient Instruction Scheduling for a Pipelined Architecture. *SIGPLAN Not.* 21, 7 (July 1986), 11–16. doi:10.1145/13310.13312

[24] Wenlei He, Julián Mestre, Sergey Pupyrev, Lei Wang, and Hongtao Yu. 2022. Profile Inference Revisited. *Proc. ACM Program. Lang.* 6, POPL, Article 52 (Jan. 2022), 24 pages. doi:10.1145/3498714

[25] Wenlei He, Hongtao Yu, Lei Wang, and Taewook Oh. 2024. Revamping Sampling-Based PGO with Context-Sensitivity and Pseudo-instrumentation. In *2024 IEEE/ACM International Symposium on Code Generation and Optimization (CGO).* 322–333. doi:10.1109/CGO57630.2024.10444807

[26] Wenju He, Maosu Zhao, Yuxin Zou, and Feng Zou. 2023. Profile Guided Optimization Transfer-Learning for OpenCL/SYCL Kernel Compilation and Runtime. In *Proceedings of the 2023 International Workshop on OpenCL* (Cambridge, United Kingdom) *(IWOCL '23)*. Association for Computing Machinery, New York, NY, USA, Article 23, 1 pages. doi:10.1145/3585341.3585359

[27] R. R. Heisch. 1994. Trace-directed Program Restructuring for AIX Executables. *IBM Journal of Research and Development* 38, 5 (1994), 595–603. doi:10.1147/rd.385.0595

[28] Michael Held and Richard M. Karp. 1971. The Traveling-salesman Problem and Minimum Spanning Trees. *Operations Research* 18, 6 (1971), 1138–1162. doi:10.1287/opre.18.6.1138

[29] Martin Hirzel. 2001. Bursty Tracing: A Framework for Low-Overhead Temporal Profiling. https://api.semanticscholar.org/CorpusID:1588330

[30] Intel Corporation. 2025. *Intel® 64 and IA-32 Architectures Software Developer's Manual.* Version latest. Intel Corporation. https://www.intel.com/content/www/us/en/developer/articles/technical/intel-sdm.html PDF manual updated Mar. 5, 2025.

[31] Thomas Kistler and Michael Franz. 2003. Continuous Program Optimization: A Case Study. *ACM Trans. Program. Lang. Syst.* 25, 4 (July 2003), 500–548. doi:10.1145/778559.778562

[32] James R. Larus. 1999. Whole Program Paths. In *Proceedings of the ACM SIGPLAN 1999 Conference on Programming Language Design and Implementation* (Atlanta, Georgia, USA) *(PLDI '99)*. Association for Computing Machinery, New York, NY, USA, 259–269. doi:10.1145/301618.301678

[33] Roy Levin, Ilan Newman, and Gadi Haber. 2008. Complementing Missing and Inaccurate Profiling Using a Minimum Cost Circulation Algorithm. In *High Performance Embedded Architectures and Compilers*, Per Stenström, Michel Dubois, Manolis Katevenis, Rajiv Gupta, and Theo Ungerer (Eds.). Springer Berlin Heidelberg, Berlin, Heidelberg, 291–304.

[34] Xinliang David Li, Raksit Ashok, and Robert Hundt. 2010. Lightweight Feedback-Directed Cross-Module Optimization. In *Proceedings of International Symposium on Code Generation and Optimization (CGO)*.

[35] Shen Lin and Brian W. Kernighan. 1973. An Effective Heuristic Algorithm for The Traveling-salesman Problem. *Operations Research* 21, 2 (1973), 498–516. doi:10.1287/opre.21.2.498

[36] Chi-Keung Luk, Robert Cohn, Robert Muth, Harish Patil, Artur Klauser, Geoff Lowney, Steven Wallace, Vijay Janapa Reddi, and Kim Hazelwood. 2005. Pin: Building Customized Program Analysis Tools with Dynamic Instrumentation. *SIGPLAN Not.* 40, 6 (June 2005), 190–200. doi:10.1145/1064978.1065034

[37] S. McFarling and J. Hennesey. 1986. Reducing The Cost of Branches. *SIGARCH Comput. Archit. News* 14, 2 (May 1986), 396–403. doi:10.1145/17356.17402

[38] David Gordon Melski and Thomas Reps. 2002. *Interprocedural Path Profiling and The Interprocedural Express-lane Transformation.* Ph. D. Dissertation. AAI3049391.

[39] Samuel Miksits, Ruimin Shi, Maya Gokhale, Jacob Wahlgren, Gabin Schieffer, and Ivy Peng. 2025. Multi-level Memory-Centric Profiling on ARM Processors with ARM SPE. In *Proceedings of the SC '24 Workshops of the International Conference on High Performance Computing, Network, Storage, and Analysis* (Atlanta, GA, USA) *(SC-W '24)*. IEEE Press, 996–1005. doi:10.1109/SCW63240.2024.00139

[40] Tipp Moseley, Alex Shye, Vijay Janapa Reddi, Dirk Grunwald, and Ramesh Peri. 2007. Shadow Profiling: Hiding Instrumentation Costs with Parallelism. In *International Symposium on Code Generation and Optimization (CGO'07)*. 198–208. doi:10.1109/CGO.2007.35

[41] Nicholas Nethercote and Julian Seward. 2007. Valgrind: a Framework for Heavyweight Dynamic Binary Instrumentation. In *Proceedings of the 28th ACM SIGPLAN Conference on Programming Language Design and Implementation* (San Diego, California, USA) *(PLDI '07)*. Association for Computing Machinery, New York, NY, USA, 89–100. doi:10.1145/1250734.1250746

[42] Maksim Panchenko, Rafael Auler, Bill Nell, and Guilherme Ottoni. 2019. BOLT: A Practical Binary Optimizer for Data Centers and Beyond. In *Proceedings of the 2019 IEEE/ACM International Symposium on Code Generation and Optimization* (Washington, DC, USA) *(CGO 2019)*. IEEE Press, 2–14.

[43] Maksim Panchenko, Rafael Auler, Laith Sakka, and Guilherme Ottoni. 2021. Lightning BOLT: Powerful, Fast, and Scalable Binary Optimization. In *Proceedings of the 30th ACM SIGPLAN International Conference on Compiler Construction* (Virtual, Republic of Korea) *(CC 2021)*. Association for Computing Machinery, New York, NY, USA, 119–130. doi:10.1145/3446804.3446843

[44] Karl Pettis and Robert C. Hansen. 1990. Profile Guided Code Positioning. *SIGPLAN Not.* 25, 6 (June 1990), 16–27. doi:10.1145/93548.93550

[45] Vinodha Ramasamy, Paul Yuan, Dehao Chen, and Robert Hundt. 2008. Feedback-Directed Optimizations in GCC with Estimated Edge Profiles from Hardware Event Sampling. In *Proceedings of GCC Summit 2008*. 87–102. http://www.capsl.udel.edu/conferences/open64/2008/Papers/113.pdf

[46] Han Shen, Krzysztof Pszeniczny, Rahman Lavaee, Snehasish Kumar, Sriraman Tallam, and Xinliang David Li. 2023. Propeller: A Profile Guided, Relinking Optimizer for Warehouse-Scale Applications. In *Proceedings of the 28th ACM*

*International Conference on Architectural Support for Programming Languages and Operating Systems, Volume 2* (Vancouver, BC, Canada) *(ASPLOS 2023)*. Association for Computing Machinery, New York, NY, USA, 617–631. doi:10.1145/3575693.3575727

[47] Michael D. Smith. 2000. Overcoming The Challenges to Feedback-directed Optimization (Keynote Talk). *SIGPLAN Not.* 35, 7 (Jan. 2000), 1–11. doi:10.1145/351403.351408

[48] Zhendong Su and Min Zhou. [n. d.]. *A Comparative Analysis of Branch Prediction Schemes*. https://www.cs.ucdavis.edu/~su/Berkeley/cs252/project.html

[49] Toshio Suganuma, Toshiaki Yasue, and Toshio Nakatani. 2006. A Region-based Compilation Technique for Dynamic Compilers. *ACM Trans. Program. Lang. Syst.* 28, 1 (Jan. 2006), 134–174. doi:10.1145/1111596.1111600

[50] Robert Endre Tarjan. 1983. *Data Structures and Network Algorithms.* Society for Industrial and Applied Mathematics. doi:10.1137/1.9781611970265 arXiv:https://epubs.siam.org/doi/pdf/10.1137/1.9781611970265

[51] Omri Traub, Stuart Schechter, and Michael D Smith. 2000. Ephemeral Instrumentation for Lightweight Program Profiling. *Unpublished technical report, Department of Electrical Engineering and Computer Science, Hardward University, Cambridge, Massachusetts* (2000).

[52] John Whaley. 2000. A Portable Sampling-based Profiler for Java Virtual Machines. In *Proceedings of the ACM 2000 Conference on Java Grande* (San Francisco, California, USA) *(JAVA '00)*. Association for Computing Machinery, New York, NY, USA, 78–87. doi:10.1145/337449.337483

[53] John Whaley. 2001. Partial Method Compilation Using Dynamic Profile Information. In *Proceedings of the 16th ACM SIGPLAN Conference on Object-Oriented Programming, Systems, Languages, and Applications* (Tampa Bay, FL, USA) *(OOPSLA '01)*. Association for Computing Machinery, New York, NY, USA, 166–179. doi:10.1145/504282.504295

[54] Baptiste Wicht, Roberto A. Vitillo, Dehao Chen, and David Levinthal. 2014. Hardware Counted Profile-Guided Optimization. *CoRR* abs/1411.6361 (2014). arXiv:1411.6361 http://arxiv.org/abs/1411.6361

[55] David Williams-King and Junfeng Yang. 2019. CodeMason: Binary-Level Profile-Guided Optimization. In *Proceedings of the 3rd ACM Workshop on Forming an Ecosystem Around Software Transformation* (London, United Kingdom) *(FEAST'19)*. Association for Computing Machinery, New York, NY, USA, 47–53. doi:10.1145/3338502.3359763

[56] Bo Wu, Mingzhou Zhou, Xipeng Shen, Yaoqing Gao, Raul Silvera, and Graham Yiu. 2013. Simple Profile Rectifications Go a Long Way. In *ECOOP 2013 – Object-Oriented Programming*, Giuseppe Castagna (Ed.). Springer Berlin Heidelberg, Berlin, Heidelberg, 654–678.

[57] Cliff Young, David S. Johnson, Michael D. Smith, and David R. Karger. 1997. Near-optimal Intraprocedural Branch Alignment. In *Proceedings of the ACM SIGPLAN 1997 Conference on Programming Language Design and Implementation* (Las Vegas, Nevada, USA) *(PLDI '97)*. Association for Computing Machinery, New York, NY, USA, 183–193. doi:10.1145/258915.258932

[58] Cliff Young and Michael D. Smith. 1999. Static Correlated Branch Prediction. *ACM Trans. Program. Lang. Syst.* 21, 5 (Sept. 1999), 1028–1075. doi:10.1145/330249.330255

[59] Mingzhou Zhou, Bo Wu, Xipeng Shen, Yaoqing Gao, and Graham Yiu. 2016. Examining and Reducing the Influence of Sampling Errors on Feedback-Driven Optimizations. *ACM Trans. Archit. Code Optim.* 13, 1, Article 6 (April 2016), 24 pages. doi:10.1145/2851502

[60] Craig Zilles and Gurindar Sohi. 2001. Execution-based Prediction Using Speculative Slices. *SIGARCH Comput. Archit. News* 29, 2 (May 2001), 2–13. doi:10.1145/384285.379246

[61] C.B. Zilles and G.S. Sohi. 2001. A Programmable Co-processor for Profiling. In *Proceedings HPCA Seventh International Symposium on High-Performance Computer Architecture*. 241–252. doi:10.1109/HPCA.2001.903267