



Modeling Dynamic (De)Allocations of Local Memory for Translation Validation

ABHISHEK ROSE, Indian Institute of Technology Delhi, India

SORAV BANSAL, Indian Institute of Technology Delhi, India

End-to-End Translation Validation is the problem of verifying the executable code generated by a compiler against the corresponding input source code for a single compilation. This becomes particularly hard in the presence of dynamically-allocated local memory where addresses of local memory may be observed by the program. In the context of validating the translation of a C procedure to executable code, a validator needs to tackle constant-length local arrays, address-taken local variables, address-taken formal parameters, variable-length local arrays, procedure-call arguments (including variadic arguments), and the `alloca()` operator. We provide an execution model, a definition of refinement, and an algorithm to soundly convert a refinement check into first-order logic queries that an off-the-shelf SMT solver can handle efficiently. In our experiments, we perform blackbox translation validation of C procedures (with up to 100+ SLOC), involving these local memory allocation constructs, against their corresponding assembly implementations (with up to 200+ instructions) generated by an optimizing compiler with complex loop and vectorizing transformations.

CCS Concepts: • **Software and its engineering** → **Formal software verification**; *Compilers*.

Additional Key Words and Phrases: Translation validation, Equivalence checking, Certified compilation

ACM Reference Format:

Abhishek Rose and Sorav Bansal. 2024. Modeling Dynamic (De)Allocations of Local Memory for Translation Validation. *Proc. ACM Program. Lang.* 8, OOPSLA1, Article 146 (April 2024), 30 pages. <https://doi.org/10.1145/3649863>

1 INTRODUCTION

Compiler bugs can be catastrophic, especially for safety-critical applications. End-to-End Translation Validation (TV for short) checks a single compilation to ascertain if the machine executable code generated by a compiler agrees with the input source program. In our work, we validate translations from *unoptimized* IR of a C program to *optimized* executable (or assembly) code, which forms an overwhelming majority of the complexity in an end-to-end compilation pipeline. In this setting, the presence of dynamic allocations and deallocations due to local variables and procedure-call arguments in the IR program presents a special challenge — in these cases, the identification and modeling of relations between a local variable (or a procedure-call argument) in IR and its stack address in assembly is often required to complete the validation proof.

Unlike IR-to-assembly, modeling dynamic local memory allocations is significantly simpler for IR-to-IR TV [Kasampalis et al. 2021; Lopes et al. 2021; Menendez et al. 2016; Namjoshi and Zuck 2013; Necula 2000; Stepp et al. 2011; Tristan et al. 2011; Zhao et al. 2012, 2013]. For example, (pseudo)register-allocation of local variables can be tackled by identifying relational invariants that equate the value contained in a local variable’s memory region (in the original program) with the value in the corresponding pseudo-register (in the transformed program) [Kang et al. 2018]. If the

Authors’ addresses: [Abhishek Rose](mailto:abhishek.rose@cse.iitd.ac.in), Indian Institute of Technology Delhi, New Delhi, India, abhishek.rose@cse.iitd.ac.in; [Sorav Bansal](mailto:sbansal@iitd.ac.in), Indian Institute of Technology Delhi, New Delhi, India, sbansal@iitd.ac.in.



This work is licensed under a Creative Commons Attribution 4.0 International License.

© 2024 Copyright held by the owner/author(s).

ACM 2475-1421/2024/4-ART146

<https://doi.org/10.1145/3649863>

address of a local variable is observable by the C program (e.g., for an address-taken local variable), we need to additionally relate the variable addresses across both programs. These address correlations can be achieved by first correlating the corresponding allocation statements in both programs (e.g., through their names) and then assuming that their return values are equal. Provenance-based syntactic pointer analyses, that show separation between distinct variables [Andersen 1994; Steensgaard 1996], thus suffice for translation validation across IR-to-IR transformations.

An IR-to-assembly transformation involves the lowering of a memory allocation (deallocation) IR instruction to a stackpointer decrement (increment) instruction in assembly. Further, the stack space in assembly is shared by multiple local variables, procedure-call arguments, and by the potential intermediate values generated by the compiler, e.g., pseudo-register spills. Provenance-based pointer analyses are thus inadequate for showing separation in assembly.

Prior work on IR-to-assembly and assembly-to-assembly TV [Churchill et al. 2019; Gupta et al. 2020; Sewell et al. 2013; Sharma et al. 2013] assumes that local variables are either absent or their addresses are not observed in the program and so they are removed through (pseudo)register-allocation. Similarly, these prior works assume that variadic parameters (and other cases of address-taken parameters) are absent in the program.

Prior work on certified compilation, embodied in CompCert [Leroy 2006], validates its own transformation passes from IR to assembly, and supports both address-taken local variables and variadic parameters. However, CompCert sidesteps the task of having to model dynamic allocations by ensuring that the generated assembly code *preallocates* the space for all local variables and procedure-call arguments at the beginning of a procedure's body. Because preallocation is not possible if the size of an allocation is not known at compile time, CompCert does not support variable-sized local variables or `alloca()`. Moreover, preallocation is prone to stack space wastage. In contrast to a certified compiler, TV needs to validate the compilation of a third-party compiler, and thus needs to support an arbitrary (potentially dynamic) allocation strategy.

Example: Consider a C and a 32-bit x86 assembly program in fig. 1. The `fib` procedure in fig. 1a accepts two integers n and m , allocates a variable-length array (VLA) v of $n+2$ elements, computes the first $m+1$ fibonacci numbers in v , calls `printf()`, and returns the m^{th} fibonacci number. Notice that for an execution free of *Undefined Behaviour* (UB), both n and m must be non-negative and m must be less than $(n+2)$. Note that the memory for local variables (v and i) and procedure-call arguments (for the call to `printf`) is allocated dynamically through the `alloc` instruction in the IR program (fig. 1b). In the assembly program (fig. 1c), memory is allocated through instructions that manipulate the stackpointer register `esp`.

If the IR program uses an address, say α , of a local variable (e.g., $\alpha \in \{p_{I1}, p_{I2}\}$) or a procedure-call argument (e.g., $\alpha \in \{p_{I7}, p_{I8}\}$) in its computation (e.g., for pointer arithmetic at lines I3 and I5, or for accessing the variadic argument at p_{I8} within `printf`), validation requires a relation between α and its corresponding stack address in assembly (e.g., $p_{I7} = \text{esp}$ at line A14).

Contributions: We formalize IR and assembly execution semantics in the presence of dynamically (de)allocated memory for local variables and procedure-call arguments, define a notion of correct translation, and provide an algorithm that converts the correctness check to first-order logic queries over bitvectors, arrays, and uninterpreted functions. Almost all production compilers (e.g., GCC) generate assembly code to dynamically allocate stack space for procedure-call arguments at the callsite, e.g., in fig. 1c, the arguments to `printf` are allocated at line A13. Ours is perhaps the first effort to enable validation of this common allocation strategy. Further, our work enables translation validation for programs with dynamically-allocated fixed-length and variable-length local variables for a wide set of allocation strategies used by a compiler including stack merging, stack reallocation (if the order of allocations is preserved), and intermittent register allocation.

```

C0: int fib(int n, int m) {
C1:   int v[n+2];
C2:   v[0]=0; v[1]=1;
C3:   for(int i=2; i<=m; i++)
C4:     v[i]=v[i-1]+v[i-2];
C5:   printf("fib(%d)=%d", m, v[m]);
C6:   return v[m];
C7: }

```

(a) C Program with VLA.

```

I0: int fib(int* n, int* m):
I1:   i=pI1=alloc 1,int,4;
I2:   v=pI2=alloc *n+2,int,4;
I3:   v[0]=0; v[1]=1; *i=2;
I4:   if(*i >_s *m) goto I7;
I5:   v[*i]=v[*i-1]+v[*i-2];
I6:   (*i)++; goto I4;
I7:   pI7=alloc 1,char*,4;
I8:   pI8=alloc 1,struct{int; int;},4;
I9:   *pI7=__S__; *pI8=*m; *(pI8+4)=v[*m];
I10:  t=call int printf(pI7, pI8);
I11:  dealloc I7;
I12:  dealloc I8;
I13:  r=v[*m];
I14:  dealloc I2;
I15:  dealloc I1;
I16:  ret r;

```

(b) (Abstracted) IR.

```

A0: fib:
A1:   push ebp; ebp = esp;
A2:   push {edi, esi, ebx}; esp = esp-12;
A3:   eax = mem4[ebp+8]; ebx = mem4[ebp+12];
A4:   esp = esp-(0xFFFFFFFF & (4*(eax+2)+15));
A4.1: vI1 = alloc_u 4,4,I1;
A4.2: alloc_s esp,4*(eax+2),4,I2;
A5:   esi = ((esp+3)/4)*4;
A6:   mem4[esi] = 0; mem4[esi+4] = 1;
A7:   if(ebx <_s 1) jmp A12;
A8:   edi = 0; edx = 1; eax = 2;
A9:   ecx = edx+edi; edi = edx; edx = ecx;
A10:  mem4[esi+4*eax] = ecx; eax = eax+1;
A11:  if(eax <_s ebx) jmp A9;
A12:  edi = mem4[esi+4*ebx]; esp = esp-4;
A13:  push {edi, ebx, __S__}; //__S__ is ptr to format string
A13.1: alloc_s esp, 4,4,I7;
A13.2: alloc_s esp+4,8,4,I8;
A14:  call int printf(<char*> esp, <struct{int; int;}> esp+4)
      GU {hp,cl,I7,I8};
A14.1: dealloc_s I7;
A14.2: dealloc_s I8;
A15:  eax = edi;
A15.1: dealloc_s I2;
A15.2: dealloc_u I1;
A16:  esp = ebp-12; pop {ebx, esi, edi, ebp};
A17:  ret;

```

(c) (Abstracted) 32-bit x86 Assembly Code.

Fig. 1. Example program with VLA and its lowerings to IR and assembly. Subscripts $_s$ and $_u$ denote signed and unsigned comparison respectively. Bold font (parts of) instructions are added by our algorithm.

2 EXECUTION SEMANTICS AND NOTION OF CORRECT TRANSLATION

We are interested in showing that an x86 assembly program \mathbb{A} is a correct translation of the unoptimized IR representation of a C program \mathbb{C} . Prior TV efforts identify a lockstep correlation between (potentially unrolled) iterations of loops in the two programs to show equivalence [Churchill et al. 2019]. These correlations can be represented through a *product program* that executes \mathbb{C} and \mathbb{A} in lockstep, using a careful choice of program path correlations, to keep the machine states of both programs related at the ends of correlated paths [Gupta et al. 2020; Zaks and Pnueli 2008].

Our TV algorithm additionally attempts to identify a lockstep correlation between the dynamic (de)allocation events and procedure-call events performed in both programs, i.e., we require the order and values of these execution events to be identical in both programs. To identify a lockstep correlation, our algorithm annotates \mathbb{A} with (de)allocation instructions and procedure-call arguments. Our key insight is to define a *refinement relation* between \mathbb{C} and \mathbb{A} through the existence of an annotation in \mathbb{A} . We also generalize the definition of a product program so it can be used to witness refinement in the presence of non-determinism due to addresses of dynamically-allocated local memory, UB, and stack overflow.

Overview through example: In \mathbb{C} , an alloc instruction returns a non-deterministic address of the newly allocated region with non-deterministic contents, e.g., in fig. 1b, the address (p_{I2}) and initial contents of VLA v allocated at I2 are non-deterministic. In fig. 1c, our algorithm annotates an alloc $_s$ instruction at A4.2 to correlate in lockstep with I2, so that p_{I2} 's determinized value is identified through its first operand (esp). An alloc $_s$ instruction allocates a contiguous address interval from the stack, starting at esp in this case, to a local variable. The second ($4*(eax+2)$), third (4), and fourth (I2) operands of alloc $_s$ specify the allocation size in bytes, required alignment, and the PC of the correlated allocation instruction in \mathbb{C} (which also identifies the local variable)

respectively. The determinized values of the initial contents of VLA v at I2 are identified to be equal to the contents of the stack region $[\text{esp}, \text{esp}+4*(\text{eax}+2)-1]$ at A4.1. A corresponding dealloc_s instruction, that correlates in lockstep with I14, is annotated at A15.1 to free the memory allocated by A4.2 (both have operand I2) and return it to stack.

A procedure call appears as an x86 call instruction and we annotate the actual arguments as its operands in \mathbb{A} . In fig. 1c, the two operands (esp and $\text{esp}+4$) annotated at A14 are the determinized values of p_{I7} and p_{I8} , as obtained through x86 calling conventions. The last annotation at A14 is the set of memory regions (e.g., G, hp, cl, \dots , as described in section 2.2.2) observable by printf in \mathbb{A} – this is equal to the set of memory regions observable by printf in \mathbb{C} , as obtained through an over-approximate points-to analysis. Annotations of alloc_s at A13.{1,2} and dealloc_s at A14.{1,2} identify the memory regions occupied by printf 's parameters during printf 's execution.

Consider the local variable i , allocated at I1, with address p_{I1} in fig. 1b. Because i 's address is never taken in the source program, a correlation of p_{I1} with its determinized value in \mathbb{A} 's stack is not necessarily required. Further, the compiler may register-allocate i in which case no stack address exists for i , e.g., i lives in eax at A8–A11 in fig. 1c. The alloc_v instruction annotated at A4.1 performs a “virtual allocation” for variable i in lockstep with I1. The first (4), second (4), and third (I1) operands of alloc_v indicate the allocation size, required alignment, and the PC of the correlated allocation in \mathbb{C} respectively. The corresponding dealloc_v instruction, annotated at A15.2, correlates in lockstep with I15. The address and initial contents of the memory allocated by alloc_v are chosen non-deterministically in \mathbb{A} , and are assumed to be equal to the address and initial contents of memory allocated by a correlated alloc in \mathbb{C} , e.g., $v_{I1} = p_{I1}$ at A4.1. A “virtually-allocated region” is never used by \mathbb{A} . We introduce the $(\text{de})\text{alloc}_{s,v}$ instructions formally in section 2.4.

Consider the memory access $v[*i]$ at I5 in fig. 1b, and assume we identify a lockstep correlation of this memory access with the assembly program's access $\text{mem}_4[\text{esi}+4*\text{eax}]$ at A10 in fig. 1c, with value relations $\text{esi}=v$ and $\text{eax}=*i$. We need to cater to the possibility where $*i >_s *n+2$ (equivalently, $\text{eax} >_s \text{mem}_4[\text{ebp}+12]+2$), which would trigger UB in \mathbb{C} , and may go out of variable bounds in stack in assembly. Our product program encodes the necessary UB semantics that allow anything to happen in assembly (including out of bound stack access) if UB is triggered in \mathbb{C} .

Finally, consider the stackpointer decrement instruction at A4 in fig. 1c. If eax (which corresponds to $*n$) is too large, this instruction at A4 may potentially overflow the stack space. Our product program encodes the assumption that an assembly program will have the necessary stack space required for execution, which is necessary to be able to validate a translation from IR to assembly.

Thus, we are interested in identifying *legal* annotations of $(\text{de})\text{alloc}_{s,v}$ instructions and operands of procedure-call instructions in \mathbb{A} , such that the execution behaviours of \mathbb{A} can be shown to refine the execution behaviours of \mathbb{C} , assuming \mathbb{A} has the required stack space for execution. We show refinement separately for each procedure C in \mathbb{C} and its corresponding implementation A in \mathbb{A} . Thereafter, a coinductive argument shows refinement for full programs \mathbb{C} and \mathbb{A} starting at the $\text{main}()$ procedure. We do not support inter-procedural transformations.

Paper organization: Sections 2.1 to 2.3 describe a procedure's execution semantics for both IR and assembly representations. Refinement, through annotations, is defined in section 2.4. Section 3 defines a product program and its associated requirements such that refinement can be witnessed, and section 4 provides an algorithm to automatically construct such a product program.

2.1 Intermediate and Assembly Representations

2.1.1 IR. The unoptimized IR used to represent \mathbb{C} is mostly a subset of LLVM – it supports all the primitive types (integer, float, code labels) and the derived types (pointer, array, struct, procedure) of LLVM. Being unoptimized, our IR does not need to support LLVM's undef and poison values,

$va_start(ap, last)$	$va_arg(ap, \tau)$	$va_copy(aq, ap)$	$va_end(ap)$
$a := va_start_ptr$ store void*, 4, a , $\langle ap \rangle$	$a := load\ void^*, 4, \langle ap \rangle$ $result := load\ \langle \tau \rangle, \langle alignof(\tau) \rangle, a$ $a' := a + \langle roundup_4(sizeof(\tau)) \rangle$ store void*, 4, a' , $\langle ap \rangle$	$a := load\ void^*, 4, \langle ap \rangle$ store void*, 4, a , $\langle aq \rangle$	store void*, 4, 0, $\langle ap \rangle$

Fig. 2. Translation of C’s variadic macros to LLVM_d instructions. $roundup_4(a)$ returns the closest multiple of 4 greater than or equal to a .

it instead treats all error conditions as UB. Syntactic conversion of C to LLVM IR entails the usual conversion of types/operators. A global variable name g or a parameter name y appearing in a C procedure body is translated to the variable’s start address in IR, denoted $\boxed{lb.g}$ and $\boxed{lb.y}$ respectively¹. A local variable declaration or a call to C’s `alloca()` operator is converted to LLVM’s `alloca` instruction, and to distinguish the two, we henceforth refer to the latter as the “alloc” instruction. Unlike LLVM, our IR also supports a `dealloc` instruction that deallocates a variable at the end of its scope — we use LLVM’s `stack{save, restore}` intrinsics (that maintain equivalent scope information for a different purpose) to introduce explicit `dealloc` instructions in our IR. Henceforth, we refer to our IR as LLVM_d (for LLVM + `dealloc`).

We discuss our logical model in the context of compilation to 32-bit x86 for the relative simplicity of the calling conventions in 32-bit mode. Like LLVM, a procedure definition in LLVM_d can only return a scalar value — aggregate return value is passed in memory. Unlike LLVM which allocates memory for a parameter only if its address is taken, LLVM_d allocates memory for all parameters — LLVM_d thus takes all parameters through pointers, e.g., both n and m are passed through pointers in fig. 1b. This makes the translation of a procedure-call from C to LLVM_d slightly more verbose, as explicit instructions to (de)allocate memory for the arguments are required. An example of this translation is shown in fig. 1 where a call to `printf` at C5 of fig. 1a translates to instructions I7-I12 in fig. 1b: the LLVM_d program performs two allocations, one for the format string, and another for the variable argument list; the latter represented as an object of “struct” type containing two ints. The call instruction takes the pointers returned by these allocations as operands.

Figure 2 shows the C-to-LLVM_d translations for variadic macros. The translation rules have template holes marked by $\langle \rangle$ for types and variables of C which are populated at the time of translation with appropriate LLVM_d entities. LLVM_d’s `va_start_ptr` instruction returns the first address of the current procedure’s variable argument list.

2.1.2 Assembly. Broadly, an assembly program \mathbb{A} consists of a code section (with a sequence of assembly instructions), a data section (with read-only and read-write global variables), and a symbol table that maps string symbols to memory addresses in code and data sections. The validator checks that the regions specified by the symbol table are well-aligned and non-overlapping, and uses it to relate a global variable or procedure in \mathbb{C} to its address or implementation in \mathbb{A} .

We assume that the OS guarantees the caller-side contract of the ABI calling conventions for the entry procedure, `main()`. For 32-bit x86, this means that at the start of program execution, the stackpointer is available in register `esp`, and the return address and input parameters (`argc`, `argv`) to `main()` are available in the stack region just above the stackpointer. For other procedure-calls, the validator verifies the adherence to calling conventions at a callsite (in the caller) and assumes adherence at procedure entry (in the callee). Heap allocation procedures like `malloc()` are left uninterpreted, and so, the only compiler-visible way to allocate (and deallocate) memory in \mathbb{A} is through the decrement (and increment) of the stackpointer stored in register `esp`.

2.1.3 Allocation and Deallocation. Allocation and deallocation instructions appear only in \mathbb{C} , and do not appear in \mathbb{A} . Let C represent a procedure in program \mathbb{C} .

¹As we will also see later, $\boxed{lb.v}$ denotes the *lower bound* of the memory addresses occupied by variable with name v .

An LLVM_d instruction “ $p_C^a: v := \text{alloc } n, \tau, \text{align}$ ” allocates a contiguous region of local memory with space for n elements of type τ aligned by align , and returns its start address in v . The PC, p_C^a , of an alloc instruction is also called an *allocation site* (denoted by z), and let the set of allocation sites in C be Z . During conversion of the C program to LLVM_d, we distinguish between allocation sites due to the declaration of a local variable (or a procedure-call argument) and allocation sites due to a call to `alloca()` – we use Z_l for the former and Z_a for the latter, so that $Z = Z_l \cup Z_a$.

The address of an allocated region is non-deterministic, but is subject to two *Well-Formedness (WF) constraints*: (1) the newly allocated memory region should be separate from all currently allocated memory regions, i.e., there should be no overlap; and (2) the address of the newly allocated memory region should be aligned by `align`.

An LLVM_d instruction “ $p_C^d: \text{dealloc } z$ ” deallocates *all* local memory regions allocated due to the execution of allocation site $z \in Z$.

2.2 Transition Graph Representation

An LLVM_d or assembly instruction may mutate the machine state, transfer control, perform I/O, or terminate the execution. We represent a C procedure, C , as a transition graph, $C = (\mathcal{N}_C, \mathcal{E}_C)$, with a finite set of nodes $\mathcal{N}_C = \{n^s = n_1, n_2, \dots, n_m\}$, and a finite set of labeled directed edges \mathcal{E}_C . A unique node n^s represents the start node or entry point of C , and every other node n_j must be reachable from n^s . A node with no outgoing edges is a *terminating node*. A variable in C is identified by its scope-resolved unique name. The machine state of C consists of the set of input parameters \vec{y} , set of temporary variables \vec{t} , and an explicit array variable M_C denoting the current state of memory. We use i_N to denote a bitvector type of size $N > 0$. M_C is of type $\text{T}(M_C) = i_{32} \rightarrow i_8$.

An assembly implementation of the C procedure C , identified through the symbol table, is the assembly procedure A . The machine state of A consists of its hardware registers \vec{regs} and memory M_A . Similarly to C , $A = (\mathcal{N}_A, \mathcal{E}_A)$ is also represented as a transition graph.

Let $P \in \{C, A\}$. In addition to the memory (data) state M_P , we also need to track the allocation state, i.e., the set of intervals of addresses that have been allocated by the procedure. We use α (potentially with a subscript) to denote a memory address of bitvector type. Let $i = [\alpha_b, \alpha_e]$ represent an *address interval* starting at α_b and ending at α_e (both inclusive), such that $\alpha_b \leq_u \alpha_e$ (where \leq_u is unsigned comparison operator). Let $[\alpha]_w$ be a shorthand for the address interval $[\alpha, \alpha + w - 1_{i_{32}}]$, where $n_{i_{32}}$ is the two’s complement representation of integer n using 32 bits.

2.2.1 Address Set. Let Σ (potentially with a sub- or superscript) represent a set of addresses, or an *address set*. An empty address set is represented by \emptyset , and an address set of contiguous addresses is an interval i . Two address sets overlap, written $\text{ov}(\Sigma_1, \Sigma_2)$, iff $\Sigma_1 \cap \Sigma_2 \neq \emptyset$. Extended to $m \geq 2$ sets, $\text{ov}(\Sigma_1, \Sigma_2, \dots, \Sigma_m) \Leftrightarrow \exists_{1 \leq j_1 < j_2 \leq m} \text{ov}(\Sigma_{j_1}, \Sigma_{j_2})$. $|\Sigma|$ represents the number of distinct addresses in Σ . For a non-empty address set, $\text{lb}(\Sigma)$ and $\text{ub}(\Sigma)$ represent the smallest and largest address respectively in Σ . $\text{comp}(\Sigma)$ represents the *complement* of Σ , so that: $\forall \alpha : (\alpha \in \Sigma) \Leftrightarrow (\alpha \notin \text{comp}(\Sigma))$.

2.2.2 Memory Regions. To support dynamic (de)allocation, an execution model needs to individually track regions of memory belonging to each variable, heap, stack, etc. We next describe the memory regions tracked by our model.

- (1) Let G be the set of names of all global variables in C . For each global variable $g \in G$, we track the memory region belonging to that variable. We use the name of a global variable $g \in G$ as its *region identifier* to identify the region belonging to g in both C and A .
- (2) For a procedure C , let Y be the set of names of formal parameters, including the variadic parameter, if present. We use the special name `vrdc` for the variadic parameter. The memory region belonging to a parameter $y \in Y$ is called y in both C and A .

- (3) The memory region allocated by allocation site $z \in Z$ is called z in C . In A , our algorithm potentially annotates allocation instructions corresponding to an allocation site z in C .
- (4) hp denotes the region belonging to the *program heap* (managed by the OS) in both C and A .
- (5) Local variables and actual arguments may be allocated by the *call chain* of a procedure (caller, caller's caller, and so on). This is denoted by region cl , or *callers' locals*, in both C and A .
- (6) In procedure A , stack memory can be allocated and deallocated through stackpointer decrement and increment. The addresses belonging to the stack frame of A (but not to a local variable $z \in Z$ or a parameter $y \in Y$) belong to the stk (stack) region in A . The stk region is absent in C .
- (7) Similarly, in A , we use cs (*callers' stack*) to identify the region that belongs to the stack space (but not to cl) of the call-chain of procedure A . cs is absent in C .
- (8) Program \mathbb{A} may use more global memory than \mathbb{C} , e.g., to store pre-computed constants to implement vectorizing transformations. Let F be the set of names of all *assembly-only global variables* in \mathbb{A} . For each $f \in F$, its memory region in A is identified by f .
- (9) The region $free$ denotes the free space, that does not belong to any of the aforementioned regions, in both C and A ,

Let $R = G \cup Y \cup Z \cup F \cup \{hp, cl, stk, cs, free\}$ represent all *region identifiers*; $B = G \cup Y \cup Z \cup \{hp, cl\}$ denote the regions in both C and A ; and $S = \{stk, cs\}$ denote the stack regions in A .

Let $G_r \subseteq G$ be the set of read-only global variables in \mathbb{C} ; and, let $G_w = G \setminus G_r$ denote the set of read-write global variables. We define $F_r \subseteq F$ and $F_w = F \setminus F_r$ analogously.

For each non-free region $r \in (R \setminus \{free\})$, the machine state of a procedure P includes a unique variable Σ_p^r that tracks region r 's address set as P executes. If Σ_p^r is a contiguous non-empty interval, we also refer to it as i_p^r . For $r \in G \cup Y \cup F \cup \{hp, cl, cs\}$, Σ_p^r remains constant throughout P 's execution. For $\vec{r} \subseteq R$, we define an expression $\Sigma_p^{\vec{r}} = \bigcup_{r \in \vec{r}} \Sigma_p^r$. Because C does not have a stack or an assembly-only global variable, $\Sigma_C^{F \cup S} = \emptyset$ holds throughout C 's execution. At any point in P 's execution, the free space can be computed as $\Sigma_p^{free} = \text{comp}(\Sigma_p^{B \cup F \cup S})$. Notice that we do not use an explicit variable to track Σ_p^{free} .

2.2.3 Ghost Variables. Our validator introduces *ghost variables* in a procedure's execution semantics, i.e., variables that were not originally present in P . We use \boxed{x} to indicate that x is a ghost variable. For each region $r \in G \cup Y \cup Z$ ($r \in F$), we introduce $\boxed{em.r}$, $\boxed{lb.r}$, and $\boxed{ub.r}$ in C (A) to track the *emptiness* (whether the region is empty), *lower bound* (smallest address), and *upper bound* (largest address) of Σ_C^r (Σ_A^r) respectively; for $r \in G \cup Y$ ($r \in F$), $\boxed{sz.r}$ tracks the size of Σ_C^r (Σ_A^r), and for $z \in Z$, $\boxed{1stSz.z}$ tracks the *size of last allocation* due to allocation-site z . Σ_p^d and Σ_p^{sw} track the set of addresses read and written by P respectively. Let $\boxed{\ast}$ be the set of all ghost variables.

2.2.4 Error Codes. Execution of C or A may terminate successfully, may never terminate, or may terminate with an error. We support two error codes to distinguish between two categories of errors: \mathcal{U} and \mathcal{W} . In C : \mathcal{U} represents an occurrence of UB, and \mathcal{W} represents a violation of a WF constraint that needs to be ensured either by the compiler or the OS (both external to the program itself). In A : \mathcal{U} represents UB or a translation error, and \mathcal{W} represents occurrence of a condition that can be assumed to never occur, e.g., if the OS ensures that it never occurs. In summary, for a procedure P , \mathcal{W} represents an error condition that P can *assume* to be absent (because the external environment ensures it), while \mathcal{U} represents an error that P must *ensure* to be absent.

2.2.5 Outside World and Observable Trace. Let Ω_P be a state of the outside world (OS/hardware) for P that supplies external inputs whenever P reads from it, and consumes external outputs generated by P . Ω_P is assumed to mutate arbitrarily but deterministically based on the values consumed or produced due to the I/O operations performed by P during execution. Let T_P be a potentially-infinite sequence of observable trace events generated by an execution of P .

2.2.6 Expressions. Let variable v and variables \vec{v} or \vec{x} be drawn from $\text{Vars} = (\vec{t}, \overrightarrow{regs}, M_p, \Sigma_p^r, \oplus)$ (for all $P \in \{C, A\}$ and for all $r \in (R \setminus \{\text{free}\})$); $e(\vec{x})$ be an expression over \vec{x} , and $E(\vec{x})$ be a list of expressions over \vec{x} . An expression $e(\vec{x})$ is a well-formed combination of constants, variables \vec{x} , and arithmetic, logical, relational, memory access, and address set operators. For memory reads and writes, select (sel for short) and store (st for short) operations are used to access and modify M_p at a given address α . Further, the sel and st operators are associated with a sz parameter: $\text{sel}_{sz}(\text{arr}, \alpha)$ returns a little-endian concatenation of sz bytes starting at α in the array arr. Similarly, $\text{st}_{sz}(\text{arr}, \alpha, \text{data})$ returns a new array that has contents identical to arr except for the sz bytes starting at α which have been replaced by data in little-endian format. To encode reads/writes to a region of memory, we define projection and updation operations.

Definition 2.1. $\pi_\Sigma(M_p)$ denotes the *projection* of M_p on addresses in Σ , i.e., if $M'_p = \pi_\Sigma(M_p)$, then $\forall \alpha \in \Sigma : \text{sel}_1(M'_p, \alpha) = \text{sel}_1(M_p, \alpha)$ and $\forall \alpha \notin \Sigma : \text{sel}_1(M'_p, \alpha) = 0$. The sentinel value 0 is used for the addresses outside Σ . We use $M_{p_1} =_\Sigma M_{p_2}$ as shorthand for $(\pi_\Sigma(M_{p_1}) = \pi_\Sigma(M_{p_2}))$.

Definition 2.2. $\text{upd}_\Sigma(M_p, M)$ denotes the *updation* of M_p on addresses in Σ using the values in M . If $M'_p = \text{upd}_\Sigma(M_p, M)$, then $M'_p =_\Sigma M$ and $M'_p =_{\text{comp}(\Sigma)} M_p$ hold.

2.2.7 Instructions. Each edge $e_p \in \mathcal{E}_p$ is labeled with one of the following *graph instructions*:

- (1) A *simultaneous assignment* of the form $\vec{v} := E(\vec{x})$. Because variables \vec{v} and \vec{x} may include M_p , an assignment suffices for encoding memory loads and stores. Similarly, because the variables may be drawn from Σ_p^z (for an allocation site z), an assignment is also used to encode the allocation of an interval $i_{\text{new}} (\Sigma_p^z := \Sigma_p^z \cup i_{\text{new}})$ and the deallocation of all addresses allocated due to z ($\Sigma_p^z := \emptyset$). Stack allocation and deallocation in A can be similarly represented as $\Sigma_A^{\text{stk}} := \Sigma_A^{\text{stk}} \cup i_{\text{new}}$ and $\Sigma_A^{\text{stk}} := \Sigma_A^{\text{stk}} \setminus i_{\text{new}}$ respectively.
- (2) A *guard* instruction, $e(\vec{x})?$, indicates that when execution reaches its head, the edge is taken iff its *edge condition* $e(\vec{x})$ evaluates to true. For every other instruction, the edge is always taken upon reaching its head, i.e., its edge condition is true. For a non-terminating node $n_p \in \mathcal{N}_p$, the guards of all edges departing from n_p must be mutually exclusive, and their disjunction must evaluate to true.
- (3) A type-parametric *choose* instruction $\theta(\vec{\tau})$. Instruction $\vec{v} := \theta(\vec{\tau})$ non-deterministically chooses values of types $\vec{\tau}$ and assigns them to variables \vec{v} , e.g., a memory with non-deterministic contents is obtained by using $\theta(i_{32} \rightarrow i_8)$.
- (4) A *read* (rd) or *write* (wr) I/O instruction. A read instruction $\vec{v} := \text{rd}(\vec{\tau})$ reads values of types $\vec{\tau}$ from the outside world into variables \vec{v} , e.g., an address set is read using $\Sigma := \text{rd}(2^{i_{32}})$. A write instruction $\text{wr}(V(E(\vec{x})))$ writes the value constructed by value constructor V using $E(\vec{x})$ to the outside world. A value constructor is defined for each type of observable event. For a procedure-call, $\text{fcall}(\rho, \vec{v}, \vec{r}, M)$ represents value constructed for a procedure-call to callee with name (or address) ρ , the actual arguments \vec{v} , callee-observable regions \vec{r} , and memory M . Similarly, $\text{ret}(E(\vec{x}))$ is a value constructed during procedure return that captures observable values computed through $E(\vec{x})$. Local (de)allocation events have their own value constructors, $\text{allocBegin}(z, w)$, $\text{allocEnd}(z, i, M)$, and $\text{dealloc}(z)$, which represent (de)allocation due to allocation site z with the associated observables w (size), i (interval), and M (memory). A read or write instruction mutates Ω_p arbitrarily based on the read and written values. Further, the data items read or written are appended to the observable trace T_p . Let $\text{read}_{\vec{\tau}}(\Omega_p)$ be an uninterpreted function that reads values of types $\vec{\tau}$ from Ω_p ; and $\text{io}(\Omega_p, \vec{v}, \text{rw})$ be an uninterpreted function that returns an updated state of Ω_p after an I/O operation of type $\text{rw} \in \{r, w\}$ (read or write) with values \vec{v} . Thus, in its explicit syntax, $\vec{v} := \text{rd}(\vec{\tau})$ translates to a sequence of instructions: $\vec{v} := \text{read}_{\vec{\tau}}(\Omega_p)$; $\Omega_p := \text{io}(\Omega_p, \vec{v}, r)$; $T_p := T_p \cdot \vec{v}$, where \cdot is the

Table 1. Definitions of operators and predicates used in translations in figs. 3 to 6

Operator	Definition
$\text{sz}(\tau)$	Returns the size (in bytes) of type τ . For example, $\text{sz}(i_{32}) = 4$, $\text{sz}(i_8^*) = 4$, and $\text{sz}(\{80 \times i_8\}) = 80$.
$\overline{T}(a)$	Returns the type τ of a where a can be a global variable, a parameter, or a register. For example, $\overline{T}(\text{eax}) = i_{32}$.
$\Delta_\tau(\text{eax}, \text{edx})$	A macro operator which derives the return value of an assembly procedure with return type τ from input registers eax and edx using the calling conventions, e.g., $\Delta_{i_8}(\text{eax}, \text{edx}) = \text{extract}_{7,0}(\text{eax})$, $\Delta_{i_{32}}(\text{eax}, \text{edx}) = \text{eax}$, $\Delta_{i_{64}}(\text{eax}, \text{edx}) = \text{concat}(\text{edx}, \text{eax})$, where $\text{extract}_{h,l}(a)$ extracts bits h down to l from a and $\text{concat}(a, b)$ returns the bitvector concatenation of a and b where b takes the less significant position.
$\nabla_\tau(v)$	Inverse of $\Delta_\tau(\text{eax}, \text{edx})$. Distributes the packed bitvector v of type τ into two bitvectors of 32 bit-width each, setting the bits not covered by v to some non-deterministic value.
$\text{ROM}_P^r(i)$	Returns a memory array containing the contents of read-only global variable named r in P . The contents are mapped at the addresses in the provided interval i .
$\text{addrSets}_F()$	Returns the address sets of the assembly-only global variables F using the symbol table in \mathbb{A} .
Predicate	Definition
$\text{aligned}_n(a)$	Bitvector a is at least n bytes aligned. Equivalent to: $a\%n = 0$, where $\%$ is remainder operator.
$\text{isAlignedIntrvl}_a(p, w)$	A w -sized sequence of addresses starting at p is aligned by a and does not wraparound. Equivalent to: $\text{aligned}_a(p) \wedge (p \leq_u p + w - 1_{i_{32}})$.
$\text{accessIsSafeC}_{\tau,a}(p, \Sigma)$	Equivalent to: $\text{isAlignedIntrvl}_{\text{align}(a)}(p, \text{sz}(\tau)) \wedge (\{p\}_{\text{sz}(\tau)} \subseteq \Sigma)$.
$\text{addrSetsAreWF}(\Sigma_P^{hp}, \Sigma_P^{cl}, \dots, i_p^g, \dots, \Sigma_P^f, \dots, i_p^y, \dots, \Sigma_P^{\text{vrdc}}, \dots, \Sigma_P^{\text{vrdc}})$	The address sets passed as parameter are well-formed with respect to C semantics. Equivalent to: $(0_{i_{32}} \notin \Sigma_P^{GUFUYU(hp, cl)}) \wedge \neg \text{ov}(\Sigma_P^{hp}, \Sigma_P^{cl}, \dots, i_p^g, \dots, \Sigma_P^f, \dots, i_p^y, \dots, \Sigma_P^{\text{vrdc}}) \wedge (\Sigma_P^{\text{vrdc}} \neq \emptyset \Rightarrow \text{isInterval}(\Sigma_P^{\text{vrdc}}) \wedge \forall r \in \text{GU}(Y^{\text{vrdc}}) \cup F(\{i_p^r\} = \text{sz}(T(r)) \wedge \text{aligned}_{\text{align}(r)}(\text{lb}(i_p^r))))$, where $\text{isInterval}(\Sigma_P^{\text{vrdc}})$ holds iff the address set Σ_P^{vrdc} is an interval, $\text{align}(r)$ returns the alignment of variable r .
$\text{intrvlInSet}(\alpha_b, \alpha_e, \Sigma)$	The pair (α_b, α_e) forms a valid interval inside the address set Σ . Equivalent to: $(\alpha_b \neq 0_{i_{32}}) \wedge (\alpha_b \leq_u \alpha_e) \wedge (\{\alpha_b, \alpha_e\} \subseteq \Sigma)$
$\text{intrvlInSet}_a(\alpha_b, \alpha_e, \Sigma)$	Equivalent to: $\text{aligned}_a(\alpha_b) \wedge \text{intrvlInSet}(\alpha_b, \alpha_e, \Sigma)$
$\text{obeyCC}(e_{\text{esp}}, \vec{\tau}, \vec{x})$	Pointers \vec{x} match the expected addresses of arguments for a procedure-call in assembly. Based on the calling conventions, obeyCC uses the value of the current stackpointer (e_{esp}) and parameter types ($\vec{\tau}$) to obtain the expected addresses of the arguments. For example, $\text{obeyCC}(\text{esp}, (i_8, i_{32}), (\text{esp}, \text{esp} + 4_{i_{32}}))$ holds.
$\text{overflow}_{mul}(a, b)$	Signed multiplication of bitvectors a, b overflows. E.g., $\text{overflow}_{mul}(2147483647_{i_{32}}, 2_{i_{32}})$ holds.
$\text{stkIsWF}(\text{esp}, \text{stk}_e, \text{cs}_e, \vec{\tau})$	The pairs $(\text{esp}, \text{stk}_e)$, $(\text{cs}_e, \text{stk}_e)$ represent well-formed intervals for initial stk and cs regions with respect to parameter types $\vec{\tau}$ and other (input) address sets in \mathbb{A} . Equivalent to: $\text{aligned}_{16}(\text{esp} + 4_{i_{32}}) \wedge (\text{esp} \leq_u \text{esp} + 4_{i_{32}}) \wedge \neg \text{ov}([\text{esp}]_{4_{i_{32}}}, \Sigma_A^{GUFUYU(hp, cl)}) \wedge \text{obeyCC}(\text{esp} + 4_{i_{32}}, \vec{\tau}, \dots, \text{lb}(i_A^y), \dots) \wedge (\text{stk}_e <_u \text{cs}_e) \wedge \neg \text{ov}([\text{stk}_e] + 1_{i_{32}}, \text{cs}_e], \Sigma_A^{GUFUYU(hp)}) \wedge \Sigma_A^{cl} \subseteq [\text{stk}_e] + 1_{i_{32}}, \text{cs}_e]$
$\text{UB}_P(\text{op}, \vec{x})$	Application of operation op of procedure P on arguments \vec{x} triggers UB. E.g., $\text{UB}_C(\text{udiv}, (1_{i_{32}}, 0_{i_{32}}))$ holds.

trace concatenation operator. Similarly, $\text{wr}(V(E(\vec{x})))$ translates to: $\Omega_P := \text{io}(\Omega_P, V(E(\vec{x})), w)$; $T_P := T_P \cdot V(E(\vec{x}))$. Henceforth, we only use the implicit syntax for brevity.

- (5) An error-free and error-indicating *halt* instruction that terminates execution. $\text{halt}(\emptyset)$ indicates termination without error, and $\text{halt}(\ast)$ indicates termination with error code $\ast \in \{\mathcal{U}, \mathcal{W}\}$. Upon termination without error, a special *exit* event is appended to observable trace T_P . Upon termination with error, the error code is appended to T_P . Thus, the destination of an edge with a *halt* instruction is a terminating node. We create a unique terminating node for an error-free exit. We also create a unique terminating node for each error code, also called an *error node*; an edge terminating at an error node is called an *error edge*. \mathcal{U}_P and \mathcal{W}_P represent error nodes in P for error \mathcal{U} and \mathcal{W} respectively. Execution transfers to an error node upon encountering the corresponding error. Let $\mathcal{N}_P^{\text{NW}} = \mathcal{N}_P \setminus \{\mathcal{U}_P, \mathcal{W}_P\}$ be the set of non-error nodes in P .

In addition to the observable trace events generated by *rd*, *wr*, and *halt* instructions, the execution of every instruction in P also appends an observable *silent trace event*, denoted \perp , to T_P . Silent trace events count the number of executed instructions as a proxy for observing the passage of time.

2.3 Translations of C and A to Their Graph Representations

Figures 3 and 4 (and figs. 5 and 6 later) present the key translation rules from LLVM_d and (abstracted) assembly instructions to graph instructions. Each rule is composed of three parts separated by a horizontal line segment: on the left is the name of the rule, above the line segment is the LLVM_d /assembly instruction, and below the line segment is the graph instructions listing. We describe the operators and predicates used in the rules in table 1. As an example, the top right

$$\begin{array}{c}
\text{(ENTRY}_C) \frac{p_C^j : \text{def } C(\vec{\tau})}{\Sigma_C^{hp}, \Sigma_C^{cl}, \dots, i_C^g, \dots, \Sigma_C^{vrdc} := \text{rd}(2^{i_{32}}, 2^{i_{32}}, \dots, 2^{i_{32}}); \\
\Sigma_C^{stk}, \Sigma_C^{cs}, \dots, \Sigma_C^f, \dots, \Sigma_C^g, \dots, \Sigma_C^{rd}, \Sigma_C^{sw} := \emptyset, \emptyset, \dots, \emptyset; \\
\text{if } (\neg \text{addrSetsAreWF}(\Sigma_C^{hp}, \Sigma_C^{cl}, \dots, i_C^g, \dots, \Sigma_C^f, \dots, i_C^g, \dots, \Sigma_C^{vrdc})) \\
\text{halt}(\mathcal{W}); \\
M_C := \theta(i_{32} \rightarrow i_8); \quad M_C := \text{upd}_{s,B}(M_C, \text{rd}(i_{32} \rightarrow i_8)); \\
\text{for } g \text{ in } G_r \{ M_C := \text{upd}_g(M_C, \text{ROM}_C^g(i_C^g)); \} \\
\text{for } z \text{ in } Z \{ \text{em.z} := \text{true}; \quad \beta_M(z) := \emptyset; \} \\
\text{for } r \text{ in } G \cup Y \cup \{hp, cl\} \{ \beta_M(r) := G \cup \{hp, cl\}; \} \\
\text{for } r \text{ in } G \cup Y \{ \text{sz.r}; \text{em.r} := |\Sigma_C^r|; (|\Sigma_C^r| = 0_{i_{32}}); \\
\text{if } (\neg \text{em.r}) \{ \text{lb.r}; \text{ub.r} := \text{lb}(\Sigma_C^r), \text{ub}(\Sigma_C^r); \} \quad \beta(\text{lb.r}) := \{r\}; \\
\} \\
\text{(ALLOC)} \frac{z : v := \text{alloc } n, \tau, a}{\text{IF } \{z \in Z_l\} \{ \text{if } (n \leq_s 0_{i_{32}} \vee \text{overflow}_{mul}(n, \text{sz}(\tau))) \text{halt}(\mathcal{W}); \} \\
\text{wr}(\text{allocBegin}(z, n * \text{sz}(\tau))); \\
\alpha_b := \theta(i_{32}); \quad \alpha_e := \alpha_b + n * \text{sz}(\tau) - 1_{i_{32}}; \\
\text{if } (\neg \text{intrvlInSet}_a(\alpha_b, \alpha_e, \Sigma_C^{\text{free}})) \text{halt}(\mathcal{W}); \\
\Sigma_C^z := \Sigma_C^z \cup [\alpha_b, \alpha_e]; \quad M_C := \text{upd}_{[\alpha_b, \alpha_e]}(M_C, \theta(i_{32} \rightarrow i_8)); \\
\text{lb.z} := \text{em.z} ? \alpha_b : \min(\text{lb.z}, \alpha_b); \quad \text{lstSz.z} := n * \text{sz}(\tau); \\
\text{ub.z} := \text{em.z} ? \alpha_e : \max(\text{ub.z}, \alpha_e); \quad \text{em.z} := \text{false}; \\
v := \alpha_b; \quad \beta(v) := \{z\}; \\
\text{wr}(\text{allocEnd}(z, [\alpha_b, \alpha_e], \pi_{[\alpha_b, \alpha_e]}(M_C))); \\
\text{(DEALLOC)} \frac{p_C^j : \text{dealloc } z}{\Sigma_C^z := \emptyset; \quad \text{em.z} := \text{true}; \\
\text{wr}(\text{dealloc}(z)); \\
\text{(VASTARTPTR)} \frac{p_C^j : p := \text{va_start_ptr}}{\text{if } (\Sigma_C^{\text{vrdc}} = \emptyset) \{ \\
p := 0_{i_{32}}; \quad \beta(p) := \emptyset; \\
\} \text{ else } \{ \\
p := \text{lb.vrdc}; \quad \beta(p) := \{\text{vrdc}\}; \\
\} \\
\text{(LOAD}_C) \frac{p_C^j : v := \text{load } \tau, a, p}{\text{if } (\neg \text{accessIsSafe}_{C,\tau,a}(p, \Sigma_C^{\beta(p)})) \text{halt}(\mathcal{W}); \\
v := \text{sel}_{\text{sz}(\tau)}(M_C, p); \\
\beta(v) := \beta_M(\beta(p)); \quad \Sigma_C^{\text{rd}} := \Sigma_C^{\text{rd}} \cup [p]_{\text{sz}(\tau)}; \\
\text{(STORE}_C) \frac{p_C^j : \text{store } \tau, a, v, p}{\text{if } (\neg \text{accessIsSafe}_{C,\tau,a}(p, \Sigma_C^{\beta(p) \setminus G_r})) \text{halt}(\mathcal{W}); \\
M_C := \text{st}_{\text{sz}(\tau)}(M_C, p, v); \\
\beta_M(\beta(p)) := \beta_M(\beta(p)) \cup \beta(v); \quad \Sigma_C^{\text{wr}} := \Sigma_C^{\text{wr}} \cup [p]_{\text{sz}(\tau)}; \\
\text{(CALLV)} \frac{p_C^j : \text{call void } \rho(\vec{\tau} \vec{x})}{\beta^* := \beta_M^* \left(\bigcup_{x \in \vec{x}} \beta(x) \cup G \cup \{hp\} \right); \\
\text{wr}(\text{fcall}(\rho, \vec{x}, \beta^*, \pi_{\Sigma_C^{\beta^*}}(M_C))); \\
M_C := \text{upd}_{\Sigma_C^{\beta^*} \setminus G_r}(M_C, \text{rd}(i_{32} \rightarrow i_8)); \\
\beta_M(\beta^* \setminus G_r) := \beta^*; \\
\text{(CALLC)} \frac{p_C^j : v := \text{call } \gamma \rho(\vec{\tau} \vec{x}) \quad \gamma \neq \text{void}}{\beta^* := \beta_M^* \left(\bigcup_{x \in \vec{x}} \beta(x) \cup G \cup \{hp\} \right); \\
\text{wr}(\text{fcall}(\rho, \vec{x}, \beta^*, \pi_{\Sigma_C^{\beta^*}}(M_C))); \\
M_C := \text{upd}_{\Sigma_C^{\beta^*} \setminus G_r}(M_C, \text{rd}(i_{32} \rightarrow i_8)); \\
v := \text{rd}(\gamma); \quad \beta(v), \beta_M(\beta^* \setminus G_r) := \beta^*, \beta^*; \\
\text{(OP)} \frac{p_C^j : v := \text{op}(\vec{x})}{\text{if } (\text{UB}_C(\text{op}, \vec{x})) \text{halt}(\mathcal{W}); \\
v := \text{op}(\vec{x}); \quad \beta(v) := \beta^{\text{op}}(\beta(\vec{x})); \\
\text{(RETV)} \frac{p_C^j : \text{ret void}}{\text{wr}(\text{ret}(\pi_{\Sigma_C^B}(M_C))); \\
\text{halt}(\emptyset); \\
\text{(RETC)} \frac{p_C^j : \text{ret } v}{\text{wr}(\text{ret}(v, \pi_{\Sigma_C^B}(M_C))); \\
\text{halt}(\emptyset); \\
\text{(ASSIGNCONST)} \frac{p_C^j : v := c}{v := c; \quad \beta(v) := \emptyset; \\
\text{(DEALLOC)} \frac{p_C^j : \text{dealloc } z}{\Sigma_C^z := \emptyset; \quad \text{em.z} := \text{true}; \\
\text{wr}(\text{dealloc}(z)); \\
\text{(VASTARTPTR)} \frac{p_C^j : p := \text{va_start_ptr}}{\text{if } (\Sigma_C^{\text{vrdc}} = \emptyset) \{ \\
p := 0_{i_{32}}; \quad \beta(p) := \emptyset; \\
\} \text{ else } \{ \\
p := \text{lb.vrdc}; \quad \beta(p) := \{\text{vrdc}\}; \\
\}
\end{array}$$

Fig. 3. Translation rules for converting LLVM_d instructions to graph instructions.

corner of fig. 3 shows the parametric (OP) rule which gives the translation of an operation using arithmetic/logical/relational operator op in LLVM_d to corresponding graph instructions. We use C-like constructs in graph instructions as syntactic sugar for brevity, e.g. ‘;’ is used for sequencing, ‘?’ is used for conditional assignment, and **if**, **else**, and **for** are used for control flow transfer. We highlight the read and write I/O instructions with a shaded background, and use **bold**, **colored** fonts for halt instructions. We use macros IF and ELSE to choose translations based on a boolean condition on the input syntax.

2.3.1 Translation of C. Figure 3 shows the translation rules for converting LLVM_d instructions to graph instructions. The (ENTRY_C) rule presents the initialization performed at the entry of a procedure C. The address sets and memory state of C are initialized using reads from the outside world Ω_C. The address sets that are read are checked for well-formedness with respect to C semantics, or else error \mathcal{W} is triggered. The ghost variables are also appropriately initialized.

The (ALLOC) and (DEALLOC) rules provide semantics for the allocation and deallocation of local memory at allocation site z — if z ∈ Z_l, n (the number of elements allocated) has additional constraints for a UB-free execution. A (de)allocation instruction generates observable traces using the wr instruction at the beginning and end of each execution of that instruction. We will later use

these traces to identify a lockstep correlation of (de)allocation events between C and A , towards validating a translation.

In (OP), an application of op may trigger UB for certain inputs, as abstracted through the $UB_C(op, \vec{x})$ operation. While there are many UBs in the C standard, we model only the ones that we have seen getting exploited by the compiler for optimization. These include the UB associated with a logical or arithmetic shift operation (second operand should be bounded by a limit which is determined by the size of the first operand), address computation (no over- and under-flow), and division operation (second operand should be non-zero). In (LOAD $_C$) and (STORE $_C$), a UB-free execution requires the dereferenced pointer p to be non-NULL ($\neq 0_{i_{32}}$ in our modeling), aligned by a , and have its access interval belong to the regions which p may point to, or p may be based on (§6.5.6p8 of the C17 standard).

To identify the regions a pointer p may point to, we define two maps: (1) $\beta : \text{Vars} \rightarrow 2^R$, so that for a (pointer) variable $x \in \text{Vars}$, $\beta(x)$ returns the set of regions x may point to; and (2) $\beta_M : R \rightarrow 2^R$, so that for a region $r \in R$, $\beta_M(r)$ returns the set of regions that some (pointer) value stored in $\pi_{\Sigma_C}^r(M_C)$ may point to. $\beta(\vec{x})$ is equivalent to $\bigcup_{x \in \vec{x}} \beta(x)$, and $\beta_M(\vec{r})$ is equivalent to $\bigcup_{r \in \vec{r}} \beta_M(r)$. Similarly, $\beta_M(\vec{r}_1) := \vec{r}_2$ is equivalent to ‘for r_1 in \vec{r}_1 { $\beta_M(r_1) := \vec{r}_2$; }’. The initialization and updation of β and β_M due to each LLVM $_d$ instruction can be seen in fig. 3. For an operation op , $\beta^{op} : 2^R \rightarrow 2^R$ represents the over-approximate abstract transfer function for $v := op(\vec{x})$, that takes as input $\beta(\vec{x})$ and returns $\beta(v)$. We use $\beta^{op}(\vec{r}) = \vec{r}$ if op is bitwise complement and unary negation. We use $\beta^{op}(\vec{r}_1, \dots, \vec{r}_m) = \bigcup_{1 \leq j \leq m} \vec{r}_j$ if op is bitvector addition, subtraction, shift, bitwise- $\{\text{and, or}\}$, extraction, or concatenation. We use $\beta^{op}(\vec{r}) = \emptyset$ if op is bitvector multiplication, division, logical, relational or any other remaining operator.

The translation of an LLVM $_d$ procedure-call is given by the rules (CALLV) and (CALLC) and involves producing non-silent observable trace events using the wr instruction for the callee name/address, arguments, and callee-accessible regions and memory state. To model return values and side-effects to the memory state due to a callee, rd instructions are used. A callee may access a memory region iff it is *transitively reachable* from a global variable $g \in G$, the heap hp , or one of the arguments $x \in \vec{x}$. The transitively reachable memory regions are over-approximately computed through a reflexive-transitive closure of β_M , denoted β_M^* .

A rd instruction clobbers the callee-observable state elements arbitrarily. Thus, if a callee procedure terminates, wr and rd instructions over-approximately model the execution of a procedure-call. Later, our definition of refinement (section 2.4) caters to the case when a callee procedure may not terminate.

2.3.2 Translation of A . The translation rules for converting assembly instructions to graph instructions are shown in fig. 4. The assembly opcodes are abstracted to an IR-like syntax for ease of exposition. For example, in (LOAD $_A$), a memory read operation is represented by a load instruction which is annotated with address p , access size w (in bytes), and required alignment a (in bytes). Similarly, in (STORE $_A$), a memory write operation is represented by a store instruction with similar operands. Both (LOAD $_A$) and (STORE $_A$) translations update the ghost address sets Σ_A^{rd} and Σ_A^{wr} , in the same manner as done in C . Exceptions like division-by-zero or memory-access errors are modeled as UB in \mathbb{A} through UB_A (rules (OP-ESP) and (OP-NESP))

(OP-ESP) shows the translation of an instruction that updates the stackpointer. Assignment to the stackpointer register esp may indicate allocation (push) or deallocation (pop) of stack space. A stackpointer assignment which corresponds to a stackpointer decrement (push) is identified through predicate $isPush(p_A^j, t_b, t_a)$ where t_b and t_a are the values of esp before and after the execution of the instruction. We use $isPush(p_A^j, t_b, t_a) \Leftrightarrow (t_b >_u t_a)$. While this choice of $isPush$ suffices for most TV settings, we show in our technical report [Rose and Bansal 2024b] that if the translation is

$$\begin{array}{c}
\text{(OP-ESP)} \frac{p_A^j : \text{esp} := \text{op}(\vec{x})}{\text{if } (\text{UB}_A(\text{op}, \vec{x})) \text{ halt}(\mathcal{U}); \\ t := \text{op}(\vec{x}); \\ \text{if } (\text{isPush}(p_A^j, \text{esp}, t)) \{ \\ \quad \text{if } (\neg \text{intrvlInSet}(t, \text{esp} - 1_{132}, \Sigma_A^{\text{free}})) \text{ halt}(\mathcal{W}); \\ \quad \Sigma_A^{\text{stk}} := \Sigma_A^{\text{stk}} \cup [t, \text{esp} - 1_{132}]; \\ \quad M_A := \text{upd}_{[t, \text{esp} - 1_{132}]}(M_A, \theta(i_{132} \rightarrow i_8)); \\ \} \text{ else if } (t \neq \text{esp}) \{ \\ \quad \text{if } (\neg \text{intrvlInSet}(\text{esp}, t - 1_{132}, \Sigma_A^{\text{stk}})) \text{ halt}(\mathcal{U}); \\ \quad \Sigma_A^{\text{stk}} := \Sigma_A^{\text{stk}} \setminus [\text{esp}, t - 1_{132}]; \\ \} \\ \text{esp} := t; \quad \boxed{\text{sp}, p_A^j} := t;
} \\
\text{(ENTRY}_A) \frac{p_A^j : \text{def } A(\vec{\tau})}{\Sigma_A^{\text{hp}}, \Sigma_A^{\text{cl}}, \dots, i_A^g, \dots, i_A^y, \dots, \Sigma_A^{\text{vrdc}} := \text{rd}(2^{132}, 2^{132}, \dots, 2^{132}); \\ \dots, \Sigma_A^f, \dots := \text{addrSets}_F(); \quad \dots, \Sigma_A^z, \dots := \dots, \emptyset, \dots; \\ \text{if } (\neg \text{addrSetsAreWF}(\Sigma_A^{\text{hp}}, \Sigma_A^{\text{cl}}, \dots, i_A^g, \dots, \Sigma_A^f, \dots, i_A^y, \dots, \Sigma_A^{\text{vrdc}})) \\ \text{halt}(\mathcal{W}); \\ M_A := \theta(i_{132} \rightarrow i_8); \quad M_A := \text{upd}_{\Sigma_B}(M_A, \text{rd}(i_{132} \rightarrow i_8)); \\ \text{for } r \text{ in } G_r \cup F_r \{ M_A := \text{upd}_{\Sigma_A}(M_A, \text{ROM}_A^r(i_A^r)); \} \\ \text{for } r \text{ in } \overline{\text{reg}}_A^s \{ r := \theta(\text{T}(r)); \} \\ \boxed{\text{stk}_e} := \Sigma_A^y \neq \emptyset ? \text{ub}(\Sigma_A^y) : \text{esp} + 3_{132}; \quad \boxed{\text{cse}} := \theta(i_{132}); \\ \text{if } (\neg \text{stkIsWF}(\text{esp}, \boxed{\text{stk}_e}, \boxed{\text{cse}}, \vec{\tau}, \Sigma_A^{\text{hp}}, \Sigma_A^{\text{cl}}, \Sigma_A^{\text{GUF}}, \dots, i_A^y, \dots, \Sigma_A^{\text{vrdc}})) \\ \text{halt}(\mathcal{W}); \\ \Sigma_A^{\text{stk}} := [\text{esp}, \boxed{\text{stk}_e}] \setminus \Sigma_A^y; \quad \Sigma_A^{\text{cs}} := [\boxed{\text{stk}_e} + 1_{132}, \boxed{\text{cse}}] \setminus \Sigma_A^{\text{cl}}; \\ \boxed{\text{sp}, \text{entry}} := \text{esp}; \quad \boxed{M}^{\text{cs}} := \pi_{\Sigma_B}(\text{M}_A); \quad \boxed{\Sigma}^{\text{rd}}, \boxed{\Sigma}^{\text{wr}} := \emptyset, \emptyset; \\ \boxed{\text{ebp}}, \boxed{\text{esi}}, \boxed{\text{edi}}, \boxed{\text{ebx}}, \boxed{\text{eip}} := \text{ebp}, \text{esi}, \text{edi}, \text{ebx}, \text{sel}_4(M_A, \text{esp}); \\ \text{for } f \text{ in } F \{ \boxed{\text{sz}, f}, \boxed{\text{em}, f}, \boxed{\text{lb}, f}, \boxed{\text{ub}, f} := |\Sigma_A^f|, |\Sigma_A^f| = 0_{132}, \text{lb}(\Sigma_A^f), \text{ub}(\Sigma_A^f); \}
} \\
\text{(LOAD}_A) \frac{p_A^j : v := \text{load } w, a, p}{\text{if } (\neg \text{isAlignedIntrvl}_a(p, w) \\ \vee \text{ov}([\boxed{p}]_w, \Sigma_A^{\text{free}})) \text{ halt}(\mathcal{U}); \\ v := \text{sel}_w(M_A, p); \\ \boxed{\Sigma}^{\text{rd}} := \boxed{\Sigma}^{\text{rd}} \cup [p]_w;
} \\
\text{(STORE}_A) \frac{p_A^j : \text{store } w, a, p, v}{\text{if } (\neg \text{isAlignedIntrvl}_a(p, w) \\ \vee \text{ov}([\boxed{p}]_w, \Sigma_A^{\text{free}} \cup G_r \cup F_r)) \\ \text{halt}(\mathcal{U}); \\ M_A := \text{st}_w(M_A, p, v); \\ \boxed{\Sigma}^{\text{wr}} := \boxed{\Sigma}^{\text{wr}} \cup [p]_w;
} \\
\text{(OP-NESP)} \frac{p_A^j : r := \text{op}(\vec{x}) \quad r \neq \text{esp}}{\text{if } (\text{UB}_A(\text{op}, \vec{x})) \text{ halt}(\mathcal{U}); \\ r := \text{op}(\vec{x});
} \\
\text{(RET}_A) \frac{p_A^j : \text{ret } \tau}{\text{if } (\boxed{\text{sp}, \text{entry}} \neq \text{esp} \\ \vee \boxed{\text{ebp}} \neq \text{ebp} \vee \boxed{\text{esi}} \neq \text{esi} \\ \vee \boxed{\text{edi}} \neq \text{edi} \vee \boxed{\text{ebx}} \neq \text{ebx} \\ \vee \boxed{\text{eip}} \neq \text{sel}_4(M_A, \text{esp}) \\ \vee \neg (\boxed{M}^{\text{cs}} := \pi_{\Sigma_B}(\text{M}_A))) \text{ halt}(\mathcal{U}); \\ \text{IF } \{\tau = \text{void}\} \{ \boxed{\text{wr}}(\pi_{\Sigma_B}(M_A)); \} \\ \text{ELSE } \{ \\ \boxed{\text{wr}}(\text{ret}(\Delta_\tau(\text{eax}, \text{edx}), \pi_{\Sigma_B}(M_A))); \\ \} \\ \text{halt}(\emptyset);
}
\end{array}$$

Fig. 4. Translation rules for converting pseudo-assembly instructions to graph instructions.

performed by an adversarial compiler, discriminating a stack push from a pop is trickier and may require external trusted guidance from the user. For a stackpointer decrement, a failure to allocate stack space, either due to wraparound or overlap with other allocated space, triggers \mathcal{W} , i.e., we expect the environment (e.g., OS) to ensure that the required stack space is available to A . For a stackpointer increment, it is a translation error if the stackpointer moves out of stack frame bounds (captured by error code \mathcal{U}). The stackpointer value at the end of an assignment instruction at PC p_A^j is saved in a ghost variable named $\boxed{\text{sp}, p_A^j}$. These ghost variables help with inference of invariants that relate a local variable's address with stack addresses (discussed in section 4.1). During push, the initial contents of the newly allocated stack region are chosen non-deterministically using θ — this admits the possibility of arbitrary clobbering of the unallocated stack region below the stackpointer due to asynchronous external interrupts, before it is allocated again.

(ENTRY $_A$) shows the initialization of state elements of procedure A . For region $r \in B$, the initialization of Σ_A^r and $\pi_{\Sigma_A}(M_A)$ is similar to (ENTRY $_C$). The address sets of all assembly-only regions $f \in F$ are initialized using A 's symbol table ($\text{addrSets}_F()$). The memory contents of a read-only global variable $r \in G_r \cup F_r$ are initialized using $\text{ROM}_A^r(i_A^r)$ (defined in table 1). The machine registers are initialized with arbitrary contents (θ) — the constraints on the esp register are checked later, and \mathcal{W} is generated if a constraint is violated. The x86 stack of an assembly procedure includes the stack frame Σ_A^{stk} of the currently executing procedure A , the parameters Σ_A^y of A , and the remaining space which includes caller-stack Σ_A^{cs} and, possibly, the locals Σ_A^{cl} defined in the call chain of A . Ghost variable $\boxed{\text{sp}, \text{entry}}$ holds the esp value at entry of A . $\boxed{\text{stk}_e}$ represents the largest address in $\Sigma_A^{Y \cup \{\text{stk}\}}$ so that at entry, $\Sigma_A^{\text{stk}} = [\boxed{\text{sp}, \text{entry}}, \boxed{\text{stk}_e}] \setminus \Sigma_A^y$. If there are no parameters, $\boxed{\text{stk}_e} = \text{esp} + 3_{132}$ represents the end of the region that holds the return address. Ghost variable

cse holds the largest address in $\Sigma_A^{\{\text{stk}, \text{cs}, \text{cl}\} \cup Y}$. At entry, due to the calling conventions, we assume (through $\text{stkIsWF}()$) that: (1) the parameters are laid out at addresses above the stackpointer as per calling conventions (obeyCC); (2) the value $\text{esp} + 4_{132}$ is 16-byte aligned; and (3) the caller stack is above A 's stack frame Σ_A^{stk} . A violation of these conditions trigger \mathscr{W} .

Upon return (rule (RET $_A$)), we require that the callee-save registers, caller stack, and the return address remain preserved – a violation of these conditions trigger \mathscr{U} . For simplicity, we only tackle scalar return values, and ignore aggregate return values that need to be passed in memory.

2.4 Observable Traces and Refinement Definition

Recall that a procedure execution yields an observable trace containing silent and non-silent events. The error code of a trace T , written $e(T)$, is either \emptyset (indicating either non-termination or error-free termination), or one of $\nu \in \{\mathscr{U}, \mathscr{W}\}$ (indicating termination with error code ν). The non-error part of a trace T , written $\tilde{e}(T)$, is T when $e(T) = \emptyset$, and T' such that $T = T' \cdot e(T)$ otherwise.

Definition 2.3. $P \downarrow_{\Omega} T$ denotes the condition that for an initial outside world Ω , the execution of a procedure P may produce an observable trace T (for some sequence of non-deterministic choices).

Definition 2.4. Traces T and T' are *stuttering equivalent*, written $T =_{st} T'$, iff they differ only by finite sequences of silent events \perp . A trace T is a *stuttering prefix* of trace T' , written $T \leq_{st} T'$, iff $(T' =_{st} T) \vee (\exists T'' : T' =_{st} (T \cdot T''))$.

Definition 2.5. $U_{\text{pre}}^{\Omega, T_A}(C)$ denotes the condition: $\exists T_C : (C \downarrow_{\Omega} T_C \cdot \mathscr{U}) \wedge (T_C \leq_{st} T_A)$.

Definition 2.6. $W_{\text{pre}}^{\Omega, T_A}(C)$ denotes the condition: $e(T_A) = \mathscr{W} \wedge (\exists T_C : (C \downarrow_{\Omega} T_C) \wedge (\tilde{e}(T_A) \leq_{st} T_C))$

Definition 2.7. $C \sqsupseteq A$, read A refines C (or C is refined by A), iff:

$$\forall \Omega : (A \downarrow_{\Omega} T_A) \Rightarrow (W_{\text{pre}}^{\Omega, T_A}(C) \vee U_{\text{pre}}^{\Omega, T_A}(C) \vee (\exists T_C : (C \downarrow_{\Omega} T_C) \wedge (T_A =_{st} T_C)))$$

The $W_{\text{pre}}^{\Omega, T_A}(C)$ and $U_{\text{pre}}^{\Omega, T_A}(C)$ conditions cater to the cases where A triggers \mathscr{W} and C triggers \mathscr{U} respectively; the constituent \leq_{st} conditions ensure that a procedure call in A has identical termination behaviour to a procedure-call in C before an error is triggered. If neither A triggers \mathscr{W} nor C triggers \mathscr{U} , $T_A =_{st} T_C$ ensures that A and C produce identical non-silent events at similar speeds. In the absence of local variables and procedure-calls in C , $C \sqsupseteq A$ implies a correct translation from C to A .

2.4.1 Refinement Definition in the Presence of Local Variables and Procedure-Calls When All Local Variables Are Allocated on the Stack in A . For each local variable (de)allocation and for each procedure-call, our execution semantics generate a wr trace event in C (fig. 3). Thus, to reason about refinement, we require correlated and equivalent trace events to be generated in A . For this, we annotate A with two types of annotations to obtain \hat{A} :

- (1) alloc_s and dealloc_s instructions are added to explicitly indicate the (de)allocation of a local variable $z \in Z$, e.g., a stack region may be marked as belonging to z through these annotations.
- (2) A procedure-call, direct or indirect, is annotated with the types and addresses of the arguments and the set of memory regions observable by the callee.

These annotations are intended to encode the correlations with the corresponding allocation, deallocation, and procedure-call events in the source procedure C . For now, we assume that the locations and values of these annotations in \hat{A} are coming from an oracle – later in section 4, we present an algorithm to identify these annotations automatically in a best-effort manner.

Figure 5 presents three new instructions in \hat{A} – alloc_s , dealloc_s , and call – and their translations to graph instructions.

$$\begin{array}{c}
\text{(ALLOC}_S\text{)} \frac{p_A^j : \text{alloc}_s e_v, e_w, a, z}{\text{wr}(\text{allocBegin}(z, e_w)); \quad v, w := e_v, e_w; \\ \text{if } (\neg \text{intrvlInSet}_a(v, v + w - 1_{32}, \Sigma_A^{stk})) \text{halt}(\mathcal{U}); \\ \Sigma_A^{stk}, \Sigma_A^z := \Sigma_A^{stk} \setminus [v]_w, \Sigma_A^z \cup [v]_w; \\ \text{wr}(\text{allocEnd}(z, [v]_w, \pi_{[v]_w}(M_A))}; \\
\text{(DEALLOC}_S\text{)} \frac{p_A^j : \text{dealloc}_s z}{\Sigma_A^z, \Sigma_A^{stk} := \emptyset, \Sigma_A^{stk} \cup \Sigma_A^z; \\ \text{wr}(\text{dealloc}(z));}
\end{array}
\quad
\begin{array}{c}
\text{(CALL}_{\hat{A}}\text{)} \frac{p_{\hat{A}}^j : \text{call } \gamma \rho(\vec{r} \vec{x}) \beta^*}{\text{if } (\neg \text{aligned}_{16}(\text{esp}) \vee \neg \text{obeyCC}(\text{esp}, \vec{r}, \vec{x})) \\ \text{halt}(\mathcal{U}); \\ \text{wr}(\text{fcall}(\rho, \vec{x}, \beta^*, \pi_{\Sigma_{\hat{A}}^{\beta^*}}(M_{\hat{A}}))); \\ M_{\hat{A}} := \text{upd}_{\Sigma_{\hat{A}}^{\beta^*} \setminus G_r}(M_{\hat{A}}, \text{rd}(i_{32} \rightarrow i_8)); \\ \text{ecx} := \theta(i_{32}); \\ \text{IF } \{\gamma = \text{void}\} \{ \text{eax, edx} := \theta(i_{32}, i_{32}); \} \\ \text{ELSE } \{ \text{eax, edx} := \nabla_{\gamma}(\text{rd}(\gamma)); \}
\end{array}$$

Fig. 5. Additional translation rules for converting pseudo-assembly instructions to graph instructions for procedures with only stack-allocated locals.

An instruction $p_A^j : \text{alloc}_s e_v, e_w, a, z$ represents the stack allocation of a local variable identified by allocation site z . e_v is the expression for start address, e_w is the expression for allocation size, and a is the required alignment of the start address. During stack allocation of a local variable (ALLOC_S), the allocated address must satisfy the required alignment and separation constraints, or else \mathcal{U} is triggered. An allocation removes an address interval from Σ_A^{stk} and adds it to Σ_A^z .

A $p_A^j : \text{dealloc}_s z$ instruction represents the deallocation of z and empties the address set Σ_A^z , adding the removed addresses to Σ_A^{stk} (DEALLOC_S).

For procedure-calls (CALL _{\hat{A}}), we annotate the call instruction in assembly to explicitly specify the start addresses of the address regions belonging to the arguments (shown as \vec{x} in fig. 5). The address region of an argument should have previously been demarcated using an alloc_s instruction. Additionally, these address regions should satisfy the constraints imposed by the calling conventions (obeyCC). The calling conventions also require the esp value to be 16-byte aligned. A procedure-call is recorded as an observable event, along with the observation of the callee name (or address), the addresses of the arguments, callee-observable regions and their memory contents. The returned values, modeled through rd($i_{32} \rightarrow i_8$) and rd(γ), include the contents of the callee-observable memory regions and the scalar values returned by the callee (in registers eax, edx). The callee additionally clobbers the caller-save registers using θ .

Definition 2.8 (Refinement in the presence of only stack-allocated locals). $C \triangleright A$ iff: $\exists \hat{A} : C \sqsupseteq \hat{A}$

$C \triangleright A$ encodes the property that it is possible to annotate A to obtain \hat{A} so that the local variable (de)allocation and procedure-call events of C and the annotated \hat{A} can be correlated in lockstep. In the presence of stack-allocated local variables and procedure-calls, $C \triangleright A$ implies a correct translation from C to A . In the absence of local variables and procedure calls, $C \triangleright A$ reduces to $C \sqsupseteq A$ with $\hat{A} = A$.

2.4.2 Capabilities and Limitations of $C \triangleright A$. $C \triangleright A$ requires that for allocations and procedure calls that reuse the same stack space, their relative order remains preserved. This requirement is sound but may be too strict for certain (arguably rare) compiler transformations that may reorder the (de)allocation instructions that reuse the same stack space. Our refinement definition admits intermittent register-allocation of (parts of) a local variable.

$C \triangleright A$ supports *merging* of multiple allocations into a single stackpointer decrement instruction. Let p_A^s be the PC of a single stackpointer decrement instruction that implements multiple allocations. Merging can be encoded by adding multiple alloc_s instructions to A , in the same order as they appear in C , to obtain \hat{A} , so that these alloc_s instructions execute only after p_A^s executes; similarly, the corresponding dealloc_s instructions must execute before a stackpointer increment instruction deallocates this stack space.

CompCert's preallocation is a special case of merging where stack space for all local variables is allocated in the assembly procedure's prologue and deallocated in the epilogue (with no reuse of

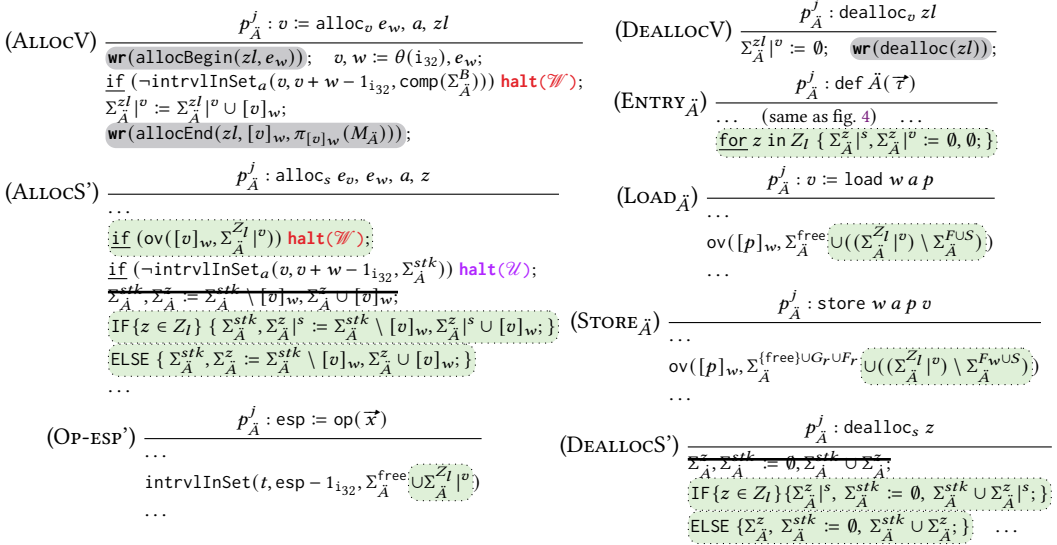


Fig. 6. Additional and revised translation rules for converting pseudo-assembly instructions to graph instructions for procedures with both stack and register allocated (or eliminated) locals.

stack space). In this case, our approach annotates A with (de)alloc_s instructions, potentially in the middle of the procedure body, such that they execute in lockstep with the (de)allocations in C .

A compiler may *reallocate* stack space by reusing the same space for two or more local variables with non-overlapping lifetimes (potentially without an intervening stackpointer increment instruction). If the relative order of (de)allocations is preserved, reallocation can be encoded by annotating \bar{A} with a dealloc_s instruction (for deallocating the first variable) immediately followed by an alloc_s instruction, such that the allocated region potentially overlaps with the previously deallocated region. Our refinement definition may not be able to cater to a translation that changes the relative order of (de)allocation instructions during reallocation.

Because our execution model observes each (de)allocation event (due to the wr instruction), a successful refinement check ensures that the allocation states of \bar{A} and C are identical at every correlated callsite. An inductive argument over \mathbb{C} and \mathbb{A} is thus used to show that the address set for region identifier cl is identical at the beginning of each correlated pair of procedures C and A (as modeled through identical reads from the outside world in (ENTRY _{P}) ($P \in \{C, A\}$) of figs. 3 and 4).

2.4.3 Refinement Definition in the Presence of Potentially Register-Allocated or Eliminated Local Variables in A . If a local variable $zl \in Z_l$ is either register-allocated or eliminated in A , there exists no stack region in A that can be associated with zl . However, recall that our execution model observes each allocation event in C through the wr instruction. Thus, for a successful refinement check, a correlated allocation event still needs to be annotated in A . We pretend that a correlated allocation occurs in A by introducing the notion of a *virtual allocation* instruction, called alloc_v, in A . Figure 6 shows the virtual (de)allocation instructions, alloc_v and dealloc_v, and the revised translations of procedure-entry and alloc_s, dealloc_s, load, store, and esp-modifying instructions. Instead of reproducing the full translations, we only show the changes with appropriate context. The additions have a highlighted background and deletions are ~~struck out~~. We update and annotate A with the translations and instructions in figs. 5 and 6 to obtain \bar{A} .

A $p_A^j : v := \text{alloc}_v e_w, a, z'$ instruction non-deterministically chooses the start address (using $\theta(i_{32})$) of a local variable z of size e_w and alignment a , performs a virtual allocation, and returns the start address in v . The chosen start address is *assumed* to satisfy the desired WF constraints,

such as separation (non-overlap) and alignment; error \mathcal{W} is triggered otherwise. Notice that this is in contrast to alloc_s where error \mathcal{U} is triggered on WF violation to indicate that it is the compiler's responsibility to ensure the satisfaction of WF constraints. Unlike a stack allocation where the compiler chooses the allocated region (and the validator identifies it through an alloc_s annotation), a virtual allocation is only a validation construct (the compiler is not involved) that is used only to enforce a lockstep correlation of allocation events. By triggering \mathcal{W} on a failure during a virtual allocation, we effectively assume that allocation through alloc_v satisfies the required WF conditions.

For simplicity, we support virtual allocations only for a variable declaration $zl \in Z_l$. Thus, we expect a call to $\text{alloca}()$ at $za \in Z_a$ to always be stack-allocated in \check{A} . In \check{A} , we replace the single variable $\Sigma_{\check{A}}^{zl}$ with two variables $\Sigma_{\check{A}}^{zl|s}$ and $\Sigma_{\check{A}}^{zl|v}$ that represent the address sets corresponding to the stack and virtual-allocations due to allocation-site zl respectively. We compute $\Sigma_{\check{A}}^{zl} = \Sigma_{\check{A}}^{zl|s} \cup \Sigma_{\check{A}}^{zl|v}$ (but we do not maintain a separate variable $\Sigma_{\check{A}}^{zl}$). We also assume that a single variable declaration zl in C may either correlate with only stack-allocations (through alloc_s) or only virtual-allocations (through alloc_v) in \check{A}^2 , i.e., $\Sigma_{\check{A}}^{zl|s} \cap \Sigma_{\check{A}}^{zl|v} = \emptyset$ holds at all times. For convenience, we define $\Sigma_{\check{A}}^{Z_l|v} = \bigcup_{zl \in Z_l} (\Sigma_{\check{A}}^{zl|v})$.

Importantly, a virtual allocation must be separate from other C allocated regions (B) but may overlap with assembly-only regions ($F \cup S$). Thus, in the revised semantics of $(\text{OP-ESP}')$, a stack push is allowed to overstep a virtually-allocated region.

The revised semantics of the alloc_s instruction (ALLOCS') assume that stack-allocated local memory is separate from virtually-allocated regions. The revised semantics of memory access instructions ($(\text{LOAD}_{\check{A}})$ and $(\text{STORE}_{\check{A}})$) enforce that a virtually-allocated region must never be accessed in \check{A} , unless it also happens to belong to the assembly-only regions ($F \cup S$).

Effectively, a lockstep correlation of virtual allocations in \check{A} with allocations in C ensures that the allocation states of both procedures always agree for regions $r \in B$.

Definition 2.9 (Refinement with stack and virtually-allocated locals). $C \ni A$ iff: $\exists \check{A} : C \sqsupseteq \check{A}$

Recall that $C \sqsupseteq \check{A}$ requires that *for all* non-deterministic choices of a virtually allocated local variable address in \check{A} (v in (ALLOCV)), there *exists* a non-deterministic choice for the correlated local variable address in C (v in (ALLOC) in fig. 3) such that: if \check{A} 's execution is well-formed (does not trigger \mathcal{W}), and C 's execution is UB-free (does not trigger \mathcal{U}), then the two allocated intervals are identical (the observable values created through allocBegin and allocEnd must be equal).

In the presence of potentially register-allocated and eliminated local variables, $C \ni A$ implies a correct translation from C to A . If all local variables are allocated in stack, $C \ni A$ reduces to $C \triangleright A$ with $\check{A} = \check{A}$. Figure 1c is an example of an annotated \check{A} .

3 WITNESSING REFINEMENT THROUGH A DETERMINIZED CROSS-PRODUCT $\check{A} \boxtimes C$

We first introduce program paths and their properties. Let $P \in \{C, \check{A}\}$. Let $e_p = (n_p \rightarrow n_p^t) \in \mathcal{E}_p$ represent an edge from node n_p to node n_p^t , both drawn from \mathcal{N}_p . A *path* ξ_p from n_p to n_p^t , written $\xi_p = n_p \rightsquigarrow n_p^t$, is a sequence of $m \geq 0$ edges $(e_p^1, e_p^2, \dots, e_p^m)$ with $\forall_{1 \leq j \leq m} : e_p^j = (n_p^{f,j} \rightarrow n_p^{t,j}) \in \mathcal{E}_p$, such that $n_p^{f,1} = n_p$, $n_p^{t,m} = n_p^t$, and $\bigwedge_{j=1}^{m-1} (n_p^{t,j} = n_p^{f,j+1})$. Nodes n_p and n_p^t are called the *source* and *sink* nodes of ξ_p respectively. Edge e_p^j (for some $1 \leq j \leq m$) is said to be present in ξ_p , written $e_p^j \in \xi_p$. An empty sequence, written ϵ , represents the *empty path*. The *path condition* of a path

²For simplicity, we do not tackle path-specializing transformations that may require, for a single variable declaration zl , a stack-allocation on one assembly path and a virtual-allocation on another. Such transformations are arguably rare.

$\xi_P = n_P \rightarrow n_P^t$, written $\text{pathcond}(\xi_P)$, is a conjunction of the edge conditions of the constituent edges. Starting at n_P , $\text{pathcond}(\xi_P)$ represents the condition that ξ_P executes to completion.

A sequence of edges corresponding to a shaded statement in the translations (figs. 3 to 6) is distinguished and identified as an *I/O path*. An I/O path must contain either a single rd or a single wr instruction. For example, the sequence of edges corresponding to “ $\text{wr}(\text{fcall}(\rho, \vec{x}, \beta^*, \pi_{(\Sigma_C^{\beta^*})}(M_C)))$ ” and “ $M_C := \text{upd}_{\Sigma_C^{\beta^*} \setminus G_r}(M_C, \text{rd}(i_{32} \rightarrow i_8))$ ” in (CALL_C) (fig. 3) refer to two separate I/O paths. A path without any rd or wr instructions is called an *I/O-free path*.

3.1 Determinized Product Graph as a Transition Graph

A product program, represented as a *determinized product graph*, also called a comparison graph or a cross-product, $X = \dot{A} \boxtimes C = (\mathcal{N}_X, \mathcal{E}_X, \mathcal{D}_X)$, is a directed multigraph with finite sets of nodes \mathcal{N}_X and edges \mathcal{E}_X , and a *deterministic choice map* \mathcal{D}_X . X is used to encode a lockstep execution of \dot{A} and C , such that $\mathcal{N}_X \subseteq \mathcal{N}_{\dot{A}} \times \mathcal{N}_C$. The start node of X is $n_X^s = (n_{\dot{A}}^s, n_C^s)$ and all nodes in \mathcal{N}_X must be reachable from n_X^s . A node $n_X = (n_{\dot{A}}, n_C)$ is an error node iff either $n_{\dot{A}}$ or n_C is an error node. $\mathcal{N}_X^{\text{DW}}$ denotes the set of non-error nodes in X , such that $n_X = (n_{\dot{A}}, n_C) \in \mathcal{N}_X^{\text{DW}} \Leftrightarrow (n_{\dot{A}} \in \mathcal{N}_{\dot{A}}^{\text{DW}} \wedge n_C \in \mathcal{N}_C^{\text{DW}})$.

Let $n_X = (n_{\dot{A}}, n_C)$ and $n_X^t = (n_{\dot{A}}^t, n_C^t)$ be nodes in \mathcal{N}_X ; let $\xi_{\dot{A}} = n_{\dot{A}} \rightarrow n_{\dot{A}}^t$ be a finite path in \dot{A} ; and let $\xi_C = n_C \rightarrow n_C^t$ be a finite path in C . Each edge, $e_X = (n_X \xrightarrow{\xi_{\dot{A}}; \xi_C} n_X^t) \in \mathcal{E}_X$, is defined as a sequential execution of $\xi_{\dot{A}}$ followed by ξ_C . The execution of e_X thus transfers control of X from n_X to n_X^t . The machine state of X is the concatenation of the machine states of \dot{A} and C . The outside world of X , written Ω_X , is a pair of the outside worlds of \dot{A} and C , i.e., $\Omega_X = (\Omega_{\dot{A}}, \Omega_C)$. Similarly, the trace generated by X , written T_X , is a pair of the traces generated by \dot{A} and C , i.e., $T_X = (T_{\dot{A}}, T_C)$.

During an execution of $e_X = (n_X \xrightarrow{\xi_{\dot{A}}; \xi_C} n_X^t) \in \mathcal{E}_X$, let $\vec{x}_{\dot{A}}$ be variables in \dot{A} just at the end of the execution of path $\xi_{\dot{A}}$ (at $n_{\dot{A}}^t$) but before the execution of path ξ_C (recall, $\xi_{\dot{A}}$ executes before ξ_C). $\mathcal{D}_X : ((\mathcal{E}_X \times \mathcal{E}_C \times \mathbb{N}) \rightarrow \text{ExprList})$, called a *deterministic choice map*, is a partial function that maps edge $e_X \in \mathcal{E}_X$, and the n^{th} (for $n \in \mathbb{N}$) occurrence of an edge ‘ $e_C^{\theta} \in \xi_C$ ’ labeled with instruction $\vec{\tau} := \theta(\vec{\tau})$ to a list of expressions $E(\vec{x}_{\dot{A}})$. The semantics of \mathcal{D}_X are such that, if $\mathcal{D}_X(e_X, e_C^{\theta}, n)$ is defined, then during an execution of e_X , an execution of the n^{th} occurrence of edge $e_C^{\theta} \in \xi_C$ labeled with $\vec{\tau} := \theta(\vec{\tau})$ is semantically equivalent to an execution of $\vec{\tau} := \mathcal{D}_X(e_X, e_C^{\theta}, n)$; otherwise, the original non-deterministic semantics of θ are used.

\mathcal{D}_X determinizes (or refines) the non-deterministic choices in C . For example, in a product graph X that correlates the programs in fig. 1b and fig. 1c, let $e_X^2 \in \mathcal{E}_X$ correlate single instructions I2 and A4.2. Let e_C^{12, θ_a} represent the edge labeled with $\alpha_b := \theta(i_{32})$ as a part of the translation of the alloc instruction at I2, as seen in (ALLOC). Then, $\mathcal{D}_X(e_X^2, e_C^{12, \theta_a}, 1) = \text{esp}$ is identified by the first operand of the annotated alloc_s instruction at A4.2. Similarly, if another edge e_C^{12, θ_m} (in the translation of alloc at I2) is labeled with $\theta(i_{32} \rightarrow i_8)$ (due to $M_C := \text{upd}_{[\alpha_b, \alpha_e]}(M_C, \theta(i_{32} \rightarrow i_8))$), then $\mathcal{D}_X(e_X^2, e_C^{12, \theta_m}, 1) = M_{\dot{A}}$, i.e., the initial contents of the newly-allocated region in C are based on the contents of the correlated uninitialized stack region in \dot{A} . Similarly, let $e_X^1 \in \mathcal{E}_X$ correlate single instructions I1 and A4.1 so that $\mathcal{D}_X(e_X^1, e_C^{11, \theta_a}, 1) = v_{I1}$ and $\mathcal{D}_X(e_X^1, e_C^{11, \theta_m}, 1) = M_{\dot{A}}$.

For a path ξ_C in C , $[\xi_C]_{\mathcal{D}_X}^{\text{EX}}$ denotes a *determinized path* that is identical to ξ_C except that: if $\mathcal{D}_X(e_X, e_C^{\theta}, n)$ is defined, then the n^{th} occurrence of edge $e_C^{\theta} \in \xi_C$, labeled with $\vec{\tau} := \theta(\vec{\tau})$, is replaced with a new edge $e_C^{\theta_n}$ labeled with $\vec{\tau} := \mathcal{D}_X(e_X, e_C^{\theta}, n)$.

Execution of a product graph X must begin at node n_X^s in an initial machine state where $\Omega_{\dot{A}} = \Omega_C$ and $T_{\dot{A}} =_{st} T_C$ hold. Thus, X is a transition graph with its execution semantics derived from the semantics of \dot{A} and C , and the map \mathcal{D}_X .

3.2 Analysis of the Determinized Product Graph

Let $X = \bar{A} \boxtimes C = (\mathcal{N}_X, \mathcal{E}_X, \mathcal{D}_X)$ be a determinized product graph. At each non-error node $n_X \in \mathcal{N}_X^{\text{DW}}$, we infer a node invariant, ϕ_{n_X} , which is a first-order logic predicate over state elements of X at node n_X that holds for all possible executions of X . A node invariant ϕ_{n_X} relates the values of state elements of C and \bar{A} that can be observed at n_X .

Definition 3.1 (Hoare Triple). Let $n_X = (n_{\bar{A}}, n_C) \in \mathcal{N}_X^{\text{DW}}$. Let $\xi_{\bar{A}} = n_{\bar{A}} \rightarrow n_{\bar{A}}^t$ and $\xi_C = n_C \rightarrow n_C^t$ be paths in \bar{A} and C . A *Hoare triple*, written $\{pre\}(\xi_{\bar{A}}; \xi_C)\{post\}$, denotes the statement: if execution starts at node n_X in state σ such that predicate $pre(\sigma)$ holds, and if paths $\xi_{\bar{A}}; \xi_C$ are executed in sequence to completion finishing in state σ' , then predicate $post(\sigma')$ holds.

Definition 3.2 (Path cover). At a node $n_X = (n_{\bar{A}}, n_C) \in \mathcal{N}_X$, for a path $\xi_{\bar{A}} = n_{\bar{A}} \rightarrow n_{\bar{A}}^t$, let $\forall_{1 \leq j \leq m} : e_X^j = n_X \xrightarrow{\xi_{\bar{A}}; \xi_C^j} n_X^{t_j}$ be all edges in \mathcal{E}_X , such that $n_X^{t_j} = (n_{\bar{A}}^t, n_C^{t_j})$. The set of edges $\{e_X^1, e_X^2, \dots, e_X^m\}$ covers path $\xi_{\bar{A}}$, written $\{e_X^1, e_X^2, \dots, e_X^m\} \langle \mathcal{D}_X, \xi_{\bar{A}} \rangle$, iff $\{\phi_{n_X}\}(\xi_{\bar{A}}; \epsilon) \{ \bigvee_{j=1}^m \text{pathcond}([\xi_C^j]_{\mathcal{D}_X}^{e_X^j}) \}$ holds.

Definition 3.3 (Path infeasibility). At a node $n_X = (n_{\bar{A}}, n_C) \in \mathcal{N}_X$, a path $\xi_{\bar{A}} = n_{\bar{A}} \rightarrow n_{\bar{A}}^t$ is *infeasible* at n_X iff $\{\phi_{n_X}\}(\xi_{\bar{A}}; \epsilon) \{\text{false}\}$ holds.

Definition 3.4 (Mutually exclusive paths). Two paths, $\xi_P^1 = n_P \rightarrow n_P^{t_1}$ and $\xi_P^2 = n_P \rightarrow n_P^{t_2}$, both originating at node n_P are *mutually-exclusive*, written $\xi_P^1 \approx \xi_P^2$, iff neither is a prefix of the other.

Definition 3.5. A *pathset* $\langle \xi \rangle_P$ is a set of pairwise mutually-exclusive paths $\langle \xi \rangle_P = \{\xi_P^1, \xi_P^2, \dots, \xi_P^m\}$ originating at the same node n_P , i.e., $\forall_{1 \leq j \leq m} : \xi_P^j = n_P \rightarrow n_P^j$ and $\forall_{1 \leq j_1 < j_2 \leq m} : (\xi_P^{j_1} \approx \xi_P^{j_2})$.

3.2.1 X Requirements. The following requirements on X help witness $C \sqsupseteq \bar{A}$:

1. (Mutex \bar{A}): For each node n_X with *all* outgoing edges $\{e_X^1, e_X^2, \dots, e_X^m\}$ such that $e_X^j = (n_X \xrightarrow{\xi_{\bar{A}}; \xi_C^j} n_X^j)$ (for $1 \leq j \leq m$), the following holds: $\forall_{1 \leq j_1, j_2 \leq m} : ((\xi_{\bar{A}}^{j_1} = \xi_{\bar{A}}^{j_2}) \vee (\xi_{\bar{A}}^{j_1} \approx \xi_{\bar{A}}^{j_2}))$.
2. (MutexC): At each node n_X , for a path $\xi_{\bar{A}}$, let $\{e_X^1, e_X^2, \dots, e_X^m\}$ be a set of *all* outgoing edges such that $e_X^j = n_X \xrightarrow{\xi_{\bar{A}}; \xi_C^j} n_X^t$ (for $1 \leq j \leq m$). Then, the set $\{\xi_C^1, \xi_C^2, \dots, \xi_C^m\}$ must be a pathset.
3. (Termination) For each non-error node $n_X = (n_{\bar{A}}, n_C) \in \mathcal{N}_X^{\text{DW}}$, $n_{\bar{A}}$ is a terminating node iff n_C is a terminating node.
4. (SingleIO): For each edge $e_X = (n_X \xrightarrow{\xi_{\bar{A}}; \xi_C} n_X^t) \in \mathcal{E}_X$, either both $\xi_{\bar{A}}$ and ξ_C are I/O paths, or both $\xi_{\bar{A}}$ and ξ_C are I/O-free.
5. (Well-formedness): If a node of the form $n_X = (_, \mathcal{W}_C)$ exists in \mathcal{N}_X , then n_X must be $(\mathcal{W}_{\bar{A}}, \mathcal{W}_C)$.
6. (Safety): If a node of the form $n_X = (\mathcal{U}_{\bar{A}}, _)$ exists in \mathcal{N}_X , then n_X must be $(\mathcal{U}_{\bar{A}}, \mathcal{U}_C)$.
7. (Similar-speed): Let $(e_X^1, e_X^2, \dots, e_X^m)$ be a cyclic path, so that $\forall_{1 \leq j \leq m} : e_X^j = (n_X^{f,j} \xrightarrow{\xi_{\bar{A}}; \xi_C^j} n_X^{t,j}) \in \mathcal{E}_X$; $n_X^{f,1} = n_X^{t,m}$; and $\bigwedge_{j=1}^{m-1} (n_X^{t,j} = n_X^{f,j+1})$. For each cyclic path, $(\neg \bigwedge_{j=1}^m (\xi_{\bar{A}}^j = \epsilon)) \wedge (\neg \bigwedge_{j=1}^m (\xi_C^j = \epsilon))$ holds.
8. (Coverage \bar{A}): For each non-error node $n_X = (n_{\bar{A}}, n_C) \in \mathcal{N}_X^{\text{DW}}$ and for each possible outgoing path $\xi_{\bar{A}}^o = n_{\bar{A}} \rightarrow n_{\bar{A}}^o$, either $\xi_{\bar{A}}^o$ is infeasible at n_X , or there exists $e_X = (n_X \xrightarrow{\xi_{\bar{A}}; \xi_C} n_X^t) \in \mathcal{E}_X$ such that either $\xi_{\bar{A}}^o$ is a prefix of $\xi_{\bar{A}}^o$ or $\xi_{\bar{A}}^o$ is a prefix of $\xi_{\bar{A}}$.
9. (CoverageC): At node n_X , for some $\xi_{\bar{A}}$, let $\{e_X^1, e_X^2, \dots, e_X^m\}$ be the set of *all* outgoing edges such that $e_X^j = n_X \xrightarrow{\xi_{\bar{A}}; \xi_C^j} (n_{\bar{A}}^t, n_C^{t_j})$ (for $1 \leq j \leq m$). Then, $\{e_X^1, e_X^2, \dots, e_X^m\} \langle \mathcal{D}_X, \xi_{\bar{A}} \rangle$ holds.
10. (Inductive): For each non-error edge $e_X = (n_X \xrightarrow{\xi_{\bar{A}}; \xi_C} n_X^t) \in \mathcal{E}_X$, $\{\phi_{n_X}\}(\xi_{\bar{A}}; [\xi_C]_{\mathcal{D}_X}^{e_X}) \{\phi_{n_X^t}\}$ holds.

11. (Equivalence): For each non-error node $n_X = (n_{\bar{A}}, n_C) \in \mathcal{N}_X^{\text{DW}}$, $\Omega_{\bar{A}} = \Omega_C$ must belong to ϕ_{n_X} .
12. (Memory Access Correspondence) or (MAC): For each edge $e_X = (n_X \xrightarrow{\xi_{\bar{A}}; \xi_C} n_X^t) \in \mathcal{E}_X$, such that $n_X^t \neq (_, \mathcal{U}_C)$, $\{\phi_{n_X} \wedge (\Sigma_{\bar{A}}^{\text{rd}} = \Sigma_C^{\text{rd}} = \emptyset)\}(\xi_{\bar{A}}; [\xi_C]_{\mathcal{D}_X}^{\text{ex}}) \{(\Sigma_{\bar{A}}^{\text{rd}} \setminus \Sigma_C^{\text{rd}}) \subseteq \Sigma_{\bar{A}}^{\text{GUF}} \cup [\text{esp}, \text{stk}_e]\}$ and $\{\phi_{n_X} \wedge (\Sigma_{\bar{A}}^{\text{wr}} = \Sigma_C^{\text{wr}} = \emptyset)\}(\xi_{\bar{A}}; [\xi_C]_{\mathcal{D}_X}^{\text{ex}}) \{(\Sigma_{\bar{A}}^{\text{wr}} \setminus \Sigma_C^{\text{wr}}) \subseteq \Sigma_{\bar{A}}^{\text{GwUFw}} \cup [\text{esp}, \text{stk}_e]\}$ hold.
13. (MemEq): For each non-error node $n_X \in \mathcal{N}_X^{\text{DW}}$, $M_{\bar{A}} =_{\Sigma_{\bar{A}}^B \setminus (\Sigma_{\bar{A}}^{Z_I})^v} M_C$ must belong to ϕ_{n_X} .

(MAC) effectively requires that for every access on path $\xi_{\bar{A}}$ to an address α belonging to region $r \in \{hp, cl\}$, there exists an access to α of the same read/write type on path $[\xi_C]_{\mathcal{D}_X}^{\text{ex}}$. This requirement allows us to soundly over-approximate the set of addresses belonging to hp and cl for a faster SMT encoding (theorem 3.7 and section 4.2.3). For (MAC) to be meaningful, $\Sigma_{\bar{A}C}^{\text{rd}}$ and $\Sigma_{\bar{A}C}^{\text{wr}}$ must not be included in X 's state elements over which a node invariant ϕ_{n_X} is inferred.

The first seven are *structural requirements* (constraints on the graph structure of X) and the remaining six are *semantic requirements* (require discharge of proof obligations). The first eleven are *soundness requirements* (required for theorem 3.6), the first twelve are *fast-encoding requirements*, and all thirteen are *search-algorithm requirements* (required for search optimizations). Excluding (Coverage \bar{A}) and (Coverage C), the remaining eleven are called *non-coverage requirements*.

THEOREM 3.6. *If there exists $X = \bar{A} \boxtimes C$ that satisfies the soundness requirements, then $C \sqsupseteq \bar{A}$ holds.*

PROOF SKETCH. (Coverage \bar{A}) and (Coverage C) ensure the coverage of \bar{A} 's and C 's traces in X . For an error-free execution of X , (Equivalence) and (Similar-speed) ensure that the generated traces are stuttering equivalent; for executions terminating in an error, (SingleIO), (Well-formedness), and (Safety) ensure that $C \sqsupseteq \bar{A}$ holds by definition. See our technical report [Rose and Bansal 2024b] for the coinductive proof. \square

3.2.2 Safety-Relaxed Semantics. Construct A' from A by using new *safety-relaxed semantics* for the assembly procedure such that: (1) a $\varphi_l = \text{ov}([p]_w, \Sigma_{\bar{A}}^{\text{free}} \cup ((\Sigma_{\bar{A}}^{Z_I})^v \setminus \Sigma_{\bar{A}}^{\text{FUS}}))$ check in (LOAD \bar{A}) in A is replaced with $\varphi'_l = \text{ov}([p]_w, (\Sigma_{\bar{A}}^{Z_I})^v \setminus (\Sigma_{\bar{A}}^{\text{F}} \cup [\text{esp}, \text{cs}_e]))$ in A' ; (2) a $\varphi_s = \text{ov}([p]_w, \Sigma_{\bar{A}}^{\{\text{free}\} \cup G_r \cup F_r} \cup ((\Sigma_{\bar{A}}^{Z_I})^v \setminus \Sigma_{\bar{A}}^{\text{FwUS}}))$ check in (STORE \bar{A}) in A is replaced with $\varphi'_s = \text{ov}([p]_w, (\Sigma_{\bar{A}}^{Z_I})^v \setminus (\Sigma_{\bar{A}}^{\text{Fw}} \cup [\text{esp}, \text{cs}_e]))$ in A' ; and (3) a $\varphi_r = \neg([\text{MCS}] =_{\Sigma_{\bar{A}}^{\text{cs}}} M_A)$ check in (RET \bar{A}) in A is replaced with $\varphi'_r = \text{false}$ in A' . Let \bar{A}' be obtained by annotating A' using instructions described in section 2.4.3. Let \bar{A} be the annotated version of A , such that the annotations made in \bar{A} and \bar{A}' are identical.

THEOREM 3.7. *Given $X' = \bar{A}' \boxtimes C$ that satisfies the fast-encoding requirements, it is possible to construct $X = \bar{A} \boxtimes C$ that also satisfies the fast-encoding requirements.*

PROOF SKETCH. Start by constructing $X = X'$. Because $\varphi'_{l,s,r} \Rightarrow \varphi_{l,s,r}$, \bar{A} may include more executions of a path of form $\xi_{\bar{A}} = n_{\bar{A}} \rightarrow \mathcal{U}_{\bar{A}}$. Add new edges to \mathcal{E}_X , where each new edge correlates $\xi_{\bar{A}}$ with some $\xi_C = n_C \rightarrow \mathcal{U}_C$. Because X' satisfies (MAC), the addition of such new edges will ensure that X satisfies (Coverage C). See our technical report [Rose and Bansal 2024b] for the proof. \square

Using theorem 3.7, hereafter, we will use only the safety-relaxed semantics of the assembly procedure. We will continue to refer to the assembly procedure with the safety-relaxed semantics as A , and the corresponding annotated procedure \bar{A} .

4 AUTOMATIC CONSTRUCTION OF A CROSS-PRODUCT

We now describe Dynamo, an algorithm that takes as input, the transition graphs corresponding to procedures C and A , and an *unroll factor* μ , and returns as output, annotated \bar{A} and product

graph $X = \ddot{A} \boxtimes C = (\mathcal{N}_X, \mathcal{E}_X, \mathcal{D}_X)$, such that all thirteen search-algorithm requirements are met. It identifies an inductive invariant network ϕ_X that maps each non-error node $n_X \in \mathcal{N}_X^{\text{DW}}$ to its node invariant ϕ_{n_X} . Given enough computational time, Dynamo is guaranteed to find the required (\ddot{A}, X) if: (a) A is a translation of C through bisimilar transformations up to a maximum unrolling of μ ; (b) for two or more allocations or procedure calls that reuse stack space in A , their relative order in C is preserved in A ; (c) the desired annotation to \ddot{A} is identifiable either through search heuristics or through compiler hints; and (d) our invariant inference procedure is able to identify the required invariant network ϕ_X that captures the compiler transformations across C and A . Dynamo constructs the solution incrementally, by relying on the property that for a non-coverage requirement to hold for fully-annotated \ddot{A} and fully-constructed X , it must also hold for partially-annotated \ddot{A} and a partially-constructed subgraph of X rooted at its entry node n_X^s .

Dynamo is presented in algorithm 1. It assumes the availability of a **chooseFrom** operator, such that $\rho \leftarrow \text{chooseFrom } \vec{\rho}$ chooses a quantity ρ from a finite set $\vec{\rho}$, such that Dynamo is able to complete the refinement proof, if such a choice exists. If the search space is limited, an exhaustive search could be used to implement **chooseFrom**. Otherwise, a counterexample-guided best-first search procedure (described later) is employed to approximate **chooseFrom**.

$\text{io}(n_P)$ evaluates to true iff n_P is either a source or sink node of an I/O path. $\text{term}(n_P)$ evaluates to true iff n_P is a terminating node. Dynamo first identifies an ordered set of nodes $Q_P \subseteq \mathcal{N}_P$, called the *cut points* in procedure P (*getCutPointsInRPO*), such that $Q_P \supseteq \{n_P : n_P \in \mathcal{N}_P \wedge (n_P = n_P^s \vee \text{io}(n_P) \vee \text{term}(n_P))\}$ and the maximum length of a path between two nodes in Q_P (not containing any other intermediate node that belongs to Q_P) is finite.

The algorithm to identify Q_P first initializes $Q_P := \{n_P : n_P \in \mathcal{N}_P \wedge (n_P = n_P^s \vee \text{io}(n_P) \vee \text{term}(n_P))\}$, and then identifies all cycles in the transition graph that do not already contain a cut point; for each such cycle, the first node belonging to that cycle in reverse postorder is added to Q_P . In fig. 1c, $Q_{\ddot{A}}$ includes constituent nodes of assembly instructions at A1, A9, A14, and `exit`, where `exit` is the destination node of the error-free `halt` instruction due to the procedure return at A17.

A *simple path* $q_P \rightarrow q_P^t$ is a path connecting two cut points $q_P, q_P^t \in Q_P$, and not containing any other cut point as an intermediate node; q_P^t is called a *cut-point successor* of q_P . By definition, a simple path must be finite. The *cutPointSuccessors()* function takes a cut point q_P and returns all its cut-point successors in reverse postorder. In our example, the cut-point successors of a node at instruction A9 are (constituent nodes of) A9, A14, $\mathcal{U}_{\ddot{A}}$, and $\mathcal{W}_{\ddot{A}}$. *getAllSimplePathsBetweenCutPoints*(q_P, q_P^t, P) returns *all* simple paths of the form $q_P \rightarrow q_P^t$, for $q_P, q_P^t \in Q_P$. Given a simple path $\xi_{\ddot{A}}$, *pathIsInfeasible*($\xi_{\ddot{A}}, q_{\ddot{A}}, \mathcal{N}_X, \phi_X$) returns true iff $\xi_{\ddot{A}}$ is infeasible at every node $n_X = (q_{\ddot{A}}, _) \in \mathcal{N}_X$; our algorithm ensures there can be at most one $n_X = (q_{\ddot{A}}, _) \in \mathcal{N}_X$ for each $q_{\ddot{A}} \in Q_{\ddot{A}}$.

correlatedPathsInCOptions(). *correlatedPathsInCOptions*($\xi_{\ddot{A}}, \dots$) identifies options for candidate pathsets $[\langle \xi \rangle_C]$, that can potentially be correlated with $\xi_{\ddot{A}} = q_{\ddot{A}} \rightarrow q_{\ddot{A}}^t$, and the **chooseFrom** operator chooses a pathset $\langle \xi \rangle_C$ from it. A path $\xi_C \in \langle \xi \rangle_C$ need not be a simple path, and can visit any node $n_C \in \mathcal{N}_C$ up to μ times. All paths in $\langle \xi \rangle_C$ must originate at a unique cut-point q_C such that $(q_{\ddot{A}}, q_C) \in \mathcal{N}_X$. By construction, there will be exactly one such $(q_{\ddot{A}}, q_C)$ in \mathcal{N}_X . Paths in $\langle \xi \rangle_C$ may have different end points however. For example, $\langle \xi \rangle_C = \{\epsilon\}$ and $\langle \xi \rangle_C = \{I3 \rightarrow I4 \rightarrow I7, I3 \rightarrow \mathcal{U}_C, I3 \rightarrow I4 \rightarrow \mathcal{U}_C\}$ may be potential candidates for $\xi_{\ddot{A}} = A9 \rightarrow A10 \rightarrow A11 \rightarrow A9$ in fig. 1.

If $q_{\ddot{A}}^t \notin \{\mathcal{U}_{\ddot{A}}, \mathcal{W}_{\ddot{A}}\}$, *correlatedPathsInCOptions()* returns candidates, where a candidate pathset $\langle \xi \rangle_C$ is a maximal set such that each path $\xi_C \in \langle \xi \rangle_C$ either (a) ends at a unique non-error destination cut-point node, say q_C^t (i.e., all paths $\xi_C \in \langle \xi \rangle_C$ ending at a non-error node end at q_C^t), or (b) ends at error node \mathcal{U}_C . This path enumeration strategy is the same as the one used in Counter [Gupta et al. 2020]; this strategy supports path specializing compiler transformations like loop peeling, unrolling, splitting, unswitching, etc., but does not support a path de-specializing transformation

Algorithm 1: Automatic construction of X .

```

1 Function Dynamo( $A, C, \mu$ )
2    $\bar{A} \leftarrow A$ ;  $N_X \leftarrow \{(n_{\bar{A}}^s, n_C^s)\}$ ;  $\mathcal{E}_X \leftarrow \{\}$ ;  $\mathcal{D}_X \leftarrow \emptyset$ ;
3    $\phi_X \leftarrow \{(n_{\bar{A}}^s, n_C^s) \mapsto (\Omega_{\bar{A}} = \Omega_C)\}$ ;  $Q_{\bar{A}} \leftarrow \text{getCutPointsInRPO}(\bar{A})$ ;
4   foreach  $q_{\bar{A}}$  in  $Q_{\bar{A}}$  do
5     foreach  $q_{\bar{A}}^t$  in  $\text{cutPointSuccessors}(q_{\bar{A}}, Q_{\bar{A}}, \bar{A})$  do
6       foreach  $\xi_{\bar{A}}$  in  $\text{getAllSimplePathsBetweenCutPoints}(q_{\bar{A}}, q_{\bar{A}}^t, \bar{A})$  do
7         if  $\text{pathIsInfeasible}(\xi_{\bar{A}}, q_{\bar{A}}, N_X, \phi_X)$  then
8           continue
9         end
10        foreach  $\xi_C$  in chooseFrom  $\text{correlatedPathsInCOptions}(\xi_{\bar{A}}, \mu, N_X, \mathcal{E}_X, \mathcal{D}_X, \phi_X, \bar{A}, C)$  do
11           $(\bar{A}, \vec{\xi}_{\bar{A}}) \leftarrow \text{chooseFrom } \text{asmAnnotOptions}(\xi_{\bar{A}}, \xi_C, N_X, \mathcal{E}_X, \mathcal{D}_X, \phi_X, \bar{A}, C)$ ;
12           $\vec{\xi}'_C \leftarrow \text{breakIntoSingleIOPaths}(\xi_C)$ ;
13          if  $\neg \text{haveSimilarStructure}(\vec{\xi}'_{\bar{A}}, \vec{\xi}'_C)$  then
14            return Failure
15          end
16          foreach  $(\xi'_{\bar{A}} = (n_{\bar{A}} \rightarrow n_{\bar{A}}^t), (\xi'_C = (n_C \rightarrow n_C^t)))$  in  $\text{zip}(\vec{\xi}'_{\bar{A}}, \vec{\xi}'_C)$  do
17             $N_X \leftarrow N_X \cup \{(n_{\bar{A}}^t, n_C^t)\}$ ; //unlike  $\mathcal{E}_X$ ,  $N_X$  may not grow on each iteration
18             $e_X \leftarrow (\xi'_{\bar{A}}; \xi'_C)$ ;  $\mathcal{E}_X \leftarrow \mathcal{E}_X \cup \{e_X\}$ ;  $\mathcal{D}_X \leftarrow \text{addDetMappings}(e_X)$ ;
19             $\phi_X \leftarrow \text{inferInvariantsAndCounterexamples}(N_X, \mathcal{E}_X, \mathcal{D}_X, \phi_X, \bar{A}, C)$ ;
20            if  $\neg \text{checkSemanticRequirementsExceptCoverage}(N_X, \mathcal{E}_X, \mathcal{D}_X, \phi_X, \bar{A}, C)$  then
21              return Failure
22            end
23          end
24        end
25      end
26    end
27  end
28  if  $\neg \text{checkCoverageRequirements}(N_X, \mathcal{E}_X, \mathcal{D}_X, \phi_X, \bar{A}, C)$  then
29    return Failure
30  end
31  return Success( $\bar{A}, N_X, \mathcal{E}_X, \mathcal{D}_X, \phi_X$ )
32 end

```

like loop re-rolling. If $q_{\bar{A}}^t = \mathcal{U}_{\bar{A}}$, $\text{correlatedPathsInCOptions}()$ returns candidates, where a candidate pathset $\langle \xi \rangle_C$ is a maximal set such that each path $\xi_C \in \langle \xi \rangle_C$ ends at \mathcal{U}_C . The algorithm identifies a correlation for a path $\xi_{\bar{A}} = q_{\bar{A}} \rightarrow \mathcal{W}_{\bar{A}}$ only after correlations for all other paths of the form $\xi_{\bar{A}}^{\mathcal{W}} = q_{\bar{A}} \rightarrow q_{\bar{A}}^{\mathcal{W}}$ (for $q_{\bar{A}}^{\mathcal{W}} \neq \mathcal{W}_{\bar{A}}$) have been identified: a pathset candidate $\langle \xi \rangle_C$ that has already been correlated with some other path $\xi_{\bar{A}}^{\mathcal{W}}$ is then prioritized for correlation with $\xi_{\bar{A}}$.

For example, in fig. 1c, for a cyclic path $\xi_{\bar{A}}$ from a node at A9 to itself, one of the candidate pathsets, $\langle \xi \rangle_C$, returned by this procedure (at $\mu = 1$) contains eleven paths originating at I4 in fig. 1b: one that cycles back to I4 and ten that terminate at \mathcal{U}_C (for each of the ten memory accesses in the path). For example, to evaluate the expression $v[*i]$, two memory loads are required, one at address i and another at $\&v[*i]$, and each such load may potentially transition to \mathcal{U}_C due to the $\text{accessIsSafeC}_{\tau,a}$ check evaluating to false in (LOAD_C) . A path that terminates at \mathcal{U}_C represents correlated transitions from node $(A9, I4)$ in X such that \bar{A} remains error-free (to end at A9) but C triggers \mathcal{U} , e.g., if the memory access $\text{mem}_4[\text{esi}+4*\text{eax}]$ in \bar{A} (corresponding to $v[*i]$ in C) overshoots the stack space corresponding to variable v but still lies within the stack region stk .

asmAnnotOptions(). For each simple path $\xi_{\bar{A}}$, and each (potentially non-simple) path ξ_C in $\langle \xi \rangle_C$ ³, *asmAnnotOptions()* enumerates the options for annotating $\xi_{\bar{A}}$ with $\text{alloc}_{s,v}$, $\text{dealloc}_{s,v}$ instructions and operands for call instructions, and the **chooseFrom** operator chooses one.

An annotation option includes the positions and the operands of the (de)allocation instructions (allocation site, alignment, address, and size). For a procedure-call, an annotation option also includes the arguments' types and values, and the set of callee-observable regions. The annotations for the callee name/address and the (de)allocations of procedure-call arguments in $\xi_{\bar{A}}$ are uniquely identified using the number and type of arguments in the candidate correlated path ξ_C using the calling conventions. Similarly, the annotation of callee-observable regions follows from the regions observable by the correlated procedure call in ξ_C .

These annotations thus update A to incrementally construct \bar{A} . If untrusted compiler hints are available, they are used to precisely identify these annotations. In a blackbox setting, where no compiler hints are available, we reduce the search space for annotations (at the cost of reduced generality) using the following three restrictions: (1) An $\text{alloc}_{s,v}$ ($\text{dealloc}_{s,v}$) annotation is annotated in $\xi_{\bar{A}}$ only if an alloc (dealloc) instruction is present in ξ_C ; (2) an $\text{alloc}_{s,v}$ ($\text{dealloc}_{s,v}$) annotation is added only after (before) an instruction that updates esp ; moreover, for alloc_s , esp is used as the local variable's address expression; (3) for a single allocation site in C , at most one $\text{alloc}_{s,v}$ instruction (but potentially multiple $\text{dealloc}_{s,v}$ instructions) is added to \bar{A} . Thus, in a blackbox setting, due to the third restriction, a refinement proof may fail if the compiler specializes a path containing a local variable allocation. Due to the second restriction, a refinement proof may fail for certain (arguably rare) types of order-preserving stack reallocation and stack merging performed by the compiler. Note that these limitations hold only for the blackbox setting.

After annotations, $\xi_{\bar{A}}$ may become a non-simple path due to the extra I/O instructions introduced by the annotations. *asmAnnotOptions* therefore additionally returns $\vec{\xi}'_{\bar{A}}$, which is a sequence of the simple paths constituting $\xi_{\bar{A}}$. The (potentially non-simple) path ξ_C is then broken into a sequence of constituent paths $\vec{\xi}'_C$ (*breakIntoSingleIOPaths*) so that each I/O path appears by itself (and not as a sub-path of a longer constituent path) in $\vec{\xi}'_C$ – this caters to the (SingleIO) requirement. A failure is returned if the sequences $\vec{\xi}'_{\bar{A}}$ and $\vec{\xi}'_C$ do not have similar structures (*haveSimilarStructure*). Let $\text{pos}(\xi, \vec{\xi})$ represent the position of path ξ in a sequence of paths $\vec{\xi}$. *haveSimilarStructure*($\vec{\xi}'_{\bar{A}}, \vec{\xi}'_C$) returns true iff $\vec{\xi}'_{\bar{A}}$ and $\vec{\xi}'_C$ are of the same size, and for paths $\xi'_C \in \vec{\xi}'_C$ and $\xi'_{\bar{A}} \in \vec{\xi}'_{\bar{A}}$, if $\text{pos}(\xi'_C, \vec{\xi}'_C) = \text{pos}(\xi'_{\bar{A}}, \vec{\xi}'_{\bar{A}})$, then either both ξ'_C and $\xi'_{\bar{A}}$ are I/O paths of same structure (i.e., they are either both reads or both writes for the same type of value) or both are I/O free.

Incremental Construction of (\bar{A}, X) . For each simple path $\xi'_{\bar{A}}$ in $\vec{\xi}'_{\bar{A}}$ enumerated in execution order, Dynamo correlates it with ξ'_C , such that $\text{pos}(\xi'_C, \vec{\xi}'_C) = \text{pos}(\xi'_{\bar{A}}, \vec{\xi}'_{\bar{A}})$ (through *zip* in algorithm 1). This candidate correlation $(\xi'_{\bar{A}}; \xi'_C)$ is added as an edge e_X to \mathcal{E}_X , adding the destination node to \mathcal{N}_X if not already present.

If ξ'_C represents a path between $\text{wr}(\text{allocBegin}(\dots))$ and $\text{wr}(\text{allocEnd}(\dots))$ for an alloc instruction in C , and $\xi'_{\bar{A}}$ is a corresponding path due to an $\text{alloc}_{s,v}$ instruction, and edges $e_C^{\theta_a}$ and $e_C^{\theta_m}$ in ξ'_C are labeled with instructions $\alpha_b := \theta(i_{32})$ and $\theta(i_{32} \rightarrow i_8)$ respectively due to (ALLOC), we add mappings $\mathcal{D}_X(e_X, e^{\theta_a}, 1) = v$ and $\mathcal{D}_X(e_X, e^{\theta_m}, 1) = M_{\bar{A}}$, where v is the address defined in $\xi'_{\bar{A}}$ due to either (ALLOCS) or (ALLOCV) (*addDetMappings* (e_X)). Notice that our algorithm only populates $\mathcal{D}_X(e_X, e^{\theta}, n)$ for $n = 1$, even though section 3.1 defines \mathcal{D}_X more generally.

³The number of paths can be exponential in procedure size, and so our implementation represents a pathset using a series-parallel digraph [Gupta et al. 2020] and annotates a pathset in \bar{A} in a single step. Similarly, a pathset in \bar{A} is correlated with a pathset in C in a single step. For easier exposition, the presented algorithm correlates each path individually.

$$\begin{array}{l}
\boxed{\text{affine}} \sum_i c_i v_i = c \quad \boxed{\text{ineqC}} \pm v \leq_s 2^c \quad \boxed{\text{ineq}} v_1 \odot v_2 \quad \boxed{\text{spOrd}} \boxed{\text{sp.p}_A^j} \leq_u (\boxed{\text{sp.p}_A^k} - v) \\
\boxed{\text{AllocEq}} \forall_{r \in B} \Sigma_C^r = \Sigma_{\bar{A}}^r \quad \boxed{\text{MemEq}} M_C = \Sigma_{\bar{A}}^B \setminus (\Sigma_{\bar{A}}^Z |^v) M_{\bar{A}} \quad \boxed{\text{zEmpty}} \{\Sigma_C^z, \Sigma_{\bar{A}}^z, \Sigma_{\bar{A}}^z |^v\} \{=, \neq\} \emptyset \\
\boxed{\text{spzBd}} \boxed{\text{em.z}} \vee (\boxed{\text{sp.p}_A^j} \odot \{\boxed{\text{lb.z}}, \boxed{\text{ub.z}}\}) \quad \boxed{\text{spzBd}'} \boxed{\text{em.z}} \vee (\boxed{\text{sp.p}_A^j} \leq_u (\boxed{\text{lb.z}} - \boxed{\text{1stSz.z}})) \\
\boxed{\text{gfySz}} \forall_{r \in GUFUY \setminus \{\text{vrdc}\}} (\boxed{\text{sz.r}} = \text{SZ}(T(r))) \quad \boxed{\text{vrdcSz}} (\boxed{\text{em.vrdc}} \Leftrightarrow \boxed{\text{sz.vrdc}} = 0) \quad \boxed{\text{Empty}} \forall_{r \in GUFUYUZ} (\Sigma_C^r = \emptyset \Leftrightarrow \boxed{\text{em.r}}) \\
\boxed{\text{gfyIntvl}} \forall_{r \in GUFUY} ((\boxed{\text{sz.r}} = 0) \vee ((\boxed{\text{lb.r}} \leq_u \boxed{\text{ub.r}}) \wedge (\boxed{\text{ub.r}} = \boxed{\text{lb.r}} + \boxed{\text{sz.r}} - 1_{i_{32}}) \wedge (\boxed{\text{lb.r}}, \boxed{\text{ub.r}}] = \Sigma_{\bar{A}}^r))) \\
\boxed{\text{zIntvl}} \boxed{\text{em.zl}} \vee ((\boxed{\text{lb.zl}} \leq_u \boxed{\text{ub.zl}}) \wedge (\boxed{\text{lb.zl}} + \boxed{\text{1stSz.zl}} - 1_{i_{32}} = \boxed{\text{ub.zl}}) \wedge (\boxed{\text{lb.zl}}, \boxed{\text{ub.zl}}] = \Sigma_{\bar{A}}^z)) \\
\boxed{\text{zaBd}} \boxed{\text{em.za}} \vee ((\boxed{\text{lb.za}} \leq_u \boxed{\text{ub.za}}) \wedge (\boxed{\text{lb.za}} + \boxed{\text{1stSz.za}} - 1_{i_{32}} \leq_u \boxed{\text{ub.za}}) \wedge (\boxed{\text{lb.za}} = \text{lb}(\Sigma_C^z) \wedge \boxed{\text{ub.za}} = \text{ub}(\Sigma_C^z))) \\
\boxed{\text{StkBd}} \Sigma_{\bar{A}}^{\{\text{stk}\}UY} \cup (\Sigma_{\bar{A}}^Z \setminus (\Sigma_{\bar{A}}^Z |^v)) = [\text{esp}, \text{stk}_e] \quad \boxed{\text{csBd}} \Sigma_{\bar{A}}^{\{\text{cs,cl}\}} = [\text{stk}_e + 1, \text{cs}_e] \\
\boxed{\text{NoOverlapC}} \neg \text{ov}(\Sigma_{\bar{A}}^{\text{hp}}, \Sigma_{\bar{A}}^{\text{cl}}, \Sigma_{\bar{A}}^{\text{vrdc}}, \dots, i_{\bar{A}}^g, \dots, i_{\bar{A}}^y, \dots, \Sigma_{\bar{A}}^z) \quad \boxed{\text{ROMC}} \forall_{r \in G_r} M_C = r_C \text{ ROM}_C^r(i_C^r) \\
\boxed{\text{NoOverlapA}} \neg \text{ov}(\Sigma_{\bar{A}}^{\{\text{hp,cl}\}UGUY}, \dots, \Sigma_{\bar{A}}^z, \dots, \Sigma_{\bar{A}}^{\text{stk}}, \Sigma_{\bar{A}}^{\text{cs}}, \Sigma_{\bar{A}}^r) \quad \boxed{\text{ROMA}} \forall_{r \in F_r} M_{\bar{A}} = r_{\bar{A}} \text{ ROM}_{\bar{A}}^r(i_{\bar{A}}^r)
\end{array}$$

Fig. 7. Predicate grammar for constructing candidate invariants. v represents a bitvector variable (registers, stack slots, and ghost variables), c represents a bitvector constant. $\odot \in \{\leq_{s,u}, <_{s,u}, >_{s,u}, \geq_{s,u}\}$.

If the destination node is not an error node, then the *inferInvariantsAndCounterexamples()* procedure updates the invariant network ϕ_X due to the addition of this new edge. The non-coverage requirements are checked after invariant inference (*checkSemanticRequirementsExceptCoverage*) and a candidate is discarded if the check fails.

When all simple paths between the cut points of \bar{A} are exhausted, the (Coverage \bar{A}) requirement must be satisfied by construction. *checkCoverageRequirements()* further checks the satisfaction of (CoverageC) before returning Success. Dynamo is sound because it returns Success only if all the thirteen search-algorithm requirements are satisfied.

The **chooseFrom** operator must attempt to maximize the chances of returning Success, even if only a fraction of the search space has been explored. Dynamo uses the counterexamples generated when a proof obligation is falsified (e.g., during invariant inference) to guide the search towards the more promising options. A counterexample is a proxy for the machine states of C and \bar{A} that may appear at a node n_X during the lockstep execution encoded by X . Thus, if at any step during the construction of X , the execution of a counterexample for a candidate partial solution (\bar{A}, X) results in the violation of a non-coverage requirement, that candidate is discarded. Further, counterexample execution opportunistically weakens the node invariants in X . Like Counter, we use the number of live registers in \bar{A} related through the current invariants in ϕ_X to rank the enumerated partial candidate solutions to implement a best-first search.

4.1 Invariant Inference

We use a counterexample-guided inference algorithm to identify node invariants [Gupta et al. 2020]. Candidate invariants at a node n_X of a partial product-graph are formed by conjuncting predicates drawn from the grammar shown in fig. 7. Apart from affine ($\boxed{\text{affine}}$) and inequality relations ($\boxed{\text{ineq}}$ and $\boxed{\text{ineqC}}$) for relating values across C and \bar{A} , the guesses attempt to equate the allocation and memory state of common regions across the two procedures ($\boxed{\text{AllocEq}}$ and $\boxed{\text{MemEq}}$).

Recall that we save stackpointer value at the boundary of a stackpointer updating instruction at PC p_A^j in ghost variable $\boxed{\text{sp.p}_A^j}$ ((OP-ESP) in fig. 4). To prove separation between different local variables, we require invariants that lower-bound the gap between two ghost variables, say $\boxed{\text{sp.p}_A^j}$ and $\boxed{\text{sp.p}_A^k}$, by some value v that depends on the allocation size operand of an `allocs` instruction ($\boxed{\text{spOrd}}$). To capture the various relations between lower bounds, upper bounds, region sizes, and $\boxed{\text{sp.p}_A^j}$, the guessing grammar includes shapes $\boxed{\text{spzBd}}$ and $\boxed{\text{spzBd}'}$ that are of the form: “either a local variable region is empty or its bounds are related to $\boxed{\text{sp.p}_A^j}$ in these possible ways”. $\boxed{\text{zEmpty}}$ tracks

the emptiness of the address-set of a local region. Together, these predicate shapes (along with $\boxed{\text{affine}}$ and $\boxed{\text{ineq}}$ relations between $\boxed{\text{sp.p}}_A^r$) enable disambiguation between stack writes involving spilled pseudo-registers and stack-allocated locals.

The predicate shapes listed below the dividing line segment in fig. 7 encode the *global invariants* that hold by construction (due to our execution semantics) at every non-error product-graph node n_X . $\boxed{\text{gfySz}}$, $\boxed{\text{vrdcSz}}$, and $\boxed{\text{gfyIntvl}}$ together encode the fact that the ghost variables associated with a region $r \in G \cup F \cup Y$ track its bounds, size, and that the address set of r is an interval. $\boxed{\text{Empty}}$ encodes that the ghost variable $\boxed{\text{em.r}}$ for $r \in G \cup F \cup Y \cup Z$ tracks the emptiness of the region r . $\boxed{\text{z1Intvl}}$ captures the property that a local variable region zl , if non-empty, must be an interval of size $\boxed{\text{1stSz.zl}}$. $\boxed{\text{zaBd}}$ captures a weaker property for a local region za (allocated using $\text{alloca}()$): if non-empty, this region must be bounded by its ghost variables and the region must be at least $\boxed{\text{1stSz.za}}$ bytes large. $\boxed{\text{StkBd}}$ encodes the invariant that the interval $[\text{esp}, \boxed{\text{stke}}]$ represents the union of the address sets of stk , regions in Y , and stack-allocated local regions ($\Sigma_{\check{A}}^Z \setminus (\Sigma_{\check{A}}^{Zl})^v$); $\boxed{\text{csBd}}$ is similarly shaped and encodes that the interval $[\boxed{\text{stke}} + 1, \boxed{\text{cse}}]$ represents the union of the address sets of regions cs and cl . $\boxed{\text{NoOverlapC}}$ encodes the disjointedness of all regions $r \in B$. $\boxed{\text{NoOverlapA}}$ encodes the disjointedness of all regions in \check{A} except virtually-allocated regions. Finally, $\boxed{\text{ROMC}}$ and $\boxed{\text{ROMA}}$ encode the preservation of memory contents of read-only regions in C and \check{A} .

A dataflow analysis [Andersen 1994] computes the possible states of $\beta()$ and $\beta_M()$ maps at each $n_C \in \mathcal{N}_C$, and the over-approximate solution is added to ϕ_{n_X} for each $n_X = (_ , n_C)$.

4.2 SMT Encoding

At a non-error node n_X , a proof obligation is represented as a first-order logic predicate over the state elements at n_X and discharged using an SMT solver. The machine states of C and \check{A} are represented using bitvectors (for a register/variable), arrays (for memory), and uninterpreted functions (for $\text{read}_{\vec{r}}(\Omega_P)$ and $\text{io}(\Omega_P, \vec{v}, r_w)$). For address sets, we encode the set-membership predicate $\alpha \in \Sigma_P^r$ for an arbitrary address α , region identifier r , and procedure $P \in \{C, \check{A}\}$. All other address set operations can be expressed in terms of the set-membership predicate. To simplify the encodings, we rely on the correct-by-construction invariants in fig. 7 and assume that ϕ_{n_X} satisfies the (Equivalence), (MAC), and (MemEq) requirements. Notice that (Equivalence) implies $\boxed{\text{AllocEq}}$.

Recall that for $z \in Z_l$, at a node $n_X \in \mathcal{N}_X$, $\Sigma_{\check{A}}^z|_s$ and $\Sigma_{\check{A}}^z|_v$ represent the address sets corresponding to the stack and virtual allocations performed in \check{A} for z . Let $Zl_s = \{z \mid z \in Z_l \wedge \Sigma_{\check{A}}^z|_s \neq \emptyset\}$ and $Zl_v = \{z \mid z \in Z_l \wedge \Sigma_{\check{A}}^z|_v \neq \emptyset\}$ represent the set of stack-allocated locals and virtually-allocated at n_X respectively. Recall that we restrict ourselves to only those compiler transformations that ensure the validity of $Zl_s \cap Zl_v = \emptyset$ at each n_X (section 2.4.3).

4.2.1 Representing Address-Sets Using Allocation State Array. Let $\mathcal{L}_P : i_{32} \rightarrow R$ be an *allocation state array* that maps an address to a region identifier in procedure P . For $r \notin Zl_v$, $\alpha \in \Sigma_P^r$ is encoded as $\text{sel}_1(\mathcal{L}_P, \alpha) = r$. Allocation of an address α to region r ($\Sigma_P^r := \Sigma_P^r \cup \{\alpha\}$) is encoded as $\mathcal{L}_P := \text{st}_1(\mathcal{L}_P, \alpha, r)$. Similarly, deallocation ($\Sigma_P^r := \Sigma_P^r \setminus \{\alpha\}$) is encoded as $\mathcal{L}_P := \text{st}_1(\mathcal{L}_P, \alpha, \text{free})$.

For $zlv \in Zl_v$, both $\alpha \in \Sigma_C^{zlv}$ and $\alpha \in \Sigma_{\check{A}}^{zlv}$ are encoded as $\text{sel}_1(\mathcal{L}_C, \alpha) = zlv$, i.e., the set-membership encodings for both procedures use \mathcal{L}_C for virtually-allocated locals (by relying on the $\boxed{\text{AllocEq}}$ invariant at n_X). In other words, $\mathcal{L}_{\check{A}}$ is not used to track the virtually-allocated locals; instead, an address belonging to a virtually allocated-region maps to one of $\{\text{free}, \text{stk}, \text{cs}\} \cup F$ regions in $\mathcal{L}_{\check{A}}$. Consequently, the (de)allocation instructions $\Sigma_{\check{A}}^{zlv}|^v := \Sigma_{\check{A}}^{zlv}|^v \cup [v]_w$ and $\Sigma_{\check{A}}^{zlv}|^v := \emptyset$ are vacuous in \check{A} , i.e., they do not change any state element in \check{A} (fig. 6).

Table 2. SMT encoding of $\alpha \in \Sigma_P^r$ for Dynamo's proof obligation O with outgoing assembly path $\xi_{\bar{A}}$.

$\alpha \in \Sigma_P^r$	Full-array encoding		Partial-interval encoding ($\Sigma_P^{Za} \neq \emptyset$)	Full-interval encoding ($\Sigma_P^{Za} = \emptyset$)
	$P = C$	$P = A$		
$r = hp$	$\text{sel}_1(\mathcal{L}_C, \alpha) = r$		$\alpha \notin (\Sigma_{\bar{A}}^{GUF} \cup Zlv_U(\xi_{\bar{A}}) \cup [SP_{min}(\xi_{\bar{A}}), \text{cse}])$	
$r = cl$			$\alpha \in [\text{stke}] + 1, \text{cse} \wedge \alpha \notin Zlv_U(\xi_{\bar{A}})$	
$r \in G \cup Zlv$				
$r \in Y \cup Z_a \cup Zls$			$\neg \text{em.r} \wedge (\text{lb.r}) \leq u \alpha \leq (\text{ub.r})$	
$r \in F$	false			
$r = cs$			$\alpha \in [\text{stke}] + 1, \text{cse} \wedge \alpha \in Zlv_U(\xi_{\bar{A}})$	
$r = stk$			$\text{sel}_1(\mathcal{L}_{\bar{A}}, \alpha) = r$	$\alpha \in [SP_{min}(\xi_{\bar{A}}), \text{stke}] \wedge \bigwedge_{r \in Y \cup Zls} (\alpha \notin \Sigma_{\bar{A}}^r)$

This encoding, based on allocation state arrays \mathcal{L}_C and $\mathcal{L}_{\bar{A}}$, is called the *full-array encoding*. The second and third columns of table 2 describe the full-array encoding for $P = C$ and $P = \bar{A}$. In the table, we use $\boxed{\text{AllocEq}}$ to replace $\text{sel}_1(\mathcal{L}_{\bar{A}}, \alpha)$ with $\text{sel}_1(\mathcal{L}_C, \alpha)$ for $r \in B$. For example, in the full-array encoding, the (MemEq) requirement $M_C =_{\Sigma_{\bar{A}}^B \setminus (\Sigma_{\bar{A}}^{Zl} \cup \sigma)} M_{\bar{A}}$ becomes $\forall \alpha : ((\text{sel}_1(\mathcal{L}_C, \alpha) \in G \cup \{hp, cl\}) \cup Y \cup Zls \cup Z_a) \Rightarrow (\text{sel}_1(M_C, \alpha) = \text{sel}_1(M_{\bar{A}}, \alpha))$.

4.2.2 Interval Encodings for $r \in G \cup F \cup Y \cup Z_l \cup \{stk\}$. We use $\boxed{\text{gfyIntvl}}$, $\boxed{\text{zlIntvl}}$, and $\boxed{\text{AllocEq}}$ invariants for a more efficient *interval encoding*: for $r \in G \cup F \cup Y \cup Z_l$, we encode $\alpha \in \Sigma_P^r$ as $\neg \text{em.r} \wedge (\text{lb.r}) \leq u \alpha \leq (\text{ub.r})$. Moreover, if there are no local variables allocated due to the $\text{alloca}()$ operator (i.e., $\Sigma_P^{Za} = \emptyset$), then all local variables are contiguous, and so, due to $\boxed{\text{StkBd}}$, the stk region can be identified as $[\text{esp}, \text{stke}] \setminus \Sigma_{\bar{A}}^{Y \cup Zls}$ – the corresponding interval encoding is shown in the right-most cell of $r = stk$ row in table 2.

4.2.3 Interval Encodings for $r \in \{hp, cl, cs\}$. Even though hp, cl, cs can be discontinuous regions in general, we over-approximate these regions to their contiguous covers to be able to soundly encode them using intervals. At a node $n_X = (n_{\bar{A}}, n_C)$, Dynamo may generate a proof obligation O of the form $\{pre\}(\xi_{\bar{A}}; [\xi_C]_{\mathcal{D}_X}^{ex})\{post\}$ – recall that path-cover and path-infeasibility conditions are also represented as Hoare triples with $\xi_C = \epsilon$. If $\xi_{\bar{A}}$ is an I/O path, its execution interacts with the outside world, and so an over-approximation of an externally-visible address set is unsound. We thus restrict our attention to an I/O-free $\xi_{\bar{A}}$ for interval encoding.

Let $n_{\bar{A}}^1, n_{\bar{A}}^2, \dots, n_{\bar{A}}^m$ be the nodes on path $\xi_{\bar{A}} = (n_{\bar{A}} \rightarrow n_{\bar{A}}^t)$, such that $n_{\bar{A}}^1 = n_{\bar{A}}$ and $n_{\bar{A}}^m = n_{\bar{A}}^t$. Let $SP_{min}(\xi_{\bar{A}})$ represent the the minimum value of esp observed at any node $n_{\bar{A}}^j$ ($1 \leq j \leq m$) visited during the execution of path $\xi_{\bar{A}}$. Similarly, let $Zlv_U(\xi_{\bar{A}})$ be the union of the values of set $\Sigma_{\bar{A}}^{Zlv}$ observed at any $n_{\bar{A}}^j$ ($1 \leq j \leq m$) visited during $\xi_{\bar{A}}$'s execution.

Let $HP(\xi_{\bar{A}}) = \text{comp}(\Sigma_{\bar{A}}^{GUF} \cup Zlv_U(\xi_{\bar{A}}) \cup [SP_{min}(\xi_{\bar{A}}), \text{cse}])$, $CL(\xi_{\bar{A}}) = [\text{stke}] + 1, \text{cse} \setminus Zlv_U(\xi_{\bar{A}})$, and $CS(\xi_{\bar{A}}) = [\text{stke}] + 1, \text{cse} \cap Zlv_U(\xi_{\bar{A}})$.

THEOREM 4.1. *Let $O = \{pre\}(\xi_{\bar{A}}; [\xi_C]_{\mathcal{D}_X}^{ex})\{post\}$ be a proof obligation generated by Dynamo. Let O' be obtained from O by strengthening precondition pre to $pre' = (pre \wedge ((\Sigma_{\bar{A}}^{hp} = HP(\xi_{\bar{A}})) \wedge (\Sigma_{\bar{A}}^{cl} = CL(\xi_{\bar{A}})) \wedge (\Sigma_{\bar{A}}^{cs} = CS(\xi_{\bar{A}})))$. If $\xi_{\bar{A}}$ is I/O-free, $O \Leftrightarrow O'$ holds.*

PROOF SKETCH. $O \Rightarrow O'$ is trivial. The proof for $O' \Rightarrow O$, available in [Rose and Bansal 2024b], relies on the limited shapes of predicates ξ that may appear in $pre, post$ – for I/O-free $\xi_{\bar{A}}$, these shapes are limited by our invariant grammar (fig. 7), and the edge conditions appearing in our execution semantics (figs. 3 to 6). The proof holds only if the safety-relaxed semantics are used. \square

Using theorem 4.1, we rewrite $\alpha \in \Sigma_P^{hp}$ to $\alpha \in HP(\xi_{\bar{A}})$, $\alpha \in \Sigma_P^{cl}$ to $\alpha \in CL(\xi_{\bar{A}})$, and $\alpha \in \Sigma_P^{cs}$ to $\alpha \in CS(\xi_{\bar{A}})$ in proof obligation O . As shown in table 2, if $\Sigma_P^{Za} = \emptyset$ holds at n_X , we encode all non-free regions using intervals (called *full-interval encoding*); else, we encode regions in $Y \cup Z_a \cup Zls \cup \{stk\}$

Table 3. Benchmarks and their programming patterns. N in $vilN$ is substituted to obtain $vil1$, $vil2$, and $vil3$. Program listings available in [Rose and Bansal 2024b].

Name	Programming pattern
ats	Address-taken local scalar <code>int ats() { int ret; foo(&ret); return ret; }</code>
atc	Address taken conditionally <code>int atc(int* p) { int x; if (!p) p = &x; foo(p); return *p }</code>
ata	Local array <code>int ata() { char ret[8]; foo(ret); return bar(ret, 0, 16); }</code>
vw1	Variadic procedure <code>int vw1(int n, ...) { va_list a; va_start(a, n); for(...){ va_arg(a,int) ... }</code>
as	GCC alloca() <code>int as(int n){...int* p=alloca(n*sizeof(n)); for(...){/*write to p*/}...</code>
vsl	VLA with loop <code>int vsl(int n){... int v[n]; for(...){/*write to v*/}...</code>
vcu	VLA conditional use <code>int vcu(int n,int k){ int a[n]; if (...) { /*rd/wr to a*/}...</code>
min	minprintf procedure from K&R [Kernighan and Ritchie 1988]
ac	alloca() conditional use <code>int ac(char*a) {...if (!a) a=alloca(n); for(...)/r/w to a*/}</code>
all	alloca() in a loop to form a linked list <code>all(..hd=NULL; for(...){..n=alloca(..);..n->nxt=hd; hd=n;} while(...){/* traverse the list starting at hd */}</code>
atail	Local array alloc. in loop <code>int atail(..){..for(...){ char a[4096]; f(a.); b(a.);...}...</code>
vilN	N VLA(s) in a loop <code>int vilN(..){..for(i=1;i<n;++i){ int v1[4*i], ... vN[4*i]; fooN(...); ..}.. }</code>
vilcc	VLA in loop with continue <code>int vilcc(..){..while(i<n){ char v[i];...if(..) continue;...}..}</code>
fib	Program from fig. 1
vilce	VLA in loop with break <code>int vilce(..){..while(i<n){ char v[i];...if(..) break;...}..}</code>
rod	A local char array initialized using string; a VLA; a for loop Available in [Rose and Bansal 2024b].

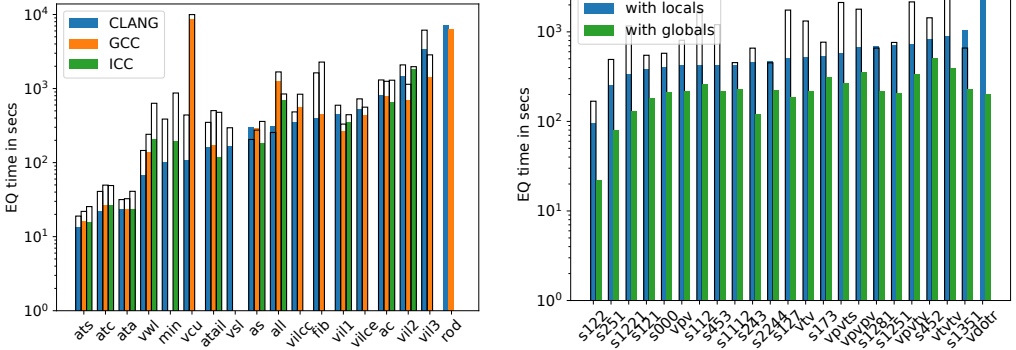
using an allocation state array, and $G \cup F \cup Zlv \cup \{hp, cl, cs\}$ using intervals (called *partial-interval encoding*).

5 EXPERIMENTS

Dynamo uses four SMT solvers running in parallel for discharging proof obligations: $z3-4.8.7$, $z3-4.8.14$, $Yices2-45e38fc$, and $cvc4-1.7$. Unless otherwise specified, we use $\mu = 64$, a timeout of ten minutes for an SMT query, and a timeout of eight hours for a refinement check.

Before checking refinement, if the address of a local variable l is never taken in C , we transform C to register-allocate l (LLVM’s `mem2reg`). This reduces the proof effort, at the cost of having to trust the pseudo-register allocation logic. `mem2reg` does not register-allocate local arrays and structs in $LLVM_d$, even though an optimizing compiler may register-allocate them in assembly – virtual allocations help validate such translations.

We first evaluate the efficacy of our implementation to handle the diverse programming patterns seen with local allocations (table 3). These include variadic procedures, VLAs allocated in loops, `alloca()` in loops, etc. Figure 8a shows the results of our experiments for these 18 programming patterns from table 3 and three compilers, namely Clang/LLVM v12.0.0, GCC v8.4.0, and ICC v2021.8.0, to generate 32-bit x86 executables at `-O3` optimization with inter-procedural analyses disabled using the compilers’ command-line flags. The X-axis lists the benchmarks and the Y-axis represents the total time taken in seconds (log scale) for a refinement check – to study the performance implications, we run a check with all three encodings for these benchmarks. The filled and empty bars represent the time taken with full-interval and partial-interval SMT encodings respectively. The figure does not show the results for the full-array encoding. A missing bar represents a failure to compute the proof. Of 54 procedure pairs, our implementation is able to check refinement for 45, 43, and 37 pairs while using full-interval, partial-interval, and full-array encodings respectively. For benchmarks where a refinement check succeeds for all encodings, the full-interval encoding performs 1.7-2.2x and 3.5-4.9x faster on average (for each compiler) than the partial-interval and full-array encodings respectively. The reasons for nine failures are: (a)



(a) Comparison of running times with full- (filled bars) and partial- (empty bars) interval encoding. (b) Comparison of running times of benchmarks with exactly same code modulo allocation.

Fig. 8. Experiments with procedures in table 3 and TSVC. Y-axis is logarithmically scaled.

Table 4. Statistics obtained by running Dynamo on procedures in the bzip2 program.

Name	SLOC	ALOC	# <i>al</i>	# <i>loop</i>	# <i>fcall</i>	D	eqT	Nodes	Edges	EXP	BT	# <i>q</i>	Avg. qT
generateMTFValues	76	144	1	6	1	2	4k	14	30	60	16	3860	0.56
recvDecodingTables	70	199	2	14	10	3	3k	38	66	102	15	5611	0.21
undoReversible-Transformation_fast	116	221	1	7	6	2	2k	21	34	43	6	2998	0.23

limitation of the blackbox annotation algorithm for one procedure-pair; (b) incompleteness of invariant inference for six procedure-pairs (e.g., requirement of non-affine invariants, choice of program variables); and (c) SMT solver timeouts for two procedure-pairs. `vilcc` and `vilce` require multiple `deallocs` instructions to be added to A for a single `dealloc` in C . An `allocv` annotation is required for the '`va_list a`' variable in the GCC and ICC compilations of `vwl` (see table 3) – while GCC and ICC register-allocate `a`, it is allocated in memory using `alloc` in $LLVM_d$ (even after `mem2reg`). The average number of best-first search backtrackings across all benchmarks is 2.8. The time spent in constructing the correct product graph forms around 70-80% of the total search time.

We next evaluate Dynamo on the TSVC suite of vectorization benchmarks with arrays and loops [Maleki et al. 2011], also used in previous work [Churchill et al. 2019; Gupta et al. 2020]. We use two versions of these benchmarks: (1) 'globals' where global variables are used for storing the output array values, and (2) 'locals' where local array variables are used for storing the output values and a procedure call is added at the end of the procedure body to print the contents of the local array variables. The compiler performs the same vectorizing transformations on both versions. Unlike globals, locals additionally requires the automatic identification of required annotations.

Figure 8b shows the execution times of Dynamo for validating the compilations produced by Clang/LLVM v12.0.0 (at -O3) for these two versions of the TSVC benchmarks. Dynamo can successfully validate these compilations. Compared to globals, refinement checks are 2.5x slower for locals (on average) due to the extra overhead of identifying the required annotations.

Our third experiment is on SPEC CPU2000's `bzip2` [Henning 2000] program compiled using Clang/LLVM v12.0.0 at three optimization levels: O1, O2, and O1-. O1- is a custom optimization level configured by us that enables all optimizations at O1 *except* (a) merging of multiple procedure calls on different paths into a single call, (b) early-CSE (common subexpression elimination), (c) loop-invariant code motion at both LLVM IR and Machine IR, (d) dead-argument elimination, (e) inter-procedural sparse conditional constant propagation, and (f) dead-code elimination of procedure calls. `bzip2` runs 2% slower with O1- than with O1; this is still 5% faster than the executable

produced by CompCert, for example. Of all 72 procedures in `bzip2`, Dynamo successfully validates the translations for 64, 60, and 54 procedures at O1-, O1, and O2 respectively at $\mu = 1$. At O1-, Dynamo takes around six CPU hours to compute refinement proofs for the 64 procedures. Dynamo times out for the remaining eight procedures, all of which are bigger than 190 ALOC.

Three of `bzip2`'s procedures for which refinement proofs are successfully computed at both O1- and O1 contain at least one local array, and table 4 presents statistics for the O1- validation experiments for these procedures. For each procedure, we show the number of source lines of code in C (SLOC), the number of assembly instructions in A (ALOC), the number of local variables ($\#_{al}$), the number of loops ($\#_{loop}$), the number of procedure calls ($\#_{call}$), and the maximum loop nest depth (D). The eqT column shows the validation times (in seconds). The Nodes and Edges columns show the number of nodes and edges in the final product graph, and BT and EXP is the number of backtrackings and the number of (partial) candidate product graphs explored by Dynamo respectively. $\#_q$ is the total number of SMT queries discharged, and Avg. qT is the average time taken by an SMT query in seconds for the refinement check.

In a separate experiment, we split the large procedures in `bzip2` into smaller procedures, so that Dynamo successfully validates the O1- compilation of the full modified `bzip2` program: the splitting disables some compiler transformations and also reduces the correlation search space.

Through our experiments, we uncovered and reported a bug in recent versions of `z3`, including `z3-4.8.14` and `z3-4.12.5`, where for an input satisfiability query Ψ , the SMT solver returns an unsound model (counterexample) that evaluates Ψ to false [z3b 2024]. When a modern SMT solver is used to validate compilations produced by a mature compiler, a bug may be found on either side.

6 RELATED WORK AND CONCLUSIONS

CoVaC [Zaks and Pnueli 2008] automatically identifies a product program that demonstrates observable equivalence for deterministic programs. Counter [Gupta et al. 2020] extends CoVaC to support path-specializing transformations, such as loop unrolling, through counterexample-guided search heuristics. We extend these prior works to support refinement between programs performing dynamic allocations with non-deterministic addresses for local variables and stack.

Recent work on bounded TV [Lee et al. 2021] models allocations through *separate blocks*, so a pointer is represented as a combination of a block-ID and an offset into a block. While this suffices for the bounded TV setting, our problem setting requires a more general representation of a dynamically-allocated variable (e.g., allocation-site) and a more general SMT encoding.

CompCert provides axiomatic semantics for memory (de)allocation in the source Clight program, and proves their preservation along the compilation pipeline [Leroy and Blazy 2008]. They restrict their proof method to CompCert's preallocation strategy for local variables, possibly to avoid the manual effort required to write mechanized proofs for a more general allocation strategy. Preallocation of local variables has also been used in prior work on TV for a verified OS kernel [Sewell et al. 2013]. Preallocation can be space inefficient and cannot support VLAs and `alloca()`. Further, TV for a third-party compiler cannot assume a particular allocation strategy.

We provide a semantic model, refinement definition, and an algorithm to determine the correctness of a third-party translation from an unoptimized high-level representation of a C program to an optimized assembly program in the presence of dynamically-allocated local memory. Our semantic model and definition of refinement require that for allocations and procedure calls that reuse stack space, their relative order is preserved in both programs. While our experiments show that this suffices in practice, a more general definition of refinement, that admits transformations that may reorder (de)allocations while reusing stack space, is perhaps a good candidate for future work.

DATA-AVAILABILITY STATEMENT

The Dynamo tool that supports section 5 is available on Zenodo [Rose and Bansal 2024a] with instructions for complete reproducibility of the presented results.

ACKNOWLEDGMENTS

We thank Shubhani Gupta for contributing towards scalability improvements of the translation validation tool. We thank Abhishek Dang for carefully reading previous drafts of the paper, and pointing out several errors, and making important suggestions that improved the paper significantly.

REFERENCES

2024. Z3 bug report for an unsound model. <https://github.com/Z3Prover/z3/issues/7132>.
- Lars Ole Andersen. 1994. *Program Analysis and Specialization for the C Programming Language*. Technical Report. Berkeley Churchill, Oded Padon, Rahul Sharma, and Alex Aiken. 2019. Semantic Program Alignment for Equivalence Checking. In *Proceedings of the 40th ACM SIGPLAN Conference on Programming Language Design and Implementation* (Phoenix, AZ, USA) (PLDI 2019). ACM, New York, NY, USA, 1027–1040. <https://doi.org/10.1145/3314221.3314596>
- Shubhani Gupta, Abhishek Rose, and Sorav Bansal. 2020. Counterexample-Guided Correlation Algorithm for Translation Validation. *Proc. ACM Program. Lang.* 4, OOPSLA, Article 221 (Nov. 2020), 29 pages. <https://doi.org/10.1145/3428289>
- John L. Henning. 2000. SPEC CPU2000: Measuring CPU performance in the new millenium. *IEEE Computer* 33, 7 (July 2000), 28–35.
- Jeehoon Kang, Yoonseung Kim, Youngju Song, Juneyoung Lee, Sanghoon Park, Mark Dongyeon Shin, Yonghyun Kim, Sungkeun Cho, Joonwon Choi, Chung-Kil Hur, and Kwangkeun Yi. 2018. Crellvm: Verified Credible Compilation for LLVM. In *Proceedings of the 39th ACM SIGPLAN Conference on Programming Language Design and Implementation* (Philadelphia, PA, USA) (PLDI 2018). ACM, New York, NY, USA, 631–645. <https://doi.org/10.1145/3192366.3192377>
- Theodoros Kasampalis, Daejun Park, Zhengyao Lin, Vikram S. Adve, and Grigore Roşu. 2021. Language-Parametric Compiler Validation with Application to LLVM. In *Proceedings of the 26th ACM International Conference on Architectural Support for Programming Languages and Operating Systems* (Virtual, USA) (ASPLOS 2021). Association for Computing Machinery, New York, NY, USA, 1004–1019.
- Brian W. Kernighan and Dennis M. Ritchie. 1988. *The C Programming Language* (2nd ed.). Prentice Hall Professional Technical Reference.
- Juneyoung Lee, Dongjoo Kim, Chung-Kil Hur, and Nuno P. Lopes. 2021. An SMT Encoding of LLVM’s Memory Model for Bounded Translation Validation. In *Computer Aided Verification*, Alexandra Silva and K. Rustan M. Leino (Eds.). Springer International Publishing, Cham, 752–776.
- Xavier Leroy. 2006. Formal certification of a compiler back-end, or: programming a compiler with a proof assistant. In *33rd ACM symposium on Principles of Programming Languages*. ACM Press, 42–54. <http://gallium.inria.fr/~xleroy/publi/compiler-certif.pdf>
- Xavier Leroy and Sandrine Blazy. 2008. Formal Verification of a C-like Memory Model and Its Uses for Verifying Program Transformations. *J. Autom. Reason.* 41, 1 (2008), 1–31. <https://doi.org/10.1007/s10817-008-9099-0>
- Nuno P. Lopes, Juneyoung Lee, Chung-Kil Hur, Zhengyang Liu, and John Regehr. 2021. Alive2: Bounded Translation Validation for LLVM. In *Proceedings of the 42nd ACM SIGPLAN International Conference on Programming Language Design and Implementation* (Virtual, Canada) (PLDI 2021). Association for Computing Machinery, New York, NY, USA, 65–79. <https://doi.org/10.1145/3453483.3454030>
- Saeed Maleki, Yaoqing Gao, Maria J. Garzarán, Tommy Wong, and David A. Padua. 2011. An Evaluation of Vectorizing Compilers. In *Proceedings of the 2011 International Conference on Parallel Architectures and Compilation Techniques (PACT ’11)*. IEEE Computer Society, Washington, DC, USA, 372–382. <https://doi.org/10.1109/PACT.2011.68>
- David Menendez, Santosh Nagarakatte, and Aarti Gupta. 2016. Alive-FP: Automated Verification of Floating Point Based Peephole Optimizations in LLVM. 317–337. https://doi.org/10.1007/978-3-662-53413-7_16
- KedarS. Namjoshi and LenoreD. Zuck. 2013. Witnessing Program Transformations. In *Static Analysis*, Francesco Logozzo and Manuel Fähndrich (Eds.). Lecture Notes in Computer Science, Vol. 7935. Springer Berlin Heidelberg, 304–323. https://doi.org/10.1007/978-3-642-38856-9_17
- George C. Necula. 2000. Translation Validation for an Optimizing Compiler. In *Proceedings of the ACM SIGPLAN 2000 Conference on Programming Language Design and Implementation* (Vancouver, British Columbia, Canada) (PLDI ’00). ACM, New York, NY, USA, 83–94. <https://doi.org/10.1145/349299.349314>
- Abhishek Rose and Sorav Bansal. 2024a. *Artifact for paper "Modeling Dynamic (De)Allocations of Local Memory for Translation Validation"*. <https://doi.org/10.5281/zenodo.10797459>

- Abhishek Rose and Sorav Bansal. 2024b. *Modeling Dynamic (De)Allocations of Local Memory for Translation Validation*. Technical Report. IIT Delhi. <https://arxiv.org/abs/2403.05302>
- Thomas Arthur Leck Sewell, Magnus O. Myreen, and Gerwin Klein. 2013. Translation Validation for a Verified OS Kernel. In *Proceedings of the 34th ACM SIGPLAN Conference on Programming Language Design and Implementation* (Seattle, Washington, USA) (*PLDI '13*). Association for Computing Machinery, New York, NY, USA, 471–482. <https://doi.org/10.1145/2491956.2462183>
- Rahul Sharma, Eric Schkufza, Berkeley Churchill, and Alex Aiken. 2013. Data-driven Equivalence Checking. In *Proceedings of the 2013 ACM SIGPLAN International Conference on Object Oriented Programming Systems Languages & Applications* (Indianapolis, Indiana, USA) (*OOPSLA '13*). ACM, New York, NY, USA, 391–406. <https://doi.org/10.1145/2509136.2509509>
- Bjarne Steensgaard. 1996. Points-to analysis in almost linear time. In *Proceedings of the 23rd ACM SIGPLAN-SIGACT symposium on Principles of programming languages*. 32–41.
- Michael Stepp, Ross Tate, and Sorin Lerner. 2011. Equality-based Translation Validator for LLVM. In *Proceedings of the 23rd International Conference on Computer Aided Verification* (Snowbird, UT) (*CAV'11*). Springer-Verlag, Berlin, Heidelberg, 737–742. <http://dl.acm.org/citation.cfm?id=2032305.2032364>
- Jean-Baptiste Tristan, Paul Govereau, and Greg Morrisett. 2011. Evaluating Value-graph Translation Validation for LLVM. In *Proceedings of the 32nd ACM SIGPLAN Conference on Programming Language Design and Implementation* (San Jose, California, USA) (*PLDI '11*). ACM, New York, NY, USA, 295–305. <https://doi.org/10.1145/1993498.1993533>
- Anna Zaks and Amir Pnueli. 2008. CoVaC: Compiler Validation by Program Analysis of the Cross-Product. In *Proceedings of the 15th International Symposium on Formal Methods* (Turku, Finland) (*FM '08*). Springer-Verlag, Berlin, Heidelberg, 35–51. https://doi.org/10.1007/978-3-540-68237-0_5
- Jianzhou Zhao, Santosh Nagarakatte, Milo M.K. Martin, and Steve Zdancewic. 2012. Formalizing the LLVM Intermediate Representation for Verified Program Transformations. In *Proceedings of the 39th Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages* (Philadelphia, PA, USA) (*POPL '12*). ACM, New York, NY, USA, 427–440. <https://doi.org/10.1145/2103656.2103709>
- Jianzhou Zhao, Santosh Nagarakatte, Milo M.K. Martin, and Steve Zdancewic. 2013. Formal Verification of SSA-based Optimizations for LLVM. In *Proceedings of the 34th ACM SIGPLAN Conference on Programming Language Design and Implementation* (Seattle, Washington, USA) (*PLDI '13*). ACM, New York, NY, USA, 175–186. <https://doi.org/10.1145/2491956.2462164>

Received 21-OCT-2023; accepted 2024-02-24