

Clang compiler frontend

Ivan Murashko

October 8, 2022

Contents

1	Introduction	5
2	Environment setup	7
2.1	Source code compilation	7
2.1.1	Getting the source code	7
2.1.2	Configuration with cmake	7
2.1.3	Build	8
3	Architecture	9
3.1	Clang driver overview	9
3.2	Clang frontend overview	10
4	Features	13
4.1	Precompiled headers	13
4.1.1	User guide	13
4.2	Modules	15
4.2.1	User guide	15
4.2.2	Implicit modules	15
4.2.3	Modules internals	17
4.3	Header-Map files	17
5	Tools	19
5.1	clang-tidy	19
5.2	clangd	19
	Index	21

Chapter 1

Introduction

The `clang` is C/C++ and ObjC compiler. It's an integral part of LLVM project. When we say about clang we can say about 2 different things. The first one is the compiler frontend i.e. the part of compiler that is responsible for parsing and semantic reasoning about the program. We also use the word `clang` when we say about the compiler itself. It's also referred as compiler driver. The driver is responsible for compiler invocation i.e. it can be considered as a manger that calls different parts of the compiler such as the compiler frontend as well as other parts that are required for successful compilation (middle-end, back-end, assembler, linker).

The book is mostly focused on the `clang` compiler frontend but it also includes some other relevant parts of LLVM that are critical for the frontend internals.

The book is separated into several parts. The first one provides basic info about LLVM project and how it can be installed. It also describes useful development tools and configurations used for LLVM code exploration later in the book.

Second part is about `clang` architecture and its place inside LLVM.

The primary part of the book is about different features of clang such as C++ modules, header maps and others.

The `clang` follows primary paradigm of LLVM - everything is a library, that allows to create a bunch of different tools. Some of them are also described: clang-tidy, clangd and others

The book also includes a lot of different examples that can be found at [\[3\]](#).

Chapter 2

Environment setup

The chapter describes basic steps to be done to setup the environment be used for future experiments with clang.

2.1 Source code compilation

we are going to compile our source code in debug mode to be suitable for future investigations with debugger.

2.1.1 Getting the source code

The clang source is a part of LLVM. You can get it with the following command

```
git clone https://github.com/llvm/llvm-project.git
cd llvm-project
```

2.1.2 Configuration with cmake

Create a build folder where the compiler and related tools will be built

```
mkdir build
cd build
```

Run configure script

```
cmake -G Ninja -DCMAKE_BUILD_TYPE=RelWithDebInfo
↪ -DLLVM_TARGETS_TO_BUILD="X86"
↪ -DLLVM_ENABLE_PROJECTS="clang;clang-tools-extra"
↪ -DLLVM_USE_LINKER=gold -DLLVM_USE_SPLIT_DWARF=ON ../llvm
```

There are several options specified:

- `-DLLVM_TARGETS_TO_BUILD="X86"` specifies exact targets to be build. It will avoid build unnecessary targets
- `LLVM_ENABLE_PROJECTS="clang;clang-tools-extra"` specifies LLVM projects that we care about
- `LLVM_USE_LINKER=gold` - uses gold linker
- `LLVM_USE_SPLIT_DWARF=ON` - splits debug information into separate files. This option saves disk space as well as memory consumption during the LLVM build. The option require compiler used for clang build to support it

For debugging purposes you might want to change the `-DCMAKE_BUILD_TYPE` into `Debug`. Thus your overall config command will look like

```
cmake -G Ninja -DCMAKE_BUILD_TYPE=Debug
↪ -DLLVM_TARGETS_TO_BUILD="X86"
↪ -DLLVM_ENABLE_PROJECTS="clang;clang-tools-extra"
↪ -DLLVM_USE_LINKER=gold -DLLVM_USE_SPLIT_DWARF=ON ../llvm
```

2.1.3 Build

The build is trivial

```
ninja clang
```

You can also run unit and end-to-end tests for the compiler with

```
ninja check-clang
```

The compiler binary can be found as `bin/clang` at the build folder.

Chapter 3

Architecture

You can find some info about clang internal architecture and relation with other LLVM components.

3.1 Clang driver overview

When we spoke about **clang** we should separate 2 things:

- driver
- compiler frontend

Both of them are called **clang** but perform different operations. The driver invokes different stages of compilation process (see fig. 3.1). The stages are standard for ordinary compiler and nothing special is there:

- Frontend: it does lexical analysis and parsing.
- Middle-end: it does different optimization on the intermediate representation (LLVM-IR) code
- Backend: Native code generation
- Assembler: Running assembler



Figure 3.1: Clang driver



Figure 3.2: Clang frontend components

- Linker: Running linker

The driver is invoked by the following command

```
clang main.cpp -o main -lstdc++
```

The driver also adds a lot of additional arguments, for instance search paths for system includes that could be platform specific. You can use `-### clang` option to print actual command line used by the driver

```
$clang -### main.cpp -o main -lstdc++
clang version 12.0.1 (Fedora 12.0.1-1.fc34)
Target: x86_64-unknown-linux-gnu
Thread model: posix
InstalledDir: /usr/bin
"/usr/bin/clang-12" "-cc1" "-triple"
↳ "x86_64-unknown-linux-gnu" "-emit-obj" "-mrelax-all" ...
```

3.2 Clang frontend overview

One may see that the clang compiler toolchain corresponds the pattern wildly described at different compiler books [4]. Despite the fact, the frontend part is quite different from a typical compiler frontend described at the books. The primary reason for this is C++ language and it's complexity. Some features (macros) can change the source code itself another (typedef) can effect on token kind. As result the relations between different frontend components can be represented as it's shown on fig 3.2

You can bypass the driver call and run the clang frontend with `-cc1` option. We are going to use it for our future experiments with one simple program:

```
int main() {
    return 0;
}
```

The first part is the **Lexer**. It's primary goal is to convert the input program into a stream of tokens. The token stream can be printed with `-dump-tokens` options as follows

```
$ clang -cc1 -dump-tokens src/simple/main.cpp
```

The output of the command is below

```
int 'int'          [StartOfLine]  Loc=<src/simple/main.cpp:1:1>
identifier 'main'  [LeadingSpace]
↳ Loc=<src/simple/main.cpp:1:5>
l_paren '('        Loc=<src/simple/main.cpp:1:9>
r_paren ')'        Loc=<src/simple/main.cpp:1:10>
l_brace '{'        [LeadingSpace] Loc=<src/simple/main.cpp:1:12>
return 'return'    [StartOfLine] [LeadingSpace]
↳ Loc=<src/simple/main.cpp:2:3>
numeric_constant '0' [LeadingSpace]
↳ Loc=<src/simple/main.cpp:2:10>
semi ';'          Loc=<src/simple/main.cpp:2:11>
r_brace '}'        [StartOfLine]  Loc=<src/simple/main.cpp:3:1>
eof ''            Loc=<src/simple/main.cpp:3:2>
```

Another part is the **Parser**. It produces AST that can be shown with the following command

```
$ clang -cc1 -ast-dump src/simple/main.cpp
```

The output of the command is below

```
TranslationUnitDecl 0x5581925de3b8 <<invalid sloc>> <invalid
↳ sloc>
|-TypedefDecl 0x5581925dec20 <<invalid sloc>> <invalid sloc>
↳ implicit __int128_t '__int128'
|  `--BuiltinType 0x5581925de980 '__int128'
|-TypedefDecl 0x5581925dec90 <<invalid sloc>> <invalid sloc>
↳ implicit __uint128_t 'unsigned __int128'
|  `--BuiltinType 0x5581925de9a0 'unsigned __int128'
|-TypedefDecl 0x5581925df008 <<invalid sloc>> <invalid sloc>
↳ implicit __NSConstantString '__NSConstantString_tag'
|  `--RecordType 0x5581925ded80 '__NSConstantString_tag'
|     `--CXXRecord 0x5581925dece8 '__NSConstantString_tag'
|-TypedefDecl 0x5581925df0a0 <<invalid sloc>> <invalid sloc>
↳ implicit __builtin_ms_va_list 'char *'
```

```

|  `~PointerType 0x5581925df060 'char *'
|    `~BuiltinType 0x5581925de460 'char'
|~TypedefDecl 0x5581926237d8 <<invalid sloc>> <invalid sloc>
|  ↪ implicit __builtin_va_list '__va_list_tag[1]'
|  `~ConstantArrayType 0x558192623780 '__va_list_tag[1]' 1
|    `~RecordType 0x5581925df190 '__va_list_tag'
|      `~CXXRecord 0x5581925df0f8 '__va_list_tag'
|~FunctionDecl 0x558192623880 <src/simple/main.cpp:1:1,
|  ↪ line:3:1> line:1:5 main 'int ()'
|    `~CompoundStmt 0x5581926239c0 <col:12, line:3:1>
|      `~ReturnStmt 0x5581926239b0 <line:2:3, col:10>
|        `~IntegerLiteral 0x558192623990 <col:10> 'int' 0

```

Chapter 4

Features

You can find some info about different clang features at the chapter

4.1 Precompiled headers

Precompiled headers or **pch** is a clang feature that was designed with the goal to improve clang frontend performance. The basic idea was to create AST for a header file and reuse the AST for some purposes.

4.1.1 User guide

Generate you pch file is simple [2]. Suppose you have a header file with name **header.h**:

```
#pragma once

void foo() {
}
```

then you can generate a pch for it with

```
clang -x c++-header header.h -o header.pch
```

the option **-x c++-header** was used there. The option says that the header file has to be treated as a c++ header file. The output file is **header.pch**.

The precompiled headers generation is not enough and you may want to start using them. Typical C++ source file that uses the header may look like

```
// test pchs

#include "header.h"

int main() {
    foo();
    return 0;
}
```

As you may see, the header is included as follows

```
...
#include "header.h"
...
```

By default clang will not use a pch at the case and you have to specify it explicitly with

```
clang -include-pch header.pch main.cpp -o main -lstdc++
```

We can check the command with debugger and it will give us

```
$ lladb ~/local/llvm-project/build/bin/clang -- -cc1
↳ -include-pch header.pch main.cpp -fsyntax-only
...
(lladb) b clang::ASTReader::ReadAST
...
(lladb) r
...
4231   llvm::SaveAndRestore<SourceLocation>
-> 4232       SetCurImportLocRAII(CurrentImportLoc, ImportLoc);
4233   llvm::SaveAndRestore<Optional<ModuleKind>>
↳   SetCurModuleKindRAII(
4234       CurrentDeserializingModuleKind, Type);
4235
(lladb) p FileName
(llvm::StringRef) $0 = (Data = "header.pch", Length = 10)
```

Note that only the first `--include-pch` option will be processed, all others will be ignored. It reflects the fact that there can be only one precompiled header for a translation unit.

4.2 Modules

Modules can be considered as a next step in evolution of precompiled headers. They also represent an parsed AST in binary form but form a DAG (tree) i.e. one module can include more than one another module ¹

4.2.1 User guide

The C++20 standard [1] introduced 2 concepts related to modules. The first one is ordinary modules described at section 10 of [1]. Another one is so call header unit that is mostly described at section 15.5. The header units can be considered as an intermediate step between ordinary headers and modules and allow to use `import` directive to import ordinary headers. The second approach was the main approach for modules implemented in clang and we will call it as **implicit modules**. The first one (primary one described at [1]) will be call **explicit modules**. We will start with the implicit modules first.

4.2.2 Implicit modules

The key point for implicit clang modules is `modulemap` file. It describes relation between different modules and interface provided by the modules. The default name for the file is `module.modulemap`. Typical content is the following

```
module header1 {
  header "header1.h"
  export *
}
```

The header paths in the modulemap file has to be ether absolute or relative on the module map file location. Thus compiler should have a chance to find them to compile.

There are 2 options to process the configuration file: explicit or implicit. The first one (explicit) assumes that you pass it via `-fmodule-map-file=<path to modulemap file>`. The second one (default) will search for modulemap files implicitly and apply them. You can turn off the behaviour with `-fno-implicit-module-maps` command line argument.

¹Compare that with precompiled header where only one precompiled header can be introduced for each compilation unit

Explicit modules

TBD

Some problems related to modules

The code that uses modules can introduce some non trivial behaviour of your program. Consider the project that consists of two headers.

header1.h:

```
#pragma once
```

```
int h1 = 1;
```

header2.h:

```
#pragma once
```

```
int h2 = 2;
```

The header1.h is included into the main.cpp

```
#include <iostream>
```

```
#include "header1.h"
```

```
int main() {
    std::cout << "Header1 value: " << h1 << std::endl;
    std::cout << "Header2 value: " << h2 << std::endl;
}
```

The code will not compile

```
clang++ -std=c++20 -fconstexpr-depth=1271242 main.cpp -o main
```

```
↪ -lstdc++
```

```
main.cpp:7:37: error: use of undeclared identifier 'h2'
```

```
    std::cout << "Header2 value: " << h2 << std::endl;
```

```
    ^
```

```
1 error generated.
```

but if you use the following module.modulemap file then it will compile with modules


```
module h1 {  
    header "header1.h"  
    export *  
    module h2 {  
        header "header2.h"  
        export *  
    }  
}
```

The example shows how the visibility scope can be leaked when modules are used in the project.

4.2.3 Modules internals

Modules are processed inside `clang::Preprocessor::HandleIncludeDirective`. There is a `clang::Preprocessor::HandleHeaderIncludeOrImport` method.

The module is loaded by `clang::CompilerInstance::loadModuleFile`. The method calls `clang::CompilerInstance::findOrCompileModuleAndReadAST`

4.3 Header-Map files

TBD

Chapter 5

Tools

There are some tools created on clang related libs

5.1 clang-tidy

TBD

5.2 clangd

TBD

Index

AST, [11](#), [13](#)

clang, [5](#)

compiler frontend, [5](#)

header unit, [15](#)

Bibliography

- [1] *C++20*. 2021.
- [2] Clang compiler user's manual, 2022.
- [3] I. Murashko. Source code examples for clang compiler frontend, 2022.
- [4] L. Torczon and K. Cooper. *Engineering A Compiler*. Elsevier Inc., 2nd edition, 2012.