

GiantSan: Efficient Operation-Level Memory Sanitization with Segment Folding

HAO LING, The Hong Kong University of Science and Technology, Hong Kong, Hong Kong HEQING HUANG^{*}, Computer Science, City University of Hong Kong, Hong Kong, Hong Kong CHENGPENG WANG, The Hong Kong University of Science and Technology, Hong Kong, Hong Kong YUANDAO CAI, The Hong Kong University of Science and Technology, Hong Kong, Hong Kong CHARLES ZHANG, The Hong Kong University of Science and Technology, Hong Kong, Hong Kong

Memory safety sanitizers, the sharp weapon for detecting invalid memory operations during execution, employ runtime metadata to model the memory and help find memory errors hidden in the programs. However, location-based methods, the most widely deployed memory sanitization methods thanks to high compatibility, face the low protection density issue: the number of bytes safeguarded by one metadata is limited. As a result, numerous memory accesses require loading excessive metadata, leading to a high runtime overhead.

To address this issue, we propose a new shadow encoding with *segment folding* to increase the protection density. Specifically, we characterize neighboring bytes with identical metadata by building novel summaries, called *folded segments*, on those bytes to reduce unnecessary metadata loadings. The new encoding uses less metadata to safeguard large memory regions with fewer instructions than existing works, speeding up memory sanitization.

We implement our designed technique as GIANTSAN. Our evaluation using the SPEC CPU 2017 benchmark shows that GIANTSAN outperforms the state-of-the-art sanitization methods with 61.37% and 41.94% less runtime overhead than ASan and ASan--, respectively. Moreover, under the same redzone setting, GIANTSAN detects 463 fewer false negative cases than ASan and ASan-- in testing the real-world project PHP.

CCS Concepts: • Security and privacy → Software security engineering.

Additional Key Words and Phrases: Sanitizers, Memory Safety

*Corresponding Author

This paper, originally titled GIANTSAN: Efficient Memory Sanitization with Segment Folding, was invited by ACM TOCS for extended publication through recommendation by the Chairs of ASPLOS 2024.

Authors' Contact Information: Hao Ling, The Hong Kong University of Science and Technology, Hong Kong, Hong Kong; e-mail: hlingab@ cse.ust.hk; Heqing Huang, Computer Science, City University of Hong Kong, Hong Kong, Hong Kong; e-mail: heqhuang@cityu.edu.hk; Chengpeng Wang, The Hong Kong University of Science and Technology, Hong Kong, Hong Kong; e-mail: cwangch@cse.ust.hk; Yuandao Cai, The Hong Kong University of Science and Technology, Hong Kong, Hong Kong; e-mail: cyaibb@cse.ust.hk; Charles Zhang, The Hong Kong University of Science and Technology, Hong Kong, Hong Kong; e-mail: ycaibb@cse.ust.hk; Charles Zhang, The Hong Kong University of Science and Technology, Hong Kong; e-mail: charlesz@ust.hk.

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 the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

© 2025 Copyright held by the owner/author(s). Publication rights licensed to ACM. ACM 1557-7333/2025/1-ART https://doi.org/10.1145/3742426

1:2 • H. Ling et al.

1 Introduction

The freedom to manipulate memory through pointers guaranteed by unsafe languages like C and C++ leads to numerous kinds of memory safety violations. As reported in the 2022 CWE Top 25 Most Dangerous Software Weaknesses [46], for instance, out-of-bounds write, out-of-bounds read, and use-after-free rank 1st, 5th, and 7th among all weaknesses, respectively. For program reliability, researchers have proposed a series of memory sanitizing techniques [2, 9, 10, 13, 21, 26, 28, 29, 34, 35, 40, 42, 49, 50] to detect invalid memory operations during the program execution.

Though tremendous efforts have been made to improve memory sanitization, most methods have limited compatibility, resulting in false negatives or low efficiency in many scenarios. Pointer-based methods, for instance, protect memory accesses with buffer bounds propagated along with pointer arithmetics. However, the propagation highly depends on program instrumentation with type information of pointers, which is not always available. It is a well-known issue [4, 9, 10, 26, 28, 34, 35, 41, 43, 45] that propagation often fails due to pointer-integer casting or uninstrumented external libraries without type information (e.g., third-party codes distributed in binary form). As a result, the pointer-based sanitizers cannot detect errors once the propagation fails.

Location-based methods stand out among the various memory sanitizers due to their high compatibility, which comes from a simpler safety model that does not rely on pointer information to maintain metadata. Specifically, each byte in the memory is assigned one of the two states, *addressable* or *non-addressable*, and a memory access is safe if the target bytes are all addressable. The addressability states are stored in a dedicated shadow memory and can be retrieved anytime, eliminating the need for instrumentation to propagate metadata. For compatibility considerations, memory sanitizers integrated into GCC [14], LLVM [27] compiler projects, and ANDROID [3] system are all location-based [30, 40, 41].

However, though location-based methods offer high compatibility and are fast in metadata maintenance [43], they are deficient in protecting memory operations ¹ involving multiple instructions, and they require excessive runtime checks compared with other methods like pointer-based solutions. Specifically, pointer-based methods safeguard memory operations by checking whether the memory region being accessed is within a safe bound. In contrast, location-based methods do not have such a bound, and they have to break operations down into instructions and check each instruction separately to ensure no non-addressable bytes are accessed. Therefore, though location-based methods save time in metadata maintenance, they incur more runtime checks, which are time-consuming.

Addressa	ble Non-Add	ressable 💋
first 4 bytes addressable	all 8 bytes addressable	non-addressable region in the heap
↑ State: 0x4	↑ State: 0x0	↑ State: 0xfa

Fig. 1. Shadow encoding in ASan, where all objects are 8-byte aligned. The addressable bytes within a segment must occupy a prefix of the segment. By default, ASan uses eight different state codes for addressable bytes within the segments and reserves the other state codes for other purposes (e.g., recording why the bytes are non-addressable).

The root cause of the excessive check issue is the low *protection density* caused by the inefficient shadow memory encoding. The protection density is the number of bytes safeguarded by one piece of metadata. Each byte in the memory has two different states: addressable or non-addressable. Technically, it requires at least one bit to distinguish the two states. Therefore, on average, location-based methods must load and decode one shadow byte for every eight memory bytes. The protection density can be slightly increased according to memory

¹In this paper, a memory operation refers to a series of instructions manipulating the memory region of the same object. For example, "memset(p, 0, 1024)" is one memory operation manipulating 1024 bytes and consists of at least 1024/8 = 128 MoV instructions related to p in a 64-bit system.

GiantSan: Efficient Operation-Level Memory Sanitization with Segment Folding • 1:3



(a) Existing location-based encoding. "good": all bytes in the segment are addressable; "bad": all bytes are nonaddressable; "part": only some bytes in the segment are addressable (partially good).



(b) Segment Folding: build a summary for "good" segments. Only one folded segment needs to be visited instead of four unfolded ones for the region [L, R).

Fig. 2. Folded segments reduce metadata loadings.

alignment: some consecutive bytes must be both addressable or non-addressable, and thus their states can be merged. However, because most objects are only guaranteed to be 8-byte aligned, and this optimization is limited to only a few neighboring bytes.

Figure 1 illustrates the shadow encoding with the low protection density in the most widely deployed sanitizer, ADDRESSSANITIZER (a.k.a. ASan) [40]. It partitions the virtual memory space into a sequence of aligned segments and employs one 8-bit integer (called the *segment state* in this paper) to encode all byte states within the segment. Segments are sized at 8 bytes so that no two objects share the same segment ². Checking a memory region containing *S* bytes requires loading $\lceil \frac{S}{8} \rceil$ segment states, which results in significant runtime overhead. For instance, checking whether a 1KB region contains a non-addressable byte requires loading 128 segment states in ASan. A past study [50] shows that ASan is about 2× slower than native execution, and about 80% of the runtime overhead comes from excessive runtime checks and metadata loadings.

This paper addresses the low protection density issue to improve the efficiency of location-based memory sanitization. Despite the various segment states, almost all segments visited during the execution are "good" segments (i.e., the segments without non-addressable bytes) because most memory operations are safe and only manipulate addressable bytes. Inspired by this observation, our key insight is to build a summary for "good" segments to help reduce segment state loadings, thus increasing the sanitizing efficiency. We call the summarizing process "segment folding".

Let us illustrate our insight with Figure 2. Figure 2a shows how existing methods work: when accessing a memory region, they need to check all segments to ensure all accessed bytes are addressable. Checking the region [L, R) involves 4 segments, and all those segments are "good" since this region is safe to access. Figure 2b shows how the segment folding works: it builds a summary of the "good" segments and uses the summary of segments to speed up the checking of the region [L, R). However, the folding is not free: storing the summary needs extra shadow memory space.

To reduce the shadow memory required to store the summary, we design the binary folding strategy: a folded segment only summarizes 2^x "good" segments for some integer *x*. In a modern 64-bit system, *x* cannot exceed 64 because the maximum object size is less than 2^{64} . As a result, six shadow bits are sufficient to record the folding degree *x*. Combined with the 8-byte alignment optimization, all the segment states and the folding degree *x* can be recorded in one 8-bit integer. As a result, the new shadow memory encoding with segment folding is compact enough to build upon the shadow memory widely adopted by existing location-based methods.

²ASan assumes all objects are 8-byte aligned, which is satisfied in most cases due to the basic assumption of heap allocation.

1:4 • H. Ling et al.

We present GIANTSAN, a dynamic memory error detector with a novel shadow encoding based on segment folding. To the best of our knowledge, GIANTSAN is the first location-based method that can safeguard a sequential region of arbitrary size in O(1) time. We evaluate GIANTSAN on SPEC CPU 2017, the industry-standard CPU-intensive benchmark suite. GIANTSAN reduces the geometric mean runtime overhead down to 43.48%, compared with 74.89% and 112.58% in the state-of-the-art location-based designs ASan-- [50] and ASan [40], respectively. The promising result indicates that GIANTSAN outperforms its competitors.

To sum up, this work makes the following contributions:

- We formulate and summarize the low protection density issue of location-based sanitizers.
- We introduce the segment folding algorithm to increase protection density significantly.
- We implement our approach as a tool named GIANTSAN and provide empirical evidence that it outperforms the state-of-the-art methods with less runtime overhead.

2 Technical Background

This section introduces fundamental knowledge about existing memory sanitizing techniques.

2.1 Existing Solutions for Memory Safety

There are two categories of memory safety violations: 1) **Spatial Errors**: access memory locations outside the allocated region of objects, and 2) **Temporal Errors**: access an object when it is not valid (e.g., unallocated or deallocated).

Although many memory safety violation detecting tools have been proposed [5–7, 9, 10, 13, 21, 24, 25, 28, 35, 38, 40, 42, 49, 50], many only provide partial memory safety guarantees. Some, like Softbound [34], Delta Pointers [25], TailCheck [16], and LFP [9, 10], only support the detection of spatial errors. In contrast, other trends of existing work, like CETS [35] and PTAuth [13], only support the detection of temporal errors.

All sanitizers need extra metadata to model the memory and validate whether one memory region can be accessed. Among the existing efforts to provide a full safety guarantee, there are two main philosophies:

- **Pointer-based**: Pointer-based methods [6, 9, 10, 13, 16, 24, 25, 28, 35, 38] model the memory from the perspective of pointers by tracking the memory region safe to access for each pointer. They encapsulate the pointer and a tag in a new pointer representation, and they use the tag as the bound for the safe region or as the index for retrieving the bound.
- Location-based: Location-based methods [5, 7, 21, 40, 42, 49, 50] model the memory from the perspective of memory bytes by recording which byte is addressable. The byte states are recorded in a compact shadow memory, and location-based methods inspect the shadow memory to check the state of each accessed byte.

The core difference between the two philosophies is *the dependence on the data type information of pointers*. Specifically, whenever pointer arithmetic creates a new pointer, pointer-based methods need to convert it into the new pointer representation and propagate the tag from the source pointer to the new pointer. Therefore, in pointer-based methods, all instructions must be aware of whether they are manipulating pointers so that the tag is propagated correctly and not misused. In contrast, memory protection in location-based methods only depends on the metadata binding to the memory address instead of pointers.

Unfortunately, the type information of pointers is not always available. For example, programs can use external libraries distributed in binary form without type information, and all values are treated as integers. Moreover, even with the source codes available, the type information of pointers may not be available since the pointer-integer casting can eliminate the type information. The casting converts pointers into integers, and later, pointers are manipulated by integer arithmetic instead of pointer arithmetic. As a result, it is challenging to distinguish between the customized pointer representation and the native integers, which might result in tag misuse or tag propagation failure [4, 9, 10, 26, 28, 34, 35, 41, 43, 45].

Analysis Method	Example	# Checks (operation-level)	# Checks (instruction-level)
Constant Propagation	p[0] + p[10] + p[20]	1	3
Predefined Semantics	memset(p, 0, N)	1	$\Theta(N)$
Loop Bound Analysis	for (auto i = 0; i < N; i++) p[i] = foo(i);	1	N
Must-alias Analysis	p[0] = 10 for (auto i : vec) p[i] = foo(i);	1 slow check + N fast checks (with bound cached)	N+1 slow checks (with nothing cached)
	p[i] = foo(i);	(with bound cached)	(with nothing cached)

Table 1. Difference between operation-level protection and instruction-level protection on the pointer *p*. The *Analysis Method* column shows the static analysis used to identify the operations in the source codes. *N* in the fourth case is the size of vec.

Once the pointer tag is lost due to propagation failure, the pointer-based methods cannot protect the pointer and all new pointers derived from it. Some efforts attempt [2, 9, 10, 24] to recover from the tag loss by obtaining a new tag based on the pointer values from dedicated data structures, e.g., shadow memory, similar to the location-based methods. However, location-based methods only require distinguishing two states of bytes with a compact shadow memory. In contrast, keeping tags to distinguish different objects requires a much larger shadow memory. Large shadow memory causes excessive memory consumption and significantly affects runtime efficiency due to a high memory footprint [40, 43].

One of the most representative efforts in tag reobtaining is the Baggy Bound Checking (BBC) [2]. To avoid large shadow memory footprints, it rounds allocation sizes up to a power of two to reduce the total variety of tags. As a result, it cannot detect errors within the rounded-up allocation size. For example, it cannot detect the out-of-bound access "p[700]" for a buffer "*char* p[600]" because the buffer is rounded up to "*char* p[1024]". Therefore, due to the tolerance of many spatial violations, BBC is less suitable for testing [2, 43].

Due to their high dependence on pointer type information, pointer-based methods are less compatible in the complicated real-world testing environment. In contrast, location-based methods are much more widely adopted because they only need to know which memory address is being accessed. That is why general-purpose compiler projects like LLVM and GCC only integrate location-based methods. However, location-based methods have their own efficiency issue, which we aim to address in this paper, discussed in the following.

2.2 Location-based checking with shadow memory

Shadow memory is a technique to monitor and maintain the states of bytes in the memory, widely used in memory safety sanitizers [5, 19, 40–42, 49, 50]. It is the most efficient data structure to implement location-based methods. Location-based methods partition the virtual memory space into fixed-sized segments and use shadow memory to record the segment state, which encodes the states of bytes within the segment. Specifically, shadow memory is an array of shadow units, each of which stores a piece of metadata. We use the notation *m* to represent the global array, *N* for the number of segments, and S_{shadow} for the size of each segment. The following is how shadow memory is declared:

ShadowUnitType m[N];

Given a memory address *p*, the state of the segment covering the address *p* can be loaded by:

 $m[(intptr_t)p/S_{shadow}]$

1:6 • H. Ling et al.

Location-based methods can only detect whether a byte is addressable but cannot guarantee that the byte belongs to the desired object. Most existing location-based methods integrate *redzones* [21, 40, 42, 49, 50] and *memory quarantine* [1, 21, 40, 42, 49, 50] to detect sophisticated memory errors. Specifically, *redzones* are non-addressable paddings between objects (for spatial error detection), and *memory quarantine* delays the re-allocation of memory regions to ensure that an object's memory region is not addressable during a particular time (for temporal error detection).

Runtime Checks. Before accessing *w* bytes starting from an address *p*, location-based methods safeguard the memory access by checking whether all *w* target bytes are addressable. The metadata indicating the addressability of bytes comes from the shadow memory. The metadata only has a limited bit width (e.g., 8 bits) to enable compact shadow memory and can not hold much information. As a result, *w* is small in existing location-based methods so that the byte states can be encoded with a limited bit width.

EXAMPLE 1. ASan [40] uses $S_{shadow} = 8$, and 8-bit signed integers as the ShadowUnitType. m[p] = 0 means the p-th segment is a "good" segment (i.e., all bytes in this segment are addressable), and m[p] = k ($1 \le k \le 7$) means the p-th segment is a k-partial segment (i.e., only the first k bytes in this segment are addressable). ASan creates one runtime check for all memory accesses with $w \le 8$:

```
1 int8_t v = m[p / 8];
2 if (v != 0 and (p & 7) + w > v) {
3 ReportError(p, w)
4 }
5 access [p, p + w)
```

The maximum allowable value of *w* determines the protection density: larger *w* means more bytes can be safeguarded by the metadata, thus resulting in fewer metadata loadings and runtime checks. However, for memory efficiency, location-based methods need to use compact shadow memory, which cannot allocate a large bit width for a piece of metadata, and inefficient shadow encoding can only employ small *w* and limits the protection density.

2.3 Problems and Challenges

In this section, we demonstrate how protection density affects sanitizing efficiency by presenting two protection principles used in different sanitizers: 1) *operation-level protection* aims to protect a memory operation consisting of multiple instructions as a whole, and 2) *instruction-level protection* safeguards each instruction separately. We discuss why operation-level protection requires a high protection density and generates fewer runtime checks. We also discuss the challenges in enabling operation-level protection in location-based methods.

A memory operation is a series of memory accesses toward the allocated region of one single object. Table 1 shows four types of commonly used runtime checks based on the semantics of memory operations, all associated with the pointer p. For example, constant propagation can tell that p[0], p[10], and p[20] are all memory accesses towards p with constant offsets. Operation-level protection safeguards all three instructions at once by testing $[\&p[0], \&p[21]) \subseteq bound(p)$. Similarly, the memset and bounded loop require only one check under operation-level protection. In contrast, the instruction-level protection checks all instructions executed separately. For example, p[0], p[10], and p[20] involve three instructions, and the instruction-level protection checks each of them separately.

Moreover, the operation-level protection can also reduce metadata loadings with caching. The operation-level protection can cache the bound of p for future memory accesses on p, as listed in the fourth case of Table 1. Once the bound of p is loaded when checking p[0] = 10, the bound can be cached in a local variable and used to check all instructions in the loop. In contrast, the instruction-level protection checks each instruction separately, and

```
1 void foo(int *p, size_t n) {
2   for (size_t i = 0; i < n; i++) {
3      p[i]++;
4   }
5 }</pre>
```

(a) A simple bounded loop, whose bound can be analyzed before optimization.

```
void foo(int *p, size_t n) {
1
2
     size_t n_vec = (n - n \% 8);
3
     size_t j = 0;
     if (n >= 8) {
4
       for (size_t i = 0; i < n_vec; i += 8) {</pre>
5
         p[i + 0] += 1;
p[i + 1] += 1;
6
7
          p[i + 2] += 1;
8
          p[i + 3] += 1;
9
10
          p[i + 4] += 1;
11
          p[i + 5] += 1;
12
          p[i + 6] += 1;
13
          p[i + 7] += 1; // to be vectorized.
14
       }
15
       i
          = n_vec;
16
17
         (; j < n; j++) p[j] += 1;
     for
  }
18
```



Fig. 3. The optimization lowers the abstraction level.

the metadata loaded can only safeguard the corresponding instruction. Caching metadata with low protection density cannot help speed up future checks because it does not contain much information.

The operation-level protection requires much fewer checks than the instruction-level protection. However, it needs to efficiently verify memory regions of arbitrary sizes, which, unfortunately, is not available in existing location-based methods, as discussed in Section 2.2. Meanwhile, operation-level protection is sensitive to the abstraction level of the code. The transformation passes from the compiler lower the abstraction level of the high-level language features in the source code, bringing the program closer to machine code. For example, a loop with a fixed bound in the source code may be unrolled or split into two loops with variant bounds during the compilation process. Building operation-level protection needs to consider the interaction with the compiler optimization.

EXAMPLE 2. Figure 3a contains a simple bounded loop. Theoretically, SCEV analysis is sufficient to infer the bound of this loop and yields the operation-level protection about the memory region [&p[0], &p[n + 1]). However, the optimization pipeline would convert the loop from a simple form into a complicated form, as shown in Figure 3b. The original single memory operation is duplicated into nine separate memory instructions, and the single loop is split into two loops. The analysis cannot handle multiple instructions across multiple loops, thus failing to infer the bound of the loop.

Summary. Existing location-based methods have the following deficiencies of the instruction-level protection, all caused by the low protection density. We attempt to address the deficiencies by increasing protection density.

- Inefficient in safeguarding large memory regions.
- Inefficient in caching history.



(a) Shadow memory technique. "good" means all bytes in the segment are addressable; "part" means partially good only the first several bytes in the segment are addressable; "bad" and "freed" represent non-addressable segments maintained by redzones and memory quarantine.



(b) "(*i*)" indicates that this segment is a folded segment combining two consecutive folded segments with the folding degree "(i - 1)". In particular, "(0)" indicates "good" segments. A segment with code "(*i*)" summarizes 2^i consecutive "good" segments.

Fig. 4. High-level comparison between GIANTSAN and existing approaches: the majority of consecutive segments can be folded and checked as a whole.

3 GIANTSAN in a Nutshell

We present GIANTSAN, a novel location-based sanitizer enabling operation-level protection. Our main observation is that most segments being visited during execution are "good" segments, so characterizing and protecting good-segment-only memory regions with a customized summary suffice in most cases. For example, in Figure 4a, the "Safe!" region requires loading 5 segment states. In contrast, in Figure 4b, GIANTSAN combines nearby "good" segments to avoid visiting "good" segments repeatedly and conducts only 2 checks. Figure 5 shows two key phases of GIANTSAN:

- The runtime support library hooks all objects' allocation and deallocation to initialize the metadata in shadow memory (Section 4.1) during the execution.
- The instrumentation system inserts checks to protect memory operations. Operation-level protection requires different instrumentation logic for consecutive region checks (Section 4.2) and history caching (Section 4.3). The construction of the operation-level protection needs to cooperate with compiler optimization (Section 4.4).

The runtime support library sets the metadata in the shadow memory. Specifically, to implement the runtime support library, we first need to design metadata modeling the memory by answering the following question: **Question 1:** How to fold segments and encode the folded segments in the shadow memory?

Solution: GIANTSAN employs the recursive binary folding strategy: two consecutive "good" segments, or two consecutive folded segments with the same size, are combined to form a new folded segment. As illustrated in Figure 4b, the (1)-folded segment combines two "good" segments, and the (2)-folded segment combines two (1)-folded segments. The folded segments summarize addressable regions, speeding up the segment checks, and only the folding degree (*i*) needs to be recorded. We discuss the details in Section 4.1.

GIANTSAN utilizes the optimized shadow memory to safeguard memory regions. To solve the deficiencies discussed in Section 2.3, we face two main questions:

Question 2: How to efficiently safeguard given memory regions with arbitrary sizes?

Solution: Safeguarding a memory region is simplified into checking whether the folding degree is large enough. More specifically, if we want to check whether *N* consecutive segments contain non-addressable bytes, we can

GiantSan: Efficient Operation-Level Memory Sanitization with Segment Folding • 1:9



Fig. 5. GIANTSAN's workflow

check whether the first and last $2^{\lfloor \log_2 N \rfloor}$ segments are folded, significantly reducing the required metadata. We place the details of locating the folded segments in Section 4.2.

Question 3: How to build a cache to speed up further checks?

Solution: GIANTSAN caches the *last* folded segment visited for a given pointer, which can be considered as a temporary bound for all accesses checked. The bound helps reduce the metadata loadings for future accesses on the same pointer. We discuss the caching algorithm in Section 4.3.

Question 4: How to cooperate with the compiler optimization?

Solution: GIANTSAN uses pseudo instrumentation to cooperate with compiler optimization. Specifically, before the optimization pipeline lowers the abstraction level, GIANTSAN scans memory operations and encodes the information about all protection tasks with a customized LLVM intrinsic, which serves as placeholders rather than actual instructions, to avoid disturbing the optimization. During the optimization pipeline, GIANTSAN refines the intrinsics (e.g., remove the ones whose corresponding memory operations are removed by the optimization). The intrinsic is materialized to sanitization instructions at the late stage of the optimization pipeline. We discuss pseudo instrumentation in Section 4.4.

4 Design

In this section, we present the design of GIANTSAN, an efficient location-based sanitizer with high protection density. Like existing location-based methods, GIANTSAN needs redzones and memory quarantine for sophisticated errors.

Figure 5 illustrates GIANTSAN's general workflow. The runtime support library hooks the object's allocation to update the shadow memory, and the instrumentation uses the shadow memory to safeguard memory regions. Sections 4.1, 4.2, and 4.3 present detailed solutions to the three questions mentioned in Section 3. Section 4.4 describes how to generate operation-level checks with pseudo instrumentation. In the end, Section 4.5 demonstrates the implementation details.

4.1 Shadow Encoding in GIANTSAN

In this section, we describe GIANTSAN's shadow memory encoding. We choose the commonly used eight-byte segment shadow memory as ASan [40]. The whole virtual memory is divided into small segments of 8 bytes, and

1:10 • H. Ling et al.

the metadata for a segment is stored in an 8-bit data type. Same as ASan, GIANTSAN ensures that all objects are 8-byte aligned, which does not make a huge difference to the memory layout because, as discussed in previous work [40], most objects in modern systems are naturally 8-byte aligned.

GIANTSAN achieves high protection density by building summaries on "good" segments, the ones containing no non-addressable bytes. The summary strategy is *binary folding*, which locates and folds consecutive 2^x "good" segments and encodes the value *x* in the shadow memory. The folded segment containing 2^x "good" segments is named as an (*x*)-*folded segment*. As illustrated in Figure 6, an *x* value in the shadow memory indicates **at least** 8×2^x and **less than** $8 \times 2^{x+1}$ consecutive bytes are addressable. In modern 64-bit systems, *x* cannot exceed 64 because the maximum object size is less than 2^{64} .

After introducing the folded segments, three categories of segment states exist: 1) the folding degree *i* for (*i*)-folded segments, 2) the value *k* for k-partial segments, which has only the first *k* bytes addressable, and 3) error codes for non-addressable segments. There are at most 64 different *i* and 7 different *k*. We use the denotation m[p] to represent the metadata stored in the *p*-th shadow byte, and m[p] is defined as follows:

DEFINITION 1 (STATE CODE). m[p] is an 8-bit unsigned integer that can store values within [0,256).

$$m[p] = \begin{cases} 64 - i, & \text{the } p\text{-th segment is an (i)-folded segmen} \\ 72 - k, & \text{the } p\text{-th segment is a } k\text{-partial segment} \\ > 72, & \text{error codes} \end{cases}$$

The monotonicity of *m* simplifies memory checks. A smaller m[p] means more consecutive addressable bytes following the *p*-th segments. Suppose that we want to check whether the *p*-th segment is a folded segment with a folding degree equal to or higher than 3. In that case, we only need to check whether $m[p] \le 64 - 3$. Any m[p] breaking the inequality indicates that there are non-addressable bytes in the memory region $[8p, 8(p+2^3))$. Checking the folding degree is the key to memory protection, which is discussed later in Section 4.2.

Though the encoding is much more complicated than existing works [2, 9, 10, 36, 40, 42, 49, 50], updating the shadow memory with the new encoding does not take extra computation. Technically, an allocated object has at most one partial segment, and all remaining segments within the allocated regions are folded. More formally, there are 2^i consecutive (*i*)-folded segments, e.g., there is one (0)-folded segment, two (1)folded segments, and four (2)-folded segments. The relative positions of the folded segments follow a simple pattern illustrated in Figure 6. Based on this pattern, GIANTSAN efficiently updates the shadow memory in linear time, the same as existing works.

4.2 Region Checking

This section introduces how to use the new shadow memory encoding to safeguard a memory region. A memory region [L, R) is safe if all except the last seg-



Fig. 6. Shadow memory encoding for an object sized 68 bytes. "(*i*)" represents an (*i*)-folded segment. "4-part" represents a partial segment with only the first 4 bytes addressable.

ment within this region are "good" segments and the first ($R \mod 8$) bytes in the last segment are addressable. GIANTSAN speeds up the "good" segment checking with folded segments. Specifically, GIANTSAN generates codes to safeguard a memory region [L, R), denoted as CI(L, R), in two steps. Let $l = \lfloor \frac{L}{8} \rfloor$, $r = \lfloor \frac{R}{8} \rfloor$:

• The *l*-th, \cdots , (r - 1)-th segments must all be "good".

GiantSan: Efficient Operation-Level Memory Sanitization with Segment Folding • 1:11



Fig. 7. Checking whether the *l*-th, *l*+1-th, \cdots , (r-1)-th segments are all "good" based on folded segments. $t = \lfloor \log_2(r-l) \rfloor$.

Algorithm 1 CI(L, R). m is the shadow memory, and L is a multiple of 8 due to the 8-	byte-alignment strategy.
1: uint8_t $v = m[\frac{L}{8}]$	$\triangleright L \equiv 0 \pmod{8}$
2: uintptr_t $u = (v \le 64) \ll (67 - v);$	
3: if $u < R - L$ then	▶ fast check
4: if $R - L \ge 8$ then	
5: if $2 * u < R - L$ then	 check folding degree
6: ReportError()	▹ of the prefix
7: end if	
8: if $m[\lfloor \frac{R-u}{8} \rfloor] \neq v$ then	▹ check folding degree
9: ReportError()	▶ <u>of the suffix</u>
10: end if	
11: end if	
12: if $m[\lfloor \frac{R-1}{8} \rfloor] > 72 - (R\&7)$ then	▹ check the partial
13: ReportError()	▹ segment at the
14: end if	⊳ end
15: end if	

• The first (*R* mod 8) bytes in the *r*-th segment are addressable.

Arbitrary N consecutive "good" segments must be a union of two $(\lfloor \log_2 N \rfloor)$ -folded segments. As illustrated in Figure 7a, if all 10 consecutive segments are "good", the first eight and the last eight "good" segments must be at least (3)-folded. Therefore, we only need to check if the folding degrees of a prefix and a suffix in the segment sequence are large enough. There are only two cases when all segments numbered from l to r - 1 are "good" $(t = \lfloor \log_2 r - l \rfloor)$:

- All segments are folded into one, and at least one (t + 1)-folded segment exists, as illustrated in Figure 7b.
- All segments are divided into two (t)-folded segments, as illustrated in Figure 7c.

An important integer trick for efficient checking is that the number of addressable bytes recorded in the *p*-th segment is $(m[p] \le 64) \ll (67 - m[p])$, where \ll is the left-shift arithmetic. The calculation result becomes 0 if m[p] does not represent a folded segment (i.e., m[p] > 64). The trick helps avoid calculating the expensive log₂ function.

Algorithm 1 shows how to safeguard the interval [L, R). It contains two stages: the *fast check* (the case in Figure 7b) and the slow check (the case in Figure 7c). The fast check is cheap and suffices to safeguard most memory regions, while the slow check handles the remaining rare cases.

• The fast check (Lines $1 \sim 3$) finds a safe region [L, L + u) without non-addressable bytes based on the folded segment recorded at $m[\frac{L}{8}]$. If [L, R) is within [L, L + u), [L, R) must be safe. According to the definition

1:12 • H. Ling et al.

of the folded segment, u covers > 50% of the addressable bytes following L; thus, u is large enough to safeguard the majority of the regions.

• The slow check (Lines $4\sim14$) ³ verifies three parts: whether non-addressable bytes exist in 1) the prefix, 2) the suffix, and 3) the last segment of [*L*, *R*). The slow check handles the case illustrated in Figure 7c, which is much more infrequent than the cases that are handled by the fast check.

This algorithm fully utilizes folded segments: folded segments summarize the majority (> 50%) of neighboring bytes in arbitrary safe regions, and the fast check efficiently safeguards any region within an existing summary. The region outside the fast check's scope is split into (at most) two folded segments and handled by the slow check, which is invoked only when the fast check fails. The slow check is also an O(1)-time algorithm with a better time complexity than existing location-based methods. Therefore, this algorithm can check a region with arbitrary size in constant time.

4.3 History Caching

History caching helps reduce metadata loadings on the same pointer. Intuitively, caching mainly speeds up memory protection within loops (the number of accesses outside loops is relatively limited). Thus, to better illustrate our method, we explain GIANTSAN's caching solution with accesses in loops.

The ideal values to be cached are the bounds of pointers since memory accesses falling within the bound do not need extra metadata. GIANTSAN can locate the bound by skipping over folded segments, as illustrated in Figure 8. The number of skipping is at most $\lceil \log_2 \frac{n}{8} \rceil$, where *n* is the size of the object, because the folding degree decreases by at least 1 after one skip and the maximum folding degree is $\lceil \log_2 \frac{n}{8} \rceil$.





Although the skipping is fast, it still takes time and is not a constant-time process. Therefore, GIANTSAN employs *on-demand skipping* to save time. Whenever GIANTSAN conducts a pointer dereference check, GIANTSAN caches the maximum valid address (called the *quasi-bound*) implied by the folded segment examined. In future dereference, the bound checks can use the quasi-bound until the dereference goes beyond the quasi-bound. GIANTSAN gets a new maximum valid address from the new folded segment visited, and reduces metadata loadings with the quasi-bound.

Figure 10 demonstrates caching logic for the memory access at Line 10 in Figure 9a. GIANTSAN creates a local variable, *ub*, as the quasi-bound for the buffer *y*. As illustrated in Figure 10, initially, the quasi-bound equals 0 because the size of the buffer is unknown. During the execution of the loop, GIANTSAN checks whether the offset j is beyond the quasi-bound (Line 4). If it goes beyond the bound, GIANTSAN checks y[j] individually (Line 5) and updates *ub* (Line 7). After the quasi-bound update, *ub* is closer to the actual bound of the region, and as discussed above, the number of *ub*'s updating is at most $\lceil \log_2 \frac{n}{8} \rceil$. Further memory accesses on *y* that fall within the quasi-bound do not need additional metadata loadings and speed up the runtime checks.

GIANTSAN also detects underflow (Lines 9-11) and temporal errors (Line 14). Technically, GIANTSAN does not create a quasi-lower bound because it is widely reported [25, 31] that the number of accesses with negative offsets is far less frequent than positive offsets. Therefore, using a dedicated CI to check underflow results in negligible cost. Moreover, the object pointed by y can be freed during the loop execution, and a final check after the loop can capture the deallocation [50].

³Codes at Lines 4, 12-14 are unnecessary if (R - L) mod 8 = 0 can be proved with static type information, e.g., reading an array of *int64_t*.

```
void foo(int **p, int N) {
1
2
3
     int *x = p[0]:
4
     int *y = p[1];
5
     for (int i = 0; i < N; i++) {</pre>
6
7
       int j = x[i];
8
9
10
       y[j] = i;
11
     }
12
     memset(x, 0, N * sizeof(int));
13
14 }
```

```
void foo(int **p, int N) {
1
2
     CI(p, p + 4);
      int *x = p[0];
3
     CI(p, p + 8);
int *y = p[1];
for (int i = 0; i < N; i++) {
4
5
6
        CI(x, x + 4 * i + 4);
7
        int j = x[i];
8
        CI(y, y + 4 * j + 4);
9
10
              = i:
        y[j]
11
     CI(x, x + 4 * N);
12
      memset(x, 0, N * sizeof(int));
13
14 }
```

(a) Source Code

(b) Check Instances (before merging)

```
void foo(int **p, int N) {
1
2
     CI(p, p + 8);
3
     int *x = p[0];
4
     int *y = p[1];
5
     CI(x, x + 4 * N);
6
     for (int i = 0; i < N; i++) {
7
8
        int j = x[i];
9
       CI(y, y + 4 * j + 4) (cached)
10
            = i;
       y[j]
11
     }
12
13
     memset(x, 0, N * sizeof(int))
14 }
```

(c) Check Instances (after merging and caching)

Fig. 9. Operation-level protection that significantly reduces runtime checks and metadata loadings

```
1 uintptr_t ub = 0;
                         i < N; i++)
2
   for (int i = 0;
                                          {
      int j = x[i];
3
      if (4 * j >= ub))
4
                              {
         CI(y, y + 4 * j + 4);
5
           [(y, y + 4 * j) >> 3];

= m[(y + 4 * j) >> 3];

h = 4 * j + (v <= 64) << (67)
6
         v
                                                  - v);
         ub
7
8
      }
9
      if (j < 0) {
         CI(y + 4 * j, y);
10
11
      3
               i;
      y[j] =
12
13 }
14 CI(y, y + ub);
```

Fig. 10. Quasi-bound instrumentation for y[j] in Figure 9a (Line 10) to reduce metadata loadings with caching.

4.4 The Construction of Operation-Level Protection

Due to the ability to handle arbitrary memory regions and history caching, GIANTSAN enables operation-level protection, significantly reducing the number of runtime checks. Some kinds of operation-level memory protection rely on program analysis (e.g., loop analysis). Lowering the abstraction level affects the performance gains from

1:14 • H. Ling et al.

the operation-level protection. However, even without any analysis, operation-level protection will not be less effective than instruction-level protection, because instruction-level protection operates at the lowest abstraction level. Building check instances based on programs with high-level abstraction makes GIANTSAN faster by fostering operation-level protection.



Fig. 11. The number of memory instructions at different stages in the optimization pipeline (based on the default extension points provided by the LLVM project).

The key challenge in establishing operation-level protection is that it sometimes clashes with compiler optimization. While the optimization pipeline in the compilation process can efficiently eliminate redundant memory operations, thereby alleviating sanitization burdens, it concurrently lowers the program's abstraction level. Typically, a compiler operates through an optimization pipeline that sequentially invokes various passes to reshape the program. Within this pipeline, the program's intermediate representation undergoes alterations, gradually converging toward machine code. For example, some passes unroll the loops to foster vectorization for the SIMD instructions. This convergence is achieved by translating the code into instructions with lower levels of abstraction.

Notice that we cannot simply instrument the program before the optimization pipeline because instrumentation instructions have side effects and fail the optimization [48]. The failure of the optimization results in low-quality binary and significantly increases the runtime overhead. Consequently, most existing sanitizers instrument the programs at the late stage of the optimization pipeline (e.g., EP_OptimizerLast) to avoid the optimization's failure. However, as Figure 11 shows, instruction duplication is a common practice in compiler optimization, and the last stage of the optimization pipeline does not contain the minimum number of memory instructions. It is essential to consider the impact of compiler optimization for sanitization carefully.

GIANTSAN aims to collect the check instances at the early stage of the optimization without failing the compilation optimization. We use *pseudo instrumentation* to separate the collection and instrumentation of CIs. The CI information is collected before the optimization is executed and is recorded by a dedicated LLVM intrinsic (called *giantsan.marker*). After the optimization, we materialize the intrinsic, i.e., replacing them with actual instrumentation instructions.

The check instances that are generated by the pseudo instrumentation need to be removed or refined during the optimization process. For example, dead code elimination removes memory operations with no side effects,

GiantSan: Efficient Operation-Level Memory Sanitization with Segment Folding • 1:15



Fig. 12. The workflow of the pseudo instrumentation and CI refinement.

and the corresponding pseudo instrumentation needs to be removed; otherwise, redundant sanitization will be introduced. The following scenarios require the removal of the intrinsics during optimization.

- Memory Operation Removal: Optimizations such as dead code elimination remove redundant memory operations. The corresponding protection should be removed.
- Memory Operation Merge: A series of instruction combination optimizations will combine memory operations that can be merged (e.g., operating on the same region of memory simultaneously). The corresponding protection should also be merged.

Since it is impossible to know if the corresponding memory operation will be removed or merged when injecting the pseudo instrumentation, we need a mechanism to reclaim the check instances during the optimization pipeline once they are removed or merged.

EXAMPLE 3. The unoptimized code in Figure 13a has two memory operations: the incrementing action inside a bounded loop and a memset idioms. Before the optimization pipeline, we generate two check instances (Lines 6 and 12). At the end of the optimization, the loop generates many memory instructions through code duplication. However, since the check instances are generated before the optimization is executed, they are not affected by the code duplication. The memset idiom will be deleted after the function bar is inlined (because it can never be triggered), and we need to delete the corresponding check instance accordingly.

Figure 12 shows the workflow of the pseudo instrumentation and CI refinement. At the early stage of the compiler optimization pipeline, GIANTSAN adds CI as LLVM intrinsic to the IR of the code. GIANTSAN first scans all memory operations in the code to collect check instances, which might contain redundant ones. The CI refinement stage will refine the check instances by removing the redundant check instances to simplify the code. For example, if two check instances for different memory operations have identical operands, they are considered redundant. The CI refinement is called multiple times because some analysis information is only available at specific stages. For example, the loop information on the LLVM SSA requires the canonical forms of loops [33], which show up at the loop optimization stage. GIANTSAN calls the CI refinement at the pass extension points provided by the LLVM projects (ModuleOptimizerEarly, CGSCCOptimizerLate, LateLoopOptimizations, LoopOptimizerEnd, ScalarOptimizerLate, VectorizerStart, OptimizerLats). These extension points are designed for customized transformation passes at different stages of the optimization pipeline.

During the compiler optimization pipeline, an associative CI removal component removes the CI when the corresponding memory operations are removed or merged. Specifically, when memory loads/stores are detached from the IR code, the corresponding CI is also removed. At the end of the IR optimization pipeline, the intrinsics are materialized into actual instructions that protect the memory operations. We discuss the CI refinement and the associative CI removal in the following.

4.4.1 *CI Refinement.* GIANTSAN first scans all instructions and memory intrinsic functions that manipulate the memory to generate the instruction-level checks. For example, there are five different codes accessing memory in

1:16 • H. Ling et al.

```
bool bar() {
                                                                  bool bar() {
1
                                                               1
2
     return false:
                                                               2
                                                                    return false:
3 }
                                                               3 }
4
                                                               4
                                                                  void foo(int *p, int *q, bool cond, int n)
5
   void foo(int *p, int *q, bool cond, int n)
                                                               5
6
  {
                                                               6
                                                                  {
7
     giantsan.marker(p, p + 4 * n);
                                                                    giantsan.marker(p, p + 4 * n);
                                                               7
     for (int i = 0; i < n; i++) {</pre>
8
                                                               8
                                                                    if (cond) {
9
                                                                      for (int i = 0; i < n; i++)
       p[i]++:
                                                               9
                                                                        p[i] += 2; // to be unrolled.
       if (cond) p[i]++;
10
                                                              10
11
                                                              11
     if (bar()) {
12
                                                              12
                                                                    else {
                                                                      for (int i = 0; i < n; i++)</pre>
13
       giantsan.marker(q, q + 4 * n);
                                                              13
                                                                        p[i]++; // to be unrolled.
14
       memset(q, 0, n * sizeof(int));
                                                              14
     }
15
                                                              15
                                                                    }
16 }
                                                              16 }
```

(a) Code before optimization

(b) Code after optimization

Fig. 13. Example about generating the check instances before the optimization.

Figure 9a. *reading* p[0] (Line 4), *reading* p[1] (Line 5), *reading* x[i] (Line 7), *writing into* y[j] (Line 8), and *memset for* x (Line 10). Figure 9b shows the checks generated in the first stage. GIANTSAN later merges checks with CI refinement; the final result is shown in Figure 9c. After the merging, only 2 checks and N cached checks are required, much fewer than the 2 + 3N checks in existing location-based methods.

Anchor-based Enhancement. Location-based methods insert *redzones* between objects to detect overflow. However, small redzones can be bypassed [19], while large redzones negatively impact memory performance. Our solution is to set a small redzone between objects and select an anchor point. When safeguarding memory accesses, GIANTSAN checks whether a redzone exists between the anchor point and the accessed location. For most memory accesses, the base pointer of a buffer is chosen as the anchor point ⁴. This optimization eliminates the trade-offs on redzone sizes and protects memory efficiently and precisely.

Take the memory access y[j] at Line 10 in Figure 9a as an example. Existing location-based sanitizers only check the region [y + 4j, y + 4j + 4) because they only protect the memory region at the instruction level. It can result in a false negative if *j* is large enough to bypass the redzone within [y, y + 4j) (if it exists). Existing methods have to enlarge the redzone size to avoid this false negative. Instead, GIANTSAN uses the base pointer *y* as the anchor point and checks the region [y, y + 4j + 4) to ensure y[j] is indeed a valid location within the same memory region as *y*. This method only requires a one-byte redzone, eliminating the need to use large redzones and significantly increasing runtime efficiency.

Check-in-Loop Promotion. Memory accesses in loops can raise multiple checks during the execution (e.g., Line 7 and Line 9 in Figure 9b). GIANTSAN runs SCEV analysis [32] to identify bounded loops and reduce runtime checks. For example, the *N* checks at Line 7 in Figure 9b are combined into one check CI(x, x + 4N). For unbounded loops, GIANTSAN employs the history caching discussed in Section 4.3.

Redundancy Elimination based on Alias. Existing efforts [9, 10, 28, 40, 50] demonstrate that sanitization tasks could be removed or merged (e.g., p[0] and p[1] in Figure 9a) to reduce the number of memory region safeguarding

⁴Some programmers would purposely employ undefined behaviors, e.g., using an out-of-bound base pointer to simulate 1-based arrays, which we consider as bugs by default. GIANTSAN can use the first dereferenced address as the anchor point to turn off the warning for the undefined behaviors.

requests if the accessed pointers are must-aliased and have dominance/post-dominance relationship. GIANTSAN adopts the LLVM's intra-procedural must-alias analysis to detect aliased checks.

Loop Independent Promotion. When the parameters of a CI inside a loop do not change with the number of loops, GIANTSAN moves the CI outside the loop. Note that a loop-independent variable is not the same as a loop invariant; a CI checks the address of a pointer, while a loop invariant examines the data contained in that address.

Full Dominance Promotion. When all the successors of a basic block contain the same CIs, they can be merged and moved to that basic block. It aims to handle complex dominance relationships. In alias-based redundancy elimination, the condition for two CIs to become redundant is the existence of dominance and the post-dominance relationship between them, which cannot handle the case in Figure 14 where two CIs are not dominated nor post-dominated by each other but can still be merged.





4.4.2 Intrinsics Property and Associative CI Removal.

The check instances generated by the pseudo instrumentation are represented by the LLVM intrinsic. The influence of LLVM Intrinics on the optimizers is mainly controlled through its memory property (no memory access / read-only / write-only / only accesses memory that is not accessible by the module being compiled). In order to minimize the impact on the optimizer, GIANTSAN sets the memory property of *giantsan.marker* to "no memory access", eliminating any potential side effect on the optimization pipeline. However, this property conflicts with dead code elimination. When the return value of the intrinsics has no users, they are recognized as dead code and removed. The logical user of *giantsan.marker* is its corresponding memory instructions, i.e., the intrinsics can only be removed when all their corresponding memory instructions have been removed, but these memory instructions do not explicitly utilize the return value of *giantsan.marker* for their computation. From the perspective of the SSA, the memory instructions are not the users of the *giantsan.marker* intrinsic.

We describe the implicit def-use relationship with *pseudo operand*. For each memory instruction, in addition to the operands required for its own computation, GIANTSAN adds an additional pseudo operand, pointing to the corresponding *giantsan.marker*. In all transformation passes involving dead code elimination, in addition to marking dead/alive according to the SSA's def-use chain, GIANTSAN also uses the pseudo operand to mark whether a *giantsan.marker* is dead or alive. A *giantsan.marker* can only be marked as dead if all of its associated memory instructions are dead. The same processing is also included in the passes associated with the instruction combination. When memory instructions are combined, the corresponding pseudo operands are also combined.

4.5 Implementation

GIANTSAN is built upon the infrastructure of ASan [40] in the LLVM Project. There are two components in the LLVM project related to memory sanitization: 1) a compilation pass that inserts runtime checks and 2) a library providing the runtime environment. Specifically, GIANTSAN modifies the framework in two aspects: the shadow memory poisoning to build folded segments and the detection logic to construct operation-level protection. **Shadow Poisoning.** GIANTSAN changes the way ASan poisons the shadow memory to build the folded segment summary. Specifically, instead of only marking the allocated region addressable (e.g., filling the shadow memory with zero values), GIANTSAN sets the folding degrees in the shadow locations of the allocated region. The other operations, e.g., redzone setting and memory unpoisoning, remain unchanged. The instrumentation is

1:18 • H. Ling et al.

implemented on top of the ASan instrumentation pass. The compilation front end controls the location of the pass in the compilation pipeline. By default, this pass is placed at the end of the optimization pipeline.

Runtime Checking. GIANTSAN changes the logic of runtime checks, prompting the instruction-level protection to the operation-level protection. ASan adds runtime protection in two ways. First, ASan employs an instrumentation pass to add runtime checks during the compilation; we modify this pass to replace ASan's runtime protection with GIANTSAN's operation-level protection. Second, ASan provides a runtime guardian function invoked before calling standard functions (e.g., strcpy). The guardian function checks contiguous regions in linear time, and we modify its implementation into GIANTSAN's constant time check.

Other implementation aspects of GIANTSAN, including shadow memory construction, shadow memory unpoisoning after object deallocation, redzone padding, and memory quarantine, are the same as the ones of ASan. Notably, the multi-thread guarantee of GIANTSAN is the same as ASan, i.e., thread-local caches are utilized to avoid locking on every call of the *malloc* and *free* functions. The collection of check instances with pseudo instrumentation starts at *EP_ModuleOptimizerEarly*, the pass extension points provided by the LLVM project.

5 Evaluation

We experimentally evaluate GIANTSAN on four questions:

- RQ1: Can GIANTSAN reduce runtime overhead?
- RQ2: What are the impacts of each optimization techniques in GIANTSAN?
- RQ3: Does pseudo instrumentation cooperate well with compiler optimization?
- RQ4: Can GIANTSAN effectively detect real bugs?

We evaluate the speed of GIANTSAN on the latest version of the industry-standard benchmark suite, SPEC CPU 2017 [44] (**RQ 1**), and conduct an ablation study to evaluate the impact of different optimizations employed by GIANTSAN with the same benchmark (**RQ 2**, **RQ 3**). We then use Juliet Test Suite [37], Magma Benchmark [22], and the Linux Flaw Project [8], the widely used vulnerability databases, to evaluate GIANTSAN's detection ability (**RQ 4**).

Configuration. GIANTSAN is built on the LLVM-12, and the experiments are conducted on a workstation with Intel(R) Xeon(R) CPU E5-2698 v3 @ 2.30GHz CPU, 128G memory (OS: ubuntu 18.04, Kernel version: 4.15.0-117-generic).

As for the sanitizer configuration, we use the default settings listed in the ASan documentation [15] for all ASan-based implementations: ASan [40], ASan-- [40], and our tool GIANTSAN, except setting *halt_on_error=false* to prevent early termination of the evaluation due to the widely-reported memory errors existing in the SPEC benchmark.

5.1 Performance Study

Setting. We use the latest version of the industry-standard CPU-intensive benchmark suite, SPEC CPU 2017 [44], to evaluate the performance improvement of GIANTSAN thoroughly. This benchmark consists of two testing modes: speed test and rate test. The speed test runs one copy of the target program to evaluate the execution time under the intensive CPU computation environment. The rate test runs multiple concurrent programs simultaneously to evaluate the throughput and performance in multi-threaded environments.

Not all programs in the benchmark are selected due to compilation issues (e.g., requiring Fortran instead of C/C++). We test projects on which at least one sanitizer can work and choose the *ref* workloads for all projects.

We choose ASan [40] (the most widely adopted location-based sanitizers) and ASan-- [50] (the state-of-the-art redundant check eliminating solution based on static analysis) as the baseline of location-based methods. We plan to use BBC [2] as the baseline of rounded-up allocation size methods, but it is not publicly available. Instead, we choose LFP [9, 10], an improved version of BBC with more variety of allocation sizes for object allocation.

Table 2. Runtime Overhead (seconds). *R* is the ratio compared to the native execution (RE: Runtime Error, CE: Compile Error). *CacheOnly* is the GIANTSAN version with history caching optimization and pseudo instrumentation, *EliminationOnly* is the one with check elimination and pseudo instrumentation, and *NoPseudo* is the one without pseudo instrumentation. The redzone sizes for location-based methods (GIANTSAN, ASan, and ASan--) are the default value (16 bytes).

	Performance Study							Ablation Study				
Programs	Native	GiantSan	R	ASan	R	ASan	R	LFP	R	CacheOnly	EliminationOnly	NoPseudo
500.perlbench_r	358	704	196.65%	822	229.61%	780	217.88%	CE	-	218.16%	224.02%	200.56%
502.gcc_r	256	713	278.52%	847	330.86%	729	284.77%	CE	-	294.14%	287.89%	278.91%
505.mcf_r	399	505	126.57%	667	167.17%	551	138.10%	602	150.88%	148.37%	144.36%	127.82%
508.namd_r	295	310	105.08%	665	225.42%	479	162.37%	675	228.81%	189.83%	172.88%	107.46%
510.parest_r	430	577	134.19%	1314	305.58%	886	206.05%	CE	-	208.84%	169.77%	136.05%
511.povray_r	426	1060	248.83%	1604	376.53%	1235	289.91%	1227	288.03%	265.02%	285.21%	250.70%
519.lbm_r	275	283	102.91%	431	156.73%	347	126.18%	554	201.45%	126.55%	126.91%	101.09%
520.omnetpp_r	343	629	183.38%	1010	294.46%	872	254.23%	532	155.10%	230.03%	240.52%	196.79%
523.xalancbmk_r	408	553	135.54%	739	181.13%	600	147.06%	418	102.45%	147.30%	150.98%	137.25%
531.deepsjeng_r	289	396	137.02%	587	203.11%	442	152.94%	595	205.88%	158.82%	163.67%	141.18%
538.imagick_r	499	659	132.06%	930	186.37%	863	172.95%	CE	-	139.08%	139.28%	136.47%
541.leela_r	456	657	144.08%	933	204.61%	808	177.19%	906	198.68%	166.01%	169.74%	145.61%
557.xz_r	362	401	110.77%	554	153.04%	488	134.81%	574	158.56%	161.33%	131.49%	114.64%
600.perlbench_s	349	723	207.16%	1113	318.91%	806	230.95%	CE	-	229.51%	233.52%	206.88%
602.gcc_s	476	603	126.68%	1341	281.72%	729	153.15%	RE		136.34%	131.93%	126.89%
605.mcf_s	788	1032	130.96%	1276	161.93%	1205	152.92%	1113	141.24%	132.74%	134.39%	134.77%
619.lbm_s	551	566	102.72%	676	122.69%	608	110.34%	535	97.10%	131.76%	135.57%	105.63%
620.omnetpp_s	323	673	208.36%	1042	322.60%	871	269.66%	518	160.37%	242.41%	256.04%	212.38%
623.xalancbmk_s	396	528	133.33%	714	180.30%	618	156.06%	417	105.30%	150.25%	156.57%	135.35%
631.deepsjeng_s	347	496	142.94%	750	216.14%	540	155.62%	705	203.17%	172.33%	174.64%	143.52%
638.imagick_s	2119	2544	120.06%	3751	177.02%	4271	201.56%	3604	170.08%	126.52%	136.81%	124.35%
641.leela_s	452	653	144.47%	1041	230.31%	816	180.53%	904	200.00%	171.02%	172.35%	148.01%
644.nab_s	1198	1319	110.10%	1915	159.85%	1480	123.54%	1464	122.20%	137.56%	139.07%	113.11%
657.xz_s	871	1049	120.44%	1323	151.89%	1342	154.08%	1240	142.37%	157.52%	147.76%	119.98%
Geometric Means.			143.48%		212.58%		174.89%		161.76%	171.32%	170.17%	146.04%

Results. The overall performance is shown in Table 2. LFP fails to build four projects *perlbench, gcc, parest,* and *imagick*. On average, GIANTSAN introduces 43.48% execution overhead on the native execution, with 61.37%, 41.94%, and 29.53% improvements over ASan, ASan--, and LFP, respectively. GIANTSAN outperforms ASan and ASan-- on all projects and is only slower than LFP on 5 out of the 24 projects. The result shows GIANTSAN has the best average performance, indicating the effectiveness of the new shadow encoding with the segment folding algorithm.

5.2 Ablation Study

This section breaks down the contributions of the three optimizations introduced in Section 4.2, Section 4.3 and Section 4.4: large region checks help eliminate unnecessary checks, history caching reduces unnecessary metadata loading, and the pseudo instrumentation collects check instances before the abstraction level is lowered.

Figure 15 demonstrates the ratio of optimized check codes in GIANTSAN by our optimizations. On average, 52.56% of the checks are optimized (30.76% eliminated and 21.80% cached). In the projects *mcf, namd*, and *lbm*, more than 80% of the checks introduced by ASan are eliminated or cached. Most of the checks in these projects are within simple loops and structure accesses with constant offsets, which our optimizations can efficiently handle. The remaining unoptimized codes include the ones that employ the fast check only and those that require the full check (i.e., fast check + slow check). GIANTSAN can remove some slow checks because memory regions of specific constant sizes (e.g., a power of 2) do not require the slow check to tackle the corner cases outside the fast check's scope. The data shows that 49.22% of the remaining unoptimized tasks only use fast checks. The result



Fig. 15. The proportion of memory instructions handled by different optimizations in GIANTSAN with ASan as the baseline. The x-labels are the project IDs. *Eliminated* are codes removed due to the check merging, and *Cached* are the ones optimized by the caching. *FastOnly* are the codes where the fast check suffices, and *FullCheck* are the ones that require both fast check and slow check.

indicates that the optimizations significantly reduce runtime checks and metadata loadings to help GIANTSAN gain high efficiency, and the fast check suffices to cover the majority of protection tasks.

The ablation study column in Table 2 shows the runtime overhead of GIANTSAN without caching, check elimination, or pseudo instrumentation, respectively. On average, compared to ASan, GIANTSAN-CacheOnly, GIANTSAN-EliminationOnly, and GIANTSAN-NoPseudo show 36.64%, 37.67%, and 59.10% improvements, respectively. Meanwhile, with either optimization enabled, GIANTSAN has comparable efficiency to ASan-- and LFP with about 70% overhead, and combining both optimizations achieves the best performance among all test configurations. GIANTSAN is faster than ASan because it supports operation-level protection with constant time region checks and history caching. Though ASan-- also uses static analysis to reduce redundant checks (it has a similar efficiency with GIANTSAN-Elimination-Only), it does not support the history cache that can further reduce runtime overhead. GIANTSAN is faster than LFP because LFP has to use extra instructions to simulate the stack due to the incomplete stack protection caused by the high memory alignment requirement. This result shows that both optimizations in GIANTSAN have significantly contributed to reducing the number of checks, and the fast check covers most of the memory protection tasks, allowing us to achieve a notable performance improvement.

5.3 The Impact of Pseudo Instrumentation

The motivation of pseudo instrumentation is to collect the information essential for operation-level sanitization early in the compilation optimization pipeline, before the abstraction level is lowered, without failing the compilation optimization. The LLVM framework provides compilation statistics (enabled by *-stats* option), which outputs the number of instructions that were successfully analyzed, simplified, or removed in specific transformation passes (e.g., the number of loops or memory instructions that were optimized). We demonstrate these statistics to observe how instrumentation affects optimizers. Figure 16 shows a report of 11 statistics directly related to performance. The results show that introducing pseudo-instrumentation has a negligible effect on the statistics. Forcing early sanitization without pseudo instrumentation results in various performance degradations. For example, early sanitization without pseudo instrumentation significantly lowers the number of vectorized loops because all sanitizer instructions without pseudo instrumentation have side effects (e.g., they may cause



Fig. 16. The compilation statistics reported by the compiler. The data shown is the ratio compared to native compilation (i.e., compilation without instrumentation).



Fig. 17. The number of check instances across various instrumentation stages.

1:22 • H. Ling et al.

```
1 void foo(int *p, int *q, bool cond, int n) {
2
     if (cond) {
        giantsan.marker(p, p + 4 * (n & 0x7FFFFF8) + 4);
3
4
        int i = 0.
        for (i = 0; i != (n & 0x7FFFFF8); i += 4) {
5
6
          p[i] += 2; p[i + 1] += 2;
7
          p[i + 2] += 2; p[i + 3] += 2;
8
        for ( ; i < n; i++) {</pre>
9
10
          giantsan.marker(p, p + 4 * i);
          p[i] += 2;
11
12
        }
13
14
     else {
15
        giantsan.marker(p, p + 4 * (n & 0x7FFFFF8) + 4);
16
        int i = 0;
17
        for (i = 0; i != (n & 0x7FFFFF8); i += 4) {
          p[i]++; p[i + 1]++;
p[i + 2]++; p[i + 3]++;
18
19
20
21
        for ( ; i < n; i++) {</pre>
22
          giantsan.marker(p, p + 4 * i);
23
          p[i]++;
24
        }
25
     }
26
  }
```

Fig. 18. The number of checks increases if the pseudo instrumentation is performed after loop vectorization. After loop fission, the original single bounded loop is divided into a bounded loop and an unbounded loop. This new structure presents challenges for SCEV analysis, resulting in redundant instance checks.

program interruptions). The results show that utilizing pseudo instrumentation to collect the check instances at the early stage of the compiler optimization pipeline has negligible influence on the optimization effect.

Instrumentation at different optimization stages results in different numbers of check instances. Figure 17 shows how the number of check instances changes when pseudo instrumentation is performed at different stages. We can see that instrumenting at EP_ModuleOptimizerEarly has the minimum number of check instances. Instrumenting at the last stage of the optimization pipeline introduces 7.29% more check instances compared to instrumenting at EP_ModuleOptimizerEarly. The reason is that the abstraction level has been significantly reduced, many memory operations have been duplicated, and only limited information can be used to infer the operation-level protection. For example, Figure 18 shows the result if the pseudo instrumentation is injected after loop vectorization for the code in Figure 13a. Compared to the early injected one (i.e., Figure 13b), Figure 18 introduces 4x more check instances because the SCEV analysis cannot handle the loop that has been fissioned.

5.4 Detectability Study

On top of the performance improvement, we also evaluate the practicalness of GIANTSAN in detecting memory errors.

Setting. We evaluate the bug detection ability on Juliet Test Suite (version 1.3) [37], Magma [22], and Linux Flaw Project [8], which are error collections widely used to evaluate the effectiveness of software assurance tools.

Juliet Test Suite contains cases that wait for an external signal (e.g., sockets), and some test cases include a randomized version (triggered with probability). We remove these cases to avoid infinitely waiting and non-deterministic results. Linux Flaw Project contains CVEs related to real-world programs, and we pick the memory-related ones, including 28 vulnerabilities from 8 programs written in C/C++. Magma [22] provides 58,969 test

Table 3. Detection capability on the Juliet Test Suite. All test cases have two versions: buggy and non-buggy versions. All
tested tools have no false-positive issues under the C/C++ standard and pass all the non-buggy tests. Therefore, only the
results for the buggy versions are presented to illustrate the false-negative issue.

CWE ID & Type	GiantSan	ASan	ASan	LFP	Total
121: Stack Buffer Overflow	1435	1435	1435	49	1439
122: Heap Buffer Overflow	1504	1504	1504	4	1504
124: Buffer Underwrite	767	767	767	767	767
126: Buffer Overread	441	441	441	352	449
127: Buffer Underread	916	916	916	916	916
416: Use After Free	393	393	393	393	393
476: NULL Pointer Dereference	288	288	288	288	288
761: Free Pointer Not at Start of Buffer	192	192	192	192	192
Total	5063	5063	5063	2088	5075
 127: Buffer Underread 416: Use After Free 476: NULL Pointer Dereference 761: Free Pointer Not at Start of Buffer Total 	916 393 288 192 5063	916 393 288 192 5063	916 393 288 192 5063	916 393 288 192 2088	916 393 288 192 5075

cases collected from its fuzzing campaign. We evaluate ASan, ASan-- and GIANTSAN on Magma to examine the effectiveness of GIANTSAN's anchor-based enhancement.

Results. Table 3 and Table 4 show the results on Juliet Test Suite and Linux Flaw Project, respectively. GIANTSAN, ASan, and ASan-- have the same results in all cases, while LFP has a significant number of false negatives in both benchmarks. LFP has many false negatives because it allocates objects with more spaces than the program requires, similar to BBC [2] discussed in Section 2.1. The cases missed by GIANTSAN, ASan, and ASan-- are potential overflow errors caused by uninitialized values. However, the uninitialized values loaded do not really trigger an overflow; thus, these tools do not generate bug reports since no overflow occurs.

For the redzone setting test, we evaluate GIANTSAN, ASan, and ASan-- on Magma, and the result is listed in Table 5. As we can see, GIANTSAN and ASan perform similarly in most projects. However, for large-scale project *PHP*, GIANTSAN reports 463 more cases than ASan and ASan-- (redzone=16) and 57 more cases than ASan and ASan-- (redzone=512). These false negatives are the POCs for *CVE-2018-14883* and are caused by the small redzone size. The result supports our conclusion in Section 4.4.1: insufficient redzone size leads to a false negative because of redzone bypassing, and GIANTSAN solves this with anchor-based enhancement.

5.5 Limitation

Because GIANTSAN only provides a single-sided summary, i.e., it summarizes segments from lower addresses to higher addresses, GIANTSAN may not effectively safeguard lower addresses given only higher addresses, causing potential efficiency deterioration in reverse traversals with unbounded loops when anchor-based enhancement is enabled.

To study this potential limitation, we conducted an additional study on Perlbench, which is a project in the SPEC CPU 2017 we used in Section 5.1. It is a program interpreter that intensively iterates the input buffer and contains different buffer iteration patterns, e.g., forward / reverse / random traversals. We evaluated the execution time to complete a traversal of the input buffer to compare the performance of GIANTSAN's history caching and ASan in different buffer traversal patterns. Each run is repeated 100 times to reduce variations, and the geometric mean is presented.

The results in Figure 19 show that GIANTSAN is 1.48x and 1.07x faster than ASan in random and forward traversals, respectively. However, due to the extra instructions to perform anchor-enhanced checks, GIANTSAN is 1.39x slower than ASan in reverse traversals. The reason is that GIANTSAN has one-sided complexity guarantees

Program	CVE ID	GiantSan	ASan	ASan	LFP
libzip	CVE-2017-12858	\checkmark	\checkmark	\checkmark	
	CVE-2017-9164	\checkmark	\checkmark	\checkmark	\checkmark
autotrace	CVE-2017-9165	\checkmark	\checkmark	\checkmark	
	CVE-2017-9166~9173	\checkmark	\checkmark	\checkmark	\checkmark
imageworsener	CVE-2017-9204~9207	\checkmark	\checkmark	\checkmark	\checkmark
lame	CVE-2015-9101	\checkmark	\checkmark	\checkmark	\checkmark
zziplib	CVE-2017-5976~5977	\checkmark	\checkmark	\checkmark	\checkmark
libtiff	CVE-2016-10270~10271	\checkmark	\checkmark	\checkmark	\checkmark
	CVE-2016-10095	\checkmark	\checkmark	\checkmark	\checkmark
potrace	CVE-2017-7263	\checkmark	\checkmark	\checkmark	$\overline{\mathbf{A}}$
mp3gain	CVE-2017-14407~14408	\checkmark	\checkmark	\checkmark	\checkmark
	CVE-2017-14409	\checkmark	\checkmark	$\overline{}$	

Table 4. Detection capability for CVEs in Linux Flaw Project.

Table 5. Detection capability in real-world projects from Magma Test Suite. rz is short for redzone.

Project (LoC)	ASan (rz=16)	ASan (rz=512)	ASan (rz=16)	ASan (rz=512)	GiantSan (rz=16)	Total
php (1.3M)	1556	1962	1556	1962	2019	3072
libpng (86K)	1881	1881	1881	1881	1881	1881
libtiff (91K)	9858	9858	9858	9858	9858	9858
libxml2 (284K)	30566	30566	30566	30566	30566	30574
openssl (535K)	46	46	46	46	46	1509
sqlite3 (367K)	1528	1528	1528	1528	1528	1528
poppler (43K)	10201	10201	10201	10201	10201	10547

with history caching, i.e., quasi-bound converges to the upper bound of the allocated region in $\lceil \log_2 \frac{n}{8} \rceil$ time; however, it does not provide time guarantees for the lower bound. Therefore, GIANTSAN is able to save time by predicting the addressability of higher addresses from lower addresses, but not vice versa.

The experimental data empirically evidence the performance difference of our approach in handling different traversal patterns, which is consistent with our theoretical justification. Fortunately, the number of reverse traversals in real-world programs is relatively limited. For example, in the real-world programs collected by the SPEC CPU 2017, only 0.39% of the buffer traversals are in reverse order. Past studies [16, 25] show that the impact of underflow is comparatively less severe than overflow. Furthermore, the SCEV optimization could eliminate the runtime checks by inferring the loop bounds, if possible.

For programs that heavily use reverse traversals, several alternatives can mitigate the efficiency deterioration. One is to remove the anchor-based enhancement in underflow detection so that GIANTSAN's detection degrades to ASan's mode (i.e., only checking the location of the access and ignoring the anchor); however, this would eliminate the superiority of GIANTSAN over ASan w.r.t. underflow detection accuracy. The second solution is to locate the lower bound before buffer reverse traversals by enumerating the folding degrees and checking whether the corresponding folded segments exist.



(a) Forward Traversal: iterate over the buffer from the lowest address to the highest address

(b) Random Traversal: iterate over the buffer in random order

(c) Reverse Traversal: iterate over the buffer from the highest address to the lowest address

Fig. 19. The time cost of GIANTSAN and ASan in three buffer traversal patterns: Forward, Random, and Reverse. The baseline *Native* is the execution time without sanitization.

Also, though GIANTSAN improves the efficiency of location-based methods, it still shares some common limitations with existing works.

Sub-object Overflow Insensitivity: GIANTSAN detects memory accesses outside objects' allocated regions but cannot detect memory safety violations related to sub-objects, which is an open question in the existing literature. The best practices in detecting sub-object overflow are pointer-based methods like Softbound+CETS [34, 35] and EffctiveSan [11]. However, they all suffer from high runtime overhead and require precise type information, which might not be always available in real-world programs.

Quarantine Bypassing: GIANTSAN detects temporal errors based on memory quarantine, but the memory quarantine can be bypassed with a small probability. It is a common issue for memory quarantine-based solutions [40, 49, 50]. In practice, the probability of bypassing the quarantine queue is low, and few related false negative reports exist.

6 Related Work

Researchers have proposed various dynamic error detectors. We further discuss existing related works about memory error detection.

Token Authentication. HWASAN [41] uses *address tagging* to replace the redzone with token authentication. A random token is attached to pointers with the *Top-Byte-Ignore* hardware support, and the token is stored in the shadow memory for memory regions. The token mismatch between pointers and memory regions results in memory errors. IntegriTag [39] performs implicit probabilistic memory access checks with the Intel® TMEMK memory encryption hardware feature, increasing the detection probability with more sparse bits compared to HWASan. StickyTags [18] replaces random tagging with persistent memory tags to provide deterministic protection, utilizing a size-class allocation strategy similar to LFP. Like GIANTSAN's anchor-based enhancement, these works mitigate the redzone dilemma. Specifically, HWASAN solves the problem that traditional location-based methods are unable to distinguish between different allocated memory regions by assigning an 8-bit identifier to each region. It propagates the identifier in a pointer-based manner and removes the need for redzones with the token-matching model.

However, it does not improve the detection efficiency of the location-based methods, where a single check only safeguards a small region (e.g., 16 bytes). Therefore, it suffers from the low protection density issue that requires excessive runtime checks to safeguard a large region, decreasing its efficiency. This efficiency decreasing issue is exactly GIANTSAN's key motivation.

1:26 • H. Ling et al.

Redzone Enhancement. Location-based solutions divide the memory into separated regions using redzones to detect sophisticated bugs. Some methods that focus on redzone enhancement aim to reduce runtime overhead with redzone poisoning or improve accuracy with adaptive redzone.

For example, in-band redzone methods [17, 20] fill the redzone with a random pattern and compare the loaded data with that pattern. If they are different, the memory access is not in the redzone and is safe. These methods reduce dedicated data structure inquiries (e.g., shadow memory), thus promoting memory locality. However, this method protects only a small region with one check and faces the same low protection density issue as other location-based methods. Similarly, it suffers from small redzone size, e.g., FloatZone [17] cannot detect CVE-2017-7263 with 16-byte in-band redzones. These two issues are what GIANTSAN addresses.

Some approaches reduce the impact of redzone sizes with adaptive settings. LBC [20] selects different redzones based on the allocated region sizes. FuZZan [23] switches between different data structures (e.g., shadow memory and RB-tree) to decrease setup/inquiry costs for various inputs in a fuzzing scenario. MEDS [19] spreads the objects evenly in the address space to increase the distance between objects as much as possible. To minimize memory consumption, MEDS uses page aliasing to allow multiple virtual pages to share the same physical page, reducing the physical memory usage.

GIANTSAN is compatible with all these redzone enhancement techniques because GIANTSAN does not impose any extra requirements on redzone settings and the contents in the redzone areas. GIANTSAN only modifies the shadow memory encoding for non-redzone areas and reduces the dependency on redzone size by modifying the runtime check logic with the selected anchors.

Pointer Tracking. Pointer-based techniques provide a memory safety guarantee by tracking the lifetime of pointers. As discussed in Section 2.1, pointer-based methods require the pointer type information to propagate tags and avoid tag misuse. The complete memory safety guarantee in pointer-based methods requires instrumenting the source codes of the whole runtime environment, which is expensive and unavailable and thus makes these methods less portable.

Traditional pointer-based solutions [4, 34, 35] require extra instructions to propagate metadata (e.g., bound) along pointer arithmetics; in contrast, location-based solutions only check pointer dereference operations, which is much fewer than pointer arithmetics. The propagation is the primary source of the pointer-based solutions' runtime overhead [43]. *Pointer tagging* is a popular solution to mitigate the overhead issue in propagation. With the proliferation of large bit-width systems (e.g., 64-bit), a single pointer structure can now represent far larger address space than a program needs, resulting in some upper spare bits in pointers. Consequently, many pointer-based methods [16, 25, 26, 28, 47] propagate metadata with the upper spare bits so that the metadata associated with pointers can be propagated automatically.

Though pointer tagging solves the efficiency problem of data propagation, it faces a new problem related to the bit width: the upper spare bits are not enough to hold the metadata. One solution is reducing the address space. For example, Delta Pointers [25] and SGXBound [26] use 32-bit address space in a 64-bit platform and record the metadata with the other 32 bits. The narrowing down of the address space makes them less suitable for programs with large memory footprints. Delta Pointers mitigate this issue by providing a trade-off between the maximum object size and the address space size. Another solution [28] is to store the metadata in a key-value database, and the pointer tag only serves as the key. Compared with the shadow memory inquiry used in location-based solutions, the key-value store takes more time to retrieve the metadata.

GIANTSAN also suffers from a bit-width limitation, i.e., a single shadow byte can only hold 256 different states. GIANTSAN solves this limitation with the on-demand inquiry. The segment folding technique in GIANTSAN can be considered as a key-value store that takes logarithmic time to index an object's bound. However, one of our key observations is that the program does not always traverse the entire allocated region, and in most cases, we only need to safeguard a subregion. This observation allows us to reduce runtime queries by looking up folding degrees on demand. The spirit of on-demand inquiry is orthogonal to the pointer-based solutions and could mitigate the bit width requirement faced by the pointer tagging technique. Integrating the on-demand inquiry spirit into pointer-based solutions is a future research direction we are going to address.

Rounded-Up Bound. Works like LFP [9, 10], RedFat [12] and BBC [2] obtain the object bound by directly fetching the bound from shadow memory. However, to enable compact shadow memory, they only support a limited set of allocation sizes to reduce the bit width for recording the bound. As a result, they overapproximate the object sizes required by the programs, leading to significant false negative issues.

BBC [2] uses the power-of-two strategy similar to GIANTSAN from a particular perspective. However, BBC uses the power-of-two spirit to *approximate* the real object bound, while GIANTSAN uses the power-of-two spirit to build *precise summaries* of addressable regions. Therefore, GIANTSAN is more precise than BBC. LFP enhances BBC by introducing more variety of allocation sizes but still has numerous false negatives, as shown in our experiments.

7 Conclusions

We present GIANTSAN, a location-based sanitizer optimizing runtime checks with segment folding. GIANTSAN summarizes segments without non-addressable bytes to increase protection density. It largely reduces 61.37% and 41.94% of the overhead introduced by ASan and ASan-- on the SPEC CPU 2017 benchmark, respectively. Furthermore, the evaluation on the PHP project demonstrates that GIANTSAN can minimize the dependence on the redzone, thus resulting in a more effective detection ability than ASan and ASan--.

8 Acknowledgements

We thank the anonymous reviewers for their valuable comments and opinions for improving this work. This work is supported by the ITS/440/18FP grant from the Hong Kong Innovation and Technology Commission and research grants from Huawei, Microsoft, and TCL. Heqing Huang is the corresponding author.

References

- Sam Ainsworth and Timothy M. Jones. 2020. MarkUs: Drop-in use-after-free prevention for low-level languages. In 2020 IEEE Symposium on Security and Privacy (SP). 578–591. doi:10.1109/SP40000.2020.00058 https://doi.org/10.1109/SP40000.2020.00058.
- [2] Periklis Akritidis, Manuel Costa, Miguel Castro, and Steven Hand. 2009. Baggy Bounds Checking: An Efficient and Backwards-Compatible Defense against out-of-Bounds Errors. In Proceedings of the 18th Conference on USENIX Security Symposium (Montreal, Canada) (SSYM'09). USENIX Association, USA, 51–66. doi:10.5555/1855768.1855772 https://www.usenix.org/legacy/events/sec09/tech/full_papers/akritidis. pdf.
- [3] Android. [n. d.]. HWASan, ASan and KASAN. https://source.android.com/docs/security/test/memory-safety/hwasan-asan-kasan.
- [4] Nathan Burow, Derrick McKee, Scott A. Carr, and Mathias Payer. 2018. CUP: Comprehensive User-Space Protection for C/C++. In Proceedings of the 2018 on Asia Conference on Computer and Communications Security (Incheon, Republic of Korea) (ASIACCS '18). Association for Computing Machinery, New York, NY, USA, 381–392. doi:10.1145/3196494.3196540 https://doi.org/10.1145/3196494. 3196540.
- [5] Microsoft Corporation. 2000. How to use Pageheap.exe in Windows XP, Windows 2000, and Windows Server 2003. https://mskb. pkisolutions.com/kb/286470.
- [6] C. Cowan. 2003. Software security for open-source systems. IEEE Security & Privacy 1, 1 (2003), 38–45. doi:10.1109/MSECP.2003.1176994 https://doi.org/10.1109/MSECP.2003.1176994.
- [7] Thurston H.Y. Dang, Petros Maniatis, and David Wagner. 2017. Oscar: A Practical Page-Permissions-Based Scheme for Thwarting Dangling Pointers. In 26th USENIX Security Symposium (USENIX Security 17). USENIX Association, Vancouver, BC, 815–832. https://www.usenix. org/conference/usenixsecurity17/technical-sessions/presentation/dang https://www.usenix.org/conference/usenixsecurity17/technical-sessions/presentation/dang https://www.usenix.org/conference/usenixsecurity17/technical-sessions/presentation/dang https://www.usenix.org/conference/usenixsecurity17/technical-sessions/presentation/dang https://www.usenix.org/conference/usenixsecurity17/technical-sessions/presentation/dang https://www.usenix.org/conference/usenixsecurity17/technical-sessions/presentation/dang https://www.usenix.org/conference/usenixsecurity17/technical-sessions/presentation/dang https://www.usenix.org/conference/usenixsecurity17/technical-sessions/presentation/dang https://www.usenix.org/conference/usenixsecurity17/technical-sessions/presentation/dang https://www.usenix.org/conference/usenixsecurity17/technical-sessions/presentation/dang https://www.usenix.org/conference/usenixsecurity17/technical-sessions/presentation/dang.
- [8] Dongliang Mu. 2017. Linux Flaw Project. https://github.com/mudongliang/LinuxFlaw.
- [9] Gregory J Duck, Roland HC Yap, and Lorenzo Cavallaro. 2017. Stack Bounds Protection with Low Fat Pointers. doi:10.14722/ndss.2017. 23287 https://doi.org/10.14722/ndss.2017.23287.

- 1:28 H. Ling et al.
- [10] Gregory J. Duck and Roland H. C. Yap. 2016. Heap Bounds Protection with Low Fat Pointers. In Proceedings of the 25th International Conference on Compiler Construction (Barcelona, Spain) (CC 2016). Association for Computing Machinery, New York, NY, USA, 132–142. doi:10.1145/2892208.2892212 https://doi.org/10.1145/2892208.2892212.
- [11] Gregory J. Duck and Roland H. C. Yap. 2018. EffectiveSan: Type and Memory Error Detection Using Dynamically Typed C/C++. In Proceedings of the 39th ACM SIGPLAN Conference on Programming Language Design and Implementation (Philadelphia, PA, USA) (PLDI 2018). Association for Computing Machinery, New York, NY, USA, 181–195. doi:10.1145/3192366.3192388 https://doi.org/10.1145/ 3192366.3192388.
- [12] Gregory J. Duck, Yuntong Zhang, and Roland H. C. Yap. 2022. Hardening binaries against more memory errors. In Proceedings of the Seventeenth European Conference on Computer Systems (Rennes, France) (EuroSys '22). Association for Computing Machinery, New York, NY, USA, 117–131. doi:10.1145/3492321.3519580
- [13] Reza Mirzazade farkhani, Mansour Ahmadi, and Long Lu. 2021. PTAuth: Temporal Memory Safety via Robust Points-to Authentication. In 30th USENIX Security Symposium (USENIX Security 21). USENIX Association, 1037–1054. https://www.usenix.org/conference/ usenixsecurity21/presentation/mirzazade https://www.usenix.org/conference/usenixsecurity21/presentation/mirzazade.
- [14] GCC. [n. d.]. The GNU Compiler Collection. https://gcc.gnu.org/.
- [15] Google. [n. d.]. AddressSanitizier Wiki. https://github.com/google/sanitizers/wiki/AddressSanitizerFlags.
- [16] Amogha Udupa Shankaranarayana Gopal, Raveendra Soori, Michael Ferdman, and Dongyoon Lee. 2023. TAILCHECK: A Lightweight Heap Overflow Detection Mechanism with Page Protection and Tagged Pointers. In 17th USENIX Symposium on Operating Systems Design and Implementation (OSDI 23). USENIX Association, Boston, MA, 535–552. https://www.usenix.org/conference/osdi23/presentation/gopal https://www.usenix.org/conference/osdi23/presentation/gopal.
- [17] Floris Gorter, Enrico Barberis, Raphael Isemann, Erik van der Kouwe, Cristiano Giuffrida, and Herbert Bos. 2023. FloatZone: Accelerating Memory Error Detection using the Floating Point Unit. In 32nd USENIX Security Symposium (USENIX Security 23). USENIX Association, Anaheim, CA, 805–822. https://www.usenix.org/conference/usenixsecurity23/presentation/gorter https://www.usenix.org/conference/ usenixsecurity23/presentation/gorter.
- [18] Floris Gorter, Taddeus Kroes, Herbert Bos, and Cristiano Giuffrida. 2024. Sticky Tags: Efficient and Deterministic Spatial Memory Error Mitigation using Persistent Memory Tags. In 2024 IEEE Symposium on Security and Privacy (SP). 4239–4257. doi:10.1109/SP54263.2024. 00263
- [19] Wookhyun Han, Byunggill Joe, Byoungyoung Lee, Chengyu Song, and Insik Shin. 2018. Enhancing Memory Error Detection for Large-Scale Applications and Fuzz Testing. doi:10.14722/ndss.2018.23312 https://doi.org/10.14722/ndss.2018.23312.
- [20] Niranjan Hasabnis, Ashish Misra, and R. Sekar. 2012. Light-Weight Bounds Checking. In Proceedings of the Tenth International Symposium on Code Generation and Optimization (San Jose, California) (CGO '12). Association for Computing Machinery, New York, NY, USA, 135–144. doi:10.1145/2259016.2259034 https://doi.org/10.1145/2259016.2259034.
- [21] Reed Hastings. 1992. Purify: Fast detection of memory leaks and access errors. In Proc. 1992 Winter USENIX Conference. 125–136. https://web.stanford.edu/class/cs343/resources/purify.pdf.
- [22] Ahmad Hazimeh, Adrian Herrera, and Mathias Payer. 2021. Magma: A Ground-Truth Fuzzing Benchmark. Proc. ACM Meas. Anal. Comput. Syst. 4, 3, Article 49 (jun 2021), 29 pages. doi:10.1145/3428334 https://doi.org/10.1145/3428334.
- [23] Yuseok Jeon, WookHyun Han, Nathan Burow, and Mathias Payer. 2020. {FuZZan}: Efficient Sanitizer Metadata Design for Fuzzing. In 2020 USENIX Annual Technical Conference (USENIX ATC 20). 249–263.
- [24] Richard W. M. Jones and Paul H. J. Kelly. 1997. Backwards-Compatible Bounds Checking for Arrays and Pointers in C Programs. In Automated and Algorithmic Debugging. https://www.doc.ic.ac.uk/~phjk/Publications/BoundsCheckingForC.pdf.
- [25] Taddeus Kroes, Koen Koning, Erik van der Kouwe, Herbert Bos, and Cristiano Giuffrida. 2018. Delta Pointers: Buffer Overflow Checks without the Checks. In Proceedings of the Thirteenth EuroSys Conference (Porto, Portugal) (EuroSys '18). Association for Computing Machinery, New York, NY, USA, Article 22, 14 pages. doi:10.1145/3190508.3190553 https://doi.org/10.1145/3190508.3190553.
- [26] 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* (Belgrade, Serbia) (*EuroSys '17*). Association for Computing Machinery, New York, NY, USA, 205–221. doi:10.1145/3064176.3064192 https: //doi.org/10.1145/3064176.3064192.
- [27] Chris Arthur Lattner. 2002. LLVM: An infrastructure for multi-stage optimization. (2002). http://llvm.org.
- [28] Yuan Li, Wende Tan, Zhizheng Lv, Songtao Yang, Mathias Payer, Ying Liu, and Chao Zhang. 2022. PACMem: Enforcing Spatial and Temporal Memory Safety via ARM Pointer Authentication. In *Proceedings of the 2022 ACM SIGSAC Conference on Computer and Communications Security* (Los Angeles, CA, USA) (CCS '22). Association for Computing Machinery, New York, NY, USA, 1901–1915. doi:10.1145/3548606.3560598 https://doi.org/10.1145/3548606.3560598.
- [29] Hao Ling, Heqing Huang, Yuandao Cai, and Charles Zhang. 2025. Efficient Fuzzing Infrastructure for Pointer-to-Object Association. ACM Trans. Softw. Eng. Methodol. (April 2025). doi:10.1145/3730580 Just Accepted, https://doi.org/10.1145/3730580.
- [30] Linux Kernel. [n. d.]. The Kernel Address Sanitizer. https://www.kernel.org/doc/html/v4.14/dev-tools/kasan.html.

- [31] David Litchfield. 2005. Buffer Underruns, DEP, ASLR and improving the Exploitation Prevention Mechanisms (XPMs) on the Windows platform. Next Generation Security Software (2005). https://research.nccgroup.com/wp-content/uploads/episerver-images/assets/ 854f87540884465e8c6930b1b2fabf9b/854f87540884465e8c6930b1b2fabf9b.pdf.
- [32] LLVM. [n. d.]. Scalar Evolution and Loop Optimization. https://llvm.org/devmtg/2009-10/ScalarEvolutionAndLoopOptimization.pdf.
- [33] LLVM Project. 2024. LLVM Loop Terminology (and Canonical Forms). https://llvm.org/docs/LoopTerminology.html.
- [34] Santosh Nagarakatte, Jianzhou Zhao, Milo M.K. Martin, and Steve Zdancewic. 2009. SoftBound: Highly Compatible and Complete Spatial Memory Safety for c. In Proceedings of the 30th ACM SIGPLAN Conference on Programming Language Design and Implementation (Dublin, Ireland) (PLDI '09). Association for Computing Machinery, New York, NY, USA, 245–258. doi:10.1145/1542476.1542504 https://doi.org/10.1145/1542476.1542504.
- [35] Santosh Nagarakatte, Jianzhou Zhao, Milo M.K. Martin, and Steve Zdancewic. 2010. CETS: Compiler Enforced Temporal Safety for C. SIGPLAN Not. 45, 8 (jun 2010), 31–40. doi:10.1145/1837855.1806657 https://doi.org/10.1145/1837855.1806657.
- [36] Nicholas Nethercote and Julian Seward. 2007. Valgrind: A Framework for Heavyweight Dynamic Binary Instrumentation. In Proceedings of the 28th ACM SIGPLAN Conference on Programming Language Design and Implementation (San Diego, California, USA) (PLDI '07). Association for Computing Machinery, New York, NY, USA, 89–100. doi:10.1145/1250734.1250746 https://doi.org/10.1145/1250734. 1250746.
- [37] NIST. 2017. Software Assurance Reference Dataset. https://samate.nist.gov/SARD/test-suites.
- [38] Olatunji Ruwase and Monica S. Lam. 2004. A Practical Dynamic Buffer Overflow Detector. In Network and Distributed System Security Symposium. https://www.ndss-symposium.org/ndss2004/practical-dynamic-buffer-overflow-detector/ https://www.ndsssymposium.org/ndss2004/practical-dynamic-buffer-overflow-detector/.
- [39] David Schrammel, Martin Unterguggenberger, Lukas Lamster, Salmin Sultana, Karanvir Grewal, Michael LeMay, David M. Durham, and Stefan Mangard. 2024. Memory Tagging using Cryptographic Integrity on Commodity x86 CPUs. In 2024 IEEE 9th European Symposium on Security and Privacy (EuroS&P). 311–326. doi:10.1109/EuroSP60621.2024.00024
- [40] Konstantin Serebryany, Derek Bruening, Alexander Potapenko, and Dmitry Vyukov. 2012. AddressSanitizer: A Fast Address Sanity Checker. In Proceedings of the 2012 USENIX Conference on Annual Technical Conference (Boston, MA) (USENIX ATC'12). USENIX Association, USA, 28. doi:10.5555/2342821.2342849 https://dl.acm.org/doi/10.5555/2342821.2342849.
- [41] Kostya Serebryany, Evgenii Stepanov, Aleksey Shlyapnikov, Vlad Tsyrklevich, and Dmitry Vyukov. 2018. Memory tagging and how it improves C/C++ memory safety. arXiv preprint arXiv:1802.09517 (2018). https://arxiv.org/abs/1802.09517.
- [42] Julian Seward and Nicholas Nethercote. 2005. Using Valgrind to Detect Undefined Value Errors with Bit-Precision. In 2005 USENIX Annual Technical Conference (USENIX ATC 05). USENIX Association, Anaheim, CA. https://www.usenix.org/conference/2005-usenixannual-technical-conference/using-valgrind-detect-undefined-value-errors-bit https://www.usenix.org/conference/2005-usenixannual-technical-conference/using-valgrind-detect-undefined-value-errors-bit.
- [43] Dokyung Song, Julian Lettner, Prabhu Rajasekaran, Yeoul Na, Stijn Volckaert, Per Larsen, and Michael Franz. 2019. SoK: Sanitizing for Security. In 2019 IEEE Symposium on Security and Privacy (SP). 1275–1295. doi:10.1109/SP.2019.00010 https://doi.org/10.1109/SP.2019.00010.
 [44] Standard Performance Evaluation Corporation. 2022. SPEC CPU® 2017. https://www.spec.org/cpu2017/.
- [45] László Szekeres, Mathias Payer, Tao Wei, and Dawn Song. 2013. SoK: Eternal War in Memory. In 2013 IEEE Symposium on Security and Privacy. 48-62. doi:10.1109/SP.2013.13 https://doi.org/10.1109/SP.2013.13.
- [46] The 2022 CWE Top 25 Team. 2022. 2022 CWE Top 25 Most Dangerous Software Weaknesses. https://cwe.mitre.org/top25/archive/2022/ 2022_cwe_top25.html.
- [47] Kui Wang, Dmitry Kasatkin, Vincent Ahlrichs, Lukas Auer, Konrad Hohentanner, Julian Horsch, and Jan-Erik Ekberg. 2024. Cherifying Linux: A Practical View on using CHERI. In Proceedings of the 17th European Workshop on Systems Security (Athens, Greece) (EuroSec '24). Association for Computing Machinery, New York, NY, USA, 15–21. doi:10.1145/3642974.3652282
- [48] Mingzhe Wang, Jie Liang, Chijin Zhou, Zhiyong Wu, Xinyi Xu, and Yu Jiang. 2022. Odin: on-demand instrumentation with on-the-fly recompilation. In *Proceedings of the 43rd ACM SIGPLAN International Conference on Programming Language Design and Implementation* (San Diego, CA, USA) (*PLDI 2022*). Association for Computing Machinery, New York, NY, USA, 1010–1024. doi:10.1145/3519939.3523428 https://doi.org/10.1145/3519939.3523428.
- [49] Jiang Zhang, Shuai Wang, Manuel Rigger, Pinjia He, and Zhendong Su. 2021. SANRAZOR: Reducing Redundant Sanitizer Checks in C/C++ Programs. In 15th USENIX Symposium on Operating Systems Design and Implementation (OSDI 21). USENIX Association, 479–494. https://www.usenix.org/conference/osdi21/presentation/zhang https://www.usenix.org/conference/osdi21/presentation/zhang.
- [50] Yuchen Zhang, Chengbin Pang, Georgios Portokalidis, Nikos Triandopoulos, and Jun Xu. 2022. Debloating Address Sanitizer. In 31st USENIX Security Symposium (USENIX Security 22). USENIX Association, Boston, MA, 4345–4363. https://www.usenix.org/conference/ usenixsecurity22/presentation/zhang-yuchen https://www.usenix.org/conference/usenixsecurity22/presentation/zhang-yuchen.

Received 24 December 2024; revised 24 December 2024; accepted 27 May 2025