

# Towards Online Code Specialization of Systems

Vaastav Anand, Deepak Garg, Antoine Kaufmann  
Max Planck Institute for Software Systems

## Abstract

Specializing low-level systems to specifics of the workload they serve and platform they are running on often significantly improves performance. However, specializing systems is difficult because of three compounding challenges: i) specialization for optimal performance requires in-depth compile-time changes; ii) the right combination of specialization choices for optimal performance is hard to predict a priori; and iii) workloads and platform details often change online. In practice, benefits of specialization are thus not attainable for many low-level systems.

To address this, we advocate for a radically different approach for performance-critical low-level systems: designing and implementing systems with and for runtime code specialization. We leverage just-in-time compilation to change systems code based on developer-specified specialization points as the system runs. The JIT runtime automatically tries out specialization choices and measures their impact on system performance, e.g. request latency or throughput, to guide the search. With Iridescent, our early prototype, we demonstrate that online specialization (i) is feasible even for low-level systems code, such as network stacks, (ii) improves system performance without the need for complex cost models, (iii) incurs low developer effort, especially compared to manual exploration. We conclude with future opportunities online system code specialization enables.

## 1 Introduction

Specializing system implementations to workload characteristics and hardware can significantly improve performance and efficiency [2, 4, 7, 9–11, 14, 15, 21, 23, 24, 28, 29, 32, 33, 35, 39]. To achieve these benefits for particular hardware and workload combination, system developers manually modify the system code and recompile to benefit from the compile-time optimizations enabled by specialization [25, 36].

System specialization comes at the cost of generality. A system heavily specialized to workload and hardware either performs poorly outside of this regime, or completely fails. As a result, developers today must carefully navigate this specialization-generalization tradeoff and optimize for the most common hardware and workload setting.

However, even selecting this design point a priori is difficult for multiple reasons. First, the workload and platform are often unknown at development time. Systems are often deployed to different hardware than they were developed

on. Second, workload and platform often change at runtime. A system might migrate to newer deployments over time. Application workloads rarely remain fixed throughout the application lifetime. Third, even for a known workload and platform, performance for the modified system is hard to predict. Modern system performance is dictated by emergent behaviors arising from the complex intra-system interactions and interactions between the system and the workload. To make matters worse, the exact same system and workload are likely to exhibit different performance on different hardware.

As a result, making optimal system code specialization choices is fundamentally an iterative process: run the system, profile, adjust the code, and repeat. Unfortunately this is also extremely laborious for developers, and thus rarely worth the cost and complexity in practice.

We propose a fundamentally different approach to system specialization for performance and efficiency: *automated online specialization*. We aim to specialize systems at runtime to the exact workload and hardware conditions where the system is deployed, *without the human developer in the loop*. To demonstrate the idea’s feasibility, we present a design proposal for a specialization toolbox that can iteratively specialize low-level systems online with runtime performance feedback based on developer-annotated possible *specialization points*. We have developed an initial prototype Iridescent, and provide early evidence for feasibility, performance improvements, and lowered developer effort. Further, we argue that this mechanism of code changes at runtime with full system performance could fundamentally change how we build systems and how we approach compiler optimization.

## 2 Specialization is Effective but Complex

We illustrate the potential and pitfalls specializations for optimizing low-level system code with a well understood algorithm: blocked/tiled matrix multiplication to optimize cache locality [5, 37, 38]. By dividing up input and output matrices into blocks and operating on these instead of full columns or rows, cache locality improves substantially. We measure different combinations of (square) input matrix sizes and block sizes 5 different processor core architectures (see rows and columns in Table 1). Blocking works as expected. We find that blocking reduces required processor cycles by  $1.15\times$  to  $28.5\times$  with the ideal block size.

**Specialization improves performance.** When trying to implement this in a system we need to choose a block size

Machine / Workload	N=1024	N=256	N=64
IceLake	32	32	32
IvyBridge	16	16	4
CoffeeLake	32	4	4
AlderLake-p	32	4	2
AlderLake-e	64	4	4

**Table 1: Optimal configurations for our block matrix multiply, across 5 hardware platforms and 3 workloads.**

to use. Unfortunately we find that the optimal block size depends on the hardware the system is running on, different cache sizes or latencies, and the workload, such as different matrix sizes. We illustrate this in Table 1 where we show the best performing block size for each combination, and find 5 different optimal block sizes depending on the datapoint. Thus choosing a specific block size is specializing the implementation to the machine and workload, reduces performance in other cases.

**Code changes are necessary.** This is a challenge since we find that the block size really has to be fixed at compile time. Because of the computationally dense nature of the algorithm, good performance really requires heavy compiler optimizations, such as loop unrolling and vectorization. And since blocking dictates the trip count of the two innermost loops, these optimizations are not feasible if the block size is a variable. Here we find that using a variable instead of a fixed block size can reduce performance by up to 6.5 $\times$ .

**Optimal choices are hard to predict.** Our results in Table 1 show that the optimal parameters are difficult to predict a priori, even for a fixed workload. Depending on the machine three different block sizes may be ideal. With dynamic workloads that change over time, or external factors such as noisy neighbor cores filling up caches, this quickly becomes completely unpredictable. This combination leaves the system developer with no good options: The best performance requires compile-time specialization, but the best choices are dynamic and unpredictable. The closest solution for simple cases is to implement multiple different versions of the code and switch between them at runtime. However, this is laborious and does not scale beyond a single parameter.

### 3 Specializing without Human in-the Loop

We enable system developers to easily build systems that perform well across a broad set of workloads and platforms with a fundamentally different approach. Instead of exclusively compiling ahead of time, we argue that this requires *online specialization* through code changes as the system runs. This enables one version of the system to automatically adapt to the workload and hardware conditions. Specializing online makes it possible to directly use system performance as feedback for choosing optimal configurations.

We propose to use just in time (JIT) compilation to be able to change the code as it runs. We combine this with developer provided annotation for potentially relevant specialization points in the code. Finally we leverage overall system performance measurements, such as request throughput, to guide an automated, empirical search for the best configuration.

#### **Lightweight developer instrumentation is sufficient.**

Finding useful specialization opportunities in low-level systems code is extremely challenging for compilers. However, developers have excellent intuition about *possible specialization points*. We thus require developers to provide lightweight annotations about potential values in the code that could be specialized to compile time constants. The key insight here, is that commonly simply turning a variable in the code into a constant, then enables a host of cascading compiler optimizations through constant propagation, dead code, elimination, unrolling, etc. Thus we can leverage minimal annotation to automatically make substantial changes to the code.

**Fallbacks for correctness are necessary.** When automatically specializing and since real systems are complex, incorrect specializations are inevitable. Crashes or incorrect behavior for these cases are not a viable outcome. To avoid this, we include guards in the specialized code to make sure the conditions hold. If a guard triggers, we fall back to the original code version for this execution, incurring a performance penalty but no other problems.

**Online measurements can replace cost models.** The combination of JIT and specialization points provides the necessary mechanisms to specialize code at runtime, but leaves the runtime with a large number of potential knobs to turn. Ahead of time compilers making similar decisions automatically typically rely on built-in cost models. Unfortunately, cost models for system performance are inevitably highly inaccurate and would result in suboptimal results. However, since we are making these decisions online as the system runs, we can simply do what human developers do as well: try out different options and measure their impact on full system performance. Online specialization thus enables fully accurate decisions for specialization choices based on measurements rather than models.

### 4 An Online Specialization Runtime

We present a design proposal for enabling online specialization. Figure 1 shows the design of Iridescent, a runtime specialization toolbox that allows developers to automatically specialize any given software based on metrics calculated at runtime. Iridescent supports exploration of the set of all possible specializations, applying a given particular specialization, and switching between different specializations at runtime to handle changing workloads.

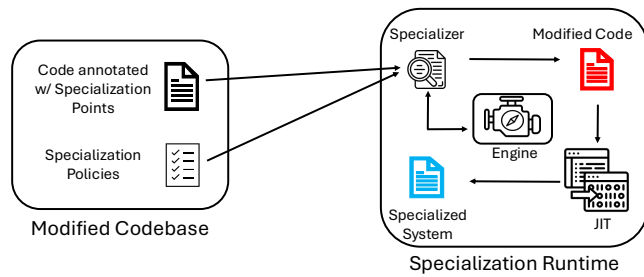


Figure 1: Iridescent design

## 4.1 System Modifications

**Source-Code Division.** Developers divide the codebase into two parts: (i) the performance-critical *handler code*; (ii) the *fixed code*. The handler code is the set of functions and their dependencies in the codebase that the developer wishes to specialize. In contrast, the fixed code entails the general-purpose part of the codebase that makes function calls into the handler code to execute incoming requests. For example, in the MMulBlockBench microbenchmark, `matrix_multiply` would comprise the performance-critical handler code. The callers of this function would comprise the fixed code.

**Specialization points.** The developer instruments the handler code in the specialization runtime to create *specialization points*. A specialization point is a variable or a parameter in the handler code that can potentially be converted into a constant during execution. Broadly, there are two types of specialization points: (i) workload knobs; and (ii) configuration knobs. Workload knobs are variables that change with potentially every request in a given workload. For example, in the MMulBlockBench microbenchmark, the block size  $s$  is the workload knob. In contrast, configuration knobs are variables that do not rely on workload data and are usually provided as input to the program as part of a command line flag or a configuration file. For example, this would be the batch size of a packet processing framework.

**Specialization Policies.** Through the specialization policies, developers control the exploration and selection of specializations. The developers initialize the specialization runtime with the target end-to-end metric and a way to measure the metric and the specific exploration policy to use to explore the specialization space of the handler code.

## 4.2 Specialization Runtime

To specialize the handler code, the developer first places the handler code in the *specialization runtime*. The specialization runtime is a wrapper around a JIT that manages the lifetime of the handler code inside the JIT.

**Specializer.** To specialize a point with a given value, the specializer specializes the handler code by replacing the variable in the code that is mapped to the specialization point

Hook	Description
<code>runtime_init</code>	Initialize runtime
<code>update_runtime</code>	Recompile handler code
<code>get_specialized_function</code>	Get address of specialized function
<code>point_specialize</code>	Specialize point to a value
<code>point_disable_spec</code>	Disable specialization
<code>point_disable_spec_check</code>	Disable specialization check
<code>point_enable_collection</code>	Enable data collection
<code>point_disable_collection</code>	Disable data collection
<code>start_exploration</code>	Start exploration
<code>end_exploration</code>	End exploration

Table 2: Runtime hooks provided by Iridescent

and fixes its value to the specific value. The handler code is then recompiled in the JIT compiler by the specialized runtime. The recompilation process allows for various compiler optimization passes to take effect that now view the previously runtime specialization point as a fixed compile time constant. This allows for the optimization passes to be more aggressive with optimizations such as constant propagation, loop unrolling, instruction vectorization, among others to produce a specialized version of the handler code.

Even though the specialization points are applied at a variable level, the code specialization takes place at the granularity of a function. That is, once a specialization point is specialized, a new version of the function containing that specialization point will be generated with the variable's value now fixed to the specialized value.

**Handling errors and side-effects.** To ensure correctness, for each specialized function, the specialization runtime inserts a runtime check called a specialization check that checks if an input value matches the current specialized value of the runtime. This check allows for the runtime to detect if a new incoming request in the workload meets the criteria to use the specialized version of the code. If this check fails, then the specialization runtime cleans up any side-effects and transfers control back to the generic version of the function. To ensure correct clean up of side effects, the developers define and register a dedicated cleanup function in the handler code for each function that could be specialized by the runtime. It is critical to note here that not all side-effects are reversible (e.g. sending a packet to a neighbor), so Iridescent performs a best-effort clean-up.

**Using the specialized code.** To use the specialized version of the handler, the fixed code must obtain the addresses of all possible specialized functions from the specialization runtime. The fixed code must use these function addresses to make function calls to the specialized version of the handler code. If no specialization has been enabled for any given

specialization point, then the function address defaults to the non-customized generic version of the function.

**Specialization Hooks.** The specialization runtime provides hooks in the fixed code that the developers can use to control the behavior of the specialization runtime. Table 2 provides a list of basic operations that are supported by the runtime. The operations include which specialization point(s) to specialize, when to specialize them, when to change the specialization, and updating the objective function to use for selecting the optimal configuration. Developers can use the hooks in the fixed code to change the behavior of the runtime.

**Exploration Engine.** As the space of possible specializations and their combinations is large, it is almost impossible for developers to know which possible specializations must be enabled in conjunction for best performance. Thus, the fixed code needs a way to try out different specializations and then settle on a set of specializations.

To achieve this, Iridescent provides an explorer engine that can be configured to explore different specializations for the various specialization points. The runtime instruments the handler code to start collecting three pieces of information for the specialization point - (i) the set of values for that point, (ii) the frequency with which each value is seen in the workload, and (iii) the number of specialization check failures. The developer can override or turn-off the data collection process by providing a fixed set of values that should be used during the exploration.

### 4.3 Prototype Details

We have implemented our current prototype in 3K lines of C++ code. Our current prototype uses the JIT compiler provided by LLVM [17] to generate the specialized code at runtime. The specializer is implemented as a set of LLVM transformation passes that operate on LLVM IR of the handler code. We are also in the process of implementing a prototype of online specialization in Go and integrating it into the microservice generation framework, Blueprint [3].

## 5 Use Cases

We showcase how we can use online specialization with Iridescent for three use-cases: (i) enabling compile-time optimizations at runtime; (ii) enabling incremental specializations at runtime; (iii) design exploration at runtime.

### 5.1 Compiler Optimizations at Runtime

Compile-time optimizations often produce more performant code. These compile time optimizations include constant propagation, loop unrolling, dead-code elimination, and use of vectorization instructions.

With Iridescent, we can enable compile time optimizations at runtime. To show the potential benefits of enabling

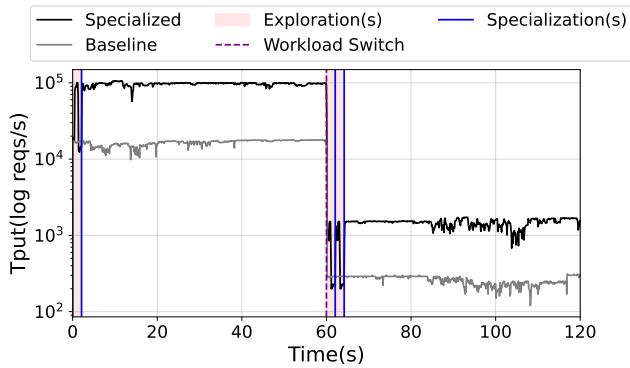
Machine (Processor)	Constant (c) (cycles/op)	Variable (v) (cycles/op)	Benefit (v/c)
IceLake	175297	284944	61%
IvyBridge	250434	661295	246%
CoffeeLake	168817	581130	348%
AlderLake-p	173350	583724	336%
AlderLake-e	206924	557572	269%

**Table 3: Impact of turning  $s$  as compile-time constant at runtime for  $N=64$  for different hardwares (Table 1)**

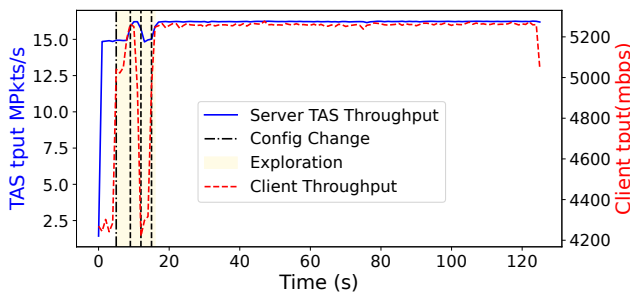
compile-time optimizations at runtime, we use the MMul-BlockBench shown in Table 1. Specifically, for the workload  $N = 64$ , we compare the performance of keeping the optimal block size,  $s$ , as a runtime parameter to that of converting it into a constant at runtime using Iridescent. For both configurations, we execute the function for a fixed number of times and measure the amount of cycles spent for executing 1 execution of the function. Table 3 shows the benefit for converting  $s$  into a compile-time constant for different hardwares. For each hardware, we get at least 50% reduction in consumed cycles, and greater than 240% reduction in consumed cycles for four of the five operating conditions.

### 5.2 Incremental Specialization

Incremental computing and has been a mature idea [26] for decades. In recent years, incremental specialization has been used for packet processing frameworks [9, 21, 27]. To show that Iridescent can easily enable incremental specialization, we extend the Network Functions from Vigor [40] and Pix [12]. Specifically, we target the lpm (longest prefix match) network function, which matches the incoming packet's destination address to one of its rules in its prefix table and then routes the packet based on the selected rule. We extend this function with Iridescent to enable a specialization similar to the hot map specialization of Morpheus [21]. In this specialization, Iridescent divides the execution into a monitoring phase and a specialization phase. In the monitoring phase, Iridescent initially instruments the lpm code to find the most common input addresses and store the calculated routes for these addresses. After a user-specified amount of time, Iridescent converts the collected information into hard-coded values at the entry point of the lpm function to return the hard-coded route if the input address is one of the hardcoded addresses. With Iridescent, we see a 9% increase in the throughput.



**Figure 2: Automatic Exploration and Specialization of MMulBlockBench for two different workloads**



**Figure 3: Automatic Exploration and Specialization of different configurations for TAS**

### 5.3 Runtime Design Exploration

Automatically searching over a design space at runtime has been successfully applied for optimizing platform configuration settings [31, 32], in our use-case we target design space exploration of the deployed system itself.

To show that Iridescent can enable automatic runtime design exploration, we first use the MMulBlockBench microbenchmark. We execute a Iridescent-enabled version of the function and a non-Iridescent enabled baseline version. In the non-baseline version, the block size  $s$  is a runtime parameter. In the Iridescent-enabled version, Iridescent explores different values of  $s$  by specializing the value in the code and converting  $s$  into a compile-time constant. We execute each version with two different workloads in succession, each lasting one minute. We measure the overall throughput (executions/s) for both configurations. Figure 2 shows how Iridescent can easily explore the different block sizes and find a performant configuration as compared to the baseline. Moreover, Iridescent can automatically detect the workload change based on the drop in throughput and restart the A/B testing process.

Next, we show runtime design exploration for TAS [16], a TCP acceleration stack. We setup our experiment with a typical server-client setting, with the server providing an echo service. The client is multi-threaded and executes an open-loop workload of 64 byte packets. Both the server and client are enabled with TAS. In our experiment, we modify the server-side TAS and measure the throughput of the server side TAS as millions of packets processed per second. With Iridescent we explore different values of the batch size in the TAS source code on the server side. Normally, the batch size is a fixed constant throughout the TAS source code. Instead of treating it as a homogeneous constant, we convert batch size into three different variables for three different locations in the source code and explore different combinations of the values of the batch sizes. Figure 3 shows how Iridescent can automatically select the best-configuration on the server side by automatically exploring the different values of the batch size for different locations in the source code.

## 6 Outlook

We believe that online optimization of systems with metric-guided exploration can open multiple research avenues. We highlight some promising future directions below.

**Empirically enabled Compiler Optimizations.** In our demonstration, we have only focused on one aspect of specialization - converting runtime parameters into compile time constants. Developers can develop custom compile-time optimization techniques that are designed specifically for leveraging the structure of their system and combining it with empirical information to optimize the system. For example, offline profile-guided prefetching instruction injection LLVM pass of APT-GET [13] can seamlessly be integrated with Iridescent to better utilize the prefetching abilities of the processor at runtime. Moreover, incorporating domain-specific techniques like offline compiler optimizations of PacketMill [9] into Iridescent represent interesting future specialization opportunities.

**Leveraging scale & ML for exploration.** Several ML-based techniques exist for tuning configurations [1, 6, 8, 13, 18–20, 27, 30, 31, 34, 35]. OPPerTune [31] provides access to a plethora of ML techniques toolbox for automatically searching through large configuration spaces. We believe Iridescent can be integrated with OPPerTune to utilize state-of-the-art exploration techniques. We believe developers can leverage the large-scale nature of modern cloud systems and explore the large specialization spaces by measuring the performance of different specializations in parallel.

**Automatic Online System Repair.** Iridescent can be extended with techniques to automatically detect performance problems and then automatically explore the design space of the system to fix the issue. For example, Iridescent can be

used to detect and fix emergent misbehaviors [22] caused by complex interactions between different components of the system.

## 7 Conclusion

We introduced online specialization to enable automatic online adaptation to workload and hardware specifics. With our proposed design of Iridescent, we enable online optimization of performance-critical systems by dynamically exploring and applying the most effective specialization. Iridescent enables developers to potentially achieve significant performance improvements without the complexity of manual specialization. We believe, Iridescent can be used as a building block for enabling custom specializations at runtime to extract maximum performance from systems.

## References

- [1] Ibrahim Umit Akgun, Ali Selman Aydin, Andrew Burford, Michael McNeill, Michael Arkhangelskiy, and Erez Zadok. Improving storage systems using machine learning. *ACM Transactions on Storage*, 19(1):1–30, 2023.
- [2] Omid Alipourfard, Hongqiang Harry Liu, Jianshu Chen, Shivaram Venkataraman, Minlan Yu, and Ming Zhang. {CherryPick}: Adaptively unearthing the best cloud configurations for big data analytics. In *14th USENIX Symposium on Networked Systems Design and Implementation (NSDI 17)*, pages 469–482, 2017.
- [3] Vaastav Anand, Deepak Garg, Antoine Kaufmann, and Jonathan Mace. Blueprint: A toolchain for highly-reconfigurable microservice applications. In *Proceedings of the 29th Symposium on Operating Systems Principles*, pages 482–497, 2023.
- [4] Sapan Bhatia, Charles Consel, A-F Le Meur, and Calton Pu. Automatic specialization of protocol stacks in operating system kernels. In *29th Annual IEEE International Conference on Local Computer Networks*, pages 152–159. IEEE, 2004.
- [5] Steve Carr and Ken Kennedy. Blocking linear algebra codes for memory hierarchies. In *PPSC*, pages 400–405. Citeseer, 1989.
- [6] Stefano Cereda, Stefano Valladares, Paolo Cremonesi, and Stefano Doni. Cgptuner: a contextual gaussian process bandit approach for the automatic tuning of it configurations under varying workload conditions. *Proceedings of the VLDB Endowment*, 14(8):1401–1413, 2021.
- [7] Bangwen Deng, Wenfei Wu, and Linhai Song. Redundant logic elimination in network functions. In *Proceedings of the Symposium on SDN Research*, pages 34–40, 2020.
- [8] Songyun Duan, Vamsidhar Thummala, and Shivnath Babu. Tuning database configuration parameters with ituned. *Proceedings of the VLDB Endowment*, 2(1):1246–1257, 2009.
- [9] Alireza Farshin, Tom Barbette, Amir Roozbeh, Gerald Q Maguire Jr, and Dejan Kostić. Packetmill: toward per-core 100-gbps networking. In *Proceedings of the 26th ACM International Conference on Architectural Support for Programming Languages and Operating Systems*, pages 1–17, 2021.
- [10] Hamid Ghasemirahni, Tom Barbette, Georgios P Katsikas, Alireza Farshin, Amir Roozbeh, Massimo Girondi, Marco Chiesa, Gerald Q Maguire Jr, and Dejan Kostić. Packet order matters! improving application performance by deliberately delaying packets. In *19th USENIX Symposium on Networked Systems Design and Implementation (NSDI 22)*, pages 807–827, 2022.
- [11] Hamid Ghasemirahni, Alireza Farshin, Dejan Kostic, and Marco Chiesa. Just-in-time packet state prefetching. *arXiv preprint arXiv:2407.04344*, 2024.
- [12] Rishabh Iyer, Katerina Argyraki, and George Candea. Performance interfaces for network functions. In *19th USENIX Symposium on Networked Systems Design and Implementation (NSDI 22)*, pages 567–584, 2022.
- [13] Saba Jamilan, Tanvir Ahmed Khan, Grant Ayers, Baris Kasikci, and Heiner Litz. Apt-get: Profile-guided timely software prefetching. In *Proceedings of the Seventeenth European Conference on Computer Systems*, pages 747–764, 2022.
- [14] Anuj Kalia, Michael Kaminsky, and David Andersen. Datacenter RPCs can be general and fast. In *16th USENIX Symposium on Networked Systems Design and Implementation*, NSDI, 2019.
- [15] Ajaykrishna Karthikeyan, Nagarajan Natarajan, Gagan Somashekar, Lei Zhao, Ranjita Bhagwan, Rodrigo Fonseca, Tatiana Racheva, and Yogesh Bansal. {SelfTune}: Tuning cluster managers. In *20th USENIX Symposium on Networked Systems Design and Implementation (NSDI 23)*, pages 1097–1114, 2023.
- [16] Antoine Kaufmann, Tim Stamler, Simon Peter, Naveen Kr. Sharma, Arvind Krishnamurthy, and Thomas Anderson. TAS: TCP acceleration as an OS service. In *14th ACM European Conference on Computer Systems*, EuroSys, 2019.
- [17] Chris Lattner and Vikram Adve. LLVM: A compilation framework for lifelong program analysis & transformation. In *9th International Symposium on Code Generation and Optimization*, CGO, 2004.
- [18] Zhao Lucis Li, Chieh-Jan Mike Liang, Wenjia He, Lianjie Zhu, Wenjun Dai, Jin Jiang, and Guangzhong Sun. Metis: Robustly tuning tail latencies of cloud systems. In *2018 USENIX Annual Technical Conference (USENIX ATC 18)*, pages 981–992, 2018.
- [19] Chen Lin, Junqing Zhuang, Jiadong Feng, Hui Li, Xuanhe Zhou, and Guoliang Li. Adaptive code learning for spark configuration tuning. In *2022 IEEE 38th International Conference on Data Engineering (ICDE)*, pages 1995–2007. IEEE, 2022.
- [20] Ashraf Mahgoub, Alexander Michaelson Medoff, Rakesh Kumar, Subrata Mitra, Ana Klimovic, Somali Chatterji, and Saurabh Bagchi. {OPTIMUSCLOUD}: Heterogeneous configuration optimization for distributed databases in the cloud. In *2020 USENIX Annual Technical Conference (USENIX ATC 20)*, pages 189–203, 2020.
- [21] Sebastiano Miano, Alireza Sanaee, Fulvio Rizzo, Gábor Rétvári, and Gianni Antichi. Domain specific run time optimization for software data planes. In *Proceedings of the 27th ACM International Conference on Architectural Support for Programming Languages and Operating Systems*, pages 1148–1164, 2022.
- [22] Jeffrey C Mogul. Emergent (mis) behavior vs. complex software systems. *ACM SIGOPS Operating Systems Review*, 40(4):293–304, 2006.
- [23] László Molnár, Gergely Pongrácz, Gábor Enyedi, Zoltán Lajos Kis, Levente Csikor, Ferenc Juhász, Attila Körösi, and Gábor Rétvári. Dataplane specialization for high-performance openflow software switching. In *Proceedings of the 2016 ACM SIGCOMM Conference*, pages 539–552, 2016.
- [24] Heng Pan, Peng He, Zhenyu Li, Pan Zhang, Junjie Wan, Yuhao Zhou, XiongChun Duan, Yu Zhang, and Gaogang Xie. Hoda: a high-performance open vswitch dataplane with multiple specialized data paths. In *Proceedings of the Nineteenth European Conference on Computer Systems*, pages 82–98, 2024.
- [25] Maksim Panchenko, Rafael Auler, Bill Nell, and Guilherme Ottoni. Bolt: a practical binary optimizer for data centers and beyond. In *2019 IEEE/ACM International Symposium on Code Generation and Optimization (CGO)*, pages 2–14. IEEE, 2019.

- [26] Ganesan Ramalingam and Thomas Reps. A categorized bibliography on incremental computation. In *Proceedings of the 20th ACM SIGPLAN-SIGACT symposium on Principles of programming languages*, pages 502–510, 1993.
- [27] Fabian Ruffy, Zhanghan Wang, Gianni Antichi, Aurojit Panda, and Anirudh Sivaraman. Incremental specialization of network programs. In *Proceedings of the 23rd ACM Workshop on Hot Topics in Networks*, pages 264–272, 2024.
- [28] Krzysztof Rządca, Paweł Findeisen, Jacek Swiderski, Przemysław Zych, Przemysław Broniek, Jarek Kusmirek, Paweł Nowak, Beata Strack, Piotr Witusowski, Steven Hand, et al. Autopilot: workload autoscaling at google. In *Proceedings of the Fifteenth European Conference on Computer Systems*, pages 1–16, 2020.
- [29] Hashim Sharif, Muhammad Abubakar, Ashish Gehani, and Fareed Zaffar. Trimmer: application specialization for code debloating. In *Proceedings of the 33rd ACM/IEEE International Conference on Automated Software Engineering*, pages 329–339, 2018.
- [30] Gagan Somashekar, Amoghavarsha Suresh, Saurabh Tyagi, Vikas Dhyani, K Donkada, Anurag Pradhan, and Anshul Gandhi. Reducing the tail latency of microservices applications via optimal configuration tuning. In *2022 IEEE International Conference on Autonomic Computing and Self-Organizing Systems (ACSOS)*, pages 111–120. IEEE, 2022.
- [31] Gagan Somashekar, Karan Tandon, Anush Kini, Chieh-Chun Chang, Petr Husak, Ranjita Bhagwan, Mayukh Das, Anshul Gandhi, and Nagarajan Natarajan. {OPPerTune}:{Post-Deployment} configuration tuning of services made easy. In *21st USENIX Symposium on Networked Systems Design and Implementation (NSDI 24)*, pages 1101–1120, 2024.
- [32] Akshitha Sriraman, Abhishek Dhanotia, and Thomas F Wenisch. Softsku: Optimizing server architectures for microservice diversity@ scale. In *Proceedings of the 46th International Symposium on Computer Architecture*, pages 513–526, 2019.
- [33] Akshitha Sriraman and Thomas F Wenisch. { $\mu$ Tune}:{Auto-Tuned} threading for {OLDI} microservices. In *13th USENIX Symposium on Operating Systems Design and Implementation (OSDI 18)*, pages 177–194, 2018.
- [34] Adith Swaminathan, Akshay Krishnamurthy, Alekh Agarwal, Miro Dudik, John Langford, Damien Jose, and Imed Zitouni. Off-policy evaluation for slate recommendation. *Advances in Neural Information Processing Systems*, 30, 2017.
- [35] Dana Van Aken, Andrew Pavlo, Geoffrey J Gordon, and Bohan Zhang. Automatic database management system tuning through large-scale machine learning. In *Proceedings of the 2017 ACM international conference on management of data*, pages 1009–1024, 2017.
- [36] Patrick Wintermeyer, Maria Apostolaki, Alexander Dietmüller, and Laurent Vanbever. P2go: P4 profile-guided optimizations. In *Proceedings of the 19th ACM Workshop on Hot Topics in Networks*, pages 146–152, 2020.
- [37] Michael Wolfe. Iteration space tiling for memory hierarchies. In *Proceedings of the Third SIAM Conference on Parallel Processing for Scientific Computing*, pages 357–361, 1987.
- [38] Michael Wolfe. More iteration space tiling. In *Proceedings of the 1989 ACM/IEEE conference on Supercomputing*, pages 655–664, 1989.
- [39] Shaofeng Wu, Qiang Su, Zhixiong Niu, and Hong Xu. Tomur: Traffic-aware performance prediction of on-nic network functions with multi-resource contention. *arXiv preprint arXiv:2405.05529*, 2024.
- [40] Arseniy Zaostrovnykh, Solal Pirelli, Rishabh Iyer, Matteo Rizzo, Luis Pedrosa, Katerina Argyraki, and George Candea. Verifying software network functions with no verification expertise. In *Proceedings of the 27th ACM Symposium on Operating Systems Principles*, pages 275–290, 2019.