

Dynamic Heap Type Inference for Program Understanding and Debugging^{*}

Marina Polishchuk

Microsoft Corporation[†]
marinapo@microsoft.com

Ben Liblit

University of Wisconsin–Madison
liblit@cs.wisc.edu

Chloë W. Schulze

Oracle Corporation[†]
chloe.schulze@oracle.com

Abstract

C programs can be difficult to debug due to lax type enforcement and low-level access to memory. We present a dynamic analysis for C that checks heap snapshots for consistency with program types. Our approach builds on ideas from physical subtyping and conservative garbage collection. We infer a program-defined type for each allocated storage location or identify “untypable” blocks that reveal heap corruption or type safety violations. The analysis exploits symbolic debug information if present, but requires no annotation or recompilation beyond a list of defined program types and allocated heap blocks. We have integrated our analysis into the GNU Debugger (gdb), and describe our initial experience using this tool with several small to medium-sized programs.

Categories and Subject Descriptors D.2.5 [Software Engineering]: Testing and Debugging—Debugging aids; D.3.2 [Programming Languages]: Language Classifications—C; D.3.3 [Programming Languages]: Language Constructs and Features—Data types and structures, Dynamic storage management; F.3.3 [Logics and Meanings of Programs]: Studies of Program Constructs—Type structure

General Terms Algorithms, Human Factors, Languages, Reliability

Keywords dynamic type inference, constraints, debugging tools, heap visualization, physical subtyping, conservative garbage collection

1. Introduction

Suppose that a programmer notices an incorrect variable value during the execution of a C program. While debugging, the programmer may try to observe values of the variable at various points during execution, either by setting a watchpoint in the debugger or by inserting print statements. However, both of these techniques may be inadequate. Debugger watchpoints can be prohibitively slow. Adding print statements may be ineffective in cases of memory

corruption, as the affected data structure may have no apparent relation to the code that corrupts it. The programmer may also run a static pointer analysis to check for erroneous memory accesses. However, typical analyses [1, 23] do not model dependencies between neighboring memory blocks. In C, many bugs are caused by buffer overruns and pointer mismanagement, so physical proximity of memory blocks is an important factor.

We have designed and implemented an automated tool to help programmers understand and debug program behavior at the physical memory level. Our tool offers the programmer a low-level, but typed, view of memory. Each allocated chunk is presented either according to its inferred overall data type or as “untypable” if no defined program type is compatible with its imposed constraints. A common scenario where this view is useful is in the case of a buffer overrun: rather than attempt to deduce what data structure lies near the corrupted location from raw memory values or printed variables, the programmer uses our tool to discover that a nearby memory block is being used as an array of a particular type. The mere presence of an array makes “buffer overrun” a good hypothesis, and the array’s type helps the programmer identify relevant code to examine for possible array bounds violations.

When memory has been corrupted, our tool may detect that values stored in one or more heap blocks do not match any feasible type for the block. In order to reason about the cause of the corruption, the programmer may want to know exactly when the heap first became corrupted. For this task, we support a binary search debugging strategy often used to find ill-behaving code: set a breakpoint at the location where a variable was last known to be correct, and another where it has the wrong value, then iteratively narrow the interval until the bad assignment is exposed. Our tool gives the programmer the ability to treat the entire heap as a memory value that is either in a good or bad state, and search for the time at which the heap was first corrupted using her usual dynamic debugging techniques, such as the binary search described above.

This paper makes the following contributions:

- We introduce the idea of a *consistent typing* for the heap at any given point during execution. Each block of allocated storage is assigned a type from among those used in the program. The type assignment satisfies a set of constraints imposed by the values stored in memory, the size of each allocated block, and a type conformance relation (\preceq) based on physical subtyping. Additional constraints may be imposed by declared program variables when debugging information is available.
- We give an algorithm that finds a consistent typing for a snapshot of the heap. When no consistent typing exists, we report the locations and causes of conflicts to the user.

The remainder of this paper is organized as follows. We establish basic terms in Section 2, and define the subtyping relation on types used for individual bytes in Section 3. Section 4 specifies the

^{*}Supported in part by NSF under grant CCR-0305387.

[†]Work performed while at the University of Wisconsin–Madison.

constraints that must be satisfied for a set of typed memory locations to be considered consistent, and Section 5 gives an algorithm to find a consistent typing for the entire heap. Section 6 presents several case studies using our tool, Section 7 places our system in context with related work, and Section 8 concludes.

2. Preliminaries

Here, we describe the scenario under which we solve the problem of finding a consistent typing for the heap and provide key definitions and notation used throughout the paper.

2.1 Definitions

A *typing* is a map $T : Store \rightarrow Types$ from each storage location ($addr \in Store$) to its corresponding type ($\tau \in Types$). Types are all those defined or used by the program, including structures, unions, pointers, arrays, and functions. Storage locations include: a fixed set of addresses holding global variables; the set of locations that hold all local variables and function arguments on the stack at the current program point, or on multiple stacks for multithreaded programs; and all memory blocks dynamically allocated since the start of the program, which we refer to as *heap storage*.

Globals, locals, and function arguments may have associated type information if symbolic debug information is present. However, memory blocks from heap storage never have associated dynamic type information. Unions are also untagged per standard C.

In the manner of conservative garbage collectors [3], we define a *valid pointer* as a block of memory whose value is in *Store* (i.e., is a storage location). A valid pointer may also point immediately after the end of a block of storage; this is a common programming idiom that is explicitly allowed by the C standard [13].

A *type constraint* restricts the types that may reside at a given storage location. Our analysis imposes constraints on the types of individual bytes of memory, termed *byte types*. Informally, a byte type indicates that a byte holds the start of some program type or any subtype thereof. Byte types may also indicate that a byte is part of the interior of a multi-byte value that starts at an earlier location.

2.2 Notation

We establish a concise notation for C types, derived from that used by Siff *et al.* [21]. An array of n elements of type τ is written $\tau[n]$, while τptr is a pointer to τ . A tuple of the form $s\{m_1, \dots, m_k\}$ denotes a `struct`, while $u\{m_1, \dots, m_k\}$ denotes an untagged union. Each m_j is a triple (τ, l, i) giving the type, name, and starting offset of one field within a structure or union. Structure fields are ordered by offset, with m_1 starting at offset zero. Union fields are unordered and all start at offset zero.

Our subtyping relation is given by \prec and its reflexive closure \preceq . We use type notation $addr : \tau$ in place of the subtyping constraint $T(addr) \preceq \tau$ when the type mapping T is evident from context.

In the discussion that follows, byte offsets within a memory block are represented in terms of addition: if *block* is the start address of some block of storage, $block + i$ denotes the address of the i^{th} byte of *block*, starting from zero. We abbreviate the address $block + 0$ as simply *block*. The predicate $validPointer(addr)$ asserts that *addr* holds the start of a valid, non-null pointer value.

3. Byte Type Lattice

Our analysis may yield multiple distinct types for the same memory location. In some cases this reveals a conflict and likely misused memory. In other cases the types may, in fact, be compatible. This section explains how we construct a lattice from the set of C program types to model type compatibility in our analysis.

The data types in Figure 1 form a running example throughout the paper; hereafter, we omit the explicit “`struct`” keyword.

```

struct Point {
    double x, y;
};

struct Shape {
    char *name;
    FILE *fptr;
};

struct Part {
    struct Point center;
    struct Shape *shape;
    struct Assembly *owner;
};

struct PartNode {
    struct Part *part;
    struct PartNode *next;
};

struct Assembly {
    struct Point center;
    struct PartNode *nodes;
    struct Assembly *owner;
};

```

Figure 1. Example data types for assembly-building program

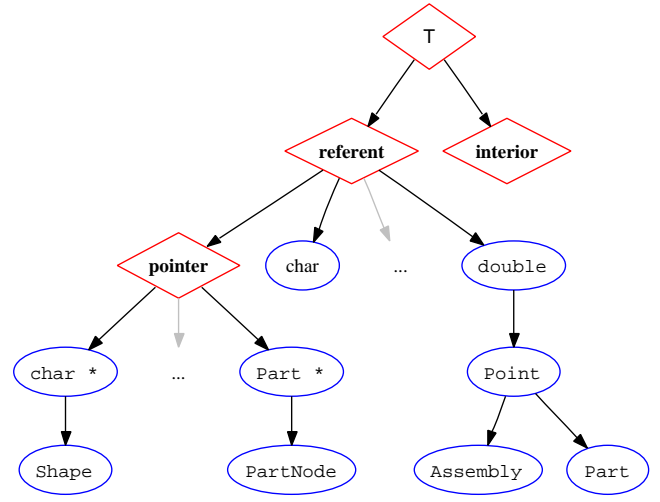


Figure 2. Byte type lattice corresponding to the data types in Figure 1

Given these types, the remainder of this section defines the subtyping relation used to construct the byte type lattice shown in Figure 2. For clarity, we omit \perp in Figure 2 and from the subtyping definitions below. \perp should be assumed as the *meet* of any two lattice elements for which no other lower bound is defined; it denotes that the corresponding two types may not be consistently stored at the same address. Treatment of function types may be found in a companion technical report [20].

3.1 Compound Types

A structure and its initial (offset 0) field have the same physical address but distinct types, as do an array and its 0^{th} element. Therefore the immediate supertype of a structure is the type of its initial field, and the immediate supertype of an array type is the type of its elements:

$$s\langle(\tau_1, l_1, 0), \dots, m_k\rangle \prec \tau_1 \quad \tau[n] \prec \tau$$

3.2 Special and Atomic Types

The four diamond-shaped nodes at the top of Figure 2 are always present. \top denotes an unconstrained memory address that may hold any type. **referent** describes all types that can be the referent of a pointer, while **interior** describes the non-initial bytes of multi-byte atomic values. Pointer targets must be **referent** subtypes and can never be **interior**. For example, if v holds an eight-byte `double`, then byte $v + 0$ has type `double` but bytes $v + 1$ through $v + 7$ all have type **interior** and may not be pointed to directly. Bitfields are also typed as **interior**, for their addresses cannot be taken.

pointer indicates that a storage location holds the 0th byte of a valid pointer, and thus is potentially consistent with any pointer type. Pointers to pointers are allowed, so **pointer** \prec **referent**. **pointer** is not the same as `void*`; the former represents all pointers, while the latter is a specific program type.

The oval nodes in the lattice correspond to actual types that may be used in a C program. Notice that the primitive atomic types are all sibling immediate subtypes of **referent**. This stipulates that two or more atomic types may not be simultaneously stored at the same address.

As a special case, we treat `void` as a zero-length type that is identical to **referent**. Although no byte should ever have type `void` in the final result, this convention allows transparent handling of aliases between `void*` and more fully typed pointers: pointers to τ and `void` may refer to the same address in our system as though `void` were a zero-length prefix of every **referent** subtype.

The subtyping relation is not extended across pointers. That is, $\tau_1 \preceq \tau_2 \not\Rightarrow \tau_1 \text{ ptr} \preceq \tau_2 \text{ ptr}$. This is the standard restriction for subtyping with updatable references, and any program that requires this form of subtyping to describe its heap must necessarily have used casts or other measures to violate type safety.

3.3 Unions

Untagged unions require special treatment, because a union may be used as any of its fields, but a consistent typing requires that every address be assigned a unique type. For each untagged union type $u\{m_1, \dots, m_k\}$, we extend the type grammar to include one tagged case $u@m_r\{m_1, \dots, m_k\}$ for each $1 \leq r \leq k$. Unions and their tagged cases adhere to the following subtyping relations:

$$\begin{aligned} u\{m_1, \dots, m_k\} &\prec \text{referent} \\ u@(\tau_r, l_r, 0)\{m_1, \dots, m_k\} &\prec u\{m_1, \dots, m_k\} \\ u@(\tau_r, l_r, 0)\{m_1, \dots, m_k\} &\prec \tau_r \end{aligned}$$

These relations forbid aliased pointers to differently-typed union fields. Each union must be used in a single consistent manner at any given point during execution. For example, a tagged union storing a `Point` can be the target of pointers to the untagged union as well as pointers to `Point` and `double`, but could not be the target of a pointer to `Part`. When two or more fields of a union have a common supertype, additional cases can be introduced that represent a subset of possible tagged cases rather than a single case. This preserves uniqueness of the lattice *meet* operation.

Note that only tagged unions and \perp have multiple immediate subtypes. The byte type lattice without these becomes a tree.

3.4 Finite Type Space

The byte type lattice contains an unbounded number of types, including arrays of arbitrary length and pointers of arbitrarily deep nesting. In practice, we only consider the following finite subset of types that are likely to be meaningful and useful for a given program:

- program-declared structures, unions, and enumerations
- tagged variants of unions or types containing unions
- arrays used by the program, e.g. `int[3][5]` if and only if at least one field or variable has exactly this type
- pointers used by the program, e.g. `int****` if and only if at least one field or variable has exactly this type
- pointers to known types up to two more levels of indirection

The number of tagged variants of union-containing types is potentially exponential. In our experience, multi-union structures and nested unions are unusual, and therefore the number of tagged variants is typically linear.

```
void main() {
    ...
    carAsm = create_assembly();
    ...
}

Assembly *create_assembly() {
1  Assembly *asm =
   malloc(sizeof(Assembly));
2  PartNode *node =
   malloc(sizeof(PartNode));

3  node->part = malloc(sizeof(Part));
4  node->next = node;
5  asm->nodes = node;

   // build part's shape and set name
6  init_part(node->part, "door", asm);
   ...
7  return asm;
}
```

Figure 3. Program that builds a simple assembly

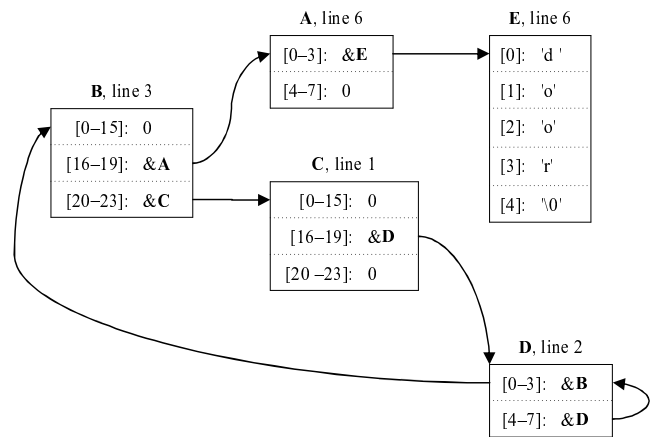


Figure 4. Allocated blocks and values after `init_part()` call in Figure 3

Additional array types are synthesized as needed during the analysis to satisfy size constraints (Section 4.2), but only using element types appearing in the original program. For example, a block of 32 bytes may be typed as `int[8]` if `int` is a known 4-byte type. We do not consider `int[2][4]` unless the corresponding element type (`int[2]`) appeared in the original program.

4. Consistency Constraints

In this section, we specify four kinds of constraints on storage locations that restrict the possible types for heap blocks. We then show how these constraints are combined to derive a consistent heap typing at one point in an example program. For simplicity, we assume a 32-bit architecture with 4-byte pointers and 8-byte doubles. Our ideas also generalize to 64-bit or other architectures.

The program excerpt in Figure 3 creates a part for a simple assembly and initializes it with its shape and owner assembly. Figure 4 shows the heap after the call to `init_part()` on line 6 of Figure 3. Blocks are labeled **A–E** arbitrarily, with the line number of each block's allocation given next to its label. Valid pointer

Block	Valid Pointers	Value-Consistent Types
A	A + 0	PartNode, Shape
B	B + 16, B + 20	Part, Assembly
C	C + 16	Part, Assembly
D	D + 0, D + 4	PartNode, Shape

Table 1. Valid pointer value and possible value-consistent types for heap in Figure 4

values are shown in a C-style syntax, and the rest of memory is assumed to be set to zero when returned by `malloc()`. Bracketed numbers indicate byte ranges within each block.

4.1 Value Constraints

Value constraints arise from concrete data values in memory at the instant the analysis is applied. They reflect the fact that some data types have limited domains that are much smaller than the set of all possible values that can fit in the allotted memory. The following value constraints, which we divide into two categories, are useful across a wide variety of C programs:

Mandatory Constraints: If a value looks like a valid pointer, we require that it be typed as a pointer (Section 5.1):

$$\text{validPointer}(addr) \Rightarrow T(addr) \preceq \text{pointer}$$

This assumes no non-pointer ever takes on a valid pointer value by chance, a strategy widely employed by conservative garbage collectors [3].

Table 1 shows valid pointer value constraints for the heap snapshot in Figure 4, and the types whose physical layouts are consistent with each constraint. An important detail is that, while locations holding zero are unconstrained, zero is consistent with either a pointer type or most primitive types. In this example, both `PartNode` and `Shape` are value-consistent with block **A** if the value at **A** + 4 is viewed as a null pointer.

Filtering Constraints: Some values may not be used to infer basic type information directly, but may help reject a hypothesized type if the value is not part of the domain of that type. We accept an `enum` as a potential byte type only if the value starting at that byte is equal to one of the defined constants for the enumerated type. Similarly, a function pointer type is consistent if its value is any word-aligned address that may contain executable code. We may also employ a character constraint in programs that manipulate ASCII text: if a block is otherwise unconstrained and ASCII character values are stored at every offset in the block, then the block should be typed as a `char` or `char[n]` rather than any other consistent primitive or primitive array type. A more detailed discussion of the rationale for the above filtering constraints may be found in our technical report [20].

Programmers may wish to define additional value constraints to reflect application-specific types and invariants. Our heap typing algorithm can accommodate arbitrary predicates that approve or reject the type proposed for a given location and value. For example, the data structure consistency specifications of Demsky *et al.* [6] could be applied as additional value constraints for selected types.

4.2 Size Constraints

The overall type for a block must fill exactly the number of bytes allocated for that block. For any address *block* which is the start of an allocated block,

$$T(\text{block}) = \tau \Rightarrow \text{sizeof}(\text{block}) = \text{sizeof}(\tau)$$

In C, dynamically allocated arrays tile multiple copies of the element type one after the other. A block holding a dynamic array

Block	Size	Size-Consistent Types
A	8	PartNode, Shape, char[8], ...
B	24	Part, Assembly, Shape[3], PartNode[3], float[6], ...
C	24	Part, Assembly, Shape[3], PartNode[3], float[6], ...
D	8	PartNode, Shape, int[2], ...
E	5	char[5]

Table 2. Size constraints and possible size-consistent types for heap in Figure 4

with *n* elements of type τ must satisfy

$$\text{sizeof}(\text{block}) = n \times \text{sizeof}(\tau)$$

for some whole number of array elements *n*.

Table 2 shows size constraints and a few illustrative size-compatible types for blocks in the example heap snapshot.

4.3 Type Constraints

Type constraints relate multiple locations, either between blocks (for pointers) or within a single block (for multi-byte structures):

- (i) If $T(addr) \preceq \text{pointer}$ and $*addr$ is within an allocated block (not one past the end), then $T(*addr) \prec \text{referent}$.
- (ii) If $T(addr) \preceq \tau$ for any atomic type τ , then $T(addr + i) = \text{interior}$ for all $1 \leq i < \text{sizeof}(\tau)$. Combined with rule (i) and the incompatibility of **interior** and **referent**, this forbids pointers into the interior of atomic values.
- (iii) For any allocated block starting at *block*, $T(\text{block}) \prec \text{referent}$. While similar to rule (i), this rule also affects leaked blocks to which nothing points.
- (iv) If $T(addr) \preceq \tau \text{ ptr}$ then $T(*addr) \preceq \tau$. Pointers and pointed-to types must be compatible modulo subtyping.
- (v) If $T(addr) \preceq s((\tau_1, l_1, i_1), \dots, (\tau_k, l_k, i_k))$, then $T(addr + i_n) = \tau_n$ for all $1 < n \leq k$. Structure fields must be compatible with the structure as a whole.
- (vi) If $T(addr) = \tau[n]$, then $T(addr + i \times \text{sizeof}(\tau)) = \tau$ for all $1 \leq i < n$. Array elements must be compatible with the array as a whole.

Implied constraints may imply additional constraints. A consistent heap typing must satisfy all transitively implied type constraints. Untagged unions induce no additional type constraints. Any tagged union type $u@(\tau, l, 0)\{\dots\}$ is also a subtype of τ and will pick up any appropriate constraints per the above rules.

4.4 Debug Constraints

If symbolic debug information is available for in-scope variables and function arguments, then this information may be added to the type map in the obvious manner. Equality constraints are appropriate here: if *x* is known to be an integer variable, then its type must be `int`, not some `int` subtype. Our main algorithm (Section 5) does not require debug information to find consistent typings, but takes debug constraints into account if present.

4.5 Example Solution

We now combine value, size, and type constraints to informally derive the consistent heap typing shown in Figure 5. Section 5 presents a systematic algorithm for finding consistent typings automatically. For clarity, we consider user-defined types before other possible matches (such as arrays of primitives). Assume that debug constraints are unavailable.

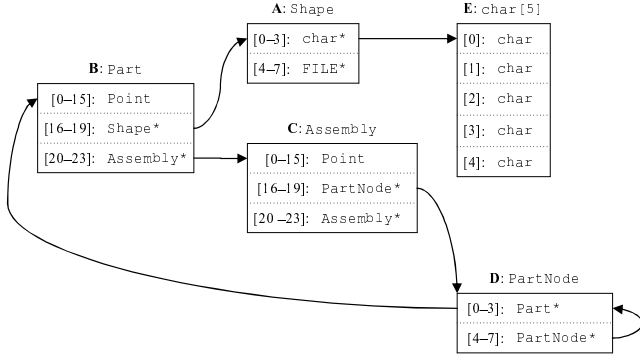


Figure 5. Fully constrained heap and derived typing for heap in Figure 4

From valid pointer value and size constraints, block **A** must have type `Shape` or `PartNode`. If `Shape` is considered first, we propagate a `char` constraint to block **E** via rule (iv).

Value and size constraints require that block **E** have type `char[5]`. Because `char[5] \preceq char`, block **E** can have type `char[5]` and still be consistent with the `char` pointer from block **A**. If we had tried `A : PartNode` first, then a conflict would have arisen; we discuss conflicts in detail in Section 5.

From size and value constraints, block **B** must have type `Part` or `Assembly`. If we try `Assembly`, then `B + 16 : PartNode*` via rule (v) and so `A : PartNode` via rule (iv), but this conflicts with `A : Shape` established earlier. Choosing `B : Part` is consistent with `A : Shape`, and also requires `C : Assembly`. This is consistent with value and size constraints on block **C**. Lastly, block **D** must have type `PartNode` due to the pointer field at `C + 16`.

5. Heap Typing Algorithm

In this section, we present an algorithm for assigning types to all storage locations, if a consistent typing is possible. Inputs to the algorithm consist of a snapshot of all values in the program’s memory; the start addresses and sizes of all allocated heap blocks; a list of defined program types; and optional symbolic debug information giving the locations, sizes, and types of in-scope variables. The output is a typing T giving consistent byte types for all allocated bytes. The byte type for the 0^{th} byte of an allocated block gives the overall type for that block. In some (but not all) cases when no globally consistent typing exists, the algorithm can identify, describe, and eliminate untypable blocks while still producing a partial typing for the remaining blocks.

The algorithm begins by assigning types to individual bytes of storage, using the values that arise in the program (e.g., valid pointers and their pointed-to locations) as a foremost source of byte type constraints (Section 5.1). Next we transitively propagate byte types from variables declared in the program (for which exact types are known), reporting any type constraint violations to the user (Section 5.2). Finally, we systematically consider possible overall types for each memory block until the typing map is fully defined or all typing alternatives for the blocks are exhausted (Section 5.3).

5.1 Pointer Constraint Gathering

We first establish valid pointer value constraints on individual bytes as described in Section 4.1. Type constraint rules (i)–(iii) induce additional constraints where appropriate. Conflicts between **referent** and **interior** may arise during this stage. For example if blocks **X** and **Y** hold valid pointers, but **Y** points to byte **X + 2**, then **X + 2** cannot simultaneously be the referent of the pointer in **Y** and the interior of the pointer in **X**. In this situation, it is difficult to know

which block is truly erroneous. We describe the conflict to the user, then mark both blocks as untypable and disregard them for the remainder of the algorithm.

After this phase, the typing constrains some bytes to be subtypes of **pointer**, **interior**, or **referent**, but many bytes remain unconstrained (\top) and no program types yet appear.

5.2 Debug Constraint Gathering

If symbolic debug information is available, debug constraints are applied next, then propagated transitively across pointers and into compound types using type constraint rules (iv)–(vi). Value and size constraints are checked where appropriate. Note, however, that size constraints are only partially enforced: a block must be at least large enough to contain the expected type, but may be larger. For example, the target of a `double*` must be at least eight bytes long, but may be longer if it is part of a larger structure, union, or array.

During this stage, conflicts may arise among debug constraints (e.g. `int` and `float` in the same location); between debug and size constraints (e.g. `int` in a two-byte block); between debug and value constraints; or between debug and **pointer**, **referent**, or **interior** constraints derived in the previous stage. If any such conflict occurs, we assume that execution has deviated from type safety and that the static type system therefore cannot be trusted to predict run-time types. We report the problem to the user and then back out all debug constraints. The remainder of the algorithm will operate using observed memory values only, without considering declared variable types. A more selective approach, which we leave as future work, would be to discard only a minimal subset of problematic debug constraints while keeping the remainder.

Barring conflicts, at the conclusion of this phase, the global typing includes program types for memory addresses that are (transitively) reached from pointers in program variables. However, these types are merely lattice upper bounds. For example, a pointer to `char` may actually point to an array of characters or to a structure with an initial `char` field. Some bytes within reachable blocks, and all bytes within unreachable blocks, still carry only the **pointer**, **referent**, and **interior** constraints added previously.

5.3 Completing the Heap Typing

Given the initial byte type constraints, we next assign an overall type to every heap block. Fully enforced size constraints ensure that the size of a block’s overall type is equal to the block size, so all allocated bytes are constrained when the typing is completed and the 0^{th} byte of each block determines the block’s overall type.

5.3.1 Typing Feasibility Check

We first verify that every heap block may be assigned at least one program type, given the initial constraints. A block that cannot be assigned any known type may be corrupt or may have been allocated in a library whose internal types are unavailable. We describe the problematic block to the user, then mark it as untypable and disregard it in the type search phase that follows.

5.3.2 Search Algorithm

The main search phase considers the possible types for each block, backtracking in the event of conflicts. When a type is verified as consistent with all current constraints on a particular block, we update the byte types for all bytes in this block to reflect the overall type, as well as propagate constraints one level forward across pointers in the block, and proceed to the next block. If no consistent type is found for some block, we backtrack to the last block that still has remaining type alternatives, and resume the search from that point. The algorithm terminates either when the last considered block is assigned a consistent type, or when all possible types for all blocks have failed, in which case no consistent typing exists.

Block	Type Considered	Outcome	Induced Constraints on Other Blocks
C	FILE	size conflict: $\text{sizeof}(\mathbf{C}) < \text{sizeof}(\text{FILE})$	
C	Part	✓	D +0: Shape
D	Shape	type conflict at D +0: $\text{meet}(\text{Shape}, \text{FILE}) = \perp$	
C	Assembly	✓	D +0: PartNode
D	PartNode	✓	B +0: Part
A	PartNode	size conflict: $\text{sizeof}(\mathbf{E}) < \text{sizeof}(\text{PartNode})$	
A	Shape	✓	E +0: char
E	char[5]	✓	

Table 3. Heap typing algorithm execution trace

Naïvely implemented, this search is geometric in the number of blocks and exponential in the number of types. However, since an entire block’s type must be a subtype of its 0th byte’s type, the search can be restricted to the sublattice below this bound. This optimization is especially effective with debug constraints enabled, as most pointers refer to the initial byte of an allocated block.

If a heap has several consistent typings, we prefer a typing that is most informative for the user. We consider larger types first, program-defined types before primitives, singletons before arrays, and fewer levels of pointer indirection before more. This order is designed to heuristically direct our solution toward more useful typings. We do not guarantee that the final heap typing is globally minimal or optimal with respect to this order, but we find that it yields good results in practice.

5.4 Example: Computing the Solution

We now illustrate the algorithm as applied to our running example, again considering the heap snapshot after the `init_part()` call on line 6 of Figure 3. To illustrate conflict handling, we modify the values shown in Figure 4 as follows: assume that bytes **B**+16..**B**+19 have been corrupted, and no longer hold a valid pointer value. All other value constraints remain as shown in Table 1.

During debug constraint collection, an inconsistency arises when a **B**:Part constraint is propagated from variable `node`, because Part requires a valid pointer or null at bytes **B**+16..**B**+19. Since a conflict is found, debug constraints are discarded. The search continues using only value constraints. Block **B** does not pass the typing feasibility check, so it is omitted from the search. Table 3 summarizes steps taken during the backtracking search for complete types. The remaining blocks are considered in their allocation order, and the types are ordered according to our sorting heuristic. The algorithm is able to recover the types of the four remaining blocks using value and size constraints only, backtracking several times throughout the search. The final heap typing for the four remaining blocks is as shown in Figure 5.

5.5 Propagation Correctness

After assigning a type to a block, we update constraints for the block itself and also propagate all pointer constraints forward across one dereference. We claim this is sufficient to ensure that no inconsistency between typed blocks is overlooked. Any two blocks *X* and *Y*, where *Y* is transitively reachable from *X*, will always be assigned consistent types regardless of the order in which they are considered. The full proof of our claim is contained in our technical report [20].

6. Evaluation

We have implemented the above algorithm within the GNU Debugger (gdb) [10], a popular symbolic debugger for C. When the program is stopped at a breakpoint, the user may type “`whatsat`

`<expr>`” to perform heap type inference and then display type-annotated memory beginning at the address computed by `<expr>`.

Implementing the *validPointer* predicate requires that the debugger probe the debuggee’s current heap allocation state. We modify the debuggee’s memory management routines to maintain a list of currently allocated blocks in a reserved global location known to the debugger; the debugger reads this list directly from the debuggee’s address space as needed. We record the start address and size of each block, plus the code address from which each block allocation call originated. `whatsat` uses the latter in diagnostic messages to report the source file, line number, and function at which each untypable block was allocated.

This extra allocation tracking uses standard hooks exposed by the GNU `libc` implementation [9] and is contained within a shared library that may be preloaded into any program one wishes to debug without recompilation or relinking. Our allocation hooks also zero-initialize newly allocated blocks. This is done to avoid spurious typing errors due to random data values in uninitialized heap memory. However, it can be useful to disable this feature in order to verify that the program under study fully initializes all of its own heap storage under normal running conditions (Section 6.3).

6.1 Visualization of Typed Memory

Following heap type inference, `whatsat` displays memory contents augmented with type information. Visualization begins at any address of the user’s choosing (e.g. “`whatsat 0x9275008`” or “`whatsat &foo[3]`”) and continues forward through raw memory under user control.

Figure 6 shows part of the type-annotated memory visualization for a program used in one of our experiments (see Section 6.2). Each line shows a capitalized hexadecimal memory address (e.g. “0X804ABC0:”), up to one word of raw memory content at that address (“0x00000001”), and an interpretation of that memory typed according to our algorithm (“`length = (int) 1`”). Indentation and field labels (“`length =`”) reflect nesting and compound types. Figure 6 shows five distinct but proximate memory blocks containing an array of structures and four function pointers. The “D” and “?” labels to the left of each address mark locations that are part of the static data section or have never been allocated, respectively. Other labels (not seen in this example) mark the stack, typed heap locations, freed heap locations, and several other categories of memory.

6.2 Schedule2

Schedule2 is a small C application from the Siemens buggy program suite [12]. Given a list of jobs and their priorities as input, the application computes and prints a schedule for running the jobs. Version 8 of `schedule2` contains a bug that causes the program to crash inside `malloc()`. A stack trace reveals that the crash is due to a bad pointer dereference: a function pointer, `__malloc_hook`, does not point to a function. `whatsat` confirms that the

```

? 0X804ABB8: 0x00000000
? 0X804ABBC: 0x00000000
      (queue [4])
      [0] = (queue)
D 0X804ABC0: 0x00000001      length = (int) 1
D 0X804ABC4: 0x0804b0d0      head = (process *) 0x804b0d0
      [1] = (queue)
D 0X804ABC8: 0x00000001      length = (int) 1
D 0X804ABCC: 0x0804b0a8      head = (process *) 0x804b0a8
      [2] = (queue)
D 0X804ABD0: 0x00000000      length = (int) 0
D 0X804ABD4: 0x00000000      head = (process *) 0x0
      [3] = (queue)
D 0X804ABD8: 0x00000000      length = (int) 0
D 0X804ABDC: 0x00000000      head = (process *) 0x0
? 0X804ABE0: 0x00000002
D 0X804ABE4: 0x0804b008      (void *(*()) 0x804b008
D 0X804ABE8: 0x00000000      (void *(*()) 0
D 0X804ABEC: 0x00749380      (void *(*()) <realloc_hook_ini>
D 0X804ABF0: 0x007493d0      (void *(*()) <memalign_hook_ini>
? 0X804ABF4: 0x0804b018
? 0X804ABF8: 0x00000000

```

Figure 6. Schedule2 global variables visualization. The types of function pointer arguments are omitted for brevity.

claimed type for `__malloc_hook` is inconsistent with its value, and therefore that debug constraints are not satisfiable. After debug constraints are discarded, `whatsat` infers that this block actually contains a `process` structure.

`__malloc_hook` is assigned from `old_malloc_hook`, which holds the same bad `process` pointer instead of a function pointer. Using `whatsat` to explore the physical memory neighborhood around `old_malloc_hook` reveals that a four-element structure array precedes `old_malloc_hook`. Figure 6 shows `whatsat`'s visualization of this area. `old_malloc_hook` appears at address `0X804ABE4`; the preceding array is clearly visible starting at address `0X804ABC0`. Observe that `old_malloc_hook` is perfectly positioned to receive an errant `process` pointer should the neighboring array overrun its bounds. Thus informed, we identify the array, the code that writes to it, and the missing bounds check that constitutes the true bug. All `whatsat` queries used in this case study completed within 0.03 seconds.

While a hardware watchpoint might also have been used to trap the bad write to `old_malloc_hook`, this would require rerunning the program and reproducing the bug. Many memory corruption bugs are difficult to reproduce on demand; we see here that `whatsat` can provide useful postmortem information on the first instance of a bug.

6.3 Exif

Exif is an open source utility for manipulating JPEG image metadata [8]. It consists of 10,375 lines of C code split into a shared library and a main driver program. We ran `exif` with a breakpoint set after `exif_loader_get_data()`, which builds an in-memory representation of a JPEG input file. We disabled zero-initialization of heap blocks to test whether `exif` performs its own initializations properly. `whatsat` identifies two untypable blocks allocated in `exif_content_add_entry()` at `exif-content.c:110`. The code in question performs a reallocation to grow an array of pointers to entry blocks:

```

entries = realloc(entries,
    sizeof (ExifEntry) * (count + 1));

```

The size calculation is incorrect. It reserves space for an array of `(count + 1)` `ExifEntry` structures, but `entries` is actually an array of *pointers* to `ExifEntry` structures. Because each `ExifEntry` is larger than a pointer, the program does not overrun this buffer. However, the extra space at the end of the array is

wasted and, because it contains uninitialized random data that may not look like valid pointers, `whatsat` determines that arrays allocated here are untypable.

We initially identified this previously unreported bug in release 0.6.9 and 0.6.10 of the `exif` driver and library. We have confirmed that it persists in the latest development snapshot as of April 15, 2006 (12,410 lines of C code). `Exif` developers have since confirmed the bug and applied our suggested fix.

After `whatsat` identifies these untypable blocks, it ignores them for the remainder of the analysis. That analysis, however, does not find a valid heap typing for `exif` in a timely manner. It is possible that no valid typing exists even though all individual blocks match at least one known type. It is also possible that a valid typing exists, but is pathologically mismatched with our heuristic search order. Improving the diagnostic capabilities of our analysis when unresolvable conflicts arise late in the search is an important area for future study.

7. Related Work

Chandra and Reps [5] and Siff *et al.* [21] introduce an alternate type system for C that allows subtyping based on the physical layouts of data structures. They describe static type checking and inference rules that test program conformance with this alternate type system. In contrast, our approach is dynamic: we examine a frozen snapshot of a running program's heap, rather than the space of all possible program heaps. This allows us to use concrete memory values and allocated block sizes to refine our analysis. As is typical for dynamic analyses, we focus on specific bugs triggered during a run without guaranteeing that all possible bugs will be detected.

The subtyping relation induced by our byte type lattice is more restrictive than the Chandra/Siff physical subtyping relation. Both allow subtyping between a structure and its first field, but we disallow more general structure prefixing or the use of `char` arrays as storage placeholders. These are merely policy choices. Our approach can use permissive Chandra/Siff subtyping or a variety of other relations with no changes to constraint collection or the core heap typing algorithm. However, not all subtyping relations are sensible in this context. For example, Cardelli's structural record subtyping relation [4], disregards field order and is therefore needlessly permissive for our scenario, where field orders are fixed.

As a dynamic heap-walking tool, our system shares some qualities with a garbage collector or leak detector, and a list of unreachable (leaked) memory blocks could easily be extracted from our analysis. Traditional garbage collectors require data structure layout information for the root set and possibly for allocated blocks as well. Conservative garbage collectors [3] relax this requirement by assuming that any location holding a valid pointer value is indeed a pointer. Our approach moves flexibly between these extremes. Ultimately, the information we recover is richer than that produced by garbage collectors: we find not only the size and embedded pointers of each allocated block, but also complete program types that are globally consistent both within and between all blocks.

Zimmermann and Zeller present strategies for extracting C heaps and displaying them to highlight key relationships [25]. Their system depends on debugger-provided type information augmented with a few C-specific heuristics also used by `whatsat`, such as pointer validity testing and dynamic array size computation. These heuristics consider only isolated blocks, though, and have no notion of global consistency. Zimmermann and Zeller acknowledge that "While such heuristics mostly make good guesses, it is safer to provide explicit disambiguation rules—either hand-crafted or inferred from the program." Dynamic heap type inference generalizes and improves upon these heuristics by defining a notion of global heap typing that considers not just local values within isolated blocks, but also the relationships between interlinked blocks. This

lets `whatsat` find globally consistent heap typings and reduces or eliminates the need for hand-crafted disambiguation rules.

The problem of heap corruption due to pointer and cast abuse is longstanding, and has inspired solutions ranging from static analysis [7, 18] to run-time instrumentation [11, 15, 16, 17, 19, 24] and the design of safer language dialects [2, 14, 22]. Our approach performs programmer-directed heap validity checks in an interactive debugging context, and does not attempt to prevent or trap errors as they occur. This allows us to be significantly less invasive: we require no changes to the C language; no recompilation or source annotation beyond a compiler-provided list of program types; no run-time instrumentation beyond a list of allocated blocks; and no dynamic type tagging or other changes to data structure layouts. Additionally, our analysis depends only on the instantaneous state of the program heap at a given moment in time: other than maintaining a list of currently allocated blocks, we do not record any trace information while the program runs.

In this sense `whatsat` can be seen as an experiment in minimalism. Rather than monitor every potentially interesting action, we ask how much information can be recovered with only the bare minimum imposition at run-time. We believe that both highly invasive and minimally invasive approaches have benefits. Exploring the extremes helps illuminate potential strategies to improve debugging tools all along the instrumentation and analysis spectrum.

8. Conclusion

Low-level programming languages sometimes require low-level debugging. However, one need not completely abandon the type system even when working with non-type-safe languages. A low-level but type-annotated view of the heap can help in debugging and more general program understanding tasks. We have presented an algorithm that infers program-defined types for memory locations. Solution consistency is defined in terms of constraints that use a novel blend of ideas from physical subtyping and conservative garbage collection. When no consistent typing exists due to heap corruption or pointer abuse, we offer focused diagnostic information to help identify the cause. Our implementation works for general C programs and requires no source annotation, no recompilation, no run-time instrumentation beyond heap allocation tracking, and no changes to physical data structure layouts. Experiences with the tool, while limited in scope, suggest that dynamic heap type inference may be a useful addition to the programmer's toolkit.

References

- [1] L. O. Andersen. *Program Analysis and Specialization for the C Programming Language*. PhD thesis, DIKU, University of Copenhagen, May 1994.
- [2] T. M. Austin, S. E. Breach, and G. S. Sohi. Efficient detection of all pointer and array access errors. In *PLDI '94: Proceedings of the ACM SIGPLAN 1994 conference on Programming language design and implementation*, pages 290–301, New York, NY, USA, 1994. ACM Press.
- [3] H.-J. Boehm and M. Weiser. Garbage collection in an uncooperative environment. *Software: Practice & Experience*, 18(9):807–820, 1988.
- [4] L. Cardelli. Structural subtyping and the notion of power type. In *POPL '88: Proceedings of the 15th ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*, pages 70–79, New York, NY, USA, 1988. ACM Press.
- [5] S. Chandra and T. W. Reps. Physical type checking for C. In *Workshop on Program Analysis For Software Tools and Engineering*, pages 66–75, 1999.
- [6] B. Demsky, M. D. Ernst, P. J. Guo, S. McCamant, J. H. Perkins, and M. Rinard. Inference and enforcement of data structure consistency specifications. In *ISSTA 2006, Proceedings of the 2006 International Symposium on Software Testing and Analysis*, Portland, ME, USA, July 18–20 2006.
- [7] D. Evans. Static detection of dynamic memory errors. In *PLDI '96: Proceedings of the ACM SIGPLAN 1996 conference on Programming language design and implementation*, pages 44–53, New York, NY, USA, 1996. ACM Press.
- [8] EXIF tag parsing library. <http://libexif.sf.net/>.
- [9] Free Software Foundation, Inc., Boston, MA, USA. *The GNU C Library*, 0.10 edition, July 6 2001.
- [10] J. Gilmore and S. Shebs. *GDB Internals*, Feb. 2004.
- [11] R. Hastings and B. Joyce. Purify: Fast detection of memory leaks and access errors. In *Proceedings of the USENIX Winter Conference*, pages 125–138, San Francisco, CA, USA, 1992. USENIX Association.
- [12] M. Hutchins, H. Foster, T. Goradia, and T. Ostrand. Experiments on the effectiveness of dataflow- and control-flow-based test adequacy criteria. In *Proceedings of the 16th International Conference on Software Engineering*, pages 191–200. IEEE Computer Society Press, May 1994.
- [13] International Organization for Standardization. *ISO/IEC 9899:1999: Programming Languages — C*. International Organization for Standardization, Geneva, Switzerland, Dec. 1999.
- [14] T. Jim, J. G. Morrisett, D. Grossman, M. W. Hicks, J. Cheney, and Y. Wang. Cyclone: A safe dialect of C. In *Proceedings of the General Track: 2002 USENIX Annual Technical Conference*, pages 275–288, Berkeley, CA, USA, 2002. USENIX Association.
- [15] R. W. M. Jones and P. H. J. Kelly. Backwards-compatible bounds checking for arrays and pointers in C programs. In *AADEBUG*, pages 13–26, 1997.
- [16] S. Kaufer, R. Lopez, and S. Pratap. Saber-C: An interpreter-based programming environment for the C language. In *Proceedings of the USENIX Summer Conference*, pages 161–171, San Francisco, CA, USA, June 1988. USENIX Association.
- [17] A. Loginov, S. H. Yong, S. Horwitz, and T. W. Reps. Debugging via run-time type checking. In *FASE '01: Proceedings of the 4th International Conference on Fundamental Approaches to Software Engineering*, pages 217–232, London, UK, 2001. Springer-Verlag.
- [18] G. C. Necula, J. Condit, M. Harren, S. McPeak, and W. Weimer. CCured: Type-safe retrofitting of legacy software. *ACM Transactions on Programming Languages and Systems*, 27(3):477–526, 2005.
- [19] N. Nethercote and J. Seward. Valgrind: A program supervision framework. *Electronic Notes in Theoretical Computer Science*, 89(2), 2003.
- [20] M. Polishchuk, B. Liblit, and C. Schulze. WHATSAT: Dynamic heap type inference for program understanding and debugging. Technical Report 1583, University of Wisconsin–Madison, Oct. 2006.
- [21] M. Siff, S. Chandra, T. Ball, K. Kunchithapadam, and T. W. Reps. Coping with type casts in C. In O. Nierstrasz and M. Lemoine, editors, *ESEC / SIGSOFT FSE*, volume 1687 of *Lecture Notes in Computer Science*, pages 180–198. Springer, 1999.
- [22] G. Smith and D. Volpano. A sound polymorphic type system for a dialect of C. *Science of Computer Programming*, 32(1-3):49–72, 1998.
- [23] B. Steensgaard. Points-to analysis in almost linear time. In *Symposium on Principles of Programming Languages*, pages 32–41, 1996.
- [24] J. L. Steffen. Adding run-time checking to the portable C compiler. *Software: Practice & Experience*, 22(4):305–316, 1992.
- [25] T. Zimmermann and A. Zeller. Visualizing memory graphs. In *Software Visualization, International Seminar Dagstuhl Castle, Germany, May 20–25, 2001, Revised Lectures*, volume 2269 of *Lecture Notes in Computer Science*, pages 191–204. Springer, May 2001.