# From parsing to sema: making sense of syntax trees

Michał Dominiak Nokia Networks griwes@griwes.info

Parser
Semantic analyzer
The curious case of C++
Optimizer
Code generator
Execution engine

## Outline

- 1. Structure of compilers and interpreters
- 2. Creating ASTs
- 3. Processing ASTs
- 4. What I tried in an interprete
- 5. What I tried in a compile
- 6 Real work

Parser
Semantic analyzer
The curious case of C++
Optimizer
Code generator

### Parser

• converts a program from source code into a structured version of it

Parser
Semantic analyzer
The curious case of C++
Optimizer
Code generator

#### Parser

- converts a program from source code into a structured version of it
- produces an AST abstract syntax tree

Parser
Semantic analyzer
The curious case of C++
Optimizer
Code generator

#### Parser

- converts a program from source code into a structured version of it
- produces an AST abstract syntax tree
- detects invalid uses of language syntax

Semantic analyzer The curious case of C++ Optimizer Code generator Execution engine

Parser

#### Parser

- converts a program from source code into a structured version of it
- produces an AST abstract syntax tree
- detects invalid uses of language syntax
- if you want to hear more watch my previous talk, "Simple hand written parsers", or reach into books

Parser

Semantic analyzer

The curious case of C++

Optimizer

Code generator

Execution engine

# Semantic analyzer

converts an AST into an annotated form

Parser

Semantic analyzer

The curious case of C++

Optimizer

Code generator

Execution engine

- converts an AST into an annotated form
- resolves names

Parser
Semantic analyzer
The curious case of C+Optimizer
Code generator
Execution engine

- converts an AST into an annotated form
- resolves names
- type checks the program

Parser

Semantic analyzer

The curious case of C++

Optimizer

Code generator

Execution engine

- converts an AST into an annotated form
- resolves names
- type checks the program
- possibly validates further aspects of the language

Parser
Semantic analyzer
The curious case of C++
Optimizer
Code generator
Execution engine

- converts an AST into an annotated form
- resolves names
- type checks the program
- possibly validates further aspects of the language
- this is the subject of this talk

Parser
Semantic analyzer
The curious case of C++
Optimizer
Code generator
Execution engine

Parser
Semantic analyzer
The curious case of C++
Optimizer
Code generator
Execution engine

## The curious case of C++

a \* b;

Parser
Semantic analyzer
The curious case of C++
Optimizer
Code generator
Execution engine

```
struct a;
a * b;
```

Parser
Semantic analyzer
The curious case of C++
Optimizer
Code generator
Execution engine

## The curious case of C++

```
int a, b;
```

a \* b;

Parser
Semantic analyzer
The curious case of C++
Optimizer
Code generator
Execution engine

#### The curious case of C++

• C++ (same as C) has an at least context-sensitive grammar (type 1 in Chomsky hierarchy)

- C++ (same as C) has an at least context-sensitive grammar (type 1 in Chomsky hierarchy)
- some argue it is even unrestricted (type 0)

Parser
Semantic analyzer
The curious case of C++
Optimizer
Code generator
Execution engine

- C++ (same as C) has an at least context-sensitive grammar (type 1 in Chomsky hierarchy)
- some argue it is even unrestricted (type 0)
- parsing C++ without doing some semantic analysis at the same time is not something you want to do

- C++ (same as C) has an at least context-sensitive grammar (type 1 in Chomsky hierarchy)
- some argue it is even unrestricted (type 0)
- parsing C++ without doing some semantic analysis at the same time is not something you want to do
- typename and template disambiguators typename T::type, T::template foo<1>, tuple.template get<3>()

Parser
Semantic analyzer
The curious case of C++
Optimizer
Code generator
Execution engine

# Optimizer

simplifies the program

Parser
Semantic analyzer
The curious case of C++
Optimizer
Code generator
Execution engine

# Optimizer

- simplifies the program
- sounds simple, doesn't it?

Parser
Semantic analyzer
The curious case of C++
Optimizer
Code generator
Execution engine

## Optimizer

- simplifies the program
- sounds simple, doesn't it?
- I don't dare entertain the idea of giving a talk on this;)

Parser
Semantic analyzer
The curious case of C++
Optimizer
Code generator
Execution engine

# Code generator

• generates the code for consumption by the next program in the pipeline

Parser
Semantic analyzer
The curious case of C++
Optimizer
Code generator
Execution engine

## Code generator

- generates the code for consumption by the next program in the pipeline
- usually some flavor of assembly, or IR (thanks, LLVM!)

Parser
Semantic analyzer
The curious case of C++
Optimizer
Code generator
Execution engine

# Code generator

- generates the code for consumption by the next program in the pipeline
- usually some flavor of assembly, or IR (thanks, LLVM!)
- seems compiler-centric, but may also be included in an interpreter (think JITs)

Parser
Semantic analyzer
The curious case of C++
Optimizer
Code generator
Execution engine

# Code generator

- generates the code for consumption by the next program in the pipeline
- usually some flavor of assembly, or IR (thanks, LLVM!)
- seems compiler-centric, but may also be included in an interpreter (think JITs)
- probably (?) not as interesting as a talk than the previous ones?

Parser
Semantic analyzer
The curious case of C++
Optimizer
Code generator
Execution engine

# Execution engine

generally found in interpreters...

Parser
Semantic analyzer
The curious case of C++
Optimizer
Code generator
Execution engine

## Execution engine

- generally found in interpreters...
- ...but ask your local C++ compiler developer about constexpr

Parser
Semantic analyzer
The curious case of C++
Optimizer
Code generator
Execution engine

## Execution engine

- generally found in interpreters...
- ...but ask your local C++ compiler developer about constexpr
- this talk touches on this topic, but not significantly

Dynamic polymorphism Static polymorphism So what?

## Outline

- 1. Structure of compilers and interpreters
- 2. Creating ASTs
- 3. Processing AST:
- 4. What I tried in an interprete
- 5. What I tried in a compile
- 6 Real work

easy to extend dynamically

- easy to extend dynamically
- use virtual dispatch for visitation

- easy to extend dynamically
- use virtual dispatch for visitation
- or use dynamic casts for visitation

```
struct visitor {
    virtual bool visit(const assignment &) {}
};
struct expression {
    virtual ~expression() = default;
    virtual bool visit(visitor & v) = 0;
};
struct assignment : expression {
    virtual bool visit(visitor & v) override {
        return visitor.visit(*this);
```

```
struct expression {
    virtual ~expression() = default;
}:
struct assignment : expression {};
struct visitor {
    std::unordered_map<std::type_index, std::function<void (expression *)>> callbacks;
    bool visit(expression * expr) {
        return callbacks.at(typeid(expr))(expr);
    template<typename T, typename F>
    void add visitor(F && f) {
        callbacks.emplace(typeid(T), [f = std::forward<F>(f)](expression * ptr) {
            f(dynamic_cast<T *>(ptr));
        });
};
```

# Static polymorphism

harder to extend, with lots of recompilation

# Static polymorphism

- harder to extend, with lots of recompilation
- type-checked for handling most (all?) cases

# Static polymorphism

- harder to extend, with lots of recompilation
- type-checked for handling most (all?) cases
- variant-based

## Static polymorphism

```
using expression = variant<</pre>
    assignment,
    binary_expression.
    function call
void analyze_expression(const expression & expr) {
    fmap(expr, make_overload_set()
        [](const assignment &) { /* handle assignment */ },
        [](const binary_expression &) { /* handle binary expression */ },
        [](const function call &) { /* handle function call */ }
    )):
```

the way the AST is structured will affect how you consume it in analyzer

- the way the AST is structured will affect how you consume it in analyzer
- making the AST polymorphic would require either:

- the way the AST is structured will affect how you consume it in analyzer
- making the AST polymorphic would require either:
  - making the parser AST aware of the analyzer types, or

- the way the AST is structured will affect how you consume it in analyzer
- making the AST polymorphic would require either:
  - making the parser AST aware of the analyzer types, or
  - using the crazy dynamic visitation scheme to spit out analyzesr tree

Multiple passes DAG of AST elements Annotating in place Generating a new structure

#### Outline

- 1. Structure of compilers and interpreters
- 2. Creating ASTs
- 3. Processing ASTs
- 4. What I tried in an interprete
- 5. What I tried in a compile
- 6 Real world

ucture of compilers and interpreters Creating ASTs Processing ASTs What I tried in an interpreter What I tried in a compiler Real world

Single pass Multiple passes DAG of AST elements Annotating in place Generating a new structure

# Single pass

go in order, once

structure of compilers and interpreters
Creating ASTs
Processing ASTs
What I tried in an interpreter
What I tried in a compiler

Single pass Multiple passes DAG of AST elements Annotating in place Generating a new structure

# Single pass

- go in order, once
- set itself on fire if anything is not resolved

ucture of compilers and interpreters Creating ASTs Processing ASTs What I tried in an interpreter What I tried in a compiler Real world

Multiple passes
DAG of AST elements
Annotating in place
Generating a new structure

# Multiple passes

go in order, multiple times

# Multiple passes

- go in order, multiple times
- do a set number of multiple passes

# Multiple passes

- go in order, multiple times
- do a set number of multiple passes, or
- keep going until everything is resolved

## Multiple passes

- go in order, multiple times
- do a set number of multiple passes, or
- keep going until everything is resolved ...or report an error once a pass makes no progress

Single pass
Multiple passes
DAG of AST elements
Annotating in place
Generating a new structure

### DAG of AST elements

• build a tree of (possibly incomplete) elements

Single passs
Multiple passes
DAG of AST elements
Annotating in place
Generating a new structure

- build a tree of (possibly incomplete) elements
- gather information about the dependencies

Multiple passes
DAG of AST elements
Annotating in place
Generating a new structure

- build a tree of (possibly incomplete) elements
- gather information about the dependencies
- execute the DAG of dependencies

- build a tree of (possibly incomplete) elements
- gather information about the dependencies
- execute the DAG of dependencies
- (shared) futures are great for this!

Multiple passes

DAG of AST elements

Annotating in place

Generating a new structure

- build a tree of (possibly incomplete) elements
- gather information about the dependencies
- execute the DAG of dependencies
- (shared) futures are great for this!
- coroutines will make them better at it

Structure of compilers and interpreters Creating ASTs Processing ASTs What I tried in an interpreter What I tried in a compiler Real world

Single pass
Multiple passes
DAG of AST elements
Annotating in place
Generating a new structure

# Annotating in place

avoids constant reallocation

Single pass
Multiple passes
DAG of AST elements
Annotating in place
Generating a new structure

## Annotating in place

- avoids constant reallocation
- possibly easier to debug, since addresses don't change

Single pass
Multiple passes
DAG of AST elements
Annotating in place
Generating a new structure

## Annotating in place

- avoids constant reallocation
- possibly easier to debug, since addresses don't change
- to an extend a necessity for C++, since you are analyzing during parsing

ructure of compilers and interpreters Creating ASTs Processing ASTs What I tried in an interpreter What I tried in a compiler Real world

Single pass
Multiple passes
DAG of AST elements
Annotating in place
Generating a new structure

## Generating a new structure

• allows for a higher degree of encapsulation between the parser and the analyzer

Structure of compilers and interpreters
Creating ASTs
Processing ASTs
What I tried in an interpreter
What I tried in a compiler

Single pass
Multiple passes
DAG of AST elements
Annotating in place
Generating a new structure

## Generating a new structure

- allows for a higher degree of encapsulation between the parser and the analyzer
- allows for more flexibility with the analyzed structure

#### Outline

- 1. Structure of compilers and interpreter
- 2. Creating ASTs
- 3. Processing ASTs
- 4. What I tried in an interpreter
- 5. What I tried in a compile
- 6. Real world

# Language

```
foo = bar
baz = "123"
fizz = buzz(bar, "123")
bar = "this is " + "text"
```

• string literals: create an element of the analyzed AST

- string literals: create an element of the analyzed AST
- names: try to resolve

- string literals: create an element of the analyzed AST
- names: try to resolve
  - if the name is known mark resolved and refer to the actual value of the name

- string literals: create an element of the analyzed AST
- names: try to resolve
  - if the name is known mark resolved and refer to the actual value of the name
  - if the name is not know save the name for later and return a "delayed expression"

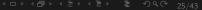
- string literals: create an element of the analyzed AST
- names: try to resolve
  - if the name is known mark resolved and refer to the actual value of the name
  - if the name is not know save the name for later and return a "delayed expression"
- addition: analyze operands

- string literals: create an element of the analyzed AST
- names: try to resolve
  - if the name is known mark resolved and refer to the actual value of the name
  - if the name is not know save the name for later and return a "delayed expression"
- addition: analyze operands
  - if both operands are resolved evaluate the operation

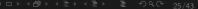
- string literals: create an element of the analyzed AST
- names: try to resolve
  - if the name is known mark resolved and refer to the actual value of the name
  - if the name is not know save the name for later and return a "delayed expression"
- addition: analyze operands
  - if both operands are resolved evaluate the operation
  - if any one of the operands are not resolved save the arguments for later and return a (different kind of a) "delayed expression"

- string literals: create an element of the analyzed AST
- names: try to resolve
  - if the name is known mark resolved and refer to the actual value of the name
  - if the name is not know save the name for later and return a "delayed expression"
- addition: analyze operands
  - if both operands are resolved evaluate the operation
  - if any one of the operands are not resolved save the arguments for later and return a (different kind of a) "delayed expression"
- type instantiation: analyze the type name, analyze the arguments

- string literals: create an element of the analyzed AST
- names: try to resolve
  - if the name is known mark resolved and refer to the actual value of the name
  - if the name is not know save the name for later and return a "delayed expression"
- addition: analyze operands
  - if both operands are resolved evaluate the operation
  - if any one of the operands are not resolved save the arguments for later and return a (different kind of a) "delayed expression"
- type instantiation: analyze the type name, analyze the arguments
  - if everything is already resolved evaluate the instantiation



- string literals: create an element of the analyzed AST
- names: try to resolve
  - if the name is known mark resolved and refer to the actual value of the name
  - if the name is not know save the name for later and return a "delayed expression"
- addition: analyze operands
  - if both operands are resolved evaluate the operation
  - if any one of the operands are not resolved save the arguments for later and return a (different kind of a) "delayed expression"
- type instantiation: analyze the type name, analyze the arguments
  - if everything is already resolved evaluate the instantiation
  - if anything is not resolved yet create a (yet another kind of a) "delayed expression"



#### Initial pass

```
struct delayed instantiation info {
    type_identifier actual_type;
    std::vector<std::shared_ptr<variable>> arguments;
struct _delayed_reference_info {
    std::vector<std::u32string> referenced_id_expression;
struct _delayed_type_info {
    std::vector<std::u32string> type_name;
    std::vector<std::shared_ptr<variable>> arguments;
struct _delayed_operation_info {
    std::shared ptr<variable> lhs:
    std::shared_ptr<variable> rhs:
    operation_type operation;
std::variant
    std::shared_ptr<variable>,
    _delayed_instantiation_info,
    _delayed_reference_info,
    _delayed_type_info,
    delayed operation info
> _state;
```

• save the count of currently unresolved expressions

- save the count of currently unresolved expressions
- iterate through the unresolved expressions, asking them to try to resolve themselves

- save the count of currently unresolved expressions
- iterate through the unresolved expressions, asking them to try to resolve themselves
- if the count of currently unresolved expressions is different than the saved one, continue

- save the count of currently unresolved expressions
- iterate through the unresolved expressions, asking them to try to resolve themselves
- if the count of currently unresolved expressions is different than the saved one, continue
- validate if all expressions are resolved

```
std::size_t previous = 0;
while (previous != ctx.unresolved.size()) {
    previous = ctx.unresolved.size();
    for (auto && u : ctx.unresolved) {
        u.first->try_resolve(ctx);
  (!ctx.unresolved.emptv()) {
    throw analysis_error{};
```

#### Outline

- 1. Structure of compilers and interpreter
- 2. Creating ASTs
- 3. Processing AST:
- 4. What I tried in an interpreter
- 5. What I tried in a compiler
- 6 Real work

# Attempt #1: mimic the parser

```
struct expression {
    range_type range;
    variant<
        literal<lexer::token_type::string>,
        literal<lexer::token_type::integer>.
        postfix_expression,
        import_expression,
        recursive_wrapper<lambda_expression>,
        recursive_wrapper<unary_expression>,
        recursive_wrapper<br/>binary_expression>
      expression_value;
};
```

#### Attempt #1: mimic the parser

```
class postfix expression { public: void analyze(): /* ... */ };
class unary_expression { public: void analyze(); /* ... */ };
class binary_expression { public: void analyze(); /* ... */ };
using expression = variant<</pre>
    literal.
    variable.
    import_expression,
    postfix_expression,
    unary_expression,
    binary_expression
```

# Attempt #1: mimic the parser

```
class postfix_expression { public: void analyze(); /* ... */ };
class unary_expression { public: void analyze(); /* ... */ };
class binary_expression { public: void analyze(); /* ... */ };
using expression = variant<</pre>
    literal.
    variable.
    import_expression,
    postfix_expression,
    unary_expression,
    binary_expression
void analyze(expression & expr) {
    fmap(expr, [](auto && v) { v.analyze(); return unit{}; });
```

- pass #1: create the objects, setup scopes so that name lookup works
- "pass" #2: go top-down analyzing dependencies of all AST nodes, and then analyzing themselves

clear inheritance hierarchy that makes repetition go away

- clear inheritance hierarchy that makes repetition go away
- expressions are statements

- clear inheritance hierarchy that makes repetition go away
- expressions are statements
- similar kinds of expressions can extend other kinds

- clear inheritance hierarchy that makes repetition go away
- expressions are statements
- similar kinds of expressions can extend other kinds
- e.g. identifier is an expression\_ref

Attempt #1: mimic the parser Attempt #2: OO-based polymorphism Making sense through dependency graphs

## Attempt #2: OO-based polymorphism

• three main base classes: statement, type, and...

• three main base classes: statement, type, and... variable

- three main base classes: statement, type, and... variable
- expressions had associated variables that represented their value

- three main base classes: statement, type, and... variable
- expressions had associated variables that represented their value
- a bunch of classes interconverting between one and the other: expression\_variable

- three main base classes: statement, type, and... variable
- expressions had associated variables that represented their value
- a bunch of classes interconverting between one and the other:
   expression\_variable, variable\_expression

- three main base classes: statement, type, and... variable
- expressions had associated variables that represented their value
- a bunch of classes interconverting between one and the other: expression\_variable, variable\_expression, expression\_ref\_variable

- three main base classes: statement, type, and... variable
- expressions had associated variables that represented their value
- a bunch of classes interconverting between one and the other: expression\_variable, variable\_expression, expression\_ref\_variable, variable\_ref\_expression

- three main base classes: statement, type, and... variable
- expressions had associated variables that represented their value
- a bunch of classes interconverting between one and the other: expression\_variable, variable\_expression, expression\_ref\_variable, variable\_ref\_expression
- got rid of variable, never been happier (should've done it much sooner)

```
class expression {
    void set_type(type *);
    type * get_type() const;
    @memoized
    virtual future<void> analyze(analysis_context &);
};
```

```
class expression {
    void set_type(type *);
    type * get_type() const;
    memoized
    virtual future<void> analyze(analysis_context &);
};
class binary_expression : public expression {
    std::unique_ptr<expression> lhs;
    std::unique_ptr<expression> rhs;
    operator op:
    std::unique_ptr<expression> call_expr;
    future<void> analyze(analysis context &) override:
};
```

```
future<void> binary_expression::analyze(analysis_context & ctx)
{
    return when_all(
        lhs->analyze(ctx),
        rhs->analyze(ctx)
)
```

```
future<void> binary_expression::analyze(analysis_context & ctx)
    return when_all(
        lhs->analyze(ctx),
        rhs->analyze(ctx)
    .then([&]{
        return resolve_overload(ctx, lhs.get(), rhs.get(), op);
    7)
    .then([&](std::unique_ptr<expression> call){
        call_expr = std::move(call);
        return call_expr->analyze(ctx):
    })
```

```
future<void> binary_expression::analyze(analysis_context & ctx)
   return when_all(
        lhs->analyze(ctx),
        rhs->analyze(ctx)
    .then([&]{
        return resolve_overload(ctx, lhs.get(), rhs.get(), op);
   7)
    .then([&](std::unique_ptr<expression> call){
        call_expr = std::move(call);
        return call_expr->analyze(ctx):
   })
    .then([&]{
        set_type(call_expr->get_type());
```

```
function factorial(i : int) -> int
{
    if (i == 1) { return 1; }
    return i * factorial(i - 1);
}
```

```
function factorial(i : int) -> int
{
    if (i == 1) { return 1; }
    return i * factorial(i - 1);
}
```

• to analyze factorial, we need to analyze its signature and its body

```
function factorial(i : int) -> int
{
    if (i == 1) { return 1; }
    return i * factorial(i - 1);
}
```

- to analyze factorial, we need to analyze its signature and its body
- to analyze the call to factorial, we need to know what its signature is

```
function factorial(i : int)
{
    if (i == 1) { return 1; }
    return i * factorial(i - 1);
}
```

- to analyze factorial, we need to analyze its signature and its body
- to analyze the call to factorial, we need to know what its signature is

```
function factorial(i : int)
{
    if (i == 1) { return 1; }
    return i * factorial(i - 1);
}
```

- to analyze factorial, we need to analyze its signature and its body
- to analyze the call to factorial, we need to know what its signature is
- now we to solve a system of equations:
   return\_type\_of(factorial) == type\_of(i \* factorial(i 1))

```
future < void > function_definition::analyze(analysis_context & ctx)
{
    return when_all(fmap(parameters, [&](auto && param){
        return param = > analyze(ctx);
    }))
```

```
future < void > function_definition::analyze(analysis_context & ctx)
{
    return when_all(fmap(parameters, [&](auto && param){
        return param > analyze(ctx);
    }))
    .then([&]{
        return return_type > analyze(ctx);
    })
```

```
future < void > function_definition::analyze(analysis_context & ctx)
{
    return when all (fmap(parameters, [&] (auto && param){
        return param->analyze(ctx);
    }))
        .then([&]{
            return return_type->analyze(ctx);
        })
        .then([&]{
            return body->analyze(ctx);
        });
```

#### Outline

- 1. Structure of compilers and interpreters
- 2. Creating AST
- 3. Processing AST
- 4. What I tried in an interprete
- 5. What I tried in a compile
- 6. Real world

#### Clang

```
Decl *Parser::ParseExportDeclaration() {
  assert(Tok.is(tok::kw_export));
  SourceLocation ExportLoc = ConsumeToken();
  ParseScope ExportScope(this, Scope::DeclScope):
  Decl *ExportDecl = Actions.ActOnStartExportDecl(getCurScope(), ExportLoc,
      Tok.is(tok::1_brace) ? Tok.getLocation() : SourceLocation());
  if (Tok.isNot(tok::1 brace)) {
    ParseExternalDeclaration(Attrs):
    return Actions.ActOnFinishExportDecl(getCurScope(), ExportDecl,
                                         SourceLocation()):
  while (/* ... SNIP ... */) {
      ParseExternalDeclaration(Attrs):
  return Actions.ActOnFinishExportDecl(getCurScope(), ExportDecl.
                                       T.getCloseLocation()):
```

Clang

#### Clang

```
Decl *Sema::ActOnStartExportDecl(Scope *S, SourceLocation ExportLoc,
                                 SourceLocation LBraceLoc) {
  ExportDecl *D = ExportDecl::Create(Context, CurContext, ExportLoc);
  if (ModuleScopes.empty() | !ModuleScopes.back().ModuleInterface)
    Diag(ExportLoc, diag::err_export_not_in_module_interface);
  if (D->isExported())
    Diag(ExportLoc, diag::err_export_within_export);
  CurContext->addDecl(D):
  PushDeclContext(S, D);
  D->setModuleOwnershipKind(Decl::ModuleOwnershipKind::VisibleWhenImported):
  return D:
```

Clang

#### Clang

```
/// Complete the definition of an export declaration.
Decl *Sema::ActOnFinishExportDecl(Scope *S, Decl *D,
                                  SourceLocation RBraceLoc) {
  auto *ED = cast<ExportDecl>(D);
  if (RBraceLoc.isValid())
    ED->setRBraceLoc(RBraceLoc):
  // FIXME: Diagnose export of internal-linkage declaration (including
  PopDeclContext():
  return D:
```

ructure of compilers and interpreters Creating ASTs Processing ASTs What I tried in an interpreter What I tried in a compiler Real world

# From parsing to sema: making sense of syntax trees

Michał Dominiak Nokia Networks griwes@griwes.info