#### Dynamic Dispatch and Duck Typing

L25: Modern Compiler Design

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 のへぐ

## Late Binding

• Static dispatch (e.g. C function calls) are jumps to specific addresses

- Object-oriented languages decouple method name from method address
- One name can map to multiple implementations (e.g. different methods for subclasses)
- Destination must be computed somehow

## Example: C++

- Mostly static language
- Methods tied to class hierarchy
- Multiple inheritance can combine class hierarchies

```
class Cls {
   virtual void method();
};
// object is an instance of Cls or a subclass of
   Cls
void function(Cls *object) {
   // Will call Cls::method or a subclass
        override
   object->method();
}
```

## Example: Objective-C

- AoT-compiled dynamic language
- Duck typing

```
@interface Cls
- (void)method;
@end
// object is an instance of any classes
void function(id object) {
    // Will call method if object is an instance
    // of a class that implements it, otherwise
    // will invoke some forwarding mechanism.
    [object method];
}
```

### Example: JavaScript

- Prototype-based dynamic object-oriented language
- Objects inherit from other objects (no classes)
- Duck typing

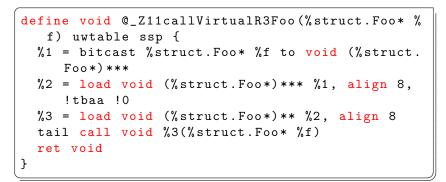
```
a.method = function() { ... };
...
// Will call method if b or an object on
// b's prototype chain provides it. No
// difference between methods and
// instance/ variables: methods are just
// instance variables containing
// closures.
b.method();
```

#### VTable-based Dispatch

- Tied to class (or interface) hierarchy
- Array of pointers (virtual function table) for method dispatch
- Method name mapped to vtable offset

```
struct Foo {
    int x;
    virtual void foo();
};
void Foo::foo() {}
void callVirtual(Foo &f) {
    f.foo();
}
void create() {
    Foo f;
    callVirtual(f);
}
```

#### Calling the method via the vtable



Call method at index 0 in vtable.

#### Creating the object

```
@_ZTV3Foo = unnamed_addr constant [3 x i8*] [
 i8* null.
 i8* bitcast ({ i8*, i8* }* @_ZTI3Foo to i8*).
 i8* bitcast (void (%struct.Foo*)*
     @ ZN3Foo3fooEv to i8*)]
define linkonce_odr void @_ZN3FooC2Ev(%struct.
   Foo* nocapture %this) {
 %1 = getelementptr inbounds %struct.Foo* %this
     , i64 0, i32 0
  store i32 (...)** bitcast
    (i8** getelementptr inbounds ([3 x i8*]*
       @_ZTV3Foo, i64 0, i64 2) to i32 (...)**),
    i32 (...)*** %1
}
```

## Devirtualisation

- Any indirect call prevents inlining
- Inlining exposes a lot of later optimisations
- If we can prove that there is only one possible callee, we can inline.
- Easy to do in JIT environments where you can *deoptimise* if you got it wrong.

• Hard to do in static compilation

#### Problems with VTable-based Dispatch

- VTable layout is per-class
- Languages with duck typing (e.g. JavaScript, Python, Objective-C) do not tie dispatch to the class hierarchy
- Dynamic languages allow methods to be added / removed dynamically
- Selectors must be more abstract than vtable offsets (e.g. globally unique integers for method names)

## Ordered Dispatch Tables

- All methods for a specific class in a sorted list
- Binary (or linear) search for lookup
- Lots of conditional branches for binary search
- Either very big dtables (one entry per class, method pair, including inherited methods) or multiple searches to look at superclasses
- Cache friendly for small dtables (entire search is in cache)

Expensive to add methods (requires lock / RCU)

# Sparse Dispatch Tables

- Tree structure, 2-3 pointer accesses + offset calculations
- Fast if in cache
- Pointer chasing is suboptimal for superscalar chips (inherently serial)

• Copy-on-write tree nodes work well for inheritance, reduce memory pressure

## Inverted Dispatch Tables

- Normal dispatch tables are a per-class (or per object) map from selector to method
- Inverted dispatch tables are a per-selector map from class (or object) to method
- If method overriding is rare, this provides smaller maps (but more of them)
- Not common, but useful in languages that encourage shallow class hierarchies

# Lookup Caching

- Method lookup can be slow or use a lot of memory (data cache)
- Caching lookups can give a performance boost
- Most object-oriented languages have a small number of classes used per callsite

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のQ@

• Have a per-callsite cache

## Callsite Categorisation

- Monomorphic: Only one method ever called
  - Huge benefit from inline caching
- Polymorphic: A small number of methods called
  - Can benefit from simple inline caching, depending on pattern

- Polymorphic inline caching (if sufficiently cheap) helps
- Megamorphic: Lots of different methods called
  - Cache usually slows things down

### Simple Cache

object.aMethod(foo);

```
static struct {
   Class cls;
   Method method;
} cache = {0, 0};
static Selector sel = compute_selector("aMethod"
   );
if (object->isa != cache->cls) {
   cache->cls = object->isa
   cache->method = method_lookup(cls, sel);
}
cache->method(object, sel, foo);
```

What's wrong with this approach?

## Simple Cache

object.aMethod(foo);

```
static struct {
   Class cls;
   Method method;
} cache = {0, 0};
static Selector sel = compute_selector("aMethod"
   );
if (object->isa != cache->cls) {
   cache->cls = object->isa
   cache->method = method_lookup(cls, sel);
}
cache->method(object, sel, foo);
```

What's wrong with this approach? Updates? Thread-safety?

# Inline caching in JITs

- Cache target can be inserted into the instruction stream
- JIT is responsible for invalidation
- Can require *deoptimisation* if a function containing the cache is on the stack

# Speculative inlining

- Lookup caching requires a mechanism to check that the lookup is still valid.
- Why not inline the expected implementation, protected by the same check?
- Essential for languages like JavaScript (lots of small methods, expensive lookups)

### Inline caching

call lookup\_fn nop

```
bne $cls, $last, fail
call method
continue:
```

▲ロト ▲帰 ト ▲ ヨ ト ▲ ヨ ト ・ ヨ ・ の Q ()

- · First call to the lookup rewrites the instruction stream
- Check jumps to code that rewrites it back

# Variation: Guarded methods

- Specialised version of each method that knows the expected class
- Jump to the lookup function replaced by call to guarded method
- Method checks receiver type and tail-calls the lookup function if it's the wrong type

# Polymorphic inline caching

```
bne $cls, $expected, cls
call method
ret
next:
bne $cls, $expected2, cls
call method
ret
```

- Branch to a jump table
- Jump table has a sequence of tests and calls
- Jump table must grow
- Too many cases can offset the speedup

#### Trace-based optimisation

- Branching is expensive
- Dynamic programming languages have lots of method calls
- Common hot code paths follow a single path
- Chain together basic blocks from different methods into a trace
- Compile with only branches leaving
- Contrast: trace vs basic block (single entry point in both, multiple exit points in a trace)

## Prototype-based Languages

- Prototype-based languages (e.g. JavaScript) don't have classes
- Any object can have methods
- Caching per class is likely to hit a lot more cases than per object

## Hidden Class Transforms

- Observation: Most objects don't have methods added to them after creation
- Create a hidden class for every constructor
- Also speed up property access by using the class to specify fixed offsets for common properties

- Create a new class whenever a property is added
- Difficult problem: Merging identical classes.

## Type specialisation

- Code paths can be optimised for specific types
- For example, elide dynamic lookup
- Common case: a+b is much faster if you know a and b are integers!

- Can use static hints, works best with dynamic profiling
- Must have fallback for when wrong

# Deoptimisation

- Disassemble existing stack frame and continue in interpreter / new JIT'd code
- Stack maps allow mapping from register / stack values to IR values
- Fall back to interpreter for new control flow
- NOPs provide places to insert new instructions
- New code paths can be created on demand
- Can be used when caches are invalidated or the first time that a cold code path is used

# LLVM: Anycall calling convention

- Used for deoptimisation
- All arguments go somewhere
- Metadata emitted to find where
- Very slow when the call is made, but no impact on register allocation
- Call is a single jump instruction, small instruction cache footprint
- Designed for slow paths, attempts not to impact fast path

## Deoptimisation example

```
JavaScript:
```

a = b + c;

Deoptimisable pseudocode:

```
if (!(is_integer(b) && is_integer(c)))
    anycall_interpreter(&a, b, c); // Function
    does not return
a = b+c;
```

### Questions?