

Modeling the Invariance of Virtual Pointers in LLVM

Piotr Padlewski
Google Research
prazek@google.com

Krzysztof Pszeniczny
Google Research
kpszeniczny@google.com

Richard Smith
Google Research
richardsmith@google.com

Abstract

Devirtualization is a compiler optimization that replaces indirect (virtual) function calls with direct calls. It is particularly effective in object-oriented languages, such as Java or C++, in which virtual methods are typically abundant.

We present a novel abstract model to express the lifetimes of C++ dynamic objects and invariance of virtual table pointers in the LLVM intermediate representation. The model and the corresponding implementation in Clang and LLVM enable full devirtualization of virtual calls whenever the dynamic type is statically known and elimination of redundant virtual table loads in other cases.

Due to the complexity of C++, this has not been achieved by any other C++ compiler so far. Although our model was designed for C++, it is also applicable to other languages that use virtual dispatch. Our benchmarks show an average of 0.8% performance improvement on real-world C++ programs, with more than 30% speedup in some cases. The implementation is already a part of the upstream LLVM/Clang and can be enabled with the `-fstrict-vtable-pointers` flag.

Keywords: devirtualization, invariant, virtual pointer, indirect call, fat pointer

1 Introduction

Devirtualization is a compiler optimization that changes virtual (dynamic) calls to direct (static) calls. The former introduce a performance penalty compared to the latter, as they cannot be inlined and are harder to speculatively execute.

Indirect calls can also serve as an attack vector if an adversary can replace the virtual pointer of an object, e.g. by performing buffer overflow.

Moreover, a recently discovered vulnerability called Spectre [13, 17] showed that the indirect branch predictor is susceptible to side-channel attacks where observable side effects could lead to private data leakage. The currently used mitigation technique (the so-called ‘retpolines’ [26]) effectively eliminates indirect branch prediction entirely.

Indirect calls also limit the extent of compiler analyses and optimizations, making it much more difficult to perform inter-procedural reasoning and thus hinder many useful intra-procedural optimizations, e.g. inlining, intra-procedural register allocation, constant propagation or function attributes inference.

Modeling the lifetime of a C++ pointer in the LLVM IR is nontrivial, because bit-wise equal C++ pointers, viz. having the same bit pattern, may nevertheless denote objects with different lifetimes. They are thus not equivalent and cannot be replaced with each other. However, this information is not reflected in the LLVM IR, because equal values are interchangeable there. The following C++ code, when translated to the LLVM IR and optimized¹, makes the issue apparent:

```

1 void test() {
2     auto* a = new A;
3     external_fun(a);
4
5     A* b = new(a) B;
6     external_fun(b);
7 }

```

```

1 define void @test() {
2     %new = call i8* @NEW(i64 8)
3     %a = bitcast i8* %1 to %struct.A*
4     call void @A_CTOR(%struct.A* %a)
5     call void @external_fun(%struct.A* %a)
6     %b = bitcast i8* %1 to %struct.B*
7     call void @B_CTOR(%struct.B* %b)
8     ; Because %a and %b are known to be equal,
9     ; the compiler decided to use %a.
10    call void @external_fun(%struct.A* %a)
11    ret void
12 }

```

Here, the so-called “placement new” was used to construct a new object in pre-allocated memory (here: in the memory pointed by `a`). As can be seen in this example, `%a` is passed to the second call of `external_fun` instead of `%b`, as the compiler figured out that they are equivalent. On the other hand, if `external_fun(b)` had been replaced with `external_fun(a)`, the behavior would be undefined per the C++ standard.

In this paper, we present a novel way to model C++ dynamic objects’ lifetimes in LLVM capable of leveraging their virtual pointer’s invariance to aid devirtualization.

¹The LLVM IR code was simplified to make it clearer. Mangled names, types or unimportant attributes will be similarly simplified in other listings.

This is a previously unresolved problem, blocking the wide deployment of devirtualization techniques.

2 Related Work

There are multiple optimizations done in hardware that speed up the execution of virtual calls. The most important one is the indirect branch predictor that makes speculative execution of indirect calls possible. Moreover, some CPU architectures provide store buffers [11, 11.10 Store Buffer] [4, 2.4.5.2 L1 DCache] [2, 8.5.1. Store buffer] [18], which temporarily hold writes before committing them to memory. When a load operation is performed, the CPU returns a value from the store buffer if it exists for the given address. Thus, after construction of an object, as long as its virtual pointer's value is in the store buffer, loading it is more efficient.

Additionally, there are also multiple software techniques used for devirtualization. Speculative devirtualization [14, 15] introduces direct calls to known implementations by guarding them with comparisons of virtual pointers or virtual function pointers. This enables inlining and relies on the branch predictor instead of the indirect call predictor, which is more accurate especially on older architectures. However, only implementations visible in the current translation unit can be speculated. An indirect call is thus still required for unknown implementations.

A further similar technique called indirect call promotion [12] uses profiling data of the program to pick the most frequent callsites and insert direct calls guarded by comparison of function pointers similarly to speculative devirtualization.

Another approach is to use information about the entire program using a family techniques called whole program optimizations (WPO), also referred as link-time optimization (LTO) [3, 7, 8, 16, 25] as it often happens during link-time. Having information about the whole program, a compiler can derive facts that might be unknown or hard to model in a single translation unit, e.g. a virtual pointer being invariant, or the definition of virtual tables, or that a function is effectively final (meaning it is not overridden by other types). Using information about the whole program, GCC in LTO mode can prune unreachable branches created by speculative devirtualization. LLVM's whole-program devirtualization [1] is also able to identify virtual functions that only return constants, and then put these constants in the vtable itself and replace calls with simple loads.

Similar optimizations can be also achieved using just-in-time compilation. However, it is rarely used in the case of C++.

Another technique is to leverage the invariance of a virtual pointer and having loaded the pointer once, reuse it for multiple virtual calls, or even devirtualize all virtual calls if the dynamic type of object is statically known. The latter holds e.g. for automatic objects in C++.

In Java, the virtual pointer is set only once in the most derived constructor. There is also no mechanism for changing the dynamic type of an existing object, and the lack of manual memory management means that no reference ever points to a dead object. Hence, one can simply inform the optimizer that virtual pointers are invariant. This technique was used in Falcon [9, 24] – an LLVM-based Java JIT compiler, resulting in 10-20% speedups².

Moreover, some LLVM front-ends decided to create a richer language-specific intermediate language used between abstract syntax tree and the LLVM IR, such as SIL (Swift), MIR (Rust), or Julia IR (Julia). Unfortunately, Clang translates its AST directly to the LLVM IR.

Previous experimental LLVM-based models [19–22] that tried to express lifetime of pointers to dynamic objects were unfortunately flawed. They failed to prevent the the compiler from being easily confused by equality of pointers to objects with different lifetimes. Hence, they could not be enabled by default, as they could lead to a miscompilation, e.g.:

```

1 void g() {
2     A *a = new A;
3     a->virt_meth();
4     A *b = new(a) B;
5     if (a == b) {
6         // Here the compiler is exposed to the fact
7         // that a == b, so it may replace the SSA
8         // value of b with a, which would result in
9         // an erroneous call to A::virt_meth.
10        b->virt_meth();
11    }
12 }

```

3 Problems Specific To C++

In C++ functions are allowed to modify any accessible memory, so the compiler needs to be very conservative with any assumptions. Moreover, compilation of separate modules – called translation units [10, 5.2 Phases of translation] – is usually independent, which makes compilation highly parallelizable, but on the other hand

²Personal communication from Philip Reames, Director of Compiler Technology at Azul Systems

limits the ability of the compiler to reason about functions from other translation units.

External functions are functions defined in a different translation unit than the currently compiled one, so that the body of a function is not visible to the compiler, unless performing whole program optimization. External functions are problematic when reasoning about virtual calls in C++ because without additional knowledge the compiler has to conservatively assume that the virtual pointer might be clobbered (overwritten).

The following code snippet demonstrates some of the shortcomings of traditional C++ devirtualization techniques that do not rely on any additional notion of virtual pointer invariance.

```

1 struct A {
2     virtual void virt_meth();
3 };
4
5 void external_fun(A *a);
6 void foo(A *a) {
7     auto *a = new A;
8     // Conservatively assumed to potentially clobber
9     // a's virtual pointer.
10    external_fun(a);
11    a->virt_meth();
12 }
13
14 void bar() {
15     auto a = new A;
16     a->virt_meth();
17     // Can be devirtualized only if the first call
18     // was inlined, as virt_meth is conservatively
19     // assumed to potentially clobber a's virtual
20     // pointer.
21     a->virt_meth();
22 }

```

In C++ objects can change their dynamic type multiple times during construction and destruction, as each constructor/destructor sets the virtual pointer to the type of the class defining them. Moreover, placement new can create a new object in the memory previously occupied by a different one. This could even happen inside a virtual call, which is allowed to call placement new on the `this` pointer:

```

1 void A::virt_meth() {
2     new(this) B;
3 }

```

Although calling placement new with `this` is allowed, one would nevertheless be unable to use an old reference or pointer to such object afterwards, as they become invalid with the end of the lifetime of the pointed-to object [10, 6.6.3 Object lifetime]. The only way to reuse these memory locations would be to use the pointer returned by placement new, or call `std::launder`, which is the standardized way of obtaining a usable pointer located at the passed address.

Consider the following C++ code:

```

1 struct A {
2     int value;
3 };
4
5 struct B : A { };
6
7 void test() {
8     auto *a = new A;
9     a->value = 0;
10    A *b = new(a)B;
11    b->value = 42;
12    if (a == b) {
13        printf("%d %d", a->value, b->value);
14    }
15 }

```

The behavior of this piece of code is undefined, because although the pointers `a` and `b` refer to the same memory location and are bit-wise equal, `a` cannot be dereferenced after calling placement new, as the latter ends the object's lifetime. Note that reading the value of a pointer to a dead object is legal, but dereferencing it is not. The latter, however, can be safely done by using `std::launder`:

```

1 if (a == b) {
2     // Fine, `a` is not referenced here.
3     printf("%d %d", std::launder(a)->value,
4             b->value);
5 }

```

The C++ rules allow the compiler to assume that whenever multiple virtual calls are performed through the same pointer, they all refer to the same dynamic type.

```

1 void multiple_calls(A *a) {
2     a->virt_meth();
3     // Because `a` is dereferenced again, one may
4     // assume that the dynamic type of `*a` did not
5     // change; otherwise the behavior would have
6     // been undefined. Hence no vtable reload is
7     // needed here, the same function is called.
8     a->virt_meth();
9 }

```

As an interesting corner case, it is nevertheless allowed to dereference a pointer after changing its dynamic type, as long as the dynamic type was changed back to the original one:

```

1 void A::virt_meth() {
2     // Changing dynamic type of 'this'.
3     auto *b = new(this) B;
4     // Calling virtual functions on 'this' is now
5     // not allowed. However, calling them on 'b'
6     // is fine.
7     b->virt_meth();
8     new(this) A;
9     // Can now legally use 'this'.
10    other_virt_meth();
11 }

```

We leverage the above semantics in order to model virtual pointers as invariant, solving the issues described in examples above, which no existing compilers can solve at the time of writing.

4 The Model

We propose an abstract model that allows us to reason about dynamic objects and their use, including object creation, performing virtual calls, and changing dynamic types of existing objects. Such a description is crucial in ensuring correctness of a translation from C++ to the LLVM IR in a way that does not inhibit optimizations.

We model pointers to objects of dynamic type as “fat pointers” – pointers that carry additional information alongside a memory location (address). In our case, we would like to think of them as pointers that also store the current dynamic type. Virtual calls will load the dynamic information from a fat pointer instead of the virtual pointer.

Creation of a fat pointer is achieved by a call to an intrinsic (a built-in function) called `launder` (as a reference to `std::launder`), that for a given address returns a fat pointer that is valid to use, as long as the given memory location is occupied by a live object. We claim that calls to `std::launder` should be translated to calls of our intrinsic, described below.

In our model, accessing a class member (field) or comparing pointers to dynamic objects requires stripping dynamic information from a fat pointer – we call this operation `strip`. Stripping is needed because those operations cannot rely on the dynamic type of object; for example, comparison between pointers should compare only the addresses and not the dynamic information.

Because a fat pointer contains information about the current dynamic type, reading the dynamic type is a pure operation that depends only on the value of the fat pointer. This means that we can model the pureness of this operation as a property of a load instruction.

4.1 Properties of the Model

- `strip` is a pure operation, i.e., its value depends only on its argument.
- `launder` is **not** a pure operation: it creates a fresh fat pointer with the current dynamic type each time it is executed.
- `strip(strip(X))` and `strip(launder(X))` are both equal to `strip(X)`. This is because the model does not care what dynamic information is stripped.
- `launder(strip(X))` and `launder(launder(X))` may be replaced with `launder(X)`. Note that this does not mean that `launder(launder(X)) == launder(X)` – `launder` is inherently nondeterministic and no two invocations of it ever return the same value.

- A pointer passed to `launder` or `strip` aliases (the same memory is accessible through different pointers) the returned pointer.

5 Modeling in LLVM

In this section, we show how to translate our model to the LLVM IR. One way to represent the fat pointer is as a struct containing an address and some additional metadata about the dynamic type, e.g. the virtual pointer. Then, every virtual call can use additional information instead of loading the virtual pointer from the object. This achieves its goals because the fat pointer does not change during the object’s lifetime, as the same SSA value is used for every virtual call on the given object. However, we observe that in such a case, the virtual pointer would be stored in two places – as a member in class and as extra data in fat pointer. Unfortunately, removing the virtual pointer from the class instances is unacceptable in a real-world C++ compiler as it breaks the established Application Binary Interface (e.g Itanium ABI [5]) – a specification describing the common behavior between the binary files (object files, executable, or a library), like the layout of the objects or calling conventions. Thus, it would not be possible to link ordinary object code with object code that uses fat pointers.

Fortunately, we can remove the additional information from the fat pointer and model the fact that loads of the virtual pointer are invariant as a property of the load instruction. This is beneficial, as it does not introduce any data redundancy and does not break the ABI. To model the fact that dynamic information is invariant for any given fat pointer, we introduce new instruction metadata – `invariant.group` – that can be attached to a `load` or `store` instruction. Instruction metadata in LLVM hints to the optimizer about special properties of an instruction. As specified by the LLVM language reference, stripping instruction metadata should not change the semantics of the program.

The presence of the `invariant.group` metadata on the instruction instructs the optimizer that every load and store to the same pointer operand can be assumed to load or store the same value. We attach this metadata on vtable loads (when performing virtual calls) and stores (in constructors and destructors).

To model the fact that the vtable contents are constant and do not change during the execution of the program, we use `invariant.load` metadata. This metadata specifies that the memory location referenced by the load contains the same value at all points in the program. Note that this requires a stronger assumption

than `invariant.group`. With `invariant.load`, multiple loads of the same function pointer from one vtable may be merged into one value.

Modeling `strip` and `launder` is done by introducing new intrinsics that take and return regular pointers, carrying imaginary meta information. The intrinsics only serve as annotations in the IR and are stripped just before the code generation phase, hence they do not introduce any runtime or object code size cost.

5.1 `llvm.strip.invariant.group`

The `strip` operation is represented by a new intrinsic function `i8* @llvm.strip.invariant.group(i8*)`, used when an invariant established by `invariant.group` metadata no longer holds, to obtain a new pointer value that does not carry the invariant information. It has the `readonly`, `speculatable`, and `nounwind` attributes [6, Function Attributes] meaning that the function is pure, can be speculated [6, Function Attributes] and that it does not throw exceptions.

5.2 `llvm.launder.invariant.group`

The `launder` operation is represented by a new intrinsic function `i8* @llvm.launder.invariant.group(i8*)`, used when an invariant established by `invariant.group` metadata no longer holds, to obtain a new pointer that carries fresh invariant group information.

It has `inaccessiblememoryonly` attribute meaning that this function can modify only memory that is not visible to the current module. In this context, we can think of it as storing the additional virtual pointer in some imaginary location. Similarly to `strip`, it has the `speculatable` and `nounwind` attributes.

The following code sample illustrates the intrinsics:

```

1 %p = alloca i8
2 store i8 42, i8* %p, !invariant.group !{}
3 call void @foo(i8* %p)
4
5 ; *%p is known not to have changed
6 %a = load i8, i8* %p, !invariant.group !{}
7 call void @foo(i8* %p)
8
9 %unknown = load i8, i8* @unknown_ptr
10
11 ; All stores to *%p write equal values,
12 ; hence we may infer that %unknown == 42
13 store i8 %unknown, i8* %p, !invariant.group !{}
14
15 call void @foo(i8* %p)
16 %q = call i8* @llvm.launder.invariant.group(i8* %p)
17 ; Cannot propagate invariant.group information
18 ; through an llvm.invariant.group.launder and
19 ; thus %d cannot be replaced with %a
20 %d = load i8, i8* %q, !invariant.group !{}

```

Here loads from the pointer `%p` can be assumed to load the same value when `invariant.group` is present, and the load of `%q` returned by the `launder` cannot.

One of the primary concerns when making the model sound was that comparison of bit-wise equal pointers, representing objects of different lifetimes, should not make it possible to make them equivalent, meaning that one pointer could be replaced with the second one.

For example, consider again the code snippet from the introduction to this paper, that demonstrated a flaw present in the previous models:

```

1 void g() {
2     A *a = new A;
3     a->virt_meth();
4     A *b = new(a) B;
5     if (a == b) {
6         // Here the compiler is exposed to fact
7         // that a == b, so it may replace the SSA
8         // value of b with a, which would result in
9         // an erroneous call to A::virt_meth.
10        b->virt_meth();
11    }
12 }

```

Under our model, the dynamic information is always stripped whenever an operation that can expose it is used, e.g. pointers about to be compared are always passed through a `strip`:

```

1 %vtable_a = load i8* %a, !invariant.group !{}
2 ; ...
3 ; if (a == b)
4 %sa = call i8* @llvm.strip.invariant.group(i8* %a)
5 %sb = call i8* @llvm.strip.invariant.group(i8* %b)
6 %bool = icmp %sa, %sb
7 br %bool, %if, %after
8 if:
9 ; a == b was lowered to %sa == %sb which does not
10 ; imply %a == %b, thus %vtable_b cannot be replaced
11 ; with %vtable_a
12 %vtable_b = load i8* %b, !invariant.group !{}

```

5.3 Emitting Launder and Strip

We apply the following rules to decide when to emit `launder` or `strip` when translating C++ to the LLVM IR:

- Constructors of derived classes need to `launder` the `this` pointer before passing it to the constructors of base classes; they may subsequently operate on the original `this` pointer.
- Likewise, destructors of derived classes must `launder` the `this` pointer before passing it to the destructors of base classes.
- Placement new and `std::launder` shall call `launder`.
- Accesses to union members must call `launder`, because the active member of the union may have changed since the last visible access.

- Whenever two pointers are to be compared they must be stripped first, because we want the result will only provide information about the address equality, not about `invariant.group` equality.
- Likewise, whenever a pointer is to be casted to an integer type, it must be stripped first.
- Whenever an integer is casted to a pointer type, the result must be laundered before it is used, because the object stored at this address could have changed since the last cast of the same integer to a pointer type.

These rules ensure that:

- no operation can leak the pointer equality to the `invariant.group` model – as every such an operation necessarily operates on stripped pointers;
- every time a pointer with a possibly unknown dynamic type is obtained, it is free of any assumptions made by the model.

5.4 Aliasing

In C++, even though pointers obtained through placement new or `std::launder` are not equivalent to the original pointer (i.e. cannot be freely substituted), they nevertheless alias the original pointer. We want to mimic this behavior in LLVM. Thus, the pointer returned by `strip` or `launder` aliases the intrinsic’s argument:

```

1 store i8 42, i8* %a
2 %b = call i8* @llvm.launder.invariant.group(i8* %a)
3 ; %b cannot be replaced with %a, but because they
4 ; alias each other, we may derive %v == 42.
5 %v = load i8, i8* %b

```

This technique of forcing aliasing on equal-but-not-equivalent pointers allows the usual memory optimizations, such as replacing redundant loads or dead store elimination through a `launder` or `strip`. The implementation simply treats `launder` and `strip` as bit-casts in the context of alias analysis.

We derive other properties based on aliasing of argument and returned value:

- An argument of a function is considered `nocapture` if the callee does not make any copies of the pointer that outlive the callee itself. If we know that a pointer returned from `launder` or `strip` is not captured, we can state that our intrinsic does not capture the argument. This is crucial, as `nocapture` is not a valid attribute for return values.
- If the argument of a `launder` or a `strip` has `nonnull` or `dereferenceable` attributes, then we can also apply it to the returned values.

5.5 Other Properties

Since the special values `undef` (unspecified bit-pattern) and `null` are unable to carry any virtual information, they can be safely propagated through both of the `invariant.group` intrinsics:

- `strip(undef) == undef == launder(undef)`
- `strip(null) == null == launder(null)`

In the case of `undef`, it is useful, as it helps to remove unreachable code, as some operations on `undef` can be fold to the `unreachable` instruction, which specifies that no execution ever reaches that point in the program. `null` is another special value, that has an unusual treatment in LLVM, e.g. an early exit caused by true comparison with `null` is considered to be a cold block (executed infrequently), which means that constant folding of `null` is important to get the same behavior from the static profile annotation pass without any modification to specially handle the new intrinsics.

The `invariant.group` metadata also make it possible to hoist loads virtual table loads out of the loops. This means that calling virtual function on an object in a loop will load virtual table and virtual function only once.

5.6 Do We Need Both Intrinsics?

With all these properties, adding calls to `strip` and `launder` related to pointer comparisons and integer-to-pointer conversions does not cause any semantic IR-level information to be lost: if any piece of information could be inferred by the optimizer about some collection of variables (e.g., that two pointers are equal), it can be inferred now about their stripped versions, no matter how many `strip` and `launder` calls have been made to obtain them in the IR. As an example, the C++ expression `ptr == std::launder(ptr)` will be optimized to true, because it is lowered to a comparison of `strip(ptr)` with `strip(launder(ptr))`, which are indeed equal under our rules.

One could argue that replacing uses of the `strip` with `launder` produces semantically correct code, although it is not strictly coherent with our model. However, having a `strip` is crucial from the optimizations point of view, as without it we lose the ability to leverage pointer equality in other optimizations. As a toy example, `strip(X) == strip(X)` can be folded to true, but it is not the case with `launder(X) == launder(X)` as `launder` is not pure.

6 Outline Constructors

External constructors pose a problem for devirtualization because after lowering to the LLVM IR it is not longer known what value is stored under a virtual pointer. This situation appears in the following code snippet:

```

1 struct C {
2     C();
3     virtual void virt_meth();
4 };
5
6 void foo() {
7     auto *c = new C;
8     // Unable to devirtualize.
9     c->virt_meth();
10 }

```

Following Padlewski [20] we use “assumption loads”, which inform the compiler about the actual dynamic type of the newly constructed object. To this end, the following three LLVM instructions are emitted after each constructor call:

1. Load the virtual pointer.
2. Compare the loaded virtual pointer with the vtable it should point to.
3. Call `@llvm.assume(i1 %b)` with the comparison result as a parameter.

The intrinsic function `@llvm.assume` takes a Boolean argument and allows the optimizer to assume that the given Boolean is true whenever the instruction is executed. In this case, it propagates equality, replacing virtual pointer loads (annotated with `!invariant.group` metadata) with the actual value. This does not introduce any program runtime cost, as `assume` is stripped before generating actual code, leaving comparison and load instructions trivially dead. The following code demonstrates the technique:

```

1 call @ctor(i8 %this_ptr)
2 %vtable = load %this_ptr, !invariant.group !{}
3 %b = icmp eq %vtable, @VTABLE_OF_C
4 call void @llvm.assume(i1 %b)

```

Note that, even in the case of assumption loads, the vtable load is marked with `invariant.group` to easily propagate the value. If the constructor is inlined, the value of the comparison can easily be determined to be true, which then allows for the instructions to be optimized away.

```

1 ; From call @ctor(i8 %this_ptr)
2 store @VTABLE_OF_C, %this_ptr, !invariant.group !{}
3 ; ...
4 ; Will be replaced with @VTABLE_OF_C.
5 %vtable = load %this_ptr, !invariant.group !{}
6 ; Optimized to true.
7 %b = icmp eq %vtable, @VTABLE_OF_C
8 ; assume(true) can be removed.
9 call void @llvm.assume(i1 %b)

```

7 Virtual Table Definition

Being able to inspect the vtable definition is necessary for full devirtualization in C++. Without it, the optimizer is only able to devirtualize vtable loads and not vfunction loads.

This is troublesome on platforms following the Itanium ABI (most of the UNIX-based platforms) as it uses the so-called “key functions” [5, 5.2.3 Virtual Tables] to emit vtable in the fewest translation units possible, which is beneficial for code size and compile time. This problem does not exist on Microsoft platforms, as the Microsoft ABI specifies that a vtable definition will be present in every translation unit that uses it.

Following the previous devirtualization models[20], we use the `available_externally` linkage [6, Linkage] already present in the LLVM IR, which provides the definition of a global variable, but only for the purposes of optimizations – the definition is stripped after optimizations leaving the variable with `external` linkage.

7.1 Emitting Virtual Inline Functions

Unfortunately, emitting a vtable definition without emitting all of the symbols referenced by it is not necessarily correct. This is because the Itanium ABI does not guarantee that symbols for inline virtual functions or ones marked with `__attribute__((visibility="hidden"))` will be exported.

Our solution is to emit vtable definitions opportunistically, viz. only when they do not reference inline virtual functions that were not emitted in that translation unit. We also implemented an option to emit all inline virtual functions with a flag `-fforce-emit-vtable` so that all definitions of all vtables can be emitted. However, this incurs non-trivial costs, as more functions need to be emitted and subsequently optimized. It also can increase the size of object files as some definitions might not be removed by the compiler, e.g. due to cyclic references (virtual table referencing a virtual function which in turn calls other virtual functions through the vtable).

8 Cross Module Optimizations

Whole program optimization is generally challenging, as modules compiled with a different set of optimizations could be optimized together. This is troublesome in our case as one C++ module could be compiled with devirtualization enabled, while another module compiled from e.g. C++, C, or Swift might not use `launder` and `strip` where the model requires it, as can be seen in the next snippet:

```

1 // Module 1: with devirtualization.
2 void with_devirt(A *a) {
3     a->foo();
4     A* other = unsafe_placement_new(a);
5     other->foo();
6 }
7
8 // Module 2: without devirtualization.
9 A* unsafe_placement_new(A *a) {
10     return new(a) B;
11 }

```

A naïve LTO between module 1 and module 2 could cause miscompilation, because our model assumes that `launder` is emitted whenever placement new is used.

A prototype approach is to disallow linking between modules with and without devirtualization. This is not ideal as it prevents performing LTO between different languages. Another possible solution would consist of stripping intrinsics and metadata from modules with devirtualization, or introducing intrinsics in modules without it.

As future work, one can alternatively perform the said stripping lazily at function-level granularity. To this end, the `invariant.group` metadata would be equipped with a tag referring to their supported devirtualization policy, to allow future extensions and special handling of other programming languages’ idiosyncracies. Moreover, functions would be annotated with a list of applicable policies. The `invariant.group` metadata and intrinsics would have no effect unless their policy is present in the policies list of the enclosing function. Otherwise, they can be safely removed. Hence, it would be possible to perform LTO across modules using `invariant.group` in different ways, e.g. a Swift module that does not need to use `launder` or `strip` with a C++ module using `invariant.group` virtual pointer loads. Whenever inlining or any other inter-procedural optimization is performed, the applicable policies list shall of course be set to the intersection of the policies lists of the optimization’s input functions.

9 Benchmarks

In previous work Padlewski [20] showed that a model that is not sound, can lead to 4 times more virtual function calls being devirtualized. Because the previous model’s unsoundness surfaces only in hand-crafted cases that should almost never appear in real code, our model achieves similar numerical results.

Our optimization was evaluated on a number of open source benchmarks using SPEC 2006, and Google internal benchmarks. In each case, the baseline was fully optimized (`-O3`), but without FDO, LTO or post-link optimizers[23]. The devirtualized build used the same

setup with an additional `-fstrict-vtable-pointers` flag.

9.1 SPEC 2006 Benchmarks

We did not see any difference on SPEC 2006 between optimized and devirtualized builds. This can be attributed to the fact that only half of benchmarks use C++, from which only 5 used virtual functions – perhaps not heavily enough to show the difference in the execution time.

9.2 Google Benchmarks

We also evaluated our changes on the standard sets of internal Google benchmarks. We observed a 0.8% improvement in the geometric mean of execution times. The results across all of the benchmarks are shown in Figure 1. Among these, a certain benchmark set heavily reliant on virtual calls showed a consistent 4.6% improvement, see Figure 2. The regressions seen in some benchmarks are most likely caused by the inliner making different decisions when exposed to more inlining opportunities. This can be possibly fixed by tuning the inliner and other passes.

Protocol Buffers is a widely-used library for serialization and deserialization of structured objects. Two sets of benchmarks for Protocol Buffers showed significant improvement, with some microbenchmarks showing over 30% improvement. The results from one of the benchmark sets are shown in Figure 3 and Figure 4.

A very large internal application (Web Search) exhibits a significant 0.65% improvement in QPS. This strongly suggests that particularly in setups without LTO and FDO, devirtualization may bring tangible benefits.

However, in the case of LTO and FDO we haven’t see improvement yet. We think that it mostly boils down to tuning other passes, because performing LTO or FDO on modules, that are proven to be faster should not regress them. We do not expect to see similar improvement with LTO and FDO because they can perform subset of optimizations using more information. It is important to note, that although some of the optimizations that we enable are theoretically possible with Whole Program Optimization, they are not feasible in practise.

We believe that better results could be seen in the languages that rely more heavily on virtual dispatch, like Java. The mentioned example of Falcon JIT [9, 24] achieved 10-20% speedup relying on unsound technique.

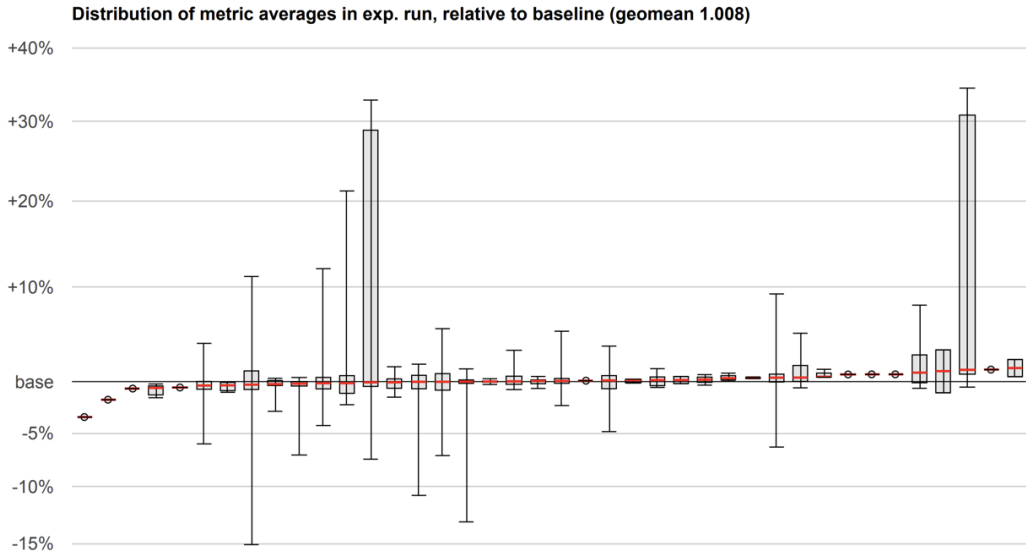


Figure 1. Results of the internal sets of benchmarks. The central lines denote the medians, grey boxes correspond to the interquartile ranges and whiskers show the minimum and maximum value.

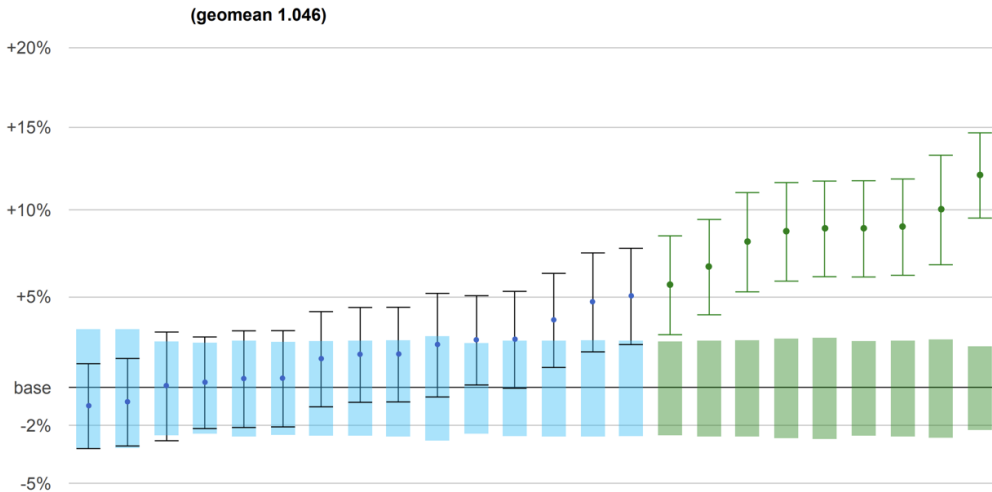


Figure 2. Average results of benchmarks within one benchmark set heavily reliant on virtual calls. 95% confidence intervals are shown: as boxes (for base) and as whiskers (for experiment).

10 Future Work

Despite significant performance improvements, our de-virtualization is disabled by default in Clang, but can be enabled with the `-fstrict-vtable-pointers` compiler flag. The last piece required to enable it by default is to implement the `supported optimizations` attribute described in Section 8: Cross Module Optimizations.

The described techniques rely heavily on the program’s behavior being defined. Some users are afraid of enabling

optimizations that can exacerbate undefined behavior as they do not know if it is present in their code. Fortunately, a tool called Undefined Behavior Sanitizer (UBSan) detects different occurrences of undefined behavior in the runtime of the program. It would be preferable to introduce new checks to UBSan that would detect places where we falsely believe the virtual pointer did not change.

One can also model the invariance of `const` fields in C and C++ objects. They too can change whenever a

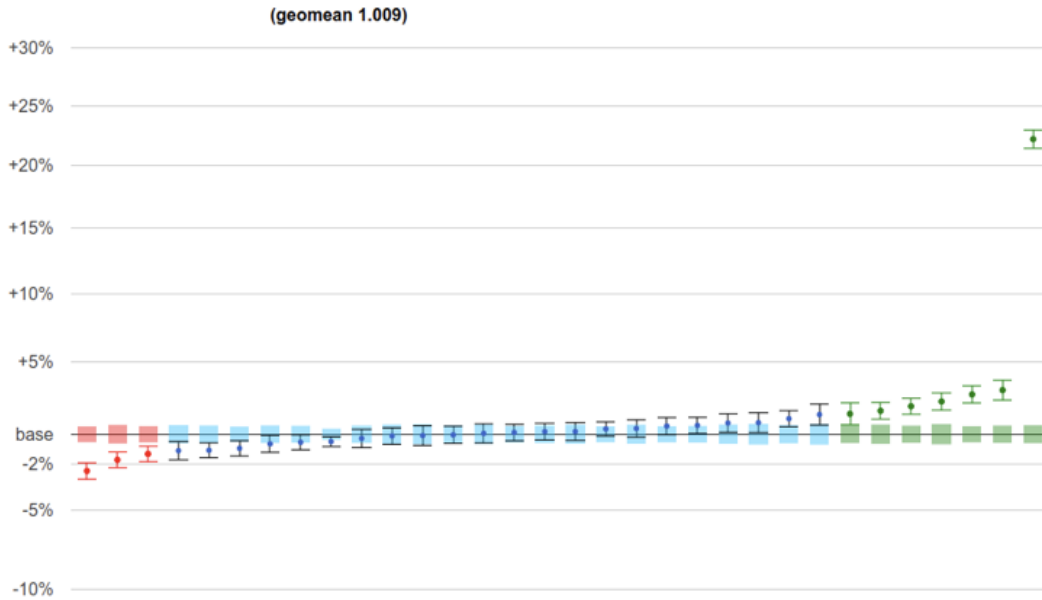


Figure 3. Average results of benchmarks within a Protocol Buffers benchmark set heavily reliant on virtual calls. 95% confidence intervals are shown: as boxes (for base) and as whiskers (for experiment).

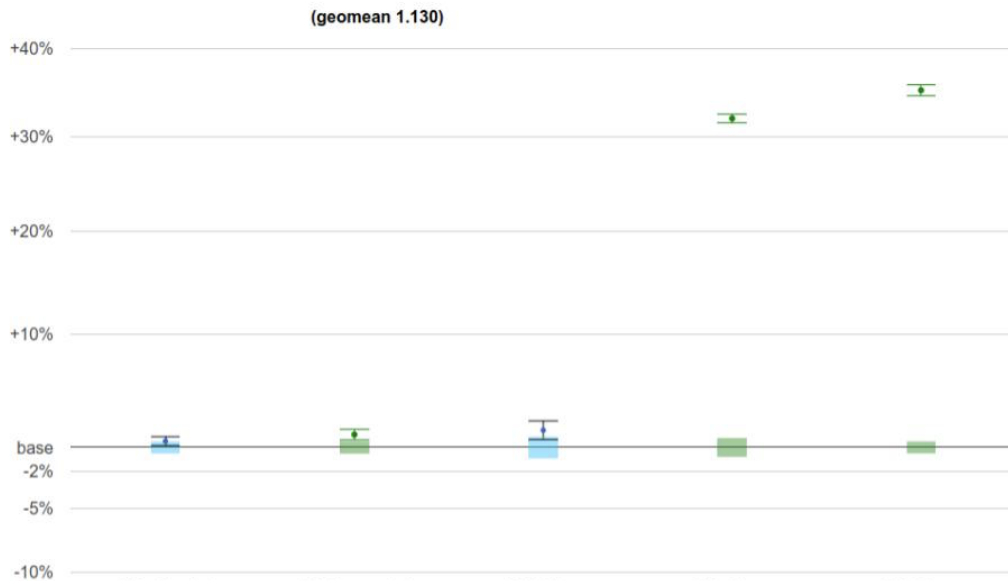


Figure 4. Average results of benchmarks within another Protocol Buffers benchmark set heavily reliant on virtual calls. 95% confidence intervals are shown: as boxes (for base) and as whiskers (for experiment).

placement new is made, so handling them necessarily requires solving the same issues as with virtual pointers. Our solution can be easily extended to this case. In the abstract model, the metadata in the fat pointer would not only store the dynamic type, but also the values of `const` members of the object. The concrete implementation in

LLVM would simply have to annotate loads and stores to `const` variables with `invariant.group`. Moreover, `launders` and `strips` would have to be emitted for every struct type instead of only those that might have a vtable.

References

- [1] Clang documentation: LTO visibility. <https://clang.llvm.org/docs/LTOVisibility.html>.
- [2] Cortex-R4 and Cortex-R4F technical reference manual. <http://infocenter.arm.com/help/index.jsp?topic=/com.arm.doc.ddi0363e/Chdcahcf.html>.
- [3] GCC link time optimization. <https://gcc.gnu.org/onlinedocs/gccint/LTO.html#LTO>.
- [4] Intel 64 and IA-32 architectures optimization reference manual. <https://software.intel.com/sites/default/files/managed/9e/bc/64-ia-32-architectures-optimization-manual.pdf>.
- [5] Itanium C++ ABI. <https://itanium-cxx-abi.github.io/cxx-abi/abi.html>.
- [6] LLVM language reference manual. <https://llvm.org/docs/LangRef.html>.
- [7] LLVM link time optimization: Design and implementation. <https://llvm.org/docs/LinkTimeOptimization.html>.
- [8] Whole program assumptions, linker plugin and symbol visibilities. <https://gcc.gnu.org/onlinedocs/gccint/WHOPR.html>.
- [9] Azul systems seizes java runtime performance lead with falcon, a new just-in-time compiler based on LLVM. https://www.azul.com/press_release/falcon-jit-compiler/, 2017.
- [10] Working draft, standard for programming language C++ . <http://www.open-std.org/jtc1/sc22/wg21/docs/papers/2017/n4659.pdf>, March 2017.
- [11] Intel 64 and IA-32 architectures software developer’s manual. <https://software.intel.com/sites/default/files/managed/39/c5/325462-sdm-vol-1-2abcd-3abcd.pdf>, 01 2019.
- [12] Ivan Baev. Profile-based indirect call promotion. <https://llvm.org/devmtg/2015-10/slides/Baev-IndirectCallPromotion.pdf>, 2015.
- [13] Jann Horn. Reading privileged memory with a side-channel. <https://googleprojectzero.blogspot.com/2018/01/reading-privileged-memory-with-side.html>, 2018.
- [14] Jan Hubicka. Devirtualization in C++, part 3 (building the type inheritance graph). <http://hubicka.blogspot.com/2014/02/devirtualization-in-c-part-3-building.html>, February 2014.
- [15] Jan Hubicka. Devirtualization in C++, part 4 (analyzing the type inheritance graph for fun and profit). <http://hubicka.blogspot.com/2014/02/devirtualization-in-c-part-4-analyzing.html>, 2014.
- [16] Teresa Johnson, Mehdi Amini, and David Xinliang Li. ThinLTO: Scalable and incremental LTO. *2017 IEEE/ACM International Symposium on Code Generation and Optimization (CGO)*, pages 111–121, 2017.
- [17] Paul Kocher, Jann Horn, Anders Fogh, Daniel Genkin, Daniel Gruss, Werner Haas, Mike Hamburg, Moritz Lipp, Stefan Mangard, Thomas Prescher, Michael Schwarz, and Yuval Yarom. Spectre attacks: Exploiting speculative execution. In *40th IEEE Symposium on Security and Privacy (S&P’19)*, 2019.
- [18] Scott Owens, Susmit Sarkar, and Peter Sewell. A better x86 memory model: X86-TSO. In *Proceedings of the 22Nd International Conference on Theorem Proving in Higher Order Logics, TPHOLs ’09*, pages 391–407, Berlin, Heidelberg, 2009. Springer-Verlag.
- [19] Piotr Padlewski. Devirtualization in LLVM. 2016 US LLVM Developers’ Meeting <http://llvm.org/devmtg/2016-11/Slides/Padlewski-DevirtualizationInLLVM.pdf>, 2016.
- [20] Piotr Padlewski. Devirtualization in LLVM. In *Proceedings Companion of the 2017 ACM SIGPLAN International Conference on Systems, Programming, Languages, and Applications: Software for Humanity, SPLASH Companion 2017*, pages 42–44, New York, NY, USA, 2017. ACM.
- [21] Piotr Padlewski. Devirtualization in LLVM and Clang. <http://blog.llvm.org/2017/03/devirtualization-in-llvm-and-clang.html>, March 2017.
- [22] Piotr Padlewski and Richard Smith. RFC: Improvements to LLVM and clang’s devirtualization support. <https://docs.google.com/document/d/1f2SGa4TIPuBGm6y6YO768GrQsA8awNfGEJSBFukLhYA/edit#heading=h.iue9731a57u3>, July 2015.
- [23] Maksim Panchenko, Rafael Auler, Bill Nell, and Guilherme Ottoni. BOLT: A practical binary optimizer for data centers and beyond. *CoRR*, abs/1807.06735, 2018.
- [24] Philip Reames. FALCON: An optimizing JAVA JIT. <https://llvm.org/devmtg/2017-10/slides/Reames-FalconKeynote.pdf>, 2017.
- [25] Mehdi Amini Teresa Johnson and David Li. ThinLTO: Scalable and incremental LTO. <http://blog.llvm.org/2016/06/thinlto-scalable-and-incremental-lto.html>, June 2016.
- [26] Paul Turner. Retpoline: a software construct for preventing branch-target-injection. <https://support.google.com/faqs/answer/7625886>, 2018.