

Occlum: Secure and Efficient Multitasking Inside a Single Enclave of Intel SGX

Youren Shen*
syr15@mails.tsinghua.edu.cn
Tsinghua University

Hongliang Tian*
tate.thl@antfin.com
Ant Financial Services Group

Yu Chen
yuchen@mail.tsinghua.edu.cn
Tsinghua University
Peng Cheng Laboratory

Kang Chen†
chenkang@tsinghua.edu.cn
Tsinghua University

Runji Wang
wrj15@mails.tsinghua.edu.cn
Tsinghua University
Ant Financial Services Group

Yi Xu‡
xu1369@purdue.edu
Purdue University
Ant Financial Services Group

Yubin Xia
xiayubin@sjtu.edu.cn
Shanghai Jiao Tong University

Shoumeng Yan
shoumeng.ysm@antfin.com
Ant Financial Services Group

Abstract

Intel Software Guard Extensions (SGX) enables user-level code to create private memory regions called *enclaves*, whose code and data are protected by the CPU from software and hardware attacks outside the enclaves. Recent work introduces library operating systems (LibOSes) to SGX so that legacy applications can run inside enclaves with few or even no modifications. As virtually any non-trivial application demands multiple processes, it is essential for LibOSes to support multitasking. However, none of the existing SGX LibOSes support multitasking both *securely* and *efficiently*.

This paper presents Occlum, a system that enables secure and efficient multitasking on SGX. We implement the LibOS processes as *SFI-Isolated Processes (SIPs)*. SFI is a software instrumentation technique for sandboxing untrusted modules (called domains). We design a novel SFI scheme named *MPX-based, Multi-Domain SFI (MMDSFI)* and leverage MMDSFI to enforce the isolation of SIPs. We also design an independent verifier to ensure the security guarantees of MMDSFI. With SIPs safely sharing the single address space of an enclave,

the LibOS can implement multitasking efficiently. The Occlum LibOS outperforms the state-of-the-art SGX LibOS on multitasking-heavy workloads by up to 6,600× on micro-benchmarks and up to 500× on application benchmarks.

CCS Concepts. • Security and privacy → Trusted computing; • Software and its engineering → Multiprocessing / multiprogramming / multitasking.

Keywords. Intel SGX, library OS, multitasking, Software Fault Isolation, Intel MPX

ACM Reference Format:

Youren Shen, Hongliang Tian, Yu Chen, Kang Chen, Runji Wang, Yi Xu, Yubin Xia, and Shoumeng Yan. 2020. Occlum: Secure and Efficient Multitasking Inside a Single Enclave of Intel SGX. In *Proceedings of the Twenty-Fifth International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS '20)*, March 16–20, 2020, Lausanne, Switzerland. ACM, New York, NY, USA, 16 pages. <https://doi.org/10.1145/3373376.3378469>

1 Introduction

Intel Software Guard Extensions (SGX) [39] is a promising trusted execution environment (TEE) technology. It enables user-level code to create private memory regions called *enclaves*, whose code and data are protected by the CPU from software attacks (e.g., a malicious OS) and hardware attacks (e.g., memory bus snooping) outside the enclaves. SGX provides a practical solution to the long-standing problem of secure computation on untrusted platforms such as public clouds. SGX developers who use Intel SGX SDK [38] are required to partition SGX-protected applications into enclave and non-enclave halves. This leads to tremendous effort to refactor legacy code for SGX. Recent work [19, 21, 55] tries to minimize the effort by introducing library operating systems (LibOSes) [30] into enclaves. With a LibOS providing system calls, legacy code can run inside enclaves with few or even no modifications.

*Both authors contributed equally to this research.

†This is the corresponding author.

‡This work was done while the author was at Tsinghua University.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

ASPLOS '20, March 16–20, 2020, Lausanne, Switzerland

© 2020 Association for Computing Machinery.

ACM ISBN 978-1-4503-7102-5/20/03...\$15.00

<https://doi.org/10.1145/3373376.3378469>

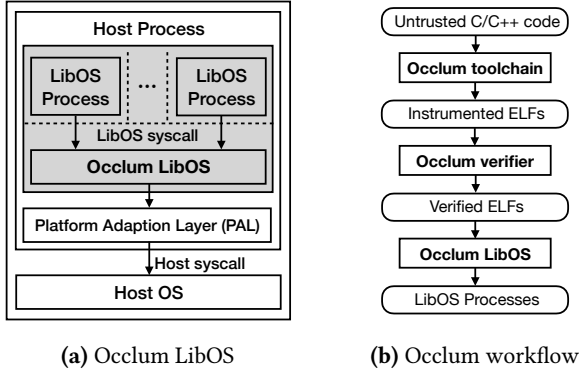


Figure 1. An overview of Occlum system, which consists of three components: the toolchain, the verifier, and the LibOS. The gray area represents an enclave.

One highly desirable feature for LibOSes, or any OSEs in general, is multitasking. Multitasking is important since virtually any non-trivial application demands more than one process. UNIX has long been known of its philosophy on the rule of composition: design programs to be connected with other programs [47]. In the modern era of cloud computing, even a single-purpose, cloud-native application in a container often requires running the main application along with some dependent services (e.g., `sshd` [14], `etcd` [48], and `fluentd` [31]). So multitasking is an indispensable feature.

However, existing SGX LibOSes cannot support multitasking both *securely* and *efficiently*. The most advanced, multitasking SGX LibOS is Graphene-SGX [55], which implements LibOS processes as *Enclave-Isolated Processes (EIPs)*. Each EIP is hosted by one instance of the LibOS inside an enclave; that is, n EIPs require n LibOS instances and n enclaves. The strong enclave-based isolation between EIPs, however, causes performance and usability issues. First, process creation is extremely expensive due to the high cost of enclave creation. Process creation on Graphene-SGX is reported to be nearly $10,000\times$ slower than that on Linux [55]. Second, inter-process communication (IPC) between EIPs is also expensive. EIPs, which are isolated completely by enclave boundaries, have to communicate with each other by transferring encrypted messages through untrusted memory. The encryption and decryption add significant overhead. Third, synchronizing between multiple LibOS instances is painful. The most notable example is the encrypted file system. As there are multiple LibOS instances, the metadata and data of the file system are thereby spread across multiple enclaves. Thus, maintaining a unified view of the file system across EIPs is difficult and inefficient. This explains why Graphene-SGX lacks a *writable*, encrypted file system. For the above reasons, it is not clear how to achieve secure and efficient multitasking in SGX LibOSes—up until now.

In this paper, we present Occlum, a system that enables *secure* and *efficient* multitasking in a LibOS for Intel SGX. Different from prior work on SGX LibOSes, we explore an opportunity of synergy between compiler techniques and LibOS design. Specifically, we propose to implement the LibOS processes as *SFI-Isolated Processes (SIPs)*, which reside alongside the LibOS in the single address space of an enclave (see Figure 1a). Software Fault Isolation (SFI) [56] is a software instrumentation technique for sandboxing untrusted modules (called *domains*). We design a novel SFI scheme named *MPX-based, Multi-Domain SFI (MMDSFI)*, which, compared with existing SFI schemes, is unique in its support of an *unlimited number* of domains without any constraints on their *addresses* and *sizes*. Thus, we can leverage MMDSFI to implement intra-enclave isolation mechanisms for SIPs, including inter-process isolation and process-LibOS isolation.

To ensure the trustworthiness of Occlum’s isolation mechanisms based on MMDSFI, we introduce the Occlum verifier, which is an independent *binary verifier* that takes as input an ELF binary and statically checks whether it is compliant with the security policies of MMDSFI. We describe in depth the design of the verifier and prove its security guarantees mathematically. By introducing the verifier, we exclude Occlum’s MMDSFI-enabled toolchain—which is large and complex—from the Trusted Computing Base (TCB) and rely only on the verifier as well as the LibOS for security.

With MMDSFI implemented by the Occlum toolchain and verified by the Occlum verifier, the Occlum LibOS can safely host multiple SIPs inside a single enclave. Since the SIPs reside inside the same address space, there are new opportunities for sharing between SIPs. As a result, the Occlum LibOS can improve both performance and usability of multitasking, achieving fast process startup, low-cost IPC, and writable encrypted file system.

We have implemented the Occlum system, which consists of the three components shown in Figure 1b: the toolchain, the verifier, and the LibOS. Our prototype implementation comprises over 20,000 lines of source code in total and is available on GitHub [10]. Experimental evaluation on CPU-intensive benchmarks shows that MMDSFI incurs an average of 36% performance overhead. Despite this overhead incurred by MMDSFI, the Occlum LibOS outperforms the state-of-the-art SGX LibOS on multitasking-heavy workloads by up to $6,600\times$ on micro-benchmarks and up to $500\times$ on application benchmarks. Furthermore, the security benchmark (RIPE [59]) shows that MMDSFI prevents all memory attacks that may break the isolation of SIPs. These results demonstrate that our SIP-based approach is secure and efficient.

This paper has four major contributions:

1. We propose *SFI-Isolated Processes (SIPs)* to realize secure and efficient multitasking on SGX LibOSes, which is radically different from the traditional approach (§3);
2. We design a novel SFI scheme named *MPX-based, Multi-Domain SFI (MMDSFI)* to enforce the isolation of SIPs (§4);

3. We describe in depth our *binary verifier* that statically checks whether an ELF binary is compliant with the security policies enforced by MMDSFI (§5) and present a security analysis against two common classes of attacks (§7);

4. We design and implement the *first SIP-based, multitasking LibOS* for Intel SGX (§6 and §8) and demonstrate its performance and security advantages with various benchmarks (§9).

2 Background and Related Work

2.1 Intel Software Guard Extensions (SGX)

The background knowledge about Intel SGX that is relevant to our discussion is summarized below.

Enclave creation. During enclave creation, the untrusted OS loads the code and data to the enclave pages, and then the enclave is marked as *initialized*. From this moment, CPU guarantees the protection of the enclave from any code outside the enclave. While an enclave is being created, its contents are cryptographically hashed to calculate the *measurement*. As this process involves a lot of cryptographic computation, it is expensive to create an enclave.

Enclave dynamic memory management. On SGX 1.0, after an enclave is initialized, enclave pages cannot be added, removed, or modified with their permissions. On SGX 2.0, this restriction has been removed by new SGX instructions. However, Intel has not yet shipped SGX 2.0 CPUs widely. So we implement Occlum on SGX 1.0 to maximize its compatibility.

Intra-enclave isolation. Currently, there is no hardware isolation mechanisms that are capable or suitable to partition an enclave into smaller security domains. Segmentation is disabled inside enclaves. Page tables are untrusted in SGX’s security model. Intel Memory Protection Keys (MPK) [26] is based on page tables. So MPK is also untrusted to SGX. This is why we turn to SFI for intra-enclave isolation.

SGX threads. Multiple SGX threads can execute inside an enclave simultaneously. The execution of an SGX thread may be interrupted by hardware exceptions, causing the CPU core of the SGX thread to exit the enclave. This is called an asynchronous enclave exit (AEX). Upon the occurrence of an AEX, the state of the CPU core is automatically stored in a prespecified, secure memory area called state save area (SSA). And later when the SGX thread is about to resume its execution inside the enclave, the state of the CPU core will be stored according to SSA.

2.2 Library OSes for Intel SGX

We compare Occlum with the state-of-the-art LibOSes for SGX—i.e., Haven [21], Scone [19], Panoply [52], and Graphene-SGX [55]—in three dimensions.

Compatibility. An SGX LibOS may be compatible with legacy applications at either binary or code level. Haven and Graphene-SGX are binary level compatible while Scone,

Panoply and Occlum are source code level compatible via cross compilers.

We argue that *the code-level compatibility is acceptable for most use cases* since it has eliminated most of the efforts required to port applications for SGX. And even a binary-compatible SGX LibOS sometimes demands recompilation so that the legacy source code can be modified to work around the hardware limitations of SGX or the system call limitations of the LibOS. Furthermore, the recompilation of source code provides the opportunity to integrate SGX-specific, compiler-based hardening techniques [42, 45].

Multitasking. Existing SGX LibOSes cannot support multitasking both securely and efficiently. Haven supports multitasking in a single-address-space architecture like Occlum, but it lacks isolation for its LibOS processes. Scone, Graphene-SGX, and Panoply support multitasking with EIPs, but suffer from performance and usability issues (§3.2).

TCB sizes. As Haven reuses the source code of Drawbridge [46], it ends up with a huge TCB. The other four LibOSes including Occlum are written from scratch and thus absent of unnecessary OS functionalities. This results in small TCBs and thus small attack surfaces.

2.3 Intel MPX and SFI

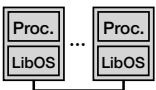
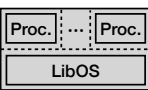
Intel Memory Protection Extensions (MPX) [37] is a set of extensions to the x86 instruction set architecture to provide bound checking at runtime. It provides four bound registers `bnd0 - bnd3`. Each bound register stores a pair of 64-bit values, one for the lower bound and the other for the upper bound. Two instructions, `bndcl` and `bndcu` are introduced to check a given address against a bound register’s lower and upper bound, respectively. If the check fails, an exception will be raised by the CPU. These four bound registers are saved when AEX occurs and are restored when the enclave resumes its execution from AEX [39]. Occlum does not use MPX’s expensive bound tables or bound table-related instructions.

Software Fault Isolation (SFI) [56] is a software instrumentation technique for sandboxing untrusted modules (called domains). Existing SFI schemes have limitations on the number, addresses, or sizes of domains to simplify their designs and minimize the overheads. For example, the well-studied SFI scheme Native Client [50] requires two 40GB unmapped memory regions around a 4GB-size domain. PittSFIeld [44] only supports at most $64 - n$ domains of n -bit address space on x86-64.

Our SFI scheme MMDSFI intends to sandbox a (potentially) large number of processes inside the limited address space of an enclave. MMDSFI aims to put no constraints on the number, addresses, or sizes of domains. We find two advantages of using MPX registers for domain bounds:

1. Bound registers can represent any address or size for a domain.
2. During thread switching, bound registers are automatically saved and stored by CPU. Thus, the maximum number

Table 1. A comparison between SIPs and EIPs.

	EIPs	SIPs
		
Process creation	Expensive	Cheap
IPC	Expensive	Cheap
Shared file systems	Read-only	Writable

of domains does not depend on the number of bound registers, but the address size of an enclave.

Existing MPX-enabled systems [29, 40] use MPX simply for reducing the overhead of SFI. MMDSFI fully leverages MPX’s advantages to achieve the flexibility on the number, addresses, and sizes of domains.

3 SFI-Isolated Processes (SIPs)

In this section, we give an overview of SIPs by first describing its threat model, then highlighting its advantages over traditional EIPs, and finally discussing its feasibility on SGX.

3.1 Threat Model and Security Goals

We assume that attackers can take full control over the hypervisor, host OS, and host applications on the target machine. SIPs are assumed to be benign (at least when loaded by the LibOS), otherwise their binaries would have been rejected by our verifier. However, SIPs may have security vulnerabilities and get compromised by attackers. In short, we assume a powerful attacker who controls both the infrastructure on which an enclave is running and some malicious SIPs inside the enclave.

To defend against such a powerful attacker, Occlum aims to isolate in-enclave SIPs. Specifically, we enforce the following two kinds of isolation: 1) *Inter-process isolation*, which protects an SIP from other SIPs; and 2) *Process-LibOS isolation*, which protects the LibOS itself from any SIP.

We assume our LibOS is implemented correctly, which is the (runtime) TCB of Occlum. Iago attacks [24] are not considered, which can be addressed by orthogonal work like Seg0 [43]. Like other SGX LibOSes, Occlum does not hide file access patterns. We do not consider denial-of-service attacks.

We do not consider side-channel attacks in this paper. Side-channel attacks have been shown to be a real threat, especially to Intel SGX [22, 25, 58]. This field is moving fast: new attacks have been kept being proposed, so have new defenses [45, 51]. We believe eventually these efforts that are independent from ours will be able to provide adequate defense against side-channel attacks. Also, covert-channel attacks are out of the scope of this paper.

3.2 Advantages of SIP

SIPs have performance and usability advantages compared to the traditional EIPs. A comparison between SIPs and EIPs is summarized in Table 1 and explained below.

Process creation. Creating a new EIP requires three steps: (1) creating a new enclave, (2) doing local attestation with other enclaves, and (3) duplicating the process state over an encrypted stream. In contrast, creating an SIP does not involve any of the three steps. Thus, SIPs are significantly cheaper to create than EIPs.

IPC. IPCs between EIPs are typically implemented by sending and receiving encrypted data through untrusted buffers outside the enclave. Yet, IPC between SIPs is simply copying data from one SIP to another, with no encryption involved.

Shared file system. In the traditional EIP-based approach, there are multiple instances of the LibOSes communicating via secure communication channels. This setup must face the challenge of data synchronization. To avoid this difficulty, traditional LibOSes like Graphene-SGX only support a read-only, encrypted file system. In contrast, all SIPs inside an enclave share the same instance of the LibOS. Thus, a writable, encrypted file system can be implemented relatively easily.

3.3 Spawn Instead of Fork

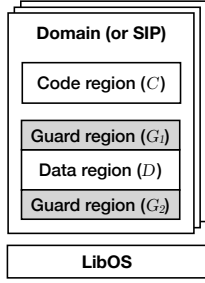
fork system call requires a forked child process to have the same address space as its parent process. This semantic is incompatible with single-address-space OSes [13], including Occlum, where all processes reside in a single address space. So, in Occlum, processes are created with the spawn system call, which is similar to libc’s `posix_spawn` and Solaris’s `spawn` [28].

Most uses of fork and exec can be readily replaced with spawn. As for those applications such as Apache, PostgreSQL, and Redis that use fork alone, the existence of their Windows ports or versions [3, 11, 15] implies that fork-based code can be rewritten to use spawn-like APIs. Our experience in porting multi-process applications for Occlum also confirms that replacing fork with spawn is not as hard as one might expect (§9.1).

Putting the compatibility issues aside, spawn is considered superior to fork. A recent paper [20] reflects on the pros and cons of fork and concludes that “we should acknowledge that fork’s continued existence as a first-class OS primitive holds back systems research, and deprecate it.

4 MPX-Based, Multi-Domain SFI (MMDSFI)

Occlum uses a novel SFI scheme named MPX-based, Multi-Domain SFI (MMDSFI) to enforce the isolation of SIPs.



(a) The memory layout

Pseudo-instructions	Equivalent Machine Instruction Sequences
mem_guard <mem>	bndcl <mem>, %bnd0 bndcu <mem>, %bnd0
cfi_label	nopl <domain_id>(%rbx,%rbx,1)
cfi_guard <reg>	movq (<reg>), <scratch_reg> bndcl <scratch_reg>, %bnd1 bndcu <scratch_reg>, %bnd1

(b) The pseudo-instructions

Before Instrumentation	After Instrumentation
movq %rdi, 8(%rcx)	mem_guard 8(%rcx) movq %rdi, 8(%rcx)
jmp %rdx	cfi_guard %rdx jmp %rdx
call %rax	cfi_guard %rax call %rax cfi_label
ret	popq %r10 cfi_guard %r10 jmp %r10

(c) Instrumentation examples

Figure 2. The design of MMDSFI. All assembly code follows the AT&T syntax. The bound register `bnd0` and `bnd1` are initialized by the LibOS to the range $[D.begin, D.end)$ and $[cfi_label, cfi_label]$, respectively. The operand `<reg>` denotes a general-purpose register; `<scratch_reg>`, a scratch register for storing intermediate values; `<mem>`, a memory address; `<domain_id>`, the ID of the domain that the code belongs to.

4.1 Overview

MMDSFI requires a memory layout shown in Figure 2a. A domain of MMDSFI has two main memory regions: the region C for the user code, and the region D for the user data. Region C and D are mapped to enclave pages with RWX and RW permissions, respectively. Region D is surrounded by two guard regions G_1 and G_2 , which are of the same size and not mapped to any enclave pages, thus triggering exceptions when accessed. These guard regions, as a common technique of SFI, are introduced to simplify the instrumentation (§4.2) and facilitate some optimizations (§4.3). Multiple domains can coexist in a single address space, and the memory ranges of these domains are exclusive to each other. All domains share a single instance of the LibOS, which is in charge of managing the domains and processing system calls from them.

As an SFI scheme, the security goal of MMDSFI is to *sandbox untrusted user code*. The untrusted code is first compiled—with the instrumentation required by MMDSFI—into a binary. Then, the binary is loaded by the LibOS into a domain: the code of the binary is loaded into the region C and the data into the region D . This domain is then run as an SIP (so actually *the two terms of domain and SIP can be used interchangeably for our discussion*). Together with the compile-time instrumentation of MMDSFI and the runtime support of the LibOS, we can guarantee that the untrusted user code is sandboxed. Specifically, the following two security policies are enforced:

The memory access policy. For any memory access instruction I in C , I must access the memory within range $[D.begin, D.end)$, where $D.begin$ and $D.end$ are addresses where the region D begins and ends.

The control transfer policy. For any control transfer instruction I in C , I must target an address within $[C.begin, C.end)$, where $C.begin$ and $C.end$ are the addresses where the region C begins and ends.

If an instruction violates the security policies above, then the instruction is *invalid*. Our goal is to guarantee that any *invalid* instruction is either *prevented* or *detected*.

4.2 Instrumentation

MMDSFI instruments the untrusted code by inserting checks for or rewriting unsafe instructions to enforce the aforementioned security policies.

Confining memory accesses. To this end, we introduce `mem_guard` pseudo-instruction, which takes a memory operand `<mem>` and checks whether `<mem>` is within the data range $[D.begin, D.end)$ of the current domain. If the check passes, it does nothing; otherwise, it triggers an exception, which can then be captured by the LibOS. As shown in Figure 2b, `mem_guard` pseudo-instruction can be implemented efficiently with two MPX bound check instructions. As the first example shown in Figure 2c, MMDSFI inserts `mem_guard` before every unsafe memory access instruction so that its memory accesses are confined within $[D.begin, D.end)$.

Confining control transfers. Before describing how MMDSFI confines control transfers, let’s take a few relevant observations. First, since x86 instructions are variable-length, a faulty control transfer could jump into the middle of (pseudo-)instructions and execute unpredictable instructions. This would completely jeopardize the validity of SFI and thus must not be allowed. Second, some instruction sequences must be treated as a whole. One such example is `mem_guard` and its guarded memory access instruction, which must be executed as a whole. Otherwise, jumping in between the two instructions would skip the `mem_guard`, thereby bypassing the memory access confinement. So, this, too, must not be allowed. Third, *indirect* control transfer instructions need runtime checks while *direct* control transfer instructions do not. A direct control transfer instruction (e.g., `jmp <imm>`, where `<imm>` denotes an immediate value) has its target address hard-coded in the instruction and thus can be verified

at compile time. Yet, an indirect control transfer instruction (e.g., `jmp <reg>`) gives its target address in a register or a memory location, which cannot be determined until runtime. Thus, indirect control transfer instructions need runtime checks.

Based on the above observations, we enforce a coarse-grained control-flow integrity (CFI) [56] in MMDSFI by introducing a pair of pseudo-instructions: `cfi_label` and `cfi_guard`, as shown in Figure 2b.

`cfi_label` has the following interesting properties:

- (1) No operation. It has no visible impact on CPU;
- (2) Alignment. The encoding has a fixed length of 8 bytes;
- (3) Nonexistence. The first 4 bytes of the encoding does not appear anywhere in the uninstrumented code, including both the untrusted user code and the trusted LibOS code;
- (4) Uniqueness. The last 4 bytes of the encoding is the unique ID of the domain where the `cfi_label` resides.

To satisfy these properties, we implement `cfi_label` with a special 8-byte nop instruction, as shown in Figure 2b. This nop is not in its most common form, so it is not supposed to be emitted by any off-the-shelf compiler. The nop is encoded in 8 bytes with the last 4 bytes being any 32-bit value of our choice. So, when loading instrumented code into a domain, the LibOS rewrites the last 4 bytes of all `cfi_labels` to the value of the ID of the current domain. This is to satisfy the uniqueness property required above.

With the properties above, we can use `cfi_labels` to “label” the valid targets of indirect control transfers. Specifically, MMDSFI inserts a `cfi_label` at every valid target address of indirect control transfers (e.g., the start address of a function). And only at these addresses can `cfi_labels` be found. Thus, the inserted `cfi_labels` can be used at runtime to verify indirect control transfers by the pseudo-instruction that is to be introduced below.

`cfi_guard` pseudo-instruction checks whether the address given in a register operand `<reg>` is a valid target of indirect control transfers. If it is a valid target, `cfi_guard` does nothing; otherwise, it raises an exception. The validity of a target address can be easily determined with `cfi_labels`. As shown in Figure 2b, `cfi_guard` is implemented in a sequence of three instructions: the first `mov` loads the value at the target address into a scratch register `<scratch_reg>`, then two following bound checks compare the value in `<scratch_reg>` with `bnd1`, which has been set by the LibOS to the range `[cfi_label, cfi_label]`. So the two bound checks are essentially a test for equality. If and only if the value at the target address equals to `cfi_label` can the target be valid. With `cfi_guards`, we now can confine indirect control transfers to target only valid addresses. See figure 2c for some typical examples of instrumentation.

Using the three pseudo-instructions, we can instrument any unsafe memory access or control transfer instruction as shown in figure 2c, thus enforcing the security policies of

Categories	All Instructions	How to Verify?
<i>Direct</i> control transfer	<code>jmp *<rel></code> <code>jcc *<rel></code> <code>loop *<rel></code> <code>loopcc *<rel></code> <code>call *<rel></code>	Check that the target specified by <rel> is not a register-based indirect control transfer instruction
<i>Register-based</i> indirect control transfer	<code>jmp *<reg></code> <code>call *<reg></code>	Check that the instruction is guarded by a <code>cfi_guard</code>
<i>Memory-based</i> indirect control transfer	<code>jmp *<mem></code> <code>call *<mem></code>	Reject
<i>Return-based</i> indirect control transfer	<code>ret</code> <code>ret <imm></code>	Reject

Figure 3. Stage 3 verifies control transfer instructions by classifying them into four categories.

MMDSFI. The correctness of MMDSFI is discussed formally in §5.

4.3 Optimizations

A naive implementation of MMDSFI instrumentation has to insert a `mem_guard` for every memory accesses. This can easily lead to an intolerable performance penalty. However, thanks to our enforcement of CFI, we can employ standard dataflow analysis [17] on the control-flow graph (CFG) to seize optimization opportunities enabled by the guard regions (§4.1). Specifically, we can use *range analysis* to track the ranges of possible values in registers before and after each instruction, and perform the following two optimizations [60]:

(1) *Redundant check elimination.* Recall that the data region D is surrounded by the guard regions G_1 and G_2 that will trigger exceptions if accessed. So, if a memory address x has been confined within the data range $[D.begin, D.end)$, then any memory address determined by the range analysis to be within $[x - G_i.size, x + G_i.size]$ is also safe and requires no check via `mem_guard`, where $G_i.size$ is the size of a guard region.

(2) *Loop check hoisting.* In a loop, if there is a `mem_guard` whose memory operand `<mem>` is increased or decreased by a constant value smaller than $G_i.size$ per loop iteration, then a new `mem_guard` can be inserted before the loop and the original, in-loop `mem_guard` can be eliminated due to the same reason in the first optimization. The net result is that the `mem_guard` will be executed only once per loop instead of once per loop iteration.

While more optimizations based on range analysis are possible, we found these two optimizations are sufficient to reduce the overhead to an acceptable level (§9.3).

5 Binary Verification

To ensure the trustworthiness of Occlum’s security properties based on MMDSFI, we introduce the Occlum verifier, which is an independent binary verifier that takes an ELF

Input: C , the code segment of an ELF binary
Output: R , the set of all reachable instructions in C

```

1  $R \leftarrow \{\}$ ;
2  $S \leftarrow$  Get all the addresses of cfi_labels by scanning  $C$  byte by byte;
3 while  $S \neq \{\}$  do
4    $addr \leftarrow$  Pop an item from  $S$ ;
5   while true do
6     if  $addr$  is not within  $C$  then
7        $\mid$  abort;
8      $instr \leftarrow$  Disassemble the instruction at  $addr$ ;
9     if  $instr$  is invalid then
10       $\mid$  abort;
11     else if  $instr \in R$  then
12       $\mid$  break;
13     else if  $instr$  overlaps with any  $i \in R$  then
14       $\mid$  abort;
15      $R \leftarrow R \cup \{instr\}$ ;
16     if  $instr$  is a direct control transfer then
17        $target\_addr \leftarrow$  Calculate the target of  $instr$ ;
18        $S \leftarrow S \cup \{target\_addr\}$ ;
19     if  $instr$  is an unconditional control transfer then
20        $\mid$  break;
21      $addr \leftarrow addr + instr.length$ ;

```

Algorithm 1: Stage 1 disassembles the binary completely and reliably, generating R the set of all reachable instructions.

Categories	Sample Instructions	How to Verify?
Scale-Index-Base (SIB)	<code>addq \$1, -8(rax,rbx,4)</code>	Check that the destination is within $[D.begin - G_i.size, D.end + G_i.size)$ according to the range analysis
Implicit register-based	<code>push \$1</code>	Same as above
RIP-relative	<code>addq %rcx, 1234(%rip)</code>	Same as above
Direct memory offset	<code>movq %rax, \$1122334455667788</code>	Reject
Vector SIB	<code>vscatterdps %zmm7, (%r12,%zmm6,4){%k3}</code>	Reject

Figure 4. Stage 4 verifies memory access instructions by classifying them into five categories. Direct memory offset, which can hard-code a 64-bit memory address in `mov`, is rejected because no fixed addresses can be assumed to be within a domain. Vector SIB is rejected because it allows one instruction to access to multiple non-contiguous memory locations.

binary as input and statically checks whether it is compliant with the security policies of MMDSFI.

By introducing the Occlum verifier, we exclude the Occlum toolchain from the TCB. This is desirable for two reasons. First, the implementation of an SFI, including MMDSFI, is error-prone since it involves dealing with the low-level details of machine code (e.g., x86 has approximately 1000 distinct instructions) and it is usually implemented upon

the huge codebase of a compiler (e.g., GCC and LLVM have millions of lines of source code). Second, even if the Occlum toolchain can be proven flawless, the Occlum verifier can still be very useful. In some interesting use cases [36, 54], the user may want to use binaries from untrusted origins. So these binaries can be generated by some arbitrary toolchains; and the Occlum toolchain being flawless is irrelevant. Thus, we must verify the binaries before loading them into enclaves.

The verifier consists of four stages: (1) complete disassembly, (2) instruction set verification, (3) control transfer verification, and (4) memory access verification. If and only if all the four stages pass, can the input binary be claimed to be compliant with MMDSFI and safe to run upon the Occlum LibOS inside an enclave.

Stage 1 - Complete disassembly. This stage aims to find out every instruction *reachable*, i.e., any instruction that may be executed as part of the untrusted user program. To this end, it disassembles the input ELF binary and generates the set of all reachable instructions, which is denoted as R . Our x86 disassembler is extended with the three pseudo-instructions introduced in the last section. So an instruction in R can be either a real x86 instruction or a pseudo-instruction, which is not distinguished for our discussion. Every instruction in R is uniquely identified by its memory address. This resulting set R is the subject of the verification to be carried out in the remaining four stages.

Disassembling an *arbitrary* binary *statically, reliably, and completely* is impossible; existing disassembly algorithms are based on heuristics, providing only best-effort results [18, 41, 49]. However, thanks to the introduction of `cfi_label`, our disassembly algorithm can give the *completely accurate* results as long as the input binary can eventually pass all the remaining stages of the verifier. (And for those binaries that are to be rejected by the verifier, it does not really matter if the disassembly results are accurate or not.)

The disassembly algorithm is shown in figure 1. Firstly, the disassembler gets all `cfi_labels` by scanning the code section of the binary byte by byte (line 2). Guaranteed by the LibOS, the user program must start its execution from `cfi_labels`. From each of these starting points (line 4), the disassembly algorithm follows the sequential execution of the program (line 21) and every direct control transfer (line 17). Note that the *indirect* control transfers can only target at `cfi_labels` (due to the control transfer policy to be verified in Stage 3), which are the starting points themselves. So the branches of indirect control transfers have been implicitly covered. As we have started from every possible starting point and followed every possible branch, we end up with R containing all reachable instructions.

Stage 2 - Instruction set verification. The goal of this stage is to make sure that R does not include any dangerous instructions that can perform privileged tasks that are meant for the LibOS. The dangerous instructions can be classified into the following three categories:

(1) SGX instructions, e.g., `eexit`, which exits the current enclave, and `emodpe/eaccept`, which extends/restricts the memory permissions of enclave pages at runtime;

(2) MPX instructions, e.g., `bndmk` and `bndmov`, which modify the values of MPX bound registers;

(3) Miscellaneous instructions, e.g., `xrstor`, which can enable or disable some CPU features including MPX, and `wrfsbase/wrgsbase`, which modifies FS/GS segment bases.

Note that these dangerous instructions are not supposed to be used in normal user programs. So forbidding them does not restrict the functionality of user programs. The implementation of this stage is quite straightforward: simply scan R to check that no such dangerous instructions exist.

Stage 3 - Control transfer verification. This stage verifies every control transfer instructions in R by classifying it into one of four categories and checking it with the criteria of its belonging category as shown in Figure 3.

If a binary passes Stage 1-3 of the verification and is loaded into a domain by the LibOS, then we can prove the lemma and theorem below.

Lemma 5.1. Any register-based indirect control transfer in the domain can only jump to the `cfi_labels` in the domain.

Proof. To prove it by contradiction, assume there exists an occurrence of register-based indirect control transfer that violates the lemma.

Let $S = \{I_1, \dots, I_n\}$ be the CPU execution trace of the domain up until the violation occurs, where I_1 is a `cfi_label`, I_n is the violating register-based indirect control transfer, and $n \geq 2$. Without loss of generality, we can assume that I_n is the first violation in the trace; that is, no $1 < j < n$ such that I_j is also an occurrence of register-based indirect control transfer that violates the lemma. Otherwise, we can simply restart our analysis on I_j instead of I_n .

Our verification in Stage 3 has checked that I_n is properly guarded with a `cfi_guard`, yet still I_n jumps to a location other than some `cfi_label` in the domain. This implies that I_{n-1} cannot be the `cfi_guard` and must be a control transfer instruction that jumps to I_n . Yet, such I_{n-1} is impossible according to the exhaustive case analysis below:

(1) I_{n-1} cannot be a *direct* control transfer since our verification in Stage 3 has checked that any direct control transfer does not jump to a register-based indirect control transfer instruction;

(2) I_{n-1} cannot be a *register-based* indirect control transfer otherwise it would be a violation that happened before I_n , which contradicts our assumption that I_n is the first;

(3) I_{n-1} cannot be a *memory-based* indirect control transfer as our verification in Stage 3 rejects it;

(4) I_{n-1} cannot be a *return-based* indirect control transfer as our verification in Stage 3 rejects it.

By contradiction, we prove the lemma. \square

Theorem 5.2. Any control transfer in the domain is compliant with the control transfer policy of MMDSFI (§4.1).

Proof. For any control transfer instruction in the domain, it must belong to one of the four categories listed in Figure 3. The first category is covered implicitly by the complete disassembly in Stage 1; the second category is implied by Lemma 5.1; the third and the fourth categories are forbidden by the verification in Stage 3. Thus, the theorem stands in all cases. \square

Stage 4 - Memory access verification. With the CFI of R verified in Stage 3, this stage first builds the CFG for R , then leverages this CFG to do the `cfi_label`-aware range analysis described in §4.3, and finally verifies every memory access instruction in R by classifying it into one of five categories and checking it with the simple and straightforward criteria of its belonging category as shown in Figure 4.

If a binary passes Stage 1-4 of the verification and is loaded into a domain by the LibOS, then we can prove the theorem below.

Theorem 5.3. Any memory access instruction in the domain is compliant with the memory access policy of MMDSFI (§4.1).

Proof. Any memory access instruction, regardless of how its destination is specified, has been either verified with the `cfi_label`-aware range analysis or rejected by our verification in Stage 4. Thus, we prove the theorem. \square

6 Library OS

In this section, we give an overview of the Occlum LibOS (see Figure 1a) with an emphasis on its unique aspects.

ELF loader. A typical program loader in a Unix-like OS performs tasks such as parsing program binaries (e.g., ELF), copying program images, and initializing CPU states. Beyond these basic tasks, the program loader in Occlum has four extra responsibilities. First, the loader checks that the ELF binaries to be loaded are verified and signed by the Occlum verifier. Second, the loader rewrites all `cfi_labels` in the program image so that the last four bytes of the `cfi_labels` are set to the ID of the domain associated with the new SIP. Third, the loader inserts into the process image a small piece of trampoline code that jumps to the entry point of the LibOS, thereby enabling the LibOS system calls. This trampoline code is the only way out of the sandbox enforced by MMDSFI. The address of this trampoline code is passed to `libc` via the auxiliary vector [34]. Fourth, the loader initializes MPX bound registers according to the memory layout of the SIP's domain.

syscall interface. LibOS system calls are just function calls, except that the users must go through the trampoline code inserted by the ELF loader. With SIPs sandboxed by MMDSFI, the only chance for faulty SIPs to compromise the LibOS is through system calls. So, at the entry point of the LibOS is a piece of carefully-written assembly, which performs sanity checks, switches between user/LibOS stack,

switches between user/LibOS thread-local storage, and eventually calls system call dispatching routine. After system call finishes, before returning to the SIP, LibOS will ensure that the return address target is a `cfi_label` of corresponding SIP.

Memory management. Occlum preallocates the available enclave pages of an MMDSFI domain during enclave initialization. For each domain, the enclave pages are allocated and set with proper permissions according to the memory layout described in §4.1. The maximum size of the code region is prespecified at compile time; so is the data region. The size of guard regions is set to 4KB. And the linker of the Occlum toolchain is aware of this 4KB gap between the code segment and the data segment when generating an ELF binary for Occlum. The maximum number of MMDSFI domains is also prespecified at compile time. Note that this preallocation of enclave pages is intended to work around the limitation of SGX 1.0 and can be avoided on SGX 2.0.

Memory mapping syscalls, e.g., `mmap` and `munmap`, can only manipulate the memory in the data region of a domain. Anonymous memory mappings are fully supported, but requires the LibOS to manually initialize the allocated memory pages to zeros. File-backed memory mappings is implemented by copying the file content to the mapped area. Shared memory mappings are impossible because SGX cannot map one EPC page to multiple virtual addresses. For the sake of security, the users are not allowed to add or remove the `X` permission of enclave pages (even on SGX 2.0).

Process management. Occlum provides the `spawn` system call, instead of `fork`, to create SIPs. SIPs are mapped one-on-one to SGX threads, which are scheduled transparently by the host OS. This frees the LibOS from the extra complexity of process scheduling. IPC between SIPs, e.g., signals, pipes, and Unix domain sockets, are implemented efficiently via shared data structures in the LibOS.

LibOS threads are treated as SIPs that happen to share resources such as virtual memory, file tables, signal handlers, etc. The synchronization between LibOS threads, e.g., `futex`, eventually relies on the host OS to sleep or wake up the corresponding SGX threads; but the semantic correctness of the synchronization primitives only relies on the LibOS.

File systems. To protect the confidentiality and integrity of persistent data, Occlum provides an encrypted file system that transparently encrypts all file data, metadata, and directories. While Intel SGX SDK contains a library named Intel SGX Protected File System [1], this library can only protect the content of individual files, not file metadata and file directories. We build upon the primitives provided by this library to implement a complete encrypted file system. Occlum also provides special file systems (e.g., `/dev/`, `/proc/`, etc.), which are completely implemented by the LibOS inside the enclave. All SIPs share a common cache for file I/O. A child SIP can inherit the open file table of its parent SIP with

minimal overhead. All SIPs have a unified view of the LibOS file systems.

Networking. Network-related operations are mostly delegated to the host OS, and the LibOS is only responsible for redirecting, bookkeeping, and sanity checks. So network I/O is not secure by default and requires user-level network encryption such as TLS for secure communication.

7 Security Analysis

With all three components of Occlum described, we can now treat the Occlum system as a whole and examine whether we have achieved our overall security goal of inter-process isolation and process-LibOS isolation. More specifically, we will consider whether two common classes of attacks—code injection and return-oriented programming (ROP)—could be used by a malicious SIP to penetrate the isolation enforced by MMDSFI and the LibOS.

Code injection attacks. SGX 1.0 does not allow enclave pages to be added, removed, or modified at runtime. To load programs into an enclave dynamically, the implementation of SGX LibOSes has to reserve a pool of enclave pages with `RWX` permissions when the enclave is launched. So the enclave is made susceptible to code injection attacks. This is a common pitfall of all existing SGX LibOSes.

Fortunately, Occlum is immune to code injection attacks. This is because 1) a SIP can only write to the data region of its domain, whose enclave pages have no executable permission (see §4.1); 2) only the LibOS can modify the executable enclave pages, e.g., when loading new verified binaries; and 3) the LibOS ensures system calls (e.g., `mmap` and `mprotect`) cannot be abused by SIPs. In short, a malicious SIP cannot inject arbitrary code to bypass the isolation mechanism.

ROP attacks. Now that a malicious SIP cannot inject new code, it can still attempt to reuse existing code gadgets for ROP attacks. ROP attacks are hard to prevent: it has been shown that static CFIs, including the coarse-grained CFI integrated with MMDSFI, cannot fully prevent ROP attacks [23, 35].

But this does not change the fact that a malicious SIP cannot break the isolation enforced by MMDSFI and the LibOS. First of all, MMDSFI does make ROP attacks harder. The coarse-grained CFI enforced by MMDSFI prevents the untrusted code in a domain to jump to arbitrary locations; in fact, it can only jump to the `cfi_labels` inside the domain. This greatly reduces the number of useful gadgets and thus the odds of success for ROP attacks. More importantly, any combination of code gadgets by ROP attacks has already been covered by the verifier (see §5). So even ROP attacks cannot violate our security policies and break the isolation.

8 Implementation

We have implemented the Occlum system, which consists of the toolchain, the verifier, and the LibOS. Our prototype

implementation comprises over 20,000 lines of source code in total (15,000 lines of code in Rust for LibOS, 3,000 lines of code in C++ for toolchain, and 2000 lines of code in Python for the verifier). The LibOS and the toolchain has been made open source on GitHub [10]. We describe briefly about how the three components are implemented.

The Occlum toolchain. The toolchain is based on LLVM 7.0 [8]. In LLVM’s x86 backend, we add two extra passes: one instruments control transfer instructions, and the other instruments memory access instructions and performs the range analysis-based optimizations.

We modify LLD so that the linker generates ELF executables that are compatible with MMDSFI. Specifically, the modified linker ensures that the generated code segments only contains code (no read-only data) and that a 4KB-gap between the code and data segments is reserved for the guard region between the code and data regions.

We modify musl libc [9] to use the system calls provided by the Occlum LibOS. The `posix_spawn` API is rewritten to use Occlum’s `spawn` instead of `vfork` and `execve`.

The Occlum verifier. The implementation of the verifier depends on two libraries. In order to decode individual x86-64 instructions correctly, we import Zydis [16], a disassembler library supporting all x86-64 instructions, including all the SGX and MPX instructions we need. In addition, we import PyVEX [53], a library that converts x86-64 instructions into its VEX intermediate representation (IR). This reduces the task of understanding the complex semantics of x86-64 instructions into manipulating the simple VEX IR instructions.

The Occlum LibOS. The LibOS is mostly written in Rust [12], a memory-safe programming language. Thanks to Rust, we can minimize the odds of low-level, memory-safety bugs in the LibOS. The LibOS is based on Intel SGX SDK [38] and Rust SGX SDK [57] for the in-enclave runtime support of C and Rust, respectively. In addition, we implement a FUSE [7]-based utility to mount and manipulate Occlum’s encrypted file system in the host environment. This simplifies the task of preparing encrypted FS images for the Occlum LibOS in development environments. Currently, the LibOS only supports loading statically-linked ELF executables; we will add support for shared libraries in the future.

9 Evaluation

In this section, we present the experimental results that answer the following questions: (1) To what extent can Occlum improve the overall performance of real-world, multi-process applications? (§9.1) (2) What about the performance of individual system calls? (§9.2) (3) Does the use of MMDSFI make user applications more secure? And how much overheads does MMDSFI incur?(§9.3)

To answer the above questions, we compare Occlum with Linux and the state-of-the-art SGX LibOS, Graphene-SGX.

The time of process creation on Graphene-SGX is sensitive to the sizes of the enclaves. For any benchmark whose result is affected by the time of process creation, we configure Graphene-SGX to use the minimal enclave size that is able to run the benchmark.

Experimental setup. We use machines with a two-core, 3.5GHz Intel Core i7 CPU (hyper-threading disabled), 32GB memory, and a 1TB SSD disk, and 1Gbps Ethernet card. Each machine supports SGX 1.0. The host runs Linux kernel 4.15. We install Intel SGX SDK v2.4.0, Rust nightly-2019-01-28, and Rust SGX SDK v1.0.6. We use the latest version of Graphene-SGX (commit f30f7e7) at the time of experiment.

9.1 Application benchmarks

We measure the performance of three widely-used applications: **Fish** (v3.0.0), a user-friendly command line shell [4]; **GCC** (v4.4.5), the GNU compiler for C-family programming languages [33]; and **Lighttpd** (v1.4.40), a fast Web server [6]. All of them are multi-process applications but have distinct workload characteristics: Fish is *process*-intensive (as every shell command is executed in a separate process), GCC is *CPU*-intensive, and Lighttpd is *I/O*-intensive.

Fish. We use the shell script provided by UnixBench [2]. The test script applies on the data from an input file a series of transformation using multiple utilities (e.g., `sort`, `od`, and `grep`), which are connected via pipes or I/O redirections. Fish does not need any code modification to run on Occlum since it uses the modern `posix_spawn`, rather than `fork`, to create processes. As shown in Figure 5a, Occlum (19.5ms) is 13.9× slower than Linux (1.4ms), but nearly 500× faster than Graphene-SGX (9.5s). Occlum is slower than Linux due to the lack of on-demand loading in enclave (see §9.2).

GCC. We use three C files of varied sizes: the first one is a “Hello World!” program (5 LOC), the other two are real-world programs (5K and 50K LOC, respectively) collected by MIT [5]. While GCC is CPU-intensive in nature, it needs to create separate processes for its preprocessor, compiler, assembler, and linker. To start these child processes, GCC originally uses `fork`; we refactor the source code to use `posix_spawn` instead, which is about 50 LOC. As shown in Figure 5b, the compilation time of the three C files ranges from 25ms to 830ms on Linux, from 9.7s to 11.7s on Graphene-SGX, and from 229ms to 3.0s on Occlum. In other words, Occlum is $3.6 \times -9.2 \times$ slower than Linux, but $3.82 \times -42 \times$ faster than Graphene-SGX.

Lighttpd. We use ApacheBench [32] on a client machine to generate HTTP requests that retrieve 10KB HTML pages from an instance of Lighttpd Web server running on Linux, Graphene-SGX, or Occlum. Both the client and server are in the same local area network. We gradually increase the concurrency of ApacheBench to simulate an increasing number of concurrent clients. We configure the master process of Lighttpd to start two worker processes, both of which inherit the listening sockets from the master and handle the

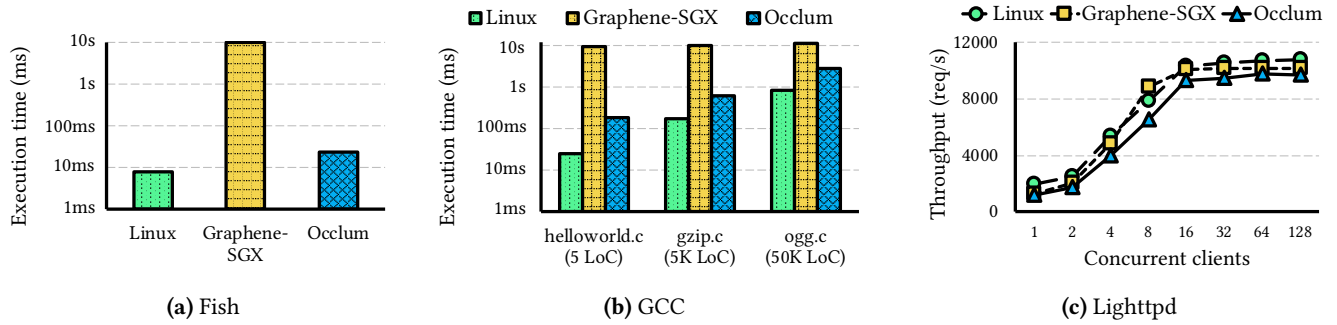


Figure 5. Application benchmarks.

connection requests to the listening sockets together. Similar to GCC, we refactor the source code to replace fork with `posix_spawn`, in about 150 LOC. The results of Lighttpd are shown in Figure 5c. The peak throughput of both Graphene-SGX (10% overhead) and Occlum (9% overhead) is slightly lower than that of Linux.

The results above show that Occlum can deliver a significant performance boost to some multi-process applications.

9.2 System call benchmarks

To better understand the results of the application benchmarks above, we measure the performance of some individual system calls. This helps us identify the sources of performance gain or loss due to the Occlum LibOS.

Process creation time. We measure the process creation time of different binaries on Linux, Graphene-SGX, and Occlum. There are three binaries of different sizes: 1) a “Hello World” program of 14KB size, 2) busybox, a collection of common UNIX utilities combined into a single small executable of 400KB size, and 3) cc1, the GCC front-end compiler of 14MB size. Process are created using `libc`’s `posix_spawn`. Occlum implements `posix_spawn` using its `spawn` system call. Linux and Graphene-SGX use `vfork` and `execve` to implement `posix_spawn`. `vfork` is more efficient than `fork` as it avoids copying the page tables on Linux and copying the process state on Graphene-SGX. As shown in 6a, Linux consumes about 170us regardless of the specific binary. For the small-sized binary (14KB), Occlum consumes 97us, which is 1.6× faster than Linux and over 6,600× faster than Graphene-SGX (0.64s). For the median-sized binary, Occlum consumes 1.7ms, which is 10.5× slower than Linux, but over 390× faster than Graphene-SGX (0.69s). For the large-sized binary, Occlum consumes 63ms, which is 13× faster than Graphene-SGX (0.89s). Linux initializes a process’s page table without actually loading all the pages with their data from disks. Thus, the process creation time on Linux is insensitive to binary sizes. Occlum has to load the entire binaries into the enclave. So the process creation time on Occlum is proportional to the sizes of binaries. Graphene-SGX has to create

a new enclave for each new process, which is very time consuming.

IPC throughput. We measure the throughput of IPC by using a common IPC method, pipe. We create two processes that are connected via a pipe. The throughput of the pipe is measured under varied buffer sizes. As shown in Figure 6b, Occlum’s throughput is on par with Linux’s throughput, which is over 3× higher than Graphene-SGX’s throughput.

File I/O throughput. To quantify the performance overhead of Occlum’s transparent file encryption, we compare Occlum’s encrypted file system with Linux’s Ext4. And we exclude Graphene-SGX from this benchmark since Graphene-SGX does not have a fully-fledged encrypted file system. The throughput of sequential file reads and writes under different buffer sizes are shown in Figure 6c and 6d. Compared with Ext4, Occlum incurs an average overhead of 39% on file reads and an average overhead of 18% on file writes.

The above results show that the Occlum LibOS greatly improves the performance of process startup and IPC, while moderately degrades the performance of file I/O in exchange for the transparent encryption of persistent data.

9.3 MMDSFI benchmarks

Security. To demonstrate the effectiveness of MMDSFI, we use RIPE [59] benchmark, which is designed to compare different defenses against attacks that exploit buffer-overflow bugs. RIPE builds up a total of 850 workable attacks, including code injection, ROP, return-to-libc, etc. We run RIPE on Graphene-SGX and Occlum for comparison.

When the stack protection is disabled by the compiler, 36 code injection, 2 ROP, and 16 return-to-libc attacks succeed on Graphene-SGX. When the stack protection is enabled, the successful attacks are reduced to 16 code injection and 12 return-to-libc. With or without the stack protection, Occlum can prevent all code injection and ROP attacks in RIPE. ROP attacks in RIPE can be fully prevented by MMDSFI as the ROP gadgets in RIPE do not start with `cfi_labels`. Return-to-libc attacks still succeed (16 without stack protection and 12 with stack protection) since `libc` functions are normal functions starting with `cfi_labels`. These results show that

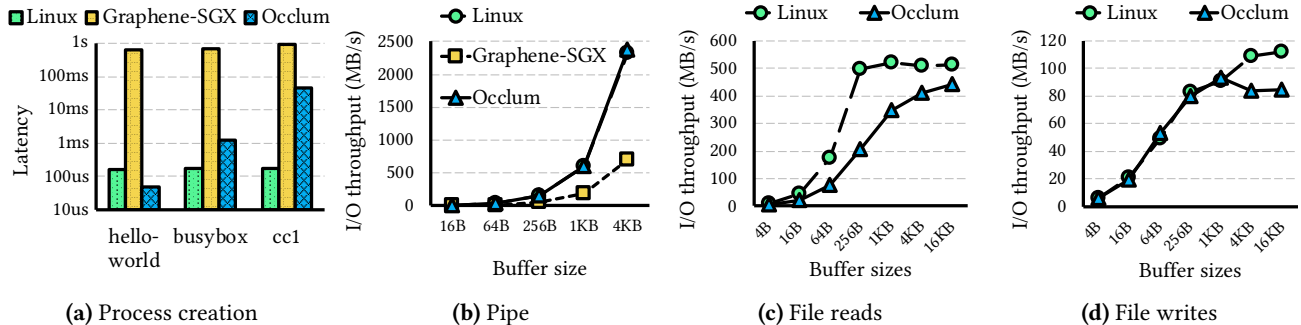
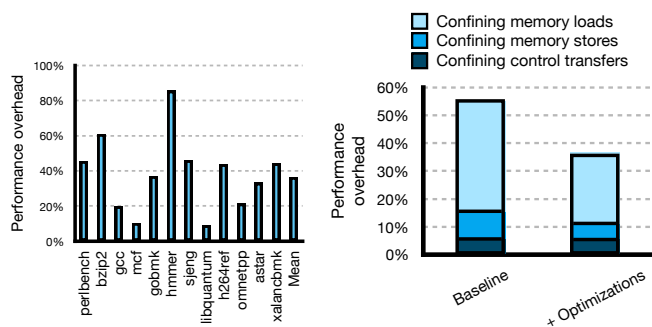


Figure 6. System call benchmarks.



(a) The overheads on SPECint2006 (b) The overhead breakdown

Figure 7. MMDSFI benchmarks

MMDSFI is effective in defending against memory-related attacks. Note that the successful attacks (e.g., return-to-libc) that are not prevented by MMDSFI do not break the isolation between SIPs.

CPU overhead. While MMDSFI is intended for the use in the Occlum LibOS, it can be potentially applied in other use cases. So it is good to know the overheads of MMDSFI alone, independent from the Occlum LibOS. We use the SPECint2006 [27] to measure the overheads of MMDSFI on CPU-intensive workloads. As shown in Figure 7a, the average overhead of MMDSFI is 36.6%.

The breakdown of CPU overheads. To further understand the overheads of MMDSFI, we break down the overheads into three sources: confining control transfers, confining memory stores, and confining memory loads. And we conduct this breakdown analysis on two implementations of MMDSFI: (1) the naive implementation that inserts `mem_guards` for all memory loads and stores; (2) the optimized implementation that leverages the range analysis-based optimization to eliminate unnecessary `mem_guards` (§4.3). Figure 7b show that the range analysis-based optimizations are effective, reducing the overhead of confining memory stores from 10.1% to 4.3% and the overhead of confining memory loads from 39.6% to 25.5%.

10 Conclusions

In this paper, we present Occlum, a system that enables secure and efficient multitasking on SGX LibOSes. We implement the LibOS processes as *SFI-Isolated Processes (SIPs)*. To this end, we propose a novel SFI scheme named *MPX-based, Multi-Domain SFI (MMDSFI)*. We also design an independent verifier to ensure the security guarantees of MMDSFI. With SIPs safely sharing the single address space of an enclave, the LibOS can implement multitasking efficiently. Experimental results show that Occlum outperforms the state-of-the-art multitasking SGX LibOS significantly.

Acknowledgments

We would like to thank ATC, SOSP and ASPLOS anonymous reviewers for their helpful comments. And we would also like to thank Junjie Mao for his valuable insights.

The authors from Tsinghua University are supported by National Key R&D Program of China (Grant No. 2017YFB0802901), National Natural Science Foundation of China (Grant No. 61772303 and 61877035), Young Scientists Fund of the National Natural Science Foundation of China (Grant No. 61802219). The author from Shanghai Jiao Tong University is sponsored by National Natural Science Foundation of China (Grant No. 61972244) and Program of Shanghai Academic/Technology Research Leader (No. 19XD1401700). Some authors are sponsored by Beijing National Research Center for Information Science and Technology (BNRist) to attend the conference.

References

- [1] 2018. *Overview of Intel Protected File System Library Using Software Guard Extensions*. <https://software.intel.com/en-us/articles/overview-of-intel-protected-file-system-library-using-software-guard-extensions>
- [2] 2019. *byte-unixbench*. <https://github.com/kdlucas/byte-unixbench>
- [3] 2019. *Download the Windows version of PostgreSQL*. <https://www.postgresql.org/download/windows/>
- [4] 2019. *fish shell*. <https://fishshell.com>
- [5] 2019. *Large single compilation-unit C programs*. <http://people.csail.mit.edu/smcc/projects/single-file-programs>
- [6] 2019. *Lighttpd, fly light*. <https://www.lighttpd.net>

- [7] 2019. *The Linux FUSE (Filesystem in Userspace) interface*. <https://github.com/libfuse/libfuse>
- [8] 2019. *The LLVM Compiler Infrastructure*. <http://llvm.org/>
- [9] 2019. *musl libc*. <https://www.musl-libc.org/>
- [10] 2019. *Occlum: A memory-safe, multi-process LibOS for Intel SGX*. <https://github.com/anonymous8923/occlum>
- [11] 2019. *Redis on Windows*. <https://github.com/ServiceStack/redis-windows>
- [12] 2019. *The Rust Programming Language*. <https://www.rust-lang.org/>
- [13] 2019. *Single Address Space Operating System*. <http://wiki.c2.com/?SingleAddressSpaceOperatingSystem>
- [14] 2019. *sshd - OpenSSH SSH daemon*. <https://linux.die.net/man/8/sshd>
- [15] 2019. *Using Apache With Microsoft Windows*. <http://structure.usc.edu/apache/windows.html>
- [16] 2019. *Zydis: Fast and lightweight x86/x86-64 disassembler library*. <https://zydis.re/>
- [17] Alfred V. Aho, Monica S. Lam, Ravi Sethi, and Jeffrey D. Ullman. 2006. *Compilers: Principles, Techniques, and Tools (2nd Edition)*. Addison-Wesley Longman Publishing Co., Inc., Boston, MA, USA.
- [18] Dennis Andriess, Xi Chen, Victor Van Der Veen, Asia Slowinska, and Herbert Bos. 2016. An In-depth Analysis of Disassembly on Full-scale x86/x64 Binaries. In *Proceedings of the 25th USENIX Conference on Security Symposium (Austin, TX, USA) (SEC'16)*. USENIX Association, Berkeley, CA, USA, 583–600. <http://dl.acm.org/citation.cfm?id=3241094.3241140>
- [19] Sergei Arnautov, Bohdan Trach, Franz Gregor, Thomas Knauth, André Martin, Christian Priebe, Joshua Lind, Divya Muthukumaran, Dan O’Keeffe, Mark Stillwell, David Goltzsche, David M. Evers, Rüdiger Kapitza, Peter R. Pietzuch, and Christof Fetzer. 2016. SCONE: Secure Linux Containers with Intel SGX. In *12th USENIX Symposium on Operating Systems Design and Implementation, OSDI 2016, Savannah, GA, USA, November 2-4, 2016*. 689–703. <https://www.usenix.org/conference/osdi16/technical-sessions/presentation/arnautov>
- [20] Andrew Baumann, Jonathan Appavoo, Orran Krieger, and Timothy Roscoe. 2019. A fork() in the road. In *Proceedings of the 18th Workshop on Hot Topics in Operating Systems, HotOS 2019, Bertinoro, Italy, May 12-15, 2019*. <https://doi.org/10.1145/3102980.3103000>
- [21] Andrew Baumann, Marcus Peinado, and Galen C. Hunt. 2015. Shielding Applications from an Untrusted Cloud with Haven. *ACM Trans. Comput. Syst.* 33, 3 (2015), 8:1–8:26. <https://doi.org/10.1145/2799647>
- [22] Jo Van Bulck, Marina Minkin, Ofir Weisse, Daniel Genkin, Baris Kasikci, Frank Piessens, Mark Silberstein, Thomas F. Wenisch, Yuval Yarom, and Raoul Strackx. 2018. Foreshadow: Extracting the Keys to the Intel SGX Kingdom with Transient Out-of-Order Execution. In *27th USENIX Security Symposium (USENIX Security 18)*. USENIX Association, Baltimore, MD, 991–1008. <https://www.usenix.org/conference/usenixsecurity18/presentation/bulck>
- [23] Nicholas Carlini, Antonio Barresi, Mathias Payer, David Wagner, and Thomas R. Gross. 2015. Control-Flow Bending: On the Effectiveness of Control-Flow Integrity. In *24th USENIX Security Symposium (USENIX Security 15)*. USENIX Association, Washington, D.C., 161–176. <https://www.usenix.org/conference/usenixsecurity15/technical-sessions/presentation/carlini>
- [24] Stephen Checkoway and Hovav Shacham. 2013. Iago Attacks: Why the System Call API is a Bad Untrusted RPC Interface. *SIGPLAN Not.* 48, 4 (March 2013), 253–264. <https://doi.org/10.1145/2499368.2451145>
- [25] Guoxing Chen, Sanchuan Chen, Yuan Xiao, Yinqian Zhang, Zhiqiang Lin, and Ten H. Lai. 2018. SgxPectre Attacks: Stealing Intel Secrets from SGX Enclaves via Speculative Execution. (2018). arXiv:arXiv:1802.09085
- [26] Jonathan Corbet. 2015. *Memory protection keys*. <https://lwn.net/Articles/643797/>
- [27] Standard Performance Evaluation Corporation. 2019. *CINT2006 (Integer Component of SPEC CPU2006)*. <https://www.spec.org/cpu2006/CINT2006/>
- [28] Casper Dik. 2018. *posix_spawn() as an actual system call*. https://blogs.oracle.com/solaris/posix_spawn-as-an-actual-system-call
- [29] Xiaowan Dong, Zhuojia Shen, John Criswell, Alan L. Cox, and Sandhya Dwarkadas. 2018. Shielding Software From Privileged Side-Channel Attacks. In *27th USENIX Security Symposium (USENIX Security 18)*. USENIX Association, Baltimore, MD, 1441–1458. <https://www.usenix.org/conference/usenixsecurity18/presentation/dong>
- [30] Dawson R. Engler, M. Frans Kaashoek, and James O’Toole. 1995. Exokernel: An Operating System Architecture for Application-Level Resource Management. In *Proceedings of the Fifteenth ACM Symposium on Operating System Principles, SOSP 1995, Copper Mountain Resort, Colorado, USA, December 3-6, 1995*. 251–266. <https://doi.org/10.1145/224056.224076>
- [31] Fluentd Project. 2019. *Fluentd is an open source data collector for unified logging layer*. <https://www.fluentd.org/>
- [32] The Apache Software Foundation. 2019. *ab - Apache HTTP server benchmarking tool*. <https://httpd.apache.org/docs/2.4/programs/ab.html>
- [33] Inc Free Software Foundation. 2019. *GCC, the GNU Compiler Collection*. <https://www.gnu.org/software/gcc>
- [34] GNU. 2019. *Auxiliary Vector*. https://www.gnu.org/software/libc/manual/html_node/Auxiliary-Vector.html
- [35] E. G  ktas, E. Athanasopoulos, H. Bos, and G. Portokalidis. 2014. Out of Control: Overcoming Control-Flow Integrity. In *2014 IEEE Symposium on Security and Privacy*. 575–589. <https://doi.org/10.1109/SP.2014.43>
- [36] Tyler Hunt, Zhiting Zhu, Yuanzhong Xu, Simon Peter, and Emmett Witchel. 2016. Ryoan: A Distributed Sandbox for Untrusted Computation on Secret Data. In *12th USENIX Symposium on Operating Systems Design and Implementation, OSDI 2016, Savannah, GA, USA, November 2-4, 2016*. 533–549. <https://www.usenix.org/conference/osdi16/technical-sessions/presentation/hunt>
- [37] Intel. 2013. *Introduction to Intel(R) Memory Protection Extensions*. <https://software.intel.com/en-us/articles/introduction-to-intel-memory-protection-extensions>
- [38] Intel. 2019. *Intel SGX for Linux*. <https://github.com/intel/linux-sgx>
- [39] Intel. 2019. *Intel(R) Software Guard Extensions SDK*. <https://software.intel.com/en-us/sgx-sdk/documentation>
- [40] Koen Koning, Xi Chen, Herbert Bos, Cristiano Giuffrida, and Elias Athanasopoulos. 2017. No Need to Hide: Protecting Safe Regions on Commodity Hardware. In *Proceedings of the Twelfth European Conference on Computer Systems, EuroSys 2017, Belgrade, Serbia, April 23-26, 2017*. 437–452. <https://doi.org/10.1145/3064176.3064217>
- [41] Christopher Kr  gel, William K. Robertson, Fredrik Valeur, and Giovanni Vigna. 2004. Static Disassembly of Obfuscated Binaries. In *Proceedings of the 13th USENIX Security Symposium, August 9-13, 2004, San Diego, CA, USA*. 255–270. <http://www.usenix.org/publications/library/proceedings/sec04/tech/kruegel.html>
- [42] Dmitrii Kuvaiskii, Oleksii Oleksenko, Sergei Arnautov, Bohdan Trach, Pramod Bhatotia, Pascal Felber, and Christof Fetzer. 2017. SGXBOUNDS: Memory Safety for Shielded Execution. In *Proceedings of the Twelfth European Conference on Computer Systems, EuroSys 2017, Belgrade, Serbia, April 23-26, 2017*. 205–221. <https://doi.org/10.1145/3064176.3064192>
- [43] Youngjin Kwon, Alan M. Dunn, Michael Z. Lee, Owen S. Hofmann, Yuanzhong Xu, and Emmett Witchel. 2016. Seg0: Pervasive Trusted Metadata for Efficiently Verified Untrusted System Services. *SIGOPS Oper. Syst. Rev.* 50, 2 (March 2016), 277–290. <https://doi.org/10.1145/2954680.2872372>
- [44] Stephen McCamant and Greg Morrisett. 2005. *Efficient, Verifiable Binary Sandboxing for a CISC Architecture*. MIT CSAIL Technical Report TR-2005-030. Massachusetts Institute of Technology, Cambridge, MA.
- [45] Oleksii Oleksenko, Bohdan Trach, Robert Krahn, Mark Silberstein, and Christof Fetzer. 2018. Varys: Protecting SGX Enclaves from Practical Side-Channel Attacks. In *2018 USENIX Annual Technical Conference*,

- USENIX ATC 2018, Boston, MA, USA, July 11-13, 2018*. 227–240.
- [46] Donald E. Porter, Silas Boyd-Wickizer, Jon Howell, Reuben Olinsky, and Galen C. Hunt. 2011. Rethinking the Library OS from the Top Down. *SIGARCH Comput. Archit. News* 39, 1 (March 2011), 291–304. <https://doi.org/10.1145/1961295.1950399>
- [47] Eric S. Raymond. 2003. *The Art of UNIX Programming*. Pearson Education.
- [48] RedHat. 2019. *etcd: a distributed, reliable key-value store for the most critical data of a distributed system*. <https://coreos.com/etcd/>
- [49] Benjamin Schwarz, Saumya K. Debray, and Gregory R. Andrews. 2002. Disassembly of Executable Code Revisited. In *9th Working Conference on Reverse Engineering (WCRE 2002), 28 October - 1 November 2002, Richmond, VA, USA*. 45–54. <https://doi.org/10.1109/WCRE.2002.1173063>
- [50] David Sehr, Robert Muth, Cliff Biffle, Victor Khimenko, Egor Pasko, Karl Schimpf, Bennet Yee, and Brad Chen. 2010. Adapting Software Fault Isolation to Contemporary CPU Architectures. In *19th USENIX Security Symposium, Washington, DC, USA, August 11-13, 2010, Proceedings*. 1–12. http://www.usenix.org/events/sec10/tech/full_papers/Sehr.pdf
- [51] Ming-Wei Shih, Sangho Lee, Taesoo Kim, and Marcus Peinado. 2017. T-SGX: Eradicating Controlled-Channel Attacks Against Enclave Programs. In *24th Annual Network and Distributed System Security Symposium, NDSS 2017, San Diego, California, USA, February 26 - March 1, 2017*. <https://www.ndss-symposium.org/ndss2017/ndss-2017-programme/t-sgx-eradicating-controlled-channel-attacks-against-enclave-programs/>
- [52] Shweta Shinde, Dat Le Tien, Shruti Tople, and Prateek Saxena. 2017. Panoply: Low-TCB Linux Applications With SGX Enclaves. In *24th Annual Network and Distributed System Security Symposium, NDSS 2017, San Diego, California, USA, February 26 - March 1, 2017*.
- [53] Yan Shoshitaishvili, Ruoyu Wang, Christophe Hauser, Christopher Kruegel, and Giovanni Vigna. 2015. Fimalice - Automatic Detection of Authentication Bypass Vulnerabilities in Binary Firmware. In *22nd Annual Network and Distributed System Security Symposium, NDSS 2015, San Diego, California, USA, February 8-11, 2015*. <https://www.ndss-symposium.org/ndss2015/fimalice-automatic-detection-authentication-bypass-vulnerabilities-binary-firmware>
- [54] Rohit Sinha, Manuel Costa, Akash Lal, Nuno P. Lopes, Sriram K. Rajamani, Sanjit A. Seshia, and Kapil Vaswani. 2016. A design and verification methodology for secure isolated regions. In *Proceedings of the 37th ACM SIGPLAN Conference on Programming Language Design and Implementation, PLDI 2016, Santa Barbara, CA, USA, June 13-17, 2016*. 665–681. <https://doi.org/10.1145/2908080.2908113>
- [55] Chia-che Tsai, Donald E. Porter, and Mona Vij. 2017. Graphene-SGX: A Practical Library OS for Unmodified Applications on SGX. In *2017 USENIX Annual Technical Conference, USENIX ATC 2017, Santa Clara, CA, USA, July 12-14, 2017*. 645–658. <https://www.usenix.org/conference/atc17/technical-sessions/presentation/tsai>
- [56] Robert Wahbe, Steven Lucco, Thomas E. Anderson, and Susan L. Graham. 1993. Efficient Software-Based Fault Isolation. In *Proceedings of the Fourteenth ACM Symposium on Operating System Principles, SOSP 1993, The Grove Park Inn and Country Club, Asheville, North Carolina, USA, December 5-8, 1993*. 203–216. <https://doi.org/10.1145/168619.168635>
- [57] Huibo Wang, Pei Wang, Yu Ding, Mingshen Sun, Yiming Jing, Ran Duan, Long Li, Yulong Zhang, Tao Wei, and Zhiqiang Lin. 2019. Towards Memory Safe Enclave Programming with Rust-SGX. In *Proceedings of the 2019 ACM SIGSAC Conference on Computer and Communications Security*. 2333–2350.
- [58] Wenhao Wang, Guoxing Chen, Xiaorui Pan, Yinqian Zhang, XiaoFeng Wang, Vincent Bindschaedler, Haixu Tang, and Carl A. Gunter. 2017. Leaky Cauldron on the Dark Land: Understanding Memory Side-Channel Hazards in SGX. In *Proceedings of the 2017 ACM SIGSAC Conference on Computer and Communications Security (Dallas, Texas, USA) (CCS '17)*. ACM, New York, NY, USA, 2421–2434. <https://doi.org/10.1145/3133956.3134038>
- [59] John Wilander, Nick Nikiforakis, Yves Younan, Mariam Kamkar, and Wouter Joosen. 2011. RIPE: Runtime Intrusion Prevention Evaluator. In *Proceedings of the 27th Annual Computer Security Applications Conference (Orlando, Florida, USA) (ACSAC '11)*. ACM, New York, NY, USA, 41–50. <https://doi.org/10.1145/2076732.2076739>
- [60] Bin Zeng, Gang Tan, and Greg Morrisett. 2011. Combining control-flow integrity and static analysis for efficient and validated data sandboxing. In *Proceedings of the 18th ACM Conference on Computer and Communications Security, CCS 2011, Chicago, Illinois, USA, October 17-21, 2011*. 29–40. <https://doi.org/10.1145/2046707.2046713>

A Artifact Appendix

A.1 Abstract

Our artifacts consist of the Occlum toolchain, the Occlum LibOS, and some scripts to set up the environment, build the projects, and run the benchmarks. The hardware requirement is Intel x86-64 CPUs with SGX and MPX support.

A.2 Artifact check-list (meta-information)

- **Compilation:** LLVM (which will be downloaded and patched by the script).
- **Transformations:** MMDSFI instrumentation implemented as LLVM passes.
- **Run-time environment:** Root access to Ubuntu Linux.
- **Hardware:** Intel x86-64 CPU with SGX and MPX extensions. We recommend i7 Kaby Lake CPUs. To test network performance, another computer is required to connect as client with 1Gbps ethernet.
- **Output:** For benchmarks, the results are printed in consoles.
- **Experiments:** Using Bash scripts.
- **How much disk space required (approximately)?:** 40GB.
- **How much time is needed to prepare workflow (approximately)?:** Half a day.
- **How much time is needed to complete experiments (approximately)?:** One hour.
- **Publicly available?:** Yes.
- **Code licenses (if publicly available)?:** BSD.
- **Archived (provide DOI)?:** 10.5281/zenodo.3565239

A.3 Description

A.3.1 How delivered. The artifacts are available on GitHub: <https://github.com/occlum/reproduce-asplos20>. We need to modify source code from other projects like LLVM. Our scripts will automatically download the original source code and then apply the patches.

Note that the artifacts provided are the archive of an early version of Occlum project. We recommend interesting users to test on the latest version of Occlum, which can be found at <https://github.com/occlum/occlum>.

A.3.2 Hardware dependencies. We recommend testing on an Intel Core i7 Kaby Lake CPU (i7-7567U was used in our test). SGX and MPX are required.

A.3.3 Software dependencies. We built and tested our system on Ubuntu 16.04 with kernel 4.15.0-65-generic.

A.4 Installation

You can get the artifacts from GitHub using the following command:

```
git clone https://github.com/occlum/reproduce-asplos20
```

A.5 Experiment workflow

The overall workflow consists of the following steps:

1. Install the dependencies;

2. Build the LibOS;
3. Build the toolchain;
4. Build and run the macro-benchmarks;
5. Build and run the micro-benchmarks.

We provide scripts for each of the steps above.

A.6 Evaluation and expected result

We will go through the entire experiment workflow by describing all the commands in each step.

A.6.1 Install the dependencies. Run the following command:

```
./prepare.sh
```

A.6.2 Build the LibOS. Run the following command:

```
./download_and_build_libos.sh.
```

This script above will download the source code of the LibOS and build it.

A.6.3 Build the toolchain. To build and install the toolchain, run the following command:

```
cd toolchain && \  
./download_and_build_toolchain.sh
```

This script will install the toolchain at `/usr/local/occlum`. To use this toolchain for building benchmarks, export the installation directory to `PATH` with the following command:
`export PATH=/usr/local/occlum/bin:$PATH`

A.6.4 Build and run the macro-benchmarks. As described in §9, there are three macro-benchmarks: fish, GCC, and lighttpd.

fish benchmark. To build the benchmark, run the following command:

```
cd apps/fish && ./download_and_build_fish.sh
```

After finishing building the benchmark, run the following command to run the benchmark:

```
cd apps/fish && ./run_fish_fish.sh
```

When the script is done, the LibOS will print the total running time of the benchmark repeating for 100 times.

GCC benchmark. To build GCC with the Occlum toolchain, run the following command:

```
cd apps/gcc && ./download_and_build_gcc.sh
```

There are three workloads, which can be tested with the following commands:

```
cd apps/gcc && \  
./run_gcc_helloworld.sh && \  
./run_gcc_5K.sh && \  
./run_gcc_50K.sh
```

lighttpd benchmark. To build lighttpd, run the following command:

```
cd apps/lighttpd && ./download_and_build_lighttpd.sh
```

Next, open `apps/lighttpd/config/lighttpd-server.conf` with your favourite text editor and set `server.bind` to the IP address of your machine.

There are two modes of the benchmark. Start a single-thread server with

```
cd apps/lighttpd &&  
./run_lighttpd_test.sh
```

or start a multi-threaded server with

```
cd apps/lighttpd &&  
./run_lighttpd_test_multithread.sh
```

After the server started, you can now test the throughput and latency by starting the http benchmark program:

```
cd apps/lighttpd &&  
./benchmark-http.sh ${SERVER_IP}:8000
```

where `SERVER_IP` is the IP address of the `lighttpd` server. This benchmark will show the latency and throughput under different number of concurrent clients.

A.6.5 Build and test micro-benchmarks. There are two micro-benchmarks: one measures the latency of process creation and the other measures the throughput of IPC.

Spawn benchmark. Run the benchmark with the following command:

```
cd bench/spawn && ./run_spawn_bench.sh
```

The result will be printed on the console.

Pipe benchmark. Run the benchmark with the following command:

```
cd bench/pipe && ./run_pipe_bench.sh
```

The result will be printed on the console.

A.7 Experiment customization

Since the LibOS and the toolchain are installed, users can build their own applications with our toolchain and run them with our LibOS.