# Targeted Testing of Compiler Optimizations via Grammar-Level Composition Styles

ZITONG ZHOU, University of California, Los Angeles (UCLA), USA

BEN LIMPANUKORN, University of California, Los Angeles (UCLA), USA

HONG JIN KANG, The University of Sydney, Australia

JIYUAN WANG, Tulane University, USA

YAOXUAN WU, University of California, Los Angeles (UCLA), USA

AKOS KISS, University of Szeged, Hungary

RENATA HODOVAN, University of Szeged, Hungary

MIRYUNG KIM, University of California, Los Angeles (UCLA), USA

Ensuring the correctness of optimizing compilers is critical yet challenging. Fuzzing offers a promising path, but existing fuzzers struggle to test optimization logic effectively. In our preliminary study, the state-of-the-art LLVM fuzzer GrayC covered only 12% of LLVM's optimization module—and left 47% of optimization passes entirely untested—after four hours of coverage-guided fuzzing. We identify two key challenges. First, most fuzzers use optimization pipelines (heuristics-based, fixed sequences of passes) as their harness. Because of the phase-ordering problem, pass ordering can enable or preempt transformations, so pipelines inevitably miss many optimization interactions; moreover, many optimizations are not scheduled, even at aggressive levels. Second, optimizations typically fire only when inputs satisfy specific structural relationships, which existing generators and mutations struggle to produce.

We propose targeted fuzzing of individual optimizations to complement pipeline-based testing. Our key idea is to exploit composition styles—structural relations over program constructs (e.g., adjacency, nesting, repetition, ordering)—that optimizations look for. We build a general-purpose, grammar-based mutational fuzzer, TARGETFUZZ, that (i) mines candidate composition styles from an optimization-relevant corpus, then (ii) rebuilds those styles inside different contexts offered by a larger, generic corpus via lightweight, synthesized mutations to exercise variations of optimization logic. TARGETFUZZ can be easily adapted to a new programming language by providing lightweight, grammar-based, construct annotations—and it automatically synthesizes mutators and crossovers to rebuild composition styles. There is no need for hand-coded generators or language-specific mutators, which is particularly useful for modular frameworks such as MLIR, whose dialect-based, rapidly evolving ecosystem makes optimizations difficult to fuzz. Our evaluation on LLVM and MLIR shows that TARGETFUZZ improves coverage by 8% and 11% and triggers optimizations 2.8× and 2.6×, compared to baseline fuzzers under the targeted fuzzing mode. We show that targeted fuzzing is complementary: it effectively tests all 37 sampled LLVM optimizations, while pipeline-fuzzing missed 12.

## 1 INTRODUCTION

Modern optimizing compilers are monolithic software, composed of layers that transform a high-level program into machine instructions. As optimizing compilers evolve, their growing complexity introduces new challenges for ensuring correctness [4, 5]. Optimizing compilers can crash, hang, or

Authors' addresses: Zitong Zhou, University of California, Los Angeles (UCLA), Los Angeles, CA, USA, zitongzhou@cs.ucla.edu; Ben Limpanukorn, University of California, Los Angeles (UCLA), Los Angeles, CA, USA, blimpan@cs.ucla.edu; Hong Jin Kang, The University of Sydney, Sydney, NSW, Australia, hongjin.kang@sydney.edu.au; Jiyuan Wang, Tulane University, New Orleans, LA, USA, wjiyuan@tulane.edu; Yaoxuan Wu, University of California, Los Angeles (UCLA), Los Angeles, CA, USA, thaddywu@cs.ucla.edu; Akos Kiss, University of Szeged, Szeged, Hungary, akiss@inf.u-szeged.hu; Renata Hodovan, University of Szeged, Szeged, Hungary, hodovan@inf.u-szeged.hu; Miryung Kim, University of California, Los Angeles (UCLA), Los Angeles, CA, USA, miryung@cs.ucla.edu.

miscompile, where bugs silently make their way into production software. In an empirical study of over 10K confirmed optimization bugs in GCC and LLVM, Zhou et al. reported that the optimizer is one of the most buggy components in the compilers [57]. Sun et al. conducted a systematic study of compiler bugs in GCC and LLVM and found that high priority tends to be assigned to optimizer bugs, most notably 30% of the bugs in GCC's inter-procedural analysis component are labeled to be the highest priority (P1) [46]. Given the importance of optimizing compilers in the software supply chain, ensuring their correctness is essential to the reliability of virtually all software.

Fuzzing or automated random testing offers a practical way to expose these hidden flaws in compilers. However, existing fuzzers are limited in their ability to thoroughly test compilers' optimization components. We performed artifact evaluation of the state-of-the-art C-language fuzzer GrayC [20]. Using its default configuration of feeding input programs to test LLVM through `clang -O3`, after 4 hours of fuzzing, 47.5% of the optimization passes were entirely untested. It only achieves total 12.12% line coverage in the transformation module of LLVM. On source files of optimization passes (excluding utilities), it has only 8.66% line coverage on average.

Most fuzzers (like GrayC) drive the compiler through aggressive, canonical optimization pipelines (e.g., `clang -O3`) in hopes of maximizing optimization behaviors. There are 2 key challenges to this approach. First, the classic phase-ordering problem [16, 48, 51]: modern optimizers apply optimization passes by heuristics, so the order of passes can enable or preempt one another. An example in LLVM is that "(fully) unrolling loops early removes opportunities for the loop vectorizer" [38]. In practice, this means that in fuzzing, many optimizations are overshadowed by others and subsequently under-tested, as the GrayC experiment demonstrates. Second, even aggressive pipelines are not *exhaustive* as they need to balance compile time and stability, so not every optimization is enabled by default [37]. In addition, the growing popularity of modular and extensible compiler frameworks exacerbates the problem. For example, MLIR [19] allows developers to build custom IR, passes and pipelines, making it even more difficult to exercise many optimizations in a single canonical pipeline. To address these two challenges, we propose to *complement whole-pipeline fuzzing (good for end-to-end interactions) with* **targeted pass fuzzing** *that exercises corner-case optimization interactions which the pipelines rarely reach.* We expect targeted fuzzing to isolate and test individual optimization passes or small groups of passes, by feeding inputs directly to them, therefore addressing these challenges. We detail the limitations of whole-pipeline fuzzing in Section 2 and targeted pass fuzzing in Section 3.

Another fundamental challenge to fuzzing compiler optimizations is that fuzzer-generated inputs cannot easily satisfy the triggering conditions expected by optimizations. This is because triggering optimizations requires further structural constraints on program inputs. For instance, optimization **Loop Fusion** (i.e., combining loop bodies) only triggers when two immediately adjacent loops share the same trip count [15]. *We observe that compiler optimizations target many structural relations (such as adjacency, repetition, and nesting) on grammar-based structures.* For example, LLVM's **Loop Flatten** [36] flattens nested loops; and **Loop Fusion** fuses adjacent loops. Our intuition traces back to the theoretical foundations of compiler optimizing transformations [11, 12]. That is, conceptually, an optimizing transformation consists of a pattern-matcher and a rewrite rule [7]. The pattern-matcher finds structural relations expressed on intermediate representations or abstract syntax trees. We propose a taxonomy of such structural relations defined at the level of grammar-based constructs, coupled with mutator synthesis to reconstruct these optimization-triggering relations on a large corpus. In this paper, we refer to these structural relations as composition styles.

For example, we define the Balanced composition that can match arithmetic expressions located in multiple branches of an `if-then-else` construct. Balanced maps to a set of composition-specific mutations, including Replicate, which, given a new program, tries to locate a similar `if-then-else` construct containing an arithmetic expression, then replicates that expression in other branches,

thus reconstructing the Balanced style in a new context. The idea is that by reconstructing these optimization-triggering relations in a variety of code contexts, we can exercise diverse variations of optimization logic.

We embody this idea of composition-based mutational fuzzing in a tool named TargetFuzz. The core functionality of TargetFuzz is that it defines a set of composition styles such as Balanced. For each composition style, a compiler developer can *register* it to match relevant *program constructs* that compiler optimizations operate on, such as if-else blocks, loops, logical and arithmetic expressions, function calls, arrays, memory references. Program constructs serve as a lightweight representation that TargetFuzz understands and operates on, unifying different programming languages. The only cost to target a programming language with TargetFuzz is to define program constructs with respect to the grammar of the language, which we estimate to take a few hours—considerably cheaper than the effort of building custom test generators.

Our contributions are as follows:

(1) We highlight the limitations of fuzzing optimizing compilers using canonical pipelines, and propose targeted pass fuzzing to complement it. Our evaluation shows that targeted fuzzing can effectively reach optimizations that are otherwise under-tested.

(2) We are the first to 1) formalize a taxonomy of structural relations (composition styles) that compiler optimizations target, and to 2) propose automatically-synthesized mutations and crossovers that rebuild these relations, for the purpose of mutational fuzzing.

(3) We implement TargetFuzz that embodies this idea of composition-based mutational fuzzing. TargetFuzz by design supports both targeted fuzzing and pipeline fuzzing. We show in evaluation that TargetFuzz improves coverage by 8%, 11% and triggers optimizations 2.8×, 2.6×, compared to baseline fuzzers under targeted fuzzing mode. TargetFuzz achieves higher coverage than all grammar-based fuzzers, and comparable coverage to the language-specific custom fuzzer GrayC, in whole-pipeline fuzzing mode. In addition, we show that targeted fuzzing is complementary to pipeline fuzzing: it effectively tests all 37 sampled LLVM optimizations, while pipeline-fuzzing missed 12.

## 2 BACKGROUND

Compiler fuzzers typically use optimization pipelines (e.g., -O3) as the test harness, feeding test programs directly into the pipeline and observing for crashes and miscompilations. This is convenient to configure, but it systematically leaves significant portions of the optimization space untested for two reasons. This motivates the need for targeted pass fuzzing.

*(Challenge 1A) Phase ordering problem.* Optimizers typically implement a large number of optimization passes (>200 for GCC, LLVM, MLIR [3–5]). Contrary to superoptimizers [14, 35, 43] that aim to enumerate a large space of optimization sequences and select the optimal sequence, mainstream, production-quality optimizers build pipelines that are fixed, empirically-proven sequences of passes, in order to speed up compilation. Critically, within these heuristics-based pipelines, the order of passes can easily enable or preempt opportunities for one another, known as optimization interactions/coupling. For example, loop unrolling may remove opportunities for subsequent loop optimizations [38]. Decades of work show that a fixed, optimal order for all programs does not exist, and searching for a good sequence is a difficult problem [16, 25–27, 42, 47, 48, 51]. Consequently, no pipeline can exercise every transformation path: some opportunities are created only after specific prior passes; others are destroyed by them. In fuzzing, this manifests as certain cold optimizations that are rarely tested under the pipeline.

*(Challenge 1B) Optimization pipelines do not include all passes.* Modern compilers design their default pipelines to balance compile time and stability rather than "run every pass."

Zitong Zhou, Ben Limpanukorn, Hong Jin Kang, Jiyuan Wang, Yaoxuan Wu, Akos Kiss, Renata Hodovan, and Miryung Kim

```
1  int sum1=0, sum2=0;
2  for(int i=0;i<10;i++){
3      sum1+=i;
4  } // adjacent loops
5  for(int j=0;j<10;j++){
6      sum2+=j;
7  }
```

```
1  int sum1=0,sum2=0;
2  for(int i=0;i<10;i++){
3      sum1+=i;
4      sum2+=i; // body of loop 2 moved here
5  }
```

Fig. 1. A test program that triggers Loop Fusion and its optimized output.

Many transformations are off by default due to profitability or maturity concerns. For example, several loop transforms have historically been disabled by default, to avoid regressions until cost models and implementations mature [37]. This challenge is exacerbated as the compiler ecosystem is shifting from monolithic optimizers toward modular, extensible infrastructures like MLIR. It organizes IR into dialects that each captures a different abstraction and employs different optimizations [19]. It is widely adopted in machine-learning and domain-specific compilers such as CIRCT [2], Triton [1], Mojo [9], and OpenXLA [10]. By design there is no single canonical sequence that covers all dialects and optimizations. Compiler fuzzers of tomorrow must adapt to this modularity.

*(Challenge 2) Optimizations expect structural relations.* Besides the 2 challenges brought by whole-pipeline fuzzing, we identified another fundamental, orthogonal challenge: optimizations are guarded by a series of strict structural and semantic conditions, making them hard to test through fuzzing. In modern compilers, conceptually, each optimization pass consists of (1) an analysis phase that checks whether structural and semantic conditions are satisfied, and (2) a transformation phase that optimizes the program. We use **Loop Fusion** in LLVM as an example. Loop Fusion is an optimization pass that merges multiple loops into one, to reduce overhead in loop control structures and achieve better instruction-level parallelism [15]. It is important in data-intensive applications and machine learning compilers. In Figure 1, the optimizer attempts to apply Loop Fusion on the input program on the left. Loop Fusion iteratively checks that each condition is satisfied. **Adjacency** requires that the two loops have no instructions in between; **Conformance** requires that the two loops have the same trip counts. There are also other semantic conditions that must be satisfied such as profitability. *If any condition is unsatisfied, the pass stops early, and the critical transformation logic is not executed.* In Figure 1, the input program satisfies all conditions, producing the optimized program with a fused loop. In this paper, we refer to such structural relationships between program constructs such as loops as composition styles. We present the formal definitions in Section 3. While useful composition styles exist in fuzzing corpora, to the best of our knowledge, no existing fuzzers explicitly leverage them to test optimizations.

*Goal.* We set out to design a general-purpose fuzzer that exploits optimization-specific composition styles to effectively test optimizing compilers. The fuzzer must be lightweight: the user only needs to supply a small corpus from where the composition styles are extracted. It must also be easily adaptable across compilers and languages (e.g., MLIR compilers), so that porting requires only supplying a grammar plus concise annotations—no need to hand-code custom generators or mutators. It must support both targeted pass fuzzing and whole-pipeline fuzzing.

## 3 APPROACH

Figure 2 gives an overview of TargetFuzz's fuzzing pipeline. In each iteration, from a target-specific corpus program, TargetFuzz extracts a composition style that may trigger an optimization. Then, it synthesizes targeted mutators and applies them to a recipient program to rebuild the composition style.
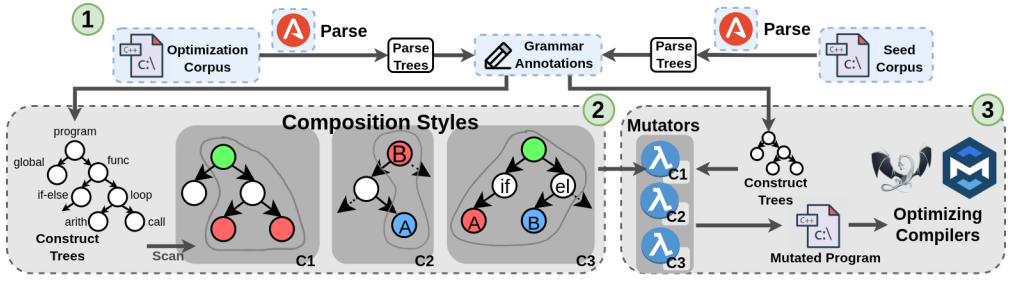
Fig. 2. Phase 1: TARGETFUZZ uses the supplied grammar to parse the optimization corpus and seed corpus into grammar parse-trees, then it uses supplied grammar annotations to translate them respectively into construct trees. Phase 2: TARGETFUZZ extracts composition styles, e.g., loops are adjacent, from the optimization construct trees - programs that contain 'breadcrumbs' of what's necessary to trigger optimizations. Phase 3: each composition style synthesizes mutators, e.g. replicating a loop. Mutators are applied on the seed construct trees to reconstruct compositions and test optimizations. Blue boxes are user inputs.

TARGETFUZZ takes *an optimization corpus, a seed corpus, a grammar, and grammar annotation as inputs.* The *optimization corpus* is where the composition styles are extracted. It consists of programs that are relevant to target optimizations. The assumption is that this corpus already has the breadcrumbs for composition styles necessary to trigger optimizations. It can be sourced from the compiler's unit tests, generated by Large Language Models (LLMs) or manually written by compiler developers. Note that TARGETFUZZ does not require the optimization corpus to be perfect, i.e., a program does not necessarily need to satisfy all conditions and trigger the optimization; it just needs to have the breadcrumbs. *The optimization corpus enables TARGETFUZZ to support two fuzzing modes in a single tool:* (1) targeted fuzzing mode for testing individual optimizations (e.g., -loop-fusion), where the optimization corpus only needs to contain a handful of programs that are relevant to the targeted optimization. This is particularly valuable for MLIR, e.g., to fuzz optimizations of the async dialect (abstractions of asynchronous execution), the optimization corpus can be async's unit tests; and (2) whole pipeline fuzzing to trigger all kinds of optimizations through an aggressive pipeline (e.g., clang -O3), when we pool optimization corpora. In evaluation, we show the two fuzzing modes are complementary to each other in terms of code coverage.

The *seed corpus* serves a different purpose. While the optimization corpus provides hints for exercising optimizations, the seed corpus provides the overall diversity of inputs. It is a generic corpus of programs that can be sourced from compiler benchmarks like SPEC CPU 2017 [17] or produced by existing program generators like CSmith [55], GrayC [20], YarpGen [32].

The *grammar* is a context-free grammar that defines the syntax of the programming language and is used to parse programs from both corpora. The *grammar annotation* (detailed in Section 3.1) extends the grammar by marking specific production rules that correspond to *program constructs* relevant to compiler optimizations (e.g., VECTOR, FUNCCALL), which may not be directly identifiable from the raw grammar rules alone. Program constructs also serve as a representation that unifies the syntax of different programming languages, enabling TARGETFUZZ to be language-agnostic.

## 3.1 Program Constructs

In Phase 1 of Figure 2, TARGETFUZZ uses an Ⓐ ANTLR context-free grammar to parse both the optimization and seed corpus [41]. Since ANTLR grammars are designed for generating parsers rather than pattern-matching, they do not always contain the production rules for semantic program constructs that compiler optimizations look for. For example, while the C grammar includes rules

```
/** C11 Spec */
grammar C;

addExpr
  : mulExpr (('+' | '-') mulExpr)*;
statement
  : labeledStatement
  | iterationStatement
  | ...;
iterationStatement
  : forStatement
  | whileStatement
  | doWhileStatement;
```
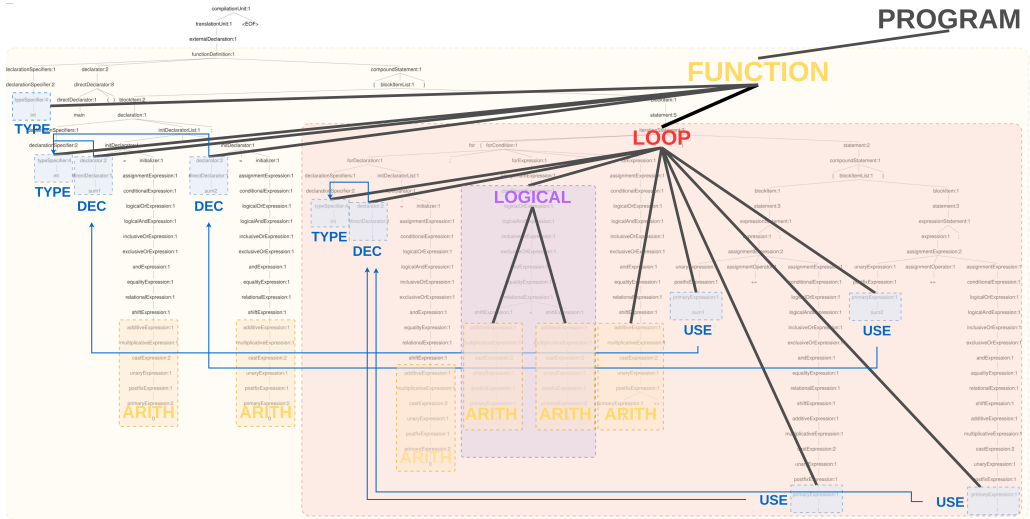
```
FOR_STMT_ = Construct(["forStatement"])
LOOPS_ = CmpdConstruct([FOR_STMT_, WHILE_STMT_, ...])
def isarith(node):
  return node.name == "addExpr" and \
    (any([n.src in ["+", "-", "*", ...] for n in
        lexers(node)]) \
    or is_numerical(node))
ARITH_EXPR_ = CustomConstruct(isarith)

def is_literal_array(node):
  ...
LTR_ARR_ = CustomConstruct(is_literal_array)
VECTOR_EXPR_ = CmpdConstruct([LTR_ARR_, ARR_ACES_])
```

(a) A snippet of the ANTLR C grammar with three production rules. '|' separates alternatives of each rule.

(b) TARGETFUZZ's Construct API enables defining custom program constructs that are relevant to compiler optimizations, through *grammar annotation*.



(c) Construct Tree of the optimized C program in Figure 1. Grammar parse-tree shown in the background. Colored boxes are nodes (constructs), bold lines are edges. Blue arrows are type-declaration-use chains.

Fig. 3. Listing 3a is C grammar. Listing 3b defines its program constructs. Listing 3c shows construct tree.

such as forStatement & ifStatement, that match for-loop and if-block, respectively, it lacks the rules for fine-grained constructs such as vectors, function calls, and memory references. The postfix expression rule postfixExpr can represent all three aforementioned constructs, but it does not distinguish them. This is problematic because compiler optimizations often operate on these fine-grained constructs.

TARGETFUZZ allows practitioners to *define custom program constructs by annotating any programming language's grammar* using its Construct API. For example, Listing 3b illustrates four program constructs: FOR_STMT_ (FORLOOP), LOOPS_ (LOOP), ARITH_EXPR_ (ARITHMETIC), and VECTOR_EXPR_ (VECTOR); the last two are not represented by rules in the grammar and need to be custom defined. FOR_STMT_ is defined as any node in the grammar parse-tree of the rule forStatement. LOOPS_ is a compound construct that is either a FOR_STMT_, WHILE_STMT_, or DO_WHILE_STMT_. ARITH_EXPR_ is a custom construct defined by the predicate function isarith. It is a addExpr grammar node (additive expressions) that is either numerical or contains arithmetic operators.

| Styles | Program Constructs Types $\langle \mathcal{T}_1, \ldots, \mathcal{T}_k \rangle$ | Context Types $\mathcal{T}_{\text{ctx}}$ | Mutation | Predicates |
|---|---|---|---|---|
| **Cousins** | $T_1 = T_2 \in \mathcal{T}$: {Loop, FuncCall, Arithmetic, Logical } | {If-Else, Loop, Func} | ALL | $k$:generation distance, $d$:token distance |
| **Nesting** | $T_1 = T_2 \in \mathcal{T}$: {Loop, If-Else } | {Loop, Func} | ALL | $d$:nesting depth |
| **Precedes** | $T_l, T_r \in \mathcal{T}$: {FuncCall, MemRef, Arithmetic} | {Loop, Func} | Ⓜ Ⓘ Ⓟ | |
| **Balanced** | $T_1 = $ If-Else; $T_2 = T_3 \in$ {Loop, FuncCall, Arithmetic, MemRef } | {If-Else, Loop, Func} | ALL | $d$:branch depth |
| **Sequence** | $T_1 = \cdots = T_n \in \mathcal{T}$: {Loop, FuncCall, Vector, Arithmetic } | {Loop, Func} | Ⓡ Ⓘ Ⓟ | $l$:seq. length |
| **Exists** | $T \in \mathcal{T}$: {Func, Loop, If-Else, FuncCall, MemRef, Vector, Arithmetic, ...} | {Vector, Loop, If-Else, Func, Program} | Ⓘ Ⓟ | $l$:min tokens |

Table 1. TargetFuzz's instantiated composition styles. For example, **Precedes** matches 2 constructs $T_l$ and $T_r$ of type FuncCall, MemRef, or Arithmetic, and a common ancestor construct (context) $T_{\text{ctx}}$ of type Loop or Func. It is equipped with three mutators: Ⓜove, Ⓘnsert, and Ⓡeplace (See Section 3.3).

Similarly, VECTOR_EXPR_ is a compound construct of LTR_ARR (literal array) and ARR_ACES_ (array access), both custom constructs defined using a predicate function.

Program constructs also serve as a lightweight, unified representation that enables TargetFuzz to be programming language-agnostic, e.g. construct ARITH_EXPR_ represents arithmetic expressions universally. Conceptually, as shown in Figure 2, TargetFuzz translates grammar parse-tree into a construct tree, where each node is a program construct. Figure 3c visualizes this step, the optimized C program in Figure 1 is parsed into a grammar parse-tree that is shown in the background, overlayed by the construct tree. Each colored box is a construct node, and bold lines are edges. Blue arrows are declaration-use chains that TargetFuzz tracks for scope well-formedness.

When adapting TargetFuzz to fuzz a new programming language $\mathcal{L}$, developers need to annotate $\mathcal{L}$'s grammar as in Figure 3b to map specific grammar parse-tree nodes to the corresponding internal constructs like ARITH_EXPR_. TargetFuzz defined common program constructs for C and MLIR, such as Loop, If-Else, FuncCall, MemRef, Vector, Arithmetic, and Logical, Jump, and many more, using 323 and 314 lines of grammar annotations, respectively. To employ TargetFuzz on a new language, practitioners need to annotate its grammar with relevant constructs, which typically means few hours of effort.

### 3.2 Composition Styles

*Intuition.* Many compiler optimizations fire only when specific structural relationships hold among program constructs (e.g., a loop nested in a loop, or sibling expressions sharing an ancestor). Our fuzzer targets such relationships directly by elevating them to first-class composition styles. A composition style specifies which types of constructs must co-occur, how they relate structurally, under what scope, and which mutations (Section 3.3) can re-create that structure.

*Definition (Composition Style).* Let $G$ be a construct tree and $N$ its constructs. Denote the type of construct $n$ as $T(n)$. A composition style is

$$C = \big(R_C, \ \langle \mathcal{T}_1, \ldots, \mathcal{T}_k \rangle, \ \mathcal{T}_{\text{ctx}}, \ p\big),$$

(i) $R_C \subseteq N^k$ is the $(k)$-ary structural relation targeted by $C$, relating $k$ constructs $\langle n_1, \ldots, n_k \rangle$. For example, we define Cousins to be a binary relation between two constructs that share a common ancestor.

(ii) $\mathcal{T}_i$ is the set of admissible types (e.g., Loop, FuncCall) of the $i$-th construct, i.e., each matched construct $n_i$'s type must satisfy $T(n_i) \in \mathcal{T}_i$. This constrains the matching of a structural relation to certain types known for exposing optimization opportunity. Some composition styles impose additional constraints on matched types: for instance, Cousins requires that
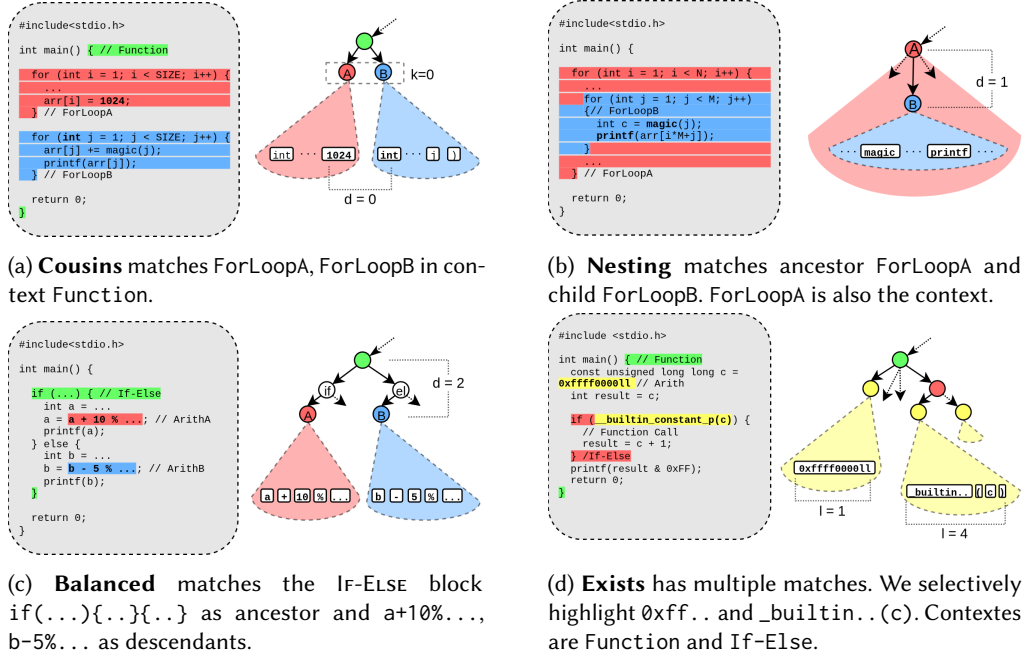
(a) **Cousins** matches ForLoopA, ForLoopB in context Function.



(b) **Nesting** matches ancestor ForLoopA and child ForLoopB. ForLoopA is also the context.



(c) **Balanced** matches the IF-ELSE block if(...){..}{..} as ancestor and a+10%..., b-5%... as descendants.



(d) **Exists** has multiple matches. We selectively highlight 0xff.. and _builtin..(c). Contextes are Function and If-Else.

Fig. 4. Examples of composition styles matching on C construct trees.

all matched constructs share the same type: $T(n_1) = T(n_2)$; while Precedes allows matched constructs of different types: $T(n_1) \neq T(n_2)$ is permitted.

(iii) $\mathcal{T}_{ctx}$ is the set of admissible types (e.g., IF-ELSE, FUNC) of the context $n_{ctx}$, i.e., it must satisfy $T(n_{ctx}) \in \mathcal{T}_{ctx}$. The context $n_{ctx}$ is the lowest common ancestor construct of all the matched constructs $\langle n_i \rangle$ and must "contain" them (defined in Section 3.3). It defines the <u>scope</u> of the structural relationship. For example, Cousins requires the context to be IF-ELSE, LOOP, or FUNC, ensuring that the matched cousins co-occur within these meaningful program structures rather than at arbitrary locations in the construct tree.

(iv) $p$ is a set of additional predicates on $n_1, \ldots, n_k$ that filter matches and constrain the structural relationship. These predicates are essential because the base structural relations are often too broad, leading to trivial or uninteresting matches. For example, without predicates, Cousins would match <u>any</u> two loops in a function since they technically share the function body as a common ancestor. To address this, each composition style defines specific predicates: Cousins bounds the generational distance between constructs to ensure meaningful proximity, while Balanced limits the branch depth explored under an IF-ELSE construct.

*Operationalization.* TARGETFUZZ implements each style $C$ with a scanning routine that programmatically defines its structural relation $R_C$:

$$C::\text{scan}(n) \rightarrow m; \qquad m : \langle C, \langle n_1, \ldots, n_k \rangle, n_{ctx}, p \rangle$$

that traverses the subtree at root node $n$ and returns a <u>match</u> $m$—a witness that the style $C$ holds at constructs $\langle n_i \rangle$ under context $n_{ctx}$ with predicates $p$ satisfied. Concretely, the returned tuple satisfies the style's relation $R_C$, the type constraints $\langle \mathcal{T}_i \rangle$ and $\mathcal{T}_{ctx}$, and the predicates $p$. For each match, the fuzzer gathers the minimal $n_{ctx}$ that encloses the declarations and uses needed for scope well-formedness. If no such tuple exists inside subtree $n$, C::scan returns no result; In Phase

2 (Figure 2), TARGETFUZZ runs C::scan on all construct trees from the optimization corpus to extract matches.

**Instantiating composition styles to extract optimization opportunities.** We extracted from *seminal works* on optimizing transformations [11, 12] and *compiler documentation* [4, 5] a set of recurring structural preconditions that commonly trigger optimizations. Then, *we instantiated six concrete composition styles* with fully specified relations, types, and predicates. Table 1 shows the registry of these concrete styles. For brevity, we highlight four representative ones and illustrate how they capture common optimization-triggering conditions.

*Cousins.* Matches two constructs of the same type—LOOP, FUNCCALL, ARITHMETIC, or LOGI-CAL—that share the same ancestor. Optional parameters bound proximity: maximum generation distance $k$ and maximum token-distance $d$. Figure 4a shows ForLoopA/ForLoopB with $k = 0$, $d = 0$. Captures redundancy from loop-level (**Loop Fusion** with the paired style in Fig. 1) to instruction-level peepholes (**InstCombine**: adjacent add/sub, shifts, casts, comparisons).

*Nesting.* Matches same-type ancestry among control constructs—nested LOOP and IF-ELSE. Optional parameter $d$ bounds nesting depth (tree distance between the matched nodes). Figure 4b illustrates a match. Targets control-flow–simplifying optimizations such as **Loop Flatten**, **Loop Interchange**, **Loop Unroll and Jam**, and **Constraint Elimination**.

*Balanced.* Matches an ancestor IF-ELSE with two descendant constructs of the same type—LOOP, FUNCCALL, ARITHMETIC, or MEMREF—appearing in different branches. Optional parameter $d$: maximum branch depth explored under the IF-ELSE. Figure 4c shows similar arithmetic in two branches, making them potential candidates for hoisting. Captures replicated work across guarded paths; common targets include **GVN** (hoisting/sinking equivalent computations) and **Jump Threading** (simplifying conditionals around calls).

*Exists.* Matches a single construct drawn from Table 1 of type: FUNC, LOOP, IF-ELSE, FUNCCALL, MEMREF, VECTOR, ARITHMETIC. Minimum lexer nodes $l$ limits the match to nodes of at least $l$ tokens. Useful to constrain fuzzing scope (e.g., scalar vs. loop) and surface magic numbers. Certain rare-typed constants are known to trigger issues in **Constant Propagation** and **Constant Hoisting**. Figure 4d highlights two matches (ARITHMETIC, FUNCCALL).

### 3.3 Mutators

Having extracted composition styles (matches) from donor programs, TARGETFUZZ synthesizes mutators to apply on programs from the seed corpus (Phase 3). By reconstructing composition styles within different program contexts, TARGETFUZZ aims to trigger *variations* of optimization logic. TARGETFUZZ implements 4 high-level program transformations: 2 mutations that rebuild structural relations using only recipient material, and 2 crossovers that transplant donor fragments into the recipient.

**Mutator.** A mutator $\mu$ is a generator with signature $\mu_m(P) \Rightarrow P'$, where $m = \langle C, \langle n_1, \ldots, n_k \rangle, n_{\text{ctx}}, p \rangle$ is a match extracted from a donor construct tree in the optimization corpus, and $P$ is a recipient construct tree in the seed corpus.

Given a match $m$ and a recipient $P$, a mutator $\mu$ performs these steps:

Ⓡ **Replicate &** Ⓜ **Move** — Mutate within the recipient (no donor construct inserted).

(1) **Partialization.** Remove one matched construct $n_i \in \langle n_1, \ldots, n_k \rangle$ from context $n_{\text{ctx}}$ to obtain a partial context $n'_{\text{ctx}}$ and an anchor set $A = \langle n_1, \ldots, n_k \rangle \setminus \{n_i\}$.

(2) **Context match.** i) Find candidate context constructs in $P$ with same type as partial context $n'_{\text{ctx}}$ and rank by structural similarity to $n'_{\text{ctx}}$ (KLR similarity [31]), then pick the most similar

recipient context $n_{\text{ctx}}^P$. ii) structurally match anchor constructs to constructs in the recipient context, $A \mapsto A^P$, that is, $(n_a \mapsto n_a^P, \dots)$. iii) Identify a <u>location</u> $\ell$ within $n_{\text{ctx}}^P$ corresponding to the missing construct $n_i$.

(3) **Transform (this distinguishes the 2 mutations).**

Ⓡ **Replicate.** With anchors $A^P$ and location $\ell$, choose some $a^P \in A^P$ and <u>clone</u> it into $\ell$ to re-create the missing counterpart. Applicable when composition style requires that all matched constructs share the same type (Table 1: Cousins; Nesting; Sequence; descendants in Balanced).

Ⓜ **Move.** Find all constructs $x$ in the <u>entire</u> recipient tree with $T(x) = T(n_i)$; rank candidates by lexical edit distance to the removed $n_i$; <u>move</u> the best-scoring $x$ to $\ell$. Applicable for both same-type styles and "heterogeneous" styles that don't require matched constructs to have the same type (Table 1: Precedes).

---

Ⓘ **Insert &** Ⓟ **Replace** − Crossover (transplant) donor constructs into the recipient

(1) **Context match.** Same as Step 2 above, except that no partialization is needed; match the full context $n_{\text{ctx}}$ and all constructs $\langle n_1, \dots, n_k \rangle$ to $n_{\text{ctx}}^P$ and $A^P = \langle n_1^P, \dots, n_k^P \rangle$.

(2) **Transform (this distinguishes the 2 crossovers).**

Ⓘ **Insert.** After matching $n_{\text{ctx}}^P$ and $A^P$, <u>insert</u> each donor construct $n_i$ next to its match $n_i^P$.

Ⓟ **Replace.** After matching $n_{\text{ctx}}^P$ and $A^P$, <u>replace</u> each $n_i^P$ with donor construct $n_i$.

*Parameterized mutation.* Edits—especially crossovers—can easily violate scope or typing (e.g., transplanting an expression $a + b$ into a context that lacks declarations for $a, b$). We therefore perform parameterized mutation [31]: every edit is re-parameterized against the recipient context so that free uses are bound to in-scope declarations and types are compatible. That is, the contexts we extract during style scanning (Section 3.2) and during context matching (Step 2 in Section 3.3) are required to satisfy a declaration–use <u>containment</u> property. We define that a context $n_{\text{ctx}}$ contains constructs $\langle n_1, \dots, n_k \rangle$, if all <u>uses</u> appearing in the constructs are resolvable to <u>declarations</u> within $n_{\text{ctx}}$ and types must be compatible. This is enabled by grammar annotations introduced in Section 3.1. Special constructs USES_, DEFS_, TYPES_ must be defined for TARGETFUZZ to identify declaration-use relations and support parametrized mutation.

When cloning an anchor $n^p \in A^P$ into $\ell$ (Replicate) or moving a construct $x$ into $\ell$ (Move), the mutator rewrites the construct's USES_ to names drawn from DEFS_ that are visible in $n_{\text{ctx}}^P$ and type-compatible per TYPES_. When inserting or replacing with donor constructs $n_i$, the mutator re-parameterizes the donor's USES_ to bind to DEFS_ visible in $n_{\text{ctx}}^P$ (and checks TYPES_). If no compatible binding exists, the edit is rejected. Parameterized mutation substantially increases edit validity rates while keeping edits local to the matched context.

*Example.* Figure 5 illustrates an end-to-end application of Ⓡ REPLICATE for the Cousins style. In the donor (left), the scanner identifies a match $m = \langle \text{Cousins}, \langle L_1, L_2 \rangle, F, p\{k = 0, d = 0\}\rangle$: two same-type loops $L_1$ (red) and $L_2$ (blue) enclosed by the function context $F$ (green). For simplicity, suppose the names SIZE, arr, magic are all declared in the function scope, i.e. $F$ "contains" the constructs $L_1, L_2$. <u>Partialization</u> removes $L_2$, retaining $L_1$ as the anchor and yielding a hole for the missing cousin. In the recipient (middle), context matching selects the most similar function $F'$, locates the anchor loop $L_1'$, and determines a position $\ell$ corresponding to the removed cousin $(F \mapsto F', L_1 \mapsto L_1', L_2 \mapsto \ell)$. REPLICATE then clones $L_1'$ into location $\ell$, producing the mutated program and, via parameterized mutation, rebinds free USES_ of aggregator sum1 within $L_1'$ to in-scope DEFS_, producing a second loop that updates sum2. Note that i is USES_ but it is declared
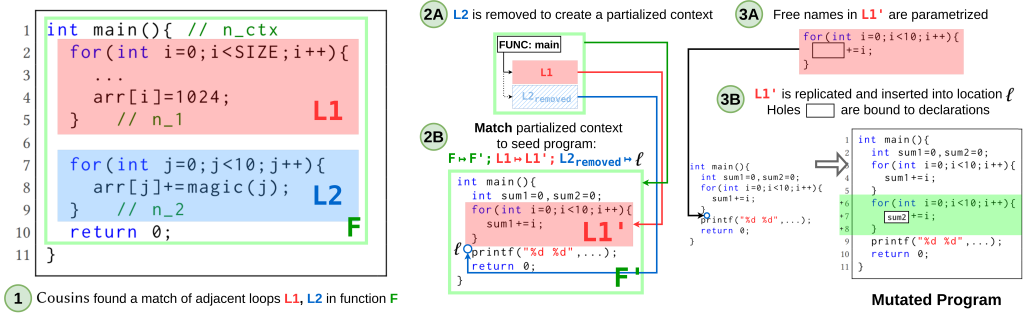
Fig. 5. A working example of TargetFuzz. First, the Cousins style matches two adjacent loops $L_1, L_2$ within function context $F$ in the donor; Partialization removes $L_2$ to create a partialized context. Then, the partialized context and its constructs are mapped to structurally similar constructs in the seed program: $F \mapsto F', L_1 \mapsto L_1', L_2$(removed) $\mapsto \ell$. Next, **Replicate** clones the matched segment $L_1'$, inserts it at location $\ell$, and substitutes the free variable name **sum1** with **sum2**.

within the loop construct itself, so there is no need to rebind it. The transformation reconstructs the Cousins relation in the recipient and exposes a canonical loop-fusion opportunity.

## 4 EVALUATION

In our study, we examine the following research questions:

**RQ1:** How effective is TargetFuzz at testing optimizations compared to other fuzzers, in terms of coverage?

**RQ2:** How does TargetFuzz's targeted fuzzing mode compare to whole pipeline fuzzing mode?

**RQ3:** What are TargetFuzz's strengths & limitations at triggering optimizations of different characteristics?

**RQ4:** What are each composition style and mutator's bug finding capability?

**RQ5:** TargetFuzz's bug findings.

We implemented TargetFuzz on top of Grammarinator in 5360 lines of code, including 323 and 314 lines of definitions of C (Figure 3b) and MLIR language constructs, respectively. Annotating grammar is a one-time effort, and we estimate that it would take around two hours, for someone experienced with grammar-based fuzzing to define the program constructs for a new language.

### 4.1 Experiment Design

*4.1.1 Compilers and optimizations under test.* We evaluate TargetFuzz on 2 optimizing compilers, namely, LLVM and MLIR. LLVM is a production-quality optimizing compiler for languages such as C/C++. It contains the most compiler optimizations, and active optimization research is usually first committed there. MLIR is a modular compiler framework that is gaining popularity as a platform for deep learning compilers. Unlike LLVM which uses a single, monolithic intermediate representation (IR), MLIR allows developers to extend its IR with custom operations, types, and optimizations, which are logically grouped into *dialects* such as vector, affine, gpu. MLIR's optimizations are mostly dialect-specific. We picked MLIR to demonstrate that TargetFuzz's approach generalizes beyond high-level imperative programming languages like C.

As explained in Section 3, TargetFuzz supports both fuzzing entire optimization pipelines and targeted fuzzing of individual optimizations. We evaluate TargetFuzz on both modes for LLVM

and only targeted fuzzing of individual optimizations for MLIR, because MLIR core dialect optimizer doesn't quite have an encompassing pipeline. For whole pipeline fuzzing of LLVM, we run `clang -O3` on TARGETFUZZ's mutated programs for 24 hours. Using this traditional evaluation setting [20, 30] ensures a fair comparison. For targeted fuzzing of individual optimizations, we run each optimization pass for one hour in RQ1 and four hours in RQ2. The shorter fuzzing budget is because individual optimization is small (mostly <10k LOC).

To "target-fuzz" individual optimizations, we need to perform sampling since technically, LLVM and MLIR each implements over 200 optimization passes. For each sampled optimization, we craft a mini "preparation" pipeline, e.g., run `-loop-simplify,-loop-rotate` before loop optimizations. We sampled 37 (LLVM) and 26 (MLIR) optimizations. We refer to them as the sampled optimizations. Our selection covers 4 LLVM optimization categories (Scalar, Loop, Instcombine, Interprocedural) and 7 MLIR core dialects: `affine`, `vector`, `arith`, `scf`, `async`, `gpu`, `memref`. We believe the diversity is enough and comparable to the literature on compiler optimization testing [54].

*4.1.2 Optimization and Seed corpus.* From the sampled optimizations, we distill the source code of each optimization pass into metadata, function declaration, and docstrings. Then we prompt GPT4o-mini to collect 100 valid programs per optimization pass as the its own optimization corpus. We collect two seed corpus of 1059 C programs and 3601 MLIR programs from the literature and compiler test suites [20, 31, 55]. In targeted fuzzing mode, we use each pass's own optimization corpus (100). In pipeline fuzzing mode, we pool all sampled optimizations' corpus (3700).

*4.1.3 Evaluation Metrics.* We use the following metrics:

- **Branch Coverage** We built the optimizing compilers (`clang`, `opt` & `mlir-opt`) with coverage instrumentation [8].
- **Optimization Trigger Throughput** is the total number successful program transformations triggered/applied in four hours. This is a standard evaluation metric in the domain of fuzzing compiler optimizations [53, 54, 56]. For example, LoopFuse's #transforms is #loops fused, and GVNHoist's #transforms is #instructions hoisted or removed.

*4.1.4 Baselines.* We evaluate TARGETFUZZ against Grammarinator [21], SynthFuzz [31], GrayC [20], and MLIRSmith [49]. To ensure fairness, we *combine* the optimization and seed corpus to use as the fuzzing corpus of Grammarinator, SynthFuzz, and GrayC. Grammarinator represents a baseline grammar-based fuzzer. SynthFuzz is a grammar-based fuzzer that synthesizes custom mutations that preserve semantic validity. We pick SynthFuzz because it is 1) novel in synthesis of grammar mutation, and 2) shown to be the state-of-the-art MLIR fuzzer [31]. GrayC is a coverage-guided, mutational fuzzer targeting the C language. It implements a set of semantic-aware, custom mutators, such as replacing an arithmetic operator with another. We pick GrayC because it's the state-of-the-art C Compiler fuzzer [20]. We also compare against GrayC's variant without coverage guidance in RQ1 as it uses the coverage as metric. MLIRSmith is a custom, generator-based fuzzer that targets MLIR's core dialects. It requires developers to manually implement custom generator for each dialects. We pick MLIRSmith as a baseline for generational MLIR fuzzers.

## 4.2 RQ1: Is TARGETFUZZ effective at testing optimizations compared to other fuzzers?

For whole optimization pipeline fuzzing, we run TARGETFUZZ and baseline fuzzers for 24 hours on the `clang -O3` pipeline, then measure branch coverage of the entire compiler (directories under `clang/` and `llvm/`). For targeted fuzzing of individual optimizations, per optimization, we run each fuzzer for one hour and measure coverage only within that specific optimization's code.
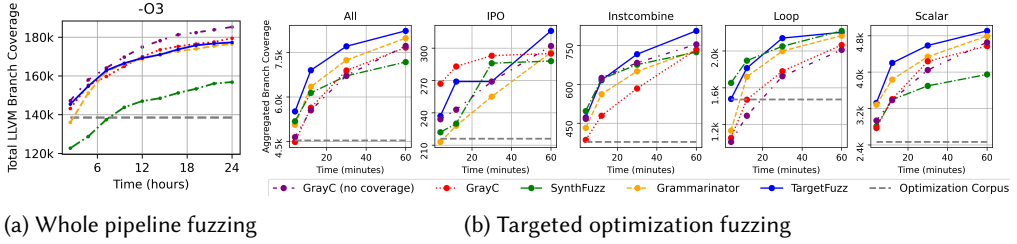
(a) Whole pipeline fuzzing                    (b) Targeted optimization fuzzing

Fig. 6. **RQ1**: (a) Coverage on entire LLVM codebase when fuzzing with -O3 pipeline for 24 hours. (b) Coverage on targeted optimizations, aggregated across 4 categories. Left to right: All optimizations, Interprocedural (IPO), InstCombine, Loop, and Scalar optimizations. In the whole pipeline fuzzing mode, TARGETFUZZ achieves lower coverage than GrayC(w/wo) coverage guidance when fuzzing the entire compiler; however, in the targeted fuzzing mode, TARGETFUZZ outperforms all baseline fuzzers on IPO, InstCombine, and Scalar families of optimizations. We show that the two fuzzing modes are complementary to each other in RQ2.
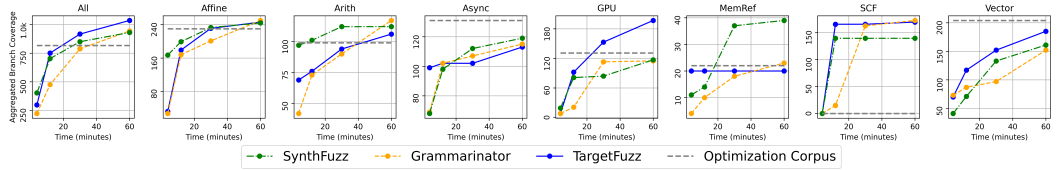


Fig. 7. **RQ1**: Coverage on targeted MLIR optimizations, summed across 7 categories. Left to right: All optimizations, Affine, Arith, Async, GPU, Memref, SCF, Vector. TARGETFUZZ achieves highest overall coverage and in dialects gpu and vector, while having comparable (<2%) coverage in affine, scf as the best baseline.

On LLVM pipeline -O3, GrayC [1] fuzzer achieves the best overall coverage, followed by TARGET-FUZZ, Grammarinator, then SynthFuzz, as shown in Figure 6a. This result aligns with expectations: GrayC is specifically designed for C/C++ programs with hand-coded mutations that operate on semantic-rich, typed ASTs to ensure high test validity. This explains why it is the most effective at testing C compiler as a whole. However, we argue that TARGETFUZZ is built with two distinctive goals: (1) unlike GrayC, TARGETFUZZ is a language-agnostic fuzzing framework that can be easily adapted to new languages (like MLIR), and (2) it is designed to be more effective at testing individual optimizations, which we will show in the rest of RQ1 and RQ2.

Notably, TARGETFUZZ achieves the highest coverage among grammar-based fuzzers. Upon a closer look, TARGETFUZZ outperforms Grammarinator a lot at the beginning of the fuzzing campaign: at the 2.5 hour mark, TARGETFUZZ achieves 7% higher coverage, but this advantage levels out in the long run. After 24 hours, TARGETFUZZ is only marginally better than Grammarinator. This is not surprising given the fact that Grammarinator's general mutations indeed *subsume* TARGETFUZZ's targeted, constrained mutators, that they eventually converge to similar coverage. However, TARGETFUZZ demonstrates an advantage in short fuzzing runs.

Under the targeted fuzzing setting, in LLVM, TARGETFUZZ achieves 3.0%, 14.7%, 7.3%, 6.5%, higher total coverage than Grammarinator, SynthFuzz, GrayC, and GrayC (without coverage), aggregating all sampled optimizations. To show generality, we further group the 37 LLVM optimizations into four categories: Scalar, Loop, Interprocedural (IPO), and InstCombine and report the aggregate coverage in Figure 6b. TARGETFUZZ achieves the highest coverage on InstCombine, IPO, and Scalar optimizations. On Loop optimizations, SynthFuzz achieves 0.7% higher coverage than TARGETFUZZ.

---

[1]GrayC wo. coverage achieves 3% higher coverage than with coverage, consistent with GrayC's own evaluation results [20]
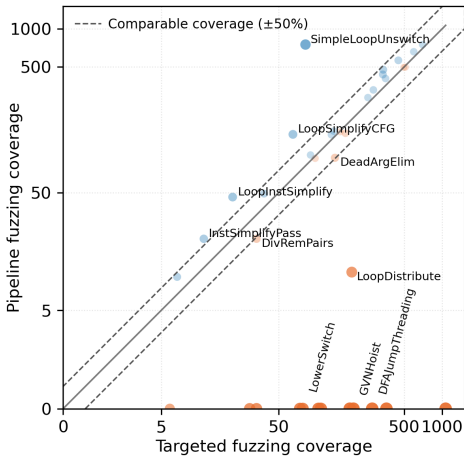
14  Ztong Zhou, Ben Limpanukorn, Hong Jin Kang, Jiyuan Wang, Yaoxuan Wu, Akos Kiss, Renata Hodovan, and Miryung Kim



Fig. 8. **RQ2**: TARGETFUZZ's coverage on sampled optimizations in targeted mode vs. pipeline mode
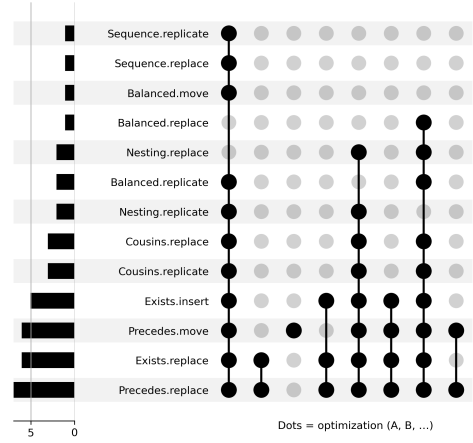
Fig. 9. **RQ4**: buggy optimizations flagged by each composition & mutator.

In MLIR targeted fuzzing, TARGETFUZZ achieves 9.7% and 11.2% higher coverage than Grammarinator and SynthFuzz, respectively. We evaluated 26 optimizations across 7 MLIR dialects. As shown in Figure 7, TARGETFUZZ achieves the highest aggregated branch coverage across all optimizations. Specifically, TARGETFUZZ achieves best performance on the `gpu` and `vector` dialects and comparable coverage on the `affine` and `async` dialects (<margin of 2%). MLIRSmith achieves near-zero coverage; we discuss the reasons for this in Section 4.4.

> TARGETFUZZ achieves the highest overall coverage, in both LLVM and MLIR, under the targeted fuzzing setting, where it performs consistently better than other language-specific fuzzers, and grammar-based fuzzers. Under the whole pipeline fuzzing setting, TARGETFUZZ is inferior to GrayC (without coverage) and perform similarly to GrayC with coverage.

### 4.3 RQ2: How does TARGETFUZZ's targeted fuzzing mode compare to whole pipeline fuzzing mode?

In RQ1, we showed that under the whole pipeline fuzzing setting, TARGETFUZZ is comparable to GrayC, but inferior to GrayC without coverage, in terms of overall coverage of the LLVM compiler codebase. This is the coverage experiment that most compiler fuzzers carry out. But this result doesn't necessarily translate to effectiveness at testing optimizations. In particular, we would like to answer the question: *Is there a need for targeted fuzzing of individual optimizations, that it complements what whole pipeline fuzzing cannot achieve?*

First, we take a closer look at coverage of LLVM's optimization module under this mode. The optimization module under `llvm/lib/Transforms` contains 300 files, including all the transformation passes and their utilities. The best-performing baseline GrayCnoCov achieves 38596 branches, while TARGETFUZZ achieves 36377 branches, a 6.1% difference. GrayCnoCov achieves higher coverage on 55 files while TARGETFUZZ achieves higher coverage on 36 files. However, *152 of the 300 files have zero coverage* from both fuzzers. This supports our hypothesis that the whole pipeline fuzzing approach leaves a lot of optimization logic untested.
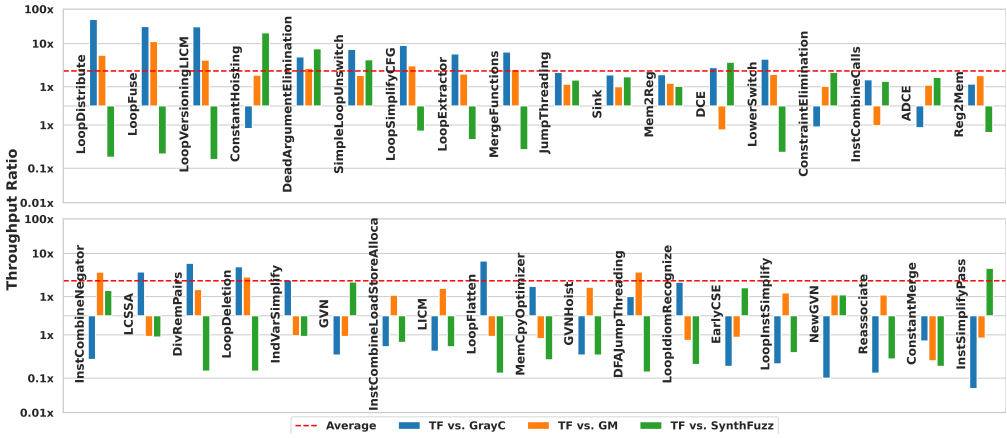
Fig. 10. **RQ3**: *LLVM* Ratio of Optimization Trigger Throughput. Each bar represents TARGETFUZZ's throughput divided by baseline's throughput. Upward bar means TARGETFUZZ outperforms the baseline, and vice versa. The red line is the average ratio: 2.8×. On average, TARGETFUZZ's trigger throughput outperforms GrayC, Grammarinator, and SynthFuzz by 5.14×, 1.74×, 1.64×.
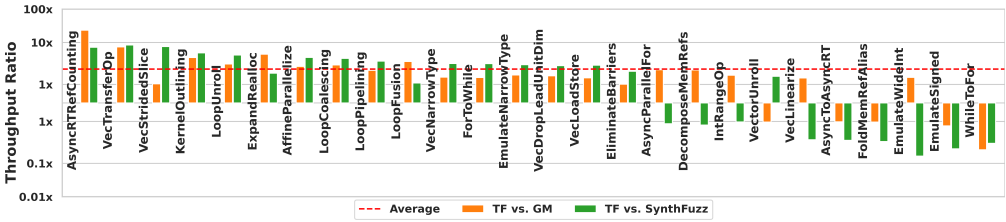


Fig. 11. **RQ3**: *MLIR* Ratio of Optimization Trigger Throughput. TARGETFUZZ's outperforms GM and SynthFuzz by 2.71×, 2.45×.

We further compare TARGETFUZZ's coverage in the two fuzzing modes. It is important to note that these two modes use different test harness (`clang -O3` vs. individual pass) and subsequently different fuzzing budget (24 hours vs. one hour). Thus, the numbers are not directly comparable, but we can still draw insights from the comparison. We show the coverage of sampled optimizations in the two fuzzing modes in Figure 8. The band around the diagonal indicates that for most optimizations, the coverage in the two modes are comparable. Notably, there are 12 passes that are entirely untested in the whole pipeline mode, but all of them are covered in targeted fuzzing mode.

> Targeted fuzzing of individual optimizations is necessary to complement whole pipeline fuzzing, as it significantly increases the number of optimizations tested.

## 4.4 RQ3: What are TARGETFUZZ's strengths and limitations at triggering optimizations of different characteristics?

To provide a more detailed assessment of TARGETFUZZ's effectiveness in targeted fuzzing, we use the optimization trigger throughput metric. Widely used in prior work [53, 54, 56], it measures how frequently a fuzzer successfully triggers optimizations. In Figure 10 and Figure 11, we show TARGETFUZZ's throughput improvements over the baseline fuzzers on all optimizations.

TARGETFUZZ outperforms Grammarinator by achieving higher trigger throughput on 26 / 37 LLVM and 22 / 26 MLIR optimizations, averaging 1.74× and 2.71× throughput of Grammarinator.

Because TARGETFUZZ is implemented on top of Grammarinator, these numbers directly reflect the gain of TARGETFUZZ's composition-based, targeted mutations over naive grammar-based fuzzing, which is less effective at exercising deep optimization logic.

TARGETFUZZ achieves 1.64× (LLVM) and 2.45× (MLIR) throughput of SynthFuzz. TARGETFUZZ outperforms SynthFuzz on 18 / 26 MLIR optimizations but SynthFuzz outperforms TARGETFUZZ on 22 / 37 LLVM optimizations. SynthFuzz is a MLIR-inspired grammar-based fuzzer, leveraging parameterized mutation to transplant grammar subtrees to a different context [31]. TARGETFUZZ's mutations are more constrained to specifically target optimizations. In addition, program constructs allow TARGETFUZZ to identify declaration, usage, and types, so TARGETFUZZ's parameterized mutation achieves higher validity (See Section 3.3). TARGETFUZZ produces C programs of 41.14% validity compared to 19.18% of SynthFuzz.

MLIRSmith triggered zero MLIR optimizations, so we omit it in Figure 11. The problem with MLIRSmith's custom generator approach is twofold. MLIRSmith's programs have an extremely low validity rate of 0.8%. This is because MLIRSmith hardcoded the semantics of MLIR operations. MLIR's constant evolution easily renders such semantic specifications obsolete. Secondly, custom generators like MLIRSmith are built to maximize extensibility, i.e., being made versatile to target many language features. They are blackbox fuzzers that are unable to take advantage of domain knowledge (e.g. the optimization corpus) about the compiler under test, which has proven to be essential in testing deep compiler logic.

TARGETFUZZ outperforms GrayC on 23 / 37 LLVM optimizations, averaging 5.14× throughput of GrayC. However, TARGETFUZZ's comparison against GrayC shows high variance: throughput ratio ranges from 50× to 0.05×. Our analysis shows that TARGETFUZZ is better at triggering more constrained optimizations: for example, it outperforms GrayC on LoopFusion, DeadArgumentElimination, MergeFunctions, by 31.2×, 4.1×, and 5.6×, respectively. GrayC performs better when the optimization is relatively "low-effort," i.e., has trivial triggering conditions. For example, GrayC outperforms TARGETFUZZ the most in InstSimplifyPass by 20×. It packages simple optimizations that remove redundant instructions, so any unoptimized C program have a high chance of triggering it. GrayC outperforms TARGETFUZZ by 10× on NewGVN (Global Value Numbering). This optimization removes equivalent expressions [6]. GrayC happens to implement a *Duplicate-Statement* mutator [20], which duplicates a statement within the same block, easily triggering NewGVN.

> TARGETFUZZ is more effective than naive grammar-based fuzzer (Grammarinator), parametrized mutators (SynthFuzz) and language-specific fuzzer (GrayC) at triggering compiler optimizations. Compared to GrayC, TARGETFUZZ is particularly effective at triggering more constrained optimizations.

### 4.5 RQ4: What are each composition style and mutator's bug-finding capabilities?

To evaluate the bug-finding capabilities of each composition style and mutator, we run TARGET-FUZZ for 6 hours per sampled LLVM optimization, focusing on miscompilations—the most severe optimization bugs. We use Alive2 [34] as the oracle to validate whether a triggered transformation is a miscompilation. For each {composition style, mutator} pair, we report # distinct optimizations produced at least one Alive2-flagged miscompilation. Results appear in Figure 9. In total, 8 optimizations are flagged for miscompilation. The most successful compositions are Precedes and Exists, and 5 optimizations are only discovered by them. It is expected that crossovers such as Exists+Replace and Exists+Insert are effective, as they are the most generic form of grammar-based mutations. However, TARGETFUZZ's more constrained composition styles and mutations are also

effective at exposing buggy optimizations. *Notably, no single style-mutator pair alone flags all 8 buggy optimizations*, proving the necessity of multiple composition styles and mutators.

> TARGETFUZZ's multiple styles and mutators are necessary to trigger a diverse set of optimization bugs.

## 4.6 Bug Findings

TARGETFUZZ found 18 previously unknown bugs in LLVM and MLIR, shown in Table 2. Among them, 10 are optimization bugs and 2 are backend bugs; 4 are miscompilations and 14 are crashes. For reference, the optimization corpus itself only exposes one bug in LLVM, proving the effectiveness of TARGETFUZZ's mutation. We used two bug-finding oracles: 1) compiler crash (LLVM & MLIR) and 2) Alive2 [34] (LLVM). We also found 2 false-positive bugs in Alive2 and both have been fixed.

| Status | LLVM | MLIR |
|---|---|---|
| Reported | 2 | 2 |
| Confirmed | 5 | 0 |
| Fixed | 2 | 5 |
| NF | 2 | 0 |
| Total | 11 | 7 |

Table 2. Bug reports; NF: not-gonna-fix.

```
1  func.func @extract_1d_constant()->(){
2    %cst = arith.constant dense<[1,2]>:
        vector<2xi32>
3    %0 = vector.extract %cst[0]: i32 fr
        vector<2xi32>
4    func.return
5  }
```

Listing 1. TARGETFUZZ program crashes **VectorLinearize**

```
1  int main() {
2    int sum1=0,sum2=0;
3    for (int i=0;i<10;i++){
4      sum1 += i;}
5    for (int i=0;i<10;i++){
6      sum2 += i*i;}
7  }
```

Listing 2. TARGETFUZZ program crashes **Loop Fuse**

We present a TARGETFUZZ-found MLIR bug that was fixed. vector.extract is an operation that extracts a $n - D$ vector at a $k - D$ position. **VectorLinearize** is an optimization pass that attempts to linearize all $n - D$ vectors to $1D$. TARGETFUZZ generated the program in Listing 1 that crashed VectorLinearize, which incorrectly assumes that $k < n$ in all vector.extract operations, i.e., it must return vectors, and attempts to linearize the vector.extract on line 3 which in fact produces a scalar value of i32. TARGETFUZZ crafted this test program using the Exists style and **Insert** mutator. First, TARGETFUZZ's identified that VECTOR constructs widely exists in the optimization corpus, matching a vector.extract operation. It then parameterizing its input argument. Upon matching the recipient context, the hole is concretized with the most structural similar node %cst. This mutation shows (1) how TARGETFUZZ uses optimization corpus to extract useful breadcrumbs (e.g., vector.extract) of the targeted optimizations, (2) parameterized mutation is critical in maintaining validty during crossover.

We present a TARGETFUZZ-found LLVM optimizer bug that was fixed. Optimization **LoopFuse** crashes on the program in Listing 2. The two loops are fusable but LoopFuser::performFusion incorrectly asserts that the loops are illegal for fusion and crashes. This bug is both a crash and a missed optimization.

## 5 RELATED WORK

While fuzzing is effective at exploring behaviors of code, a drawback is that it may only execute shallow code. As such, it will struggle to test deep code, such as code related to optimization.
**Grammar-based fuzzing.** Grammar-based fuzzers, including Grammarinator [21], Nautilus [13], LangFuzz [22], Skyfire [50], PolyGlot [18], and SynthFuzz [31] take an input grammar to generate syntactically correct programs and mutate on parse trees. Generator-based fuzzers [40] often rely on human effort to encode language-specific semantics. Even if semantically-valid inputs are generated, they do not necessarily satisfy the constraints for reaching deep code, such as

optimizations. TARGETFUZZ, instead, takes advantage of composition styles, which allows it to attain higher coverage of compiler optimization code.

**LLM-based fuzzing.** Fuzz4All [52] was the seminal work to leverage LLMs to both generate and mutate tests programs for fuzzing compiler. MetaMut [39] automatically synthesizes semantics-aware mutators using LLMs through an elaborate, multi-stage prompting strategy. TARGETFUZZ doesn't rely on LLMs for mutation. It only asks LLMs for an initial corpus of optimization tests, from where the mutations are derived through composition-styles and applied on a larger seed corpus at scale. This approach significantly reduces the number of slow and expensive LLM queries. In evaluation, TARGETFUZZ's LLM-generated optimization corpus only exposes 1 bug; TARGETFUZZ's bug findings come from its composition-style based mutations.

**Semantic-aware fuzzing.** Some compiler fuzzers aim to more effectively test deep code by ensuring semantic well-formedness. MLIRSmith [49] encodes the semantics of MLIR dialects but cannot keep up with its constant evolution. YarpGen [32] generates C/C++ programs that avoid undefined behaviors by static analysis and conservatively replacing unsafe operations. YarpGen v2 [33] added UB-free loop generation to target loop optimizations. GrayC [20] hand-coded a small set of mutators on semantic-rich, typed-AST to guarantee program well-formedness. Ensuring semantic well-formedness is an orthogonal problem; TARGETFUZZ focuses on satisfying the structural constraints necessary for testing compiler optimizations. Our evaluation shows that structural constraints is easier to achieve, i.e., without the need for custom, semantic generators/mutators, and are more important than well-formesness in testing optimizations.

In theory, such constraints can be specified explicitly using ISLa [44]: a specification language for generating structured inputs. However, doing so one-at-a-time for each optimization, is prohibitively labor-intensive. TARGETFUZZ bypasses this problem by automatically extracting composition styles from the optimization corpus, essentially loosely-overapproximating these structural constraints.

The closest idea to TARGETFUZZ is Creal [30], which injects functions of real-world code snippets into randomly generated programs (seed corpus) to increase test diversity. Creal essentially implements one "function injection" mutator: replace expressions in a seed program with calls to extracted real-world functions, but it ensures the 1) semantic equivalence and 2) well-formedness of the mutation, through a mix of profiling and static analysis. TARGETFUZZ's program transformations are more general, e.g. function injection can simply be achieved by the Exists composition and the Replace mutator. However, TARGETFUZZ does not aim for semantic equivalence or well-formedness. TARGETFUZZ and Creal share a similar goal of breaking reliance on predefined mutators, but TAR-GETFUZZ focuses on leveraging structural relations to test optimization logic. Creal only generates C programs, while TARGETFUZZ is language-agnostic. In addition, using Creal requires a large function dataset (50k+ extracted from Github in its evaluation). TARGETFUZZ only needs a small optimization corpus (100 programs) to derive its mutators.

**Fuzzing compiler optimizations.** There exists a group of compiler fuzzers that are dedicated to, or are especially suitable for fuzzing optimizations. WhiteFox [54] also leverages LLMs and use optimization-triggering tests for few-shot prompting. Thus, they steer testing towards compiler optimization logic. Optimuzz [28] uses directed fuzzing to test newly-updated optimization code, while using translation validation for detecting miscompilations [34]. It relies on some basic mutators that preserve the control-flow, limiting its exploration space. MopFuzzer [53] crafted several mutators that invoke optimizations, e.g, replacing an expression with a trivial function call that returns itself, in hopes of invoking function inlining. MopFuzzer is built for 13 specific optimizations in JVM's JIT compiler, while TARGETFUZZ is more general and can be applied to any optimizations in any compiler. CTOS [24], to the best of our knowledge, is the only optimization fuzzer other than TARGETFUZZ that also explicitly emphasizes and addresses the phase-ordering problem. CTOS uses vector embeddings of test programs and optimization sequences to cluster and

prune similar instances, increasing test diversity. TARGETFUZZ mines and reconstructs grammar-based composition styles to directly test each specific compiler optimization with a targeted pipeline harness. The EMI family [23, 29, 45] of compiler fuzzers are also applicable to fuzzing optimizations by performing semantics-preserving transformation to obtain free miscompilation-oracles.

## 6 CONCLUSION

In summary, this paper is the first to repurpose grammar-based fuzzing with composition styles to carry out targeted fuzzing of compiler optimizations. We propose TARGETFUZZ, a general-purpose fuzzer that can be adapted to fuzz different compilers and programming languages at a low cost. Our evaluation shows that the success of composition-based mutators generalizes to a wide range of optimizations in LLVM and MLIR and that targeted fuzzing complements the existing compiler fuzzing practice.

## REFERENCES

[1] 2021. Introducing Triton: Open-source GPU programming for neural networks. https://openai.com/index/triton/. Accessed 28-08-2025.

[2] 2024. CIRCT: Circuit IR Compilers and Tools. https://circt.llvm.org/. Accessed 28-08-2024.

[3] 2024. GCC Optimize Options. https://gcc.gnu.org/onlinedocs/gcc/Optimize-Options.html. Accessed 28-08-2024.

[4] 2024. LLVM pass registry. https://github.com/llvm/llvm-project/blob/main/llvm/lib/Passes/PassRegistry.def. Accessed 28-08-2024.

[5] 2024. MLIR Passes. https://mlir.llvm.org/docs/Passes/#general-transformation-passes. Accessed 28-08-2024.

[6] 2024. NewGVN.cpp. https://github.com/llvm/llvm-project/blob/main/llvm/lib/Transforms/Scalar/NewGVN.cpp. Accessed 28-08-2024.

[7] 2024. Pattern Rewriting : Generic DAG-to-DAG Rewriting. https://mlir.llvm.org/docs/PatternRewriter/. Accessed 25-02-2025.

[8] 2024. Source-based Code Coverage - Clang 20.0.0git documentation. https://clang.llvm.org/docs/SourceBasedCodeCoverage.html. Accessed 28-08-2024.

[9] 2025. Mojo: Powerful CPU+GPU Programming. https://www.modular.com/mojo. Accessed 28-08-2025.

[10] 2025. OpenXLA Project. https://openxla.org/. Accessed 28-08-2024.

[11] Alfred V. Aho, Ravi Sethi, and J. D. Ullman. 1970. A formal approach to code optimization. SIGPLAN Not. 5, 7 (July 1970), 86–100. doi:10.1145/390013.808486

[12] Frances E Allen and John Cocke. 1971. A Catalogue of Optimization Transformations. https://www.clear.rice.edu/comp512/Lectures/Papers/1971-allen-catalog.pdf. Accessed 25-02-2025.

[13] Cornelius Aschermann, Tommaso Frassetto, Thorsten Holz, Patrick Jauernig, Ahmad-Reza Sadeghi, and Daniel Teuchert. 2019. NAUTILUS: Fishing for Deep Bugs with Grammars. In NDSS.

[14] Sorav Bansal and Alex Aiken. 2006. Automatic generation of peephole superoptimizers. In Proceedings of the 12th International Conference on Architectural Support for Programming Languages and Operating Systems (San Jose, California, USA) (ASPLOS XII). Association for Computing Machinery, New York, NY, USA, 394–403. doi:10.1145/1168857.1168906

[15] Kit Barton, Johannes Doerfert, Hal Finkel, and Michael Kruse. 2019. Loop Fusion, Loop Distribution and Their Place in the Loop Optimization Pipeline. https://llvm.org/devmtg/2019-04/slides/TechTalk-Barton-Loop_fusion_loop_distribution_and_their_place_in_the_loop_optimization_pipeline.pdf. Accessed 25-02-2025.

[16] M. E. Benitez and J. W. Davidson. 1988. A portable global optimizer and linker. SIGPLAN Not. 23, 7 (June 1988), 329–338. doi:10.1145/960116.54023

[17] James Bucek, Klaus-Dieter Lange, and Jóakim v. Kistowski. 2018. SPEC CPU2017: Next-Generation Compute Benchmark. In Companion of the 2018 ACM/SPEC International Conference on Performance Engineering (Berlin, Germany) (ICPE '18). Association for Computing Machinery, New York, NY, USA, 41–42. doi:10.1145/3185768.3185771

[18] Yongheng Chen, Rui Zhong, Hong Hu, Hangfan Zhang, Yupeng Yang, Dinghao Wu, and Wenke Lee. 2021. One Engine to Fuzz 'em All: Generic Language Processor Testing with Semantic Validation. In 2021 IEEE Symposium on Security and Privacy (SP). 642–658. doi:10.1109/SP40001.2021.00071

[19] LLVM Contributors. 2023. MLIR Language Reference. https://mlir.llvm.org/docs/LangRef/

[20] Karine Even-Mendoza, Arindam Sharma, Alastair F. Donaldson, and Cristian Cadar. 2023. GrayC: Greybox Fuzzing of Compilers and Analysers for C. In Proceedings of the 32nd ACM SIGSOFT International Symposium on Software Testing and Analysis (<conf-loc>, <city>Seattle</city>, <state>WA</state>, <country>USA</country>, </conf-loc>) (ISSTA 2023). Association for Computing Machinery, New York, NY, USA, 1219–1231. doi:10.1145/3597926.3598130

[21] Renáta Hodován, Ákos Kiss, and Tibor Gyimóthy. 2018. Grammarinator: a grammar-based open source fuzzer. In Proceedings of the 9th ACM SIGSOFT International Workshop on Automating TEST Case Design, Selection, and Evaluation (Lake Buena Vista, FL, USA) (A-TEST 2018). Association for Computing Machinery, New York, NY, USA, 45–48. doi:10.1145/3278186.3278193

[22] Christian Holler, Kim Herzig, and Andreas Zeller. 2012. Fuzzing with code fragments. In 21st USENIX Security Symposium (USENIX Security 12). 445–458.

[23] Bo Jiang, Xiaoyan Wang, W. K. Chan, T. H. Tse, Na Li, Yongfeng Yin, and Zhenyu Zhang. 2020. CUDAsmith: A Fuzzer for CUDA Compilers. In 2020 IEEE 44th Annual Computers, Software, and Applications Conference (COMPSAC). 861–871. doi:10.1109/COMPSAC48688.2020.0-156

[24] He Jiang, Zhide Zhou, Zhilei Ren, Jingxuan Zhang, and Xiaochen Li. 2022. CTOS: Compiler Testing for Optimization Sequences of LLVM. IEEE Transactions on Software Engineering 48, 7 (2022), 2339–2358. doi:10.1109/TSE.2021.3058671

[25] Prasad Kulkarni, Stephen Hines, Jason Hiser, David Whalley, Jack Davidson, and Douglas Jones. 2004. Fast searches for effective optimization phase sequences. In Proceedings of the ACM SIGPLAN 2004 Conference on Programming Language Design and Implementation (Washington DC, USA) (PLDI '04). Association for Computing Machinery, New York, NY, USA, 171–182. doi:10.1145/996841.996863

[26] Prasad Kulkarni, Wankang Zhao, Hwashin Moon, Kyunghwan Cho, David Whalley, Jack Davidson, Mark Bailey, Yunheung Paek, and Kyle Gallivan. 2003. Finding effective optimization phase sequences. SIGPLAN Not. 38, 7 (June 2003), 12–23. doi:10.1145/780731.780735

[27] Sameer Kulkarni and John Cavazos. 2012. Mitigating the compiler optimization phase-ordering problem using machine learning. SIGPLAN Not. 47, 10 (Oct. 2012), 147–162. doi:10.1145/2398857.2384628

[28] Jaeseong Kwon, Bongjun Jang, Juneyoung Lee, and Kihong Heo. 2025. Optimization-Directed Compiler Fuzzing for Continuous Translation Validation. Proc. ACM Program. Lang. 9, PLDI, Article 172 (June 2025), 24 pages. doi:10.1145/3729275

[29] Vu Le, Mehrdad Afshari, and Zhendong Su. 2014. Compiler validation via equivalence modulo inputs. ACM Sigplan Notices 49, 6 (2014), 216–226.

[30] Shaohua Li, Theodoros Theodoridis, and Zhendong Su. 2024. Boosting Compiler Testing by Injecting Real-World Code. Proc. ACM Program. Lang. 8, PLDI, Article 156 (June 2024), 23 pages. doi:10.1145/3656386

[31] Ben Limpanukorn, Jiyuan Wang, Hong Jin Kang, Eric Zitong Zhou, and Miryung Kim. 2024. Fuzzing MLIR Compilers with Custom Mutation Synthesis. arXiv:2404.16947 [cs.SE] https://arxiv.org/abs/2404.16947

[32] Vsevolod Livinskii, Dmitry Babokin, and John Regehr. 2020. Random testing for C and C++ compilers with YARPGen. Proc. ACM Program. Lang. 4, OOPSLA, Article 196 (Nov. 2020), 25 pages. doi:10.1145/3428264

[33] Vsevolod Livinskii, Dmitry Babokin, and John Regehr. 2023. Fuzzing Loop Optimizations in Compilers for C++ and Data-Parallel Languages. Proc. ACM Program. Lang. 7, PLDI, Article 181 (June 2023), 22 pages. doi:10.1145/3591295

[34] 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. doi:10.1145/3453483.3454030

[35] Henry Massalin. 1987. Superoptimizer: a look at the smallest program. SIGPLAN Not. 22, 10 (Oct. 1987), 122–126. doi:10.1145/36205.36194

[36] Sjoerd Meijer. 2018. [LoopFlatten] Add a loop-flattening pass. https://reviews.llvm.org/D42365. Accessed 25-02-2025.

[37] Sjoerd Meijer. 2021. [llvm-dev] (Loop) passes disabled by default. Google Groups. https://groups.google.com/g/llvm-dev/c/y8PMF9ITkcI Accessed: 2025-09-08.

[38] Sjoerd Meijer. 2021. [LoopUnroll] Don't unroll before vectorisation. LLVM Phabricator. https://reviews.llvm.org/D102748 Accessed: 2025-09-08.

[39] Xianfei Ou, Cong Li, Yanyan Jiang, and Chang Xu. 2025. The Mutators Reloaded: Fuzzing Compilers with Large Language Model Generated Mutation Operators. In Proceedings of the 29th ACM International Conference on Architectural Support for Programming Languages and Operating Systems, Volume 4 (Hilton La Jolla Torrey Pines, La Jolla, CA, USA) (ASPLOS '24). Association for Computing Machinery, New York, NY, USA, 298–312. doi:10.1145/3622781.3674171

[40] Rohan Padhye, Caroline Lemieux, and Koushik Sen. 2019. JQF: Coverage-guided property-based testing in Java. In Proceedings of the 28th ACM SIGSOFT International Symposium on Software Testing and Analysis. 398–401.

[41] Terence Parr. [n. d.]. ANTLR (ANother Tool for Language Recognition). https://www.antlr.org/index.html. Accessed 25-02-2025.

[42] Suresh Purini and Lakshya Jain. 2013. Finding good optimization sequences covering program space. ACM Trans. Archit. Code Optim. 9, 4, Article 56 (Jan. 2013), 23 pages. doi:10.1145/2400682.2400715

[43] Eric Schkufza, Rahul Sharma, and Alex Aiken. 2013. Stochastic superoptimization. In Proceedings of the Eighteenth International Conference on Architectural Support for Programming Languages and Operating Systems (Houston,

Texas, USA) (ASPLOS '13). Association for Computing Machinery, New York, NY, USA, 305–316. doi:10.1145/2451116.2451150

[44] Dominic Steinhöfel and Andreas Zeller. 2022. Input invariants. In Proceedings of the 30th ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering (Singapore, Singapore) (ESEC/FSE 2022). Association for Computing Machinery, New York, NY, USA, 583–594. doi:10.1145/3540250.3549139

[45] Chengnian Sun, Vu Le, and Zhendong Su. 2016. Finding compiler bugs via live code mutation. SIGPLAN Not. 51, 10 (Oct. 2016), 849–863. doi:10.1145/3022671.2984038

[46] Chengnian Sun, Vu Le, Qirun Zhang, and Zhendong Su. 2016. Toward understanding compiler bugs in GCC and LLVM. In Proceedings of the 25th International Symposium on Software Testing and Analysis (Saarbrücken, Germany) (ISSTA 2016). Association for Computing Machinery, New York, NY, USA, 294–305. doi:10.1145/2931037.2931074

[47] S. Triantafyllis, M. Vachharajani, N. Vachharajani, and D.I. August. 2003. Compiler optimization-space exploration. In International Symposium on Code Generation and Optimization, 2003. CGO 2003. 204–215. doi:10.1109/CGO.2003.1191546

[48] Steven R. Vegdahl. 1982. Phase coupling and constant generation in an optimizing microcode compiler. SIGMICRO Newsl. 13, 4 (Oct. 1982), 125–133. doi:10.1145/1014194.800942

[49] Haoyu Wang, Junjie Chen, Chuyue Xie, Shuang Liu, Zan Wang, Qingchao Shen, and Yingquan Zhao. 2023. MLIR-Smith: Random Program Generation for Fuzzing MLIR Compiler Infrastructure. In 2023 38th IEEE/ACM International Conference on Automated Software Engineering (ASE). 1555–1566. doi:10.1109/ASE56229.2023.00120

[50] Junjie Wang, Bihuan Chen, Lei Wei, and Yang Liu. 2017. Skyfire: Data-Driven Seed Generation for Fuzzing. In 2017 IEEE Symposium on Security and Privacy (SP). 579–594. doi:10.1109/SP.2017.23

[51] Deborah L. Whitfield and Mary Lou Soffa. 1997. An approach for exploring code improving transformations. ACM Trans. Program. Lang. Syst. 19, 6 (Nov. 1997), 1053–1084. doi:10.1145/267959.267960

[52] C. Xia, M. Paltenghi, J. Tian, M. Pradel, and L. Zhang. 2024. Fuzz4All: Universal Fuzzing with Large Language Models. In 2024 IEEE/ACM 46th International Conference on Software Engineering (ICSE). IEEE Computer Society, Los Alamitos, CA, USA, 910–910. https://doi.ieeecomputersociety.org/

[53] Zifan Xie, Ming Wen, Shiyu Qiu, and Hai Jin. 2025. Validating JVM Compilers via Maximizing Optimization Interactions. In Proceedings of the 29th ACM International Conference on Architectural Support for Programming Languages and Operating Systems, Volume 4 (Hilton La Jolla Torrey Pines, La Jolla, CA, USA) (ASPLOS '24). Association for Computing Machinery, New York, NY, USA, 345–360. doi:10.1145/3622781.3674188

[54] Chenyuan Yang, Yinlin Deng, Runyu Lu, Jiayi Yao, Jiawei Liu, Reyhaneh Jabbarvand, and Lingming Zhang. 2024. WhiteFox: White-Box Compiler Fuzzing Empowered by Large Language Models. Proc. ACM Program. Lang. 8, OOPSLA2, Article 296 (Oct. 2024), 27 pages. doi:10.1145/3689736

[55] Xuejun Yang, Yang Chen, Eric Eide, and John Regehr. 2011. Finding and understanding bugs in C compilers. In Proceedings of the 32nd ACM SIGPLAN Conference on Programming Language Design and Implementation (San Jose, California, USA) (PLDI '11). Association for Computing Machinery, New York, NY, USA, 283–294. doi:10.1145/1993498.1993532

[56] Chijin Zhou, Bingzhou Qian, Gwihwan Go, Quan Zhang, Shanshan Li, and Yu Jiang. 2024. PolyJuice: Detecting Mis-compilation Bugs in Tensor Compilers with Equality Saturation Based Rewriting. Proc. ACM Program. Lang. 8, OOPSLA2, Article 317 (Oct. 2024), 27 pages. doi:10.1145/3689757

[57] Zhide Zhou, Zhilei Ren, Guojun Gao, and He Jiang. 2021. An empirical study of optimization bugs in GCC and LLVM. Journal of Systems and Software 174 (2021), 110884. doi:10.1016/j.jss.2020.110884