

# Clang compiler frontend and Clang tools

Understand internals of a top-rated C/C++ compiler frontend and create your own tools <sup>1</sup>

**Ivan Murashko**  
August 13, 2023

**Packt>**

BIRMINGHAM—MUMBAI

---

<sup>1</sup>This hasn't been finalized yet

# Clang compiler frontend and Clang tools

Copyright © 2022 Packt Publishing

*All rights reserved.* No part of this book may be reproduced, stored in a retrieval system, or transmitted in any form or by any means, without the prior written permission of the publisher, except in the case of brief quotations embedded in critical articles or reviews.

Every effort has been made in the preparation of this book to ensure the accuracy of the information presented. However, the information contained in this book is sold without warranty, either express or implied. Neither the authors, nor Packt Publishing or its dealers and distributors, will be held liable for any damages caused or alleged to have been caused directly or indirectly by this book. Packt Publishing has endeavored to provide trademark information about all of the companies and products mentioned in this book by the appropriate use of capitals. However, Packt Publishing cannot guarantee the accuracy of this information.

**Group Product Manager:** TBD

**Publishing Product Manager:** Kunal Sawant

**Content Development Editor:** Rosal Colaco

**Technical Editor:** TBD

**Copy Editor:** Safis Editing

**Project Coordinator:** Manisha Singh

**Proofreader:** Safis Editing

**Indexer:** Subalakshmi Govindhan

**Production Designer:** Sinhayna Bais

**Marketing Coordinator:** TBD

First published: March 2023

Production reference: 1260522

Published by Packt Publishing Ltd.

Livery Place

35 Livery Street

Birmingham

B3 2PB, UK.

ISBN 978-1-80461-383-2

[www.packt.com](http://www.packt.com)

# Contributors

## About the author

**Ivan Murashko** is a C++ software engineer with a PhD in physics from Peter the Great St. Petersburg Polytechnic University in Russia. His scientific interests include quantum optics, quantum information theory, and quantum computation. Ivan has more than 20 years of experience with C/C++ programming and has worked on projects involving cryptography, search engines, VoIP, and compilers. In addition to these areas, he also has an interest in functional languages such as Haskell and Scala, as well as category theory and logic. Since 2020, Ivan has been working with LLVM, focusing on the clang compiler frontend and code analysis and modification tools.

TBD<sup>2</sup>

## About the reviewer

**TBD** has been developing software for nearly 20 years, of which Python has always been a major element. He's been a core developer of the Python language since 2010. Tal holds a B.Sc. in Math & Physics from Tel Aviv University. He likes hiking, computer games, philosophical sci-fi, and spending time with his family.

For the past eight years, Tal has been developing educational technology, first at *Compedia* where he built a group developing VR and AR education apps, and later at the startup *FullProof* of which he was a co-founder.

Tal currently works at *Rhino Health*, a startup working to enable development and use of medical AI models with patient data from across the globe while preserving patient privacy.<sup>3</sup>

## Share Your Thoughts

Once you've read **Clang compiler frontend and Clang tools**, we'd love to hear your thoughts! Scan the QR code below to go straight to the Amazon review page for this book and share your feedback.

---

<sup>2</sup>This is an example of dedication, please edit it.

<sup>3</sup>TBD



<https://packt.link/r/1-804-61383-5>

Your review is important to us and the tech community and will help us make sure we're delivering excellent quality content.

Replace **NameOfTheProduct** with your title name.

Replace the dummy QR Code image with your product QR code image.

Replace *ISBN10P* with the 10-digit *ISBN-10P* from EPIC.

This page should be added after the Front Matter.

# Table Of Contents

|          |  |           |
|----------|--|-----------|
| <b>1</b> | <b>Environment setup</b>                             | <b>15</b> |
| 1.1      | Technical requirements . . . . .                     | 15        |
| 1.1.1    | CMake as build configuration tool . . . . .          | 16        |
| 1.1.2    | Ninja as build tool . . . . .                        | 16        |
| 1.2      | Getting to know LLVM . . . . .                       | 17        |
| 1.2.1    | Short LLVM history . . . . .                         | 17        |
| 1.2.2    | OS support . . . . .                                 | 18        |
| 1.2.2.1  | Linux . . . . .                                      | 18        |
| 1.2.2.2  | Darwin . . . . .                                     | 18        |
| 1.2.2.3  | Windows . . . . .                                    | 18        |
| 1.2.3    | LLVM/clang project structure . . . . .               | 19        |
| 1.3      | Source code compilation . . . . .                    | 21        |
| 1.3.1    | Configuration with CMake . . . . .                   | 21        |
| 1.3.2    | Build . . . . .                                      | 23        |
| 1.3.3    | The LLVM debugger, its build and usage . . . . .     | 24        |
| 1.4      | Test project: syntax check with clang tool . . . . . | 26        |
| 1.5      | Summary . . . . .                                    | 31        |
| 1.6      | Further reading . . . . .                            | 32        |
| <b>2</b> | <b>Clang architecture</b>                            | <b>33</b> |
| 2.1      | Getting started with compilers . . . . .             | 33        |

|          |                                      |           |
|----------|--------------------------------------|-----------|
| 2.1.1    | Explore compiler workflow            | 34        |
| 2.1.2    | Frontend                             | 35        |
| 2.1.2.1  | Lexer                                | 38        |
| 2.1.2.2  | Parser                               | 38        |
| 2.1.2.3  | Codegen                              | 41        |
| 2.2      | Clang driver overview                | 42        |
| 2.2.1    | Example program                      | 42        |
| 2.2.2    | Compilation phases                   | 43        |
| 2.2.3    | Tools execution                      | 43        |
| 2.2.4    | Combine all together                 | 46        |
| 2.2.5    | Debugging clang                      | 48        |
| 2.3      | Clang frontend overview              | 51        |
| 2.3.1    | Frontend action                      | 51        |
| 2.3.2    | Preprocessor                         | 53        |
| 2.3.3    | Parser and Sema                      | 55        |
| 2.4      | Summary                              | 61        |
| 2.5      | Further reading                      | 61        |
| <b>3</b> | <b>Clang AST</b>                     | <b>63</b> |
| 3.1      | AST                                  | 63        |
| 3.1.1    | Statements                           | 64        |
| 3.1.2    | Declarations                         | 65        |
| 3.1.3    | Types                                | 65        |
| 3.2      | AST traversal                        | 67        |
| 3.2.1    | DeclVisitor test tool                | 67        |
| 3.2.2    | Visitor implementation               | 73        |
| 3.3      | Recursive AST Visitor                | 76        |
| 3.4      | AST matchers                         | 80        |
| 3.5      | Explore clang AST with clang-query   | 85        |
| 3.6      | Processing AST in the case of errors | 86        |
| 3.7      | Summary                              | 88        |
| 3.8      | Further reading                      | 88        |

|          |   |            |
|----------|---|------------|
| <b>4</b> | <b>Basic libraries and tools</b>                    | <b>89</b>  |
| 4.1      | LLVM coding style . . . . .                         | 89         |
| 4.2      | Basic libs . . . . .                                | 89         |
| 4.3      | TableGen . . . . .                                  | 90         |
| 4.4      | LLVM test framework . . . . .                       | 90         |
| 4.5      | Summary . . . . .                                   | 90         |
| 4.6      | Further reading . . . . .                           | 90         |
| <b>5</b> | <b>clang-tidy linter framework</b>                  | <b>93</b>  |
| 5.1      | clang-tidy overview and usage examples . . . . .    | 93         |
| 5.2      | clang-tidy internal design . . . . .                | 93         |
| 5.3      | Custom clang-tidy check . . . . .                   | 93         |
| 5.4      | Results in the case of compilation errors . . . . . | 94         |
| <b>6</b> | <b>Advanced code analysis</b>                       | <b>95</b>  |
| 6.1      | Usage cases . . . . .                               | 95         |
| 6.2      | CFG and life time analysis . . . . .                | 95         |
| 6.3      | Custom CFG check . . . . .                          | 95         |
| <b>7</b> | <b>Refactoring tools</b>                            | <b>97</b>  |
| 7.1      | Code modification and clang-format . . . . .        | 97         |
| 7.2      | Custom code modification tool . . . . .             | 97         |
| <b>8</b> | <b>IDE support and code navigation</b>              | <b>99</b>  |
| 8.1      | VSCode and LSP . . . . .                            | 99         |
| 8.2      | clangd internals . . . . .                          | 99         |
| 8.3      | Custom extension for LSP . . . . .                  | 100        |
| <b>9</b> | <b>Features</b>                                     | <b>103</b> |
| 9.1      | Precompiled headers . . . . .                       | 103        |
| 9.1.1    | User guide . . . . .                                | 103        |
| 9.2      | Modules . . . . .                                   | 105        |
| 9.2.1    | User guide . . . . .                                | 105        |

|         |  |     |
|---------|--|-----|
| 9.2.2   | Implicit modules . . . . .                 | 106 |
| 9.2.2.1 | Explicit modules . . . . .                 | 106 |
| 9.2.2.2 | Some problems related to modules . . . . . | 106 |
| 9.2.3   | Modules internals . . . . .                | 108 |
| 9.3     | Header-Map files . . . . .                 | 108 |



# Preface

The Clang is a C/C++ and Objective-C compiler that is an integral part of the LLVM (Low Level Virtual Machine) project. When we talk about Clang, we can refer to two different things. The first one is the compiler frontend, which is the part of the compiler responsible for parsing and performing semantic reasoning about the program. We also use the word Clang to refer to the compiler itself, which is also referred to as the compiler driver. The driver is responsible for invoking the compiler, which can be thought of as a manager that calls different parts of the compiler such as the compiler frontend and other parts necessary 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 LLVM project evolves very fast, and some of its parts may be completely rewritten between different revisions. We will use a specific version of LLVM in the book - version 16.x, which was first released in March 2023 [[Community, 2022c](#)].

The Clang is compiler for C family of languages. Thus it supports such languages as C, Objective-C, C++ and Objective-C++. We are going mostly focus on C++ realisation. That assumes that we will refer C++ standard very often. Despite the fact that LLVM uses C++17 [[for Standardization, 2017](#)] for implementation it implements the latest version of standard and we will use C++20 version of standard [[for Standardization, 2020](#)] for references.

## Who this book is for

This book is intended for experienced C++ software engineers who have no prior experience with compiler design, but who want to gain this knowledge and put it into practice. It may also be useful for engineers who want to learn about how Clang works, as well as its specific features such as performance improvements and modularity, which enables the creation of powerful custom compiler tools

## What this book covers

The book is divided into two parts. The first one provides basic information about the LLVM project and how it can be installed. It also describes useful development tools and configurations used for exploring LLVM code later in the book. The internal Clang architecture is the next main topic in the first part of the book. Knowledge about the Clang internals and its place inside LLVM is essential for any development related to Clang. The Clang is also very good example of well designed software that can be used as a sample of good design pattern.

The final topic in the first part is compilation performance, particularly how it can be improved. We describe several Clang features that may significantly improve compilation speed, such as C++ modules, header maps, and others.

The Clang follows the primary paradigm of LLVM - everything is a library - which allows the creation of a

variety of different tools. The second part of the book is about such tools. We discuss clang-tidy, a powerful framework for creating lint checks. We examine simple checks based on AST (abstract syntax tree) matching, as well as more powerful ones based on advanced techniques like CFG (control flow graph). The list of tools is not limited to code analysis, but also includes refactoring tools and IDE support.

## Download the example code files

The code bundle for the book is also hosted on GitHub at <https://github.com/PacktPublishing/Clang-Compiler-Frontend>. In case there's an update to the code, it will be updated on the existing GitHub repository.

We also have other code bundles from our rich catalog of books and videos available at <https://github.com/PacktPublishing/>. Check them out!

## Download the color images

We also provide a PDF file that has color images of the screenshots/diagrams used in this book. You can download it here: [https://static.packt-cdn.com/downloads/9781801071109\\_ColorImages.pdf](https://static.packt-cdn.com/downloads/9781801071109_ColorImages.pdf).

## Conventions used

There are a number of text conventions used throughout this book.

**CodeInText:** Indicates code words in text, database table names, folder names, filenames, file extensions, pathnames, dummy URLs, and user input. Here is an example: "Any attempt to run the code that has such issues will immediately cause the interpreter to fail, raising a `SyntaxError` exception."

A block of code is set as follows:

---

```
1 int main() {  
2     return 0;  
3 }
```

---

Any command-line input or output is written as follows:

```
$ python3 script.py
```

Some code examples will be representing input of shells. You can recognize them by specific prompt characters:

- »> for interactive Python shell
- \$ for Bash shell (macOS and Linux)
- > for CMD or PowerShell (Windows)

Warnings or important notes appear like this.

**Important note**

Warnings or important notes appear like this.

Tips and tricks appear like this.

**Tips or tricks**

Appear like this.

## Get in touch

Feedback from our readers is always welcome.

**General feedback:** If you have questions about any aspect of this book, mention the book title in the subject of your message and email us at [customercare@packtpub.com](mailto:customercare@packtpub.com).

**Errata:** Although we have taken every care to ensure the accuracy of our content, mistakes do happen. If you have found a mistake in this book, we would be grateful if you would report this to us. Please visit [www.packt.com/submit-errata](http://www.packt.com/submit-errata), selecting your book, clicking on the Errata Submission Form link, and entering the details.

**Piracy:** If you come across any illegal copies of our works in any form on the Internet, we would be grateful if you would provide us with the location address or website name. Please contact us at [copyright@packt.com](mailto:copyright@packt.com) with a link to the material.

**If you are interested in becoming an author:** If there is a topic that you have expertise in and you are interested in either writing or contributing to a book, please visit [authors.packtpub.com](http://authors.packtpub.com).

## Reviews

Please leave a review. Once you have read and used this book, why not leave a review on the site that you purchased it from? Potential readers can then see and use your unbiased opinion to make purchase decisions, we at Packt can understand what you think about our products, and our authors can see your feedback on their book. Thank you!

For more information about Packt, please visit [www.packt.com](http://www.packt.com).

# Part 1:

# Clang setup and architecture

You can find some info about LLVM internal architecture and how clang fits into it. There is also description how to install and build required clang and clang-tools, description for basic LLVM libraries and tools used across LLVM project and essential for clang development. You can find description for some clang features and their internal implementation.



# Environment setup

In this chapter we will discuss the basic steps to be done to setup the environment for future experiments with Clang . The setup is appropriate for Unix-based systems such as Linux and Mac OS (Darwin). In addition, the reader will get important information on how to download, configure, and build the LLVM source code. We will continue with a short session that shows you how to build and use the LLVM debugger (LLDB ), which will be used as the primary tool for code investigation throughout the book. Finally, we will finish with a simple clang tool that can check C/C++ files for compilation errors. We will use LLDB for a simple debug session for the created tool and clang internal.

We will cover the following topics:

- Technical requirements
- Prerequisite
- Getting to know LLVM
- Source code compilation
- How to create a custom clang tool

## 1.1 Technical requirements

Downloading and building LLVM code is very easy and does not require any paid tools. You will require the following:

- Unix based OS (Linux, Darwin)
- Command line git
- Build tools: CMake and Ninja

We will use debugger as the source investigation tool. LLVM has its own debugger - LLDB. We will build it as our first tool build from LLVM monorepo.

Any build process consists of two steps. The first one is the project configuration and the last one is the build itself. LLVM uses CMake as build configuration tool. It also can use wide range of different build tools such as Unix Makefiles, Ninja and so on. It can also generate project files for popular IDEs such as Visual Studio and XCode. We are going to use Ninja as the build tool as soon as it provides more faster build process and the most LLVM developers use it [[Community, 2023d](#)]

## 1.1.1 CMake as build configuration tool

CMake is an open-source, cross-platform build system generator. It has been used as the primary build system for LLVM since version 3.3, which was released in 2013.

Before LLVM began using CMake, it used autoconf, a tool that generates a configure script that can be used to build and install software on a wide range of Unix-like systems. However, autoconf has several limitations, such as being difficult to use and maintain, and having poor support for cross-platform builds. CMake was chosen as an alternative to autoconf because it addresses these limitations and is easier to use and maintain.

In addition to being used as the build system for LLVM, CMake is also used for many other software projects, including Qt [[Wikipedia contributors, 2022c](#)], OpenCV [[Wikipedia contributors, 2022b](#)], Google Test [[Wikipedia contributors, 2022a](#)], and others.

## 1.1.2 Ninja as build tool

Ninja is a small build system with a focus on speed. It is designed to be used in conjunction with a build generator, such as CMake, which generates a build file that describes the build rules for a project.

One of the main advantages of Ninja is its speed. It is able to execute builds much faster than other build systems, such as Unix Makefiles, by only rebuilding the minimum set of files necessary to complete the build. This is because it keeps track of the dependencies between build targets and only rebuilds targets that are out of date.

Additionally, Ninja is simple and easy to use. It has a small and straightforward command-line interface, and the build files it uses are simple text files that are easy to read and understand.

Overall, Ninja is a good choice for build systems when speed is a concern, and when a simple and easy-to-use tool is desired.



One of the most useful Ninja option is `-j`. This option allows you to specify the number of commands to be run in parallel. You may want to specify the number depending on the hardware you are using.

Our next goal is to download LLVM code and investigate the project structure. We also need to set up the necessary utilities for the build process and establish the environment for our future experiments with LLVM code. This will ensure that we have the tools and dependencies in place to proceed with our work efficiently.

## 1.2 Getting to know LLVM

Let's begin by covering some foundational information about LLVM, including the project history as well as its structure.

### 1.2.1 Short LLVM history

The Clang compiler is a part of LLVM project. The project was started in 2000 by Chris Lattner and Vikram Adve as their project at the University of Illinois at Urbana–Champaign [\[Lattner and Adve, 2004\]](#).

LLVM was originally designed to be a next-generation code generation infrastructure that could be used to build optimizing compilers for many programming languages. However, it has since evolved into a full-featured platform that can be used to build a wide variety of tools, including debuggers, profilers, and static analysis tools.

LLVM has been widely adopted in the software industry and is used by many companies and organizations to build a variety of tools and applications. It is also used in academic research and teaching, and has inspired the development of similar projects in other fields.

The project received an additional boost when Apple hired Chris Lattner in 2005 and formed a team to work on LLVM. LLVM became an integral part of the development tools created by Apple (XCode).

Initially, GCC (GNU Compile Collection) was used as the C/C++ frontend for LLVM. But that had some disadvantages, one of them was related to GNU General Public License (GPL) that prevented the frontend usage at some proprietary projects. Another disadvantage was the limited support for the Objective-C language in GCC at the time, which was important for Apple. The clang project was started by Chris Lattner in 2006 to address the issues.

Clang was originally designed as a unified parser for the C family of languages, including C, Objective-C, C++, and Objective-C++. This unification was intended to simplify maintenance by using a single frontend implementation for multiple languages, rather than maintaining multiple implementations for each language.

The project became successful very quickly. One of the primary reasons for the success of Clang and LLVM was their modularity. Everything in LLVM is a library, including Clang. It opened the opportunity to create a lot of amazing tools based on Clang and LLVM, such as clang-tidy and clangd, which will be covered later in the book (see [Chapter 5, clang-tidy linter framework](#) and [Chapter 8, IDE support and code navi-](#)

gation).

LLVM and Clang have a very clear architecture and are written in C++. That makes possible to investigate and use it by any C++ developer. As result we can see huge community created around LLVM and extremely fast grows of its usage.

## 1.2.2 OS support

We are planning to focus on OS for personal computers here, such as Linux, Darwin and Windows. From other side the Clang usage is not limited by the personal computers but can also be used to compile code for mobile platforms such as iOS and different embedded systems.

### 1.2.2.1 Linux

The GCC (GNU Compiler Collection) is the default set of dev tools on Linux, especially `gcc (g++)` is the default C/C++ compiler. The Clang can also be used to compile source code on Linux. Moreover it mimics to `gcc` and supports most of its options. From the other side LLVM support might be limited for some GNU tools, for instance GNU Emacs does not support LLDB as debugger. But despite the fact, the Linux is the most suitable OS for LLVM development and investigation, thus we will mainly use this OS for future examples.

### 1.2.2.2 Darwin

The Clang is considered as the main build tool for Darwin. The entire build infrastructure is based on LLVM, and Clang is the default C/C++ compiler. The developer tools, such as the debugger (LLDB), also come from LLVM. You can get the primary developer utilities from XCode, which are LLVM-based. However, you may need to install additional command-line tools, such as CMake and Ninja, either as separate packages or through package systems like MacPorts or Homebrew. For example, you can get CMake using Homebrew as follows

```
$ brew install cmake
```

or for MacPorts:

```
$ sudo port install cmake
```

### 1.2.2.3 Windows

On Windows, Clang can be used as a command-line compiler or as part of a larger development environment such as Visual Studio. Clang on Windows includes support for the Microsoft Visual C++ (MSVC) ABI, so you can use Clang to compile programs that use the Microsoft C runtime library (CRT) and the C++ Standard Library (STL). Clang also supports many of the same language features as GCC, so it can be used as a drop-in replacement for GCC on Windows in many cases.

It's worth mentioning `clang-cl` [Community, "2023"], it is a command-line compiler driver for Clang that is designed to be used as a drop-in replacement for the Microsoft Visual C++ (MSVC) compiler, `cl.exe`. It was introduced as part of the Clang compiler, and is created to be used with the LLVM toolchain.

Like `cl.exe`, `clang-cl` is designed to be used as part of the build process for Windows programs, and it supports many of the same command-line options as the MSVC compiler. It can be used to compile C, C++, and Objective-C code on Windows, and it can also be used to link object files and libraries to create executable programs or dynamic-link libraries (DLLs).

*The development process for Windows is different from that of Unix-like systems, which would require additional specifics that might make the book material quite complicated. To avoid this complexity, our primary goal is to focus on Unix-based systems, such as Linux and Darwin, and we will omit Windows-specific examples in the book.*

### 1.2.3 LLVM/clang project structure

The Clang source is a part of LLVM monorepo (monolithic repository). LLVM started to use monorepo in 2019 as a part of its transition to git [Community, 2019] (TBD - verify the statement). The decision was driven by number of factors such as better code reuse, improved efficiency and collaboration. Thus you can find all LLVM projects in one place. As seen in Preface, we will be using LLVM version 16.x in this book. The following command will allow you to download it:

```
$ git clone https://github.com/llvm/llvm-project.git -b release/16.x
$ cd llvm-project
```

The most important parts of the `llvm-project` that will be used in the book are shown in fig. 1.1. There are

- `lld` - the LLVM linker tool. You may want to use it as a replacement for standard linker tools such as GNU `ld`.
- `llvm` - common libraries for LLVM project
- `clang` - the clang driver and frontend
- `clang-tools-extra` - there are different clang tools that will be covered at the second part of the book

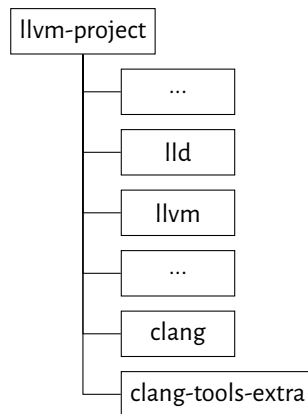


Figure 1.1: LLVM project tree: there are projects that are important for the book: *lldb* (LLDB debugger), *llvm* (some common libs are located here), *clang* (C/C++ compiler frontend and driver) and *clang-tools-extra* (different clang tools such as *clang-tidy* and *clangd*)

Most projects have the same structure shown in fig. 1.2.

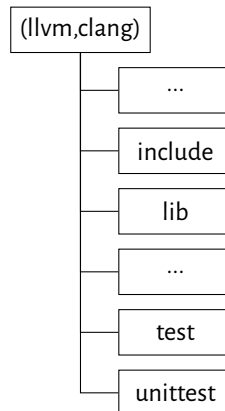


Figure 1.2: Typical LLVM project structure: the primary source code is located at *include* and *lib* folders. The *test* folder keeps LIT end-to-end tests and *unittest* folder keeps tests written for google test framework

LLVM projects, such as *clang* or *llvm*, typically contain two primary folders: *include* and *lib*. The *include* folder contains the project interfaces (header files), while the *lib* folder contains the implementation. Each LLVM project has a variety of different tests, which can be divided into two primary groups: unit tests located in the 'unittests' folder and implemented using the Google Test framework, and end-to-end tests implemented using the LLVM Integrated Tester (LIT) framework. You can get more info

about LLVM/Clang testing in [Section 4.4, LLVM test framework](#).

The most important projects for us are `clang` and `clang-tools-extra`. The `clang` folder contains the frontend and driver.<sup>1</sup> For instance, the lexer implementation is located in the `clang/lib/Lex` folder. You can also see the `clang/test` folder, which contains end-to-end tests, and the `clang/unittest` folder, which contains unit tests for the frontend and driver.

Another important folder is `clang-tools-extra`. It contains some tools based on clang Abstract Syntax Tree (AST). There are:

- `clang-tools-extra/clangd` - a language server that provides navigation info for IDEs such as VSCode.
- `clang-tools-extra/clang-tidy` - a powerful lint framework with several hundreds of different checks.
- `clang-tools-extra/clang-format` - a code formatting tool.

After obtaining the source code and setting up build tools, we are ready to compile the LLVM source code.

## 1.3 Source code compilation

We are compiling our source code in debug mode to make it suitable for future investigations with a debugger. We are using LLDB as the debugger. We will start with an overview of the build process and finish with a concrete example of the build for LLDB.

### 1.3.1 Configuration with CMake

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

```
$ mkdir build
$ cd build
```

The minimal configuration command looks like

```
$ cmake -DCMAKE_BUILD_TYPE=Debug ../llvm
```

The command requires the build type to be specified (e.g. `Debug` in our case) as well as the primary argument that points to a folder with the build configuration file. The configuration file is stored as `CMakeLists.txt` and is located in the `llvm` folder, which explains the `../llvm` argument usage. The command generates `Makefile` located at the build folder, thus you can use simple `make` command to start the build process.

We will use more advanced configuration commands in the book. The command is shown in [fig. 1.3](#).

---

<sup>1</sup>the compiler driver is used to run different stages of compilation (parsing, optimisation, link etc.). You can get more info about it at [Section 2.2, Clang driver overview](#)

---

```
cmake -G Ninja -DCMAKE_BUILD_TYPE=Debug -DCMAKE_INSTALL_PREFIX=./install
→ -DLLVM_TARGETS_TO_BUILD="X86"
→ -DLLVM_ENABLE_PROJECTS="lldb;clang;clang-tools-extra"
→ -DLLVM_USE_SPLIT_DWARF=ON ../llvm
```

---

*Figure 1.3: CMake: basic configuration used for the build at the book*

There are several LLVM/cmake options specified:

- `-DCMAKE_BUILD_TYPE=Debug` sets the build mode. The build with debug info will be created. There is a primary build configuration for Clang internals investigations
- `-DCMAKE_INSTALL_PREFIX=./install` specifies the installation folder
- `-DLLVM_TARGETS_TO_BUILD="X86"` sets exact targets to be build. It will avoid build unnecessary targets
- `-DLLVM_ENABLE_PROJECTS="lldb;clang;clang-tools-extra"` specifies LLVM projects that we care about
- `-DLLVM_USE_SPLIT_DWARF=ON` - splits debug information into separate files. This option saves disk space as well as memory consumption during the LLVM build.

We used `-DLLVM_USE_SPLIT_DWARF=ON` to save some space on the disc. For instance the clang build with the option enabled takes 20Gb<sup>2</sup> space but it takes 31Gb space with the option disabled. Note that the option requires compiler used for clang build to support it. You might also notice that we create the build for one specific architecture - "X86". This option also saved some space for us because otherwise all supported architecture will be built and the required space will increase from 20Gb to 27Gb.

You can save more space if you will use dynamic libraries instead of static ones. `-DBUILD_SHARED_LIBS=ON` will build each LLVM component as shared library. The used space will be 14Gb at the case and the overall config command will look like this

---

```
cmake -G Ninja -DCMAKE_BUILD_TYPE=Debug -DCMAKE_INSTALL_PREFIX=./install
→ -DLLVM_TARGETS_TO_BUILD="X86"
→ -DLLVM_ENABLE_PROJECTS="lldb;clang;clang-tools-extra"
→ -DLLVM_USE_SPLIT_DWARF=ON -DBUILD_SHARED_LIBS=ON ../llvm
```

---

<sup>2</sup>The configuration specified at fig. 1.3 was used for the build. The build was run as `ninja clang`

For performance purposes on Linux you might want to use `gold` linker instead of default one. The `gold` linker is an alternative to the GNU Linker that was developed as part of the GNU Binary Utilities (binutils) package. It is designed to be faster and more efficient than the GNU Linker, especially when linking large projects. One way it achieves this is by using a more efficient algorithm for symbol resolution and a more compact file format for the resulting executable. It can be enabled with `-DLLVM_USE_LINKER=gold` option. The result configuration command will look like.

---

```
cmake -G Ninja -DCMAKE_BUILD_TYPE=Debug -DCMAKE_INSTALL_PREFIX=./install
→ -DLLVM_TARGETS_TO_BUILD="X86"
→ -DLLVM_ENABLE_PROJECTS="lldb;clang;clang-tools-extra"
→ -DLLVM_USE_LINKER=gold -DLLVM_USE_SPLIT_DWARF=ON -DBUILD_SHARED_LIBS=ON
→ ../llvm
```

---

The debug build can be very slow and you may want to consider an alternative. A good compromise between debuggability and performance is the release build with debug information. To obtain this build, you can change the `CMAKE_BUILD_TYPE` flag to `RelWithDebInfo` in your overall configuration command. The command will then look like this:

---

```
cmake -G Ninja -DCMAKE_BUILD_TYPE=RelWithDebInfo
→ -DCMAKE_INSTALL_PREFIX=./install -DLLVM_TARGETS_TO_BUILD="X86"
→ -DLLVM_ENABLE_PROJECTS="lldb;clang;clang-tools-extra"
→ -DLLVM_USE_SPLIT_DWARF=ON ../llvm
```

---

The following table keeps the list of some popular options [[Community, 2023](#)].

## 1.3.2 Build

We need to call Ninja to build the desired project. The command for Clang build will look like

```
$ 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` in the `build` folder.

You can also install the binaries into the folder specified with `-DCMAKE_INSTALL_PREFIX` option. It can be done as follows

```
$ ninja install clang
```

| Option                | Description   |
|-----------------------|---|
| CMAKE_BUILD_TYPE      | Specifies build configuration.<br>Possible values are Release   Debug   RelWithDebInfo   MinSizeRel.<br>The Release and RelWithDebInfo are optimized for performance, MinSizeRel is optimized for size. |
| CMAKE_INSTALL_PREFIX  | Installation prefix   |
| CMAKE_C, CXX_FLAGS    | Extra C/C++ flags be used for compilation   |
| CMAKE_C, CXX_COMPILER | C/C++ compiler be used for compilation.<br>You might want to specify a non-default compiler to use some options that are not available or not supported by the default compiler                         |
| LLVM_ENABLE_PROJECTS  | The projects to be enabled. We will use clang; clang-tools-extra  |
| LLVM_USE_LINKER       | Specifies the linker be used.<br>There are several options that include gold and lld  |

*Table 1.1: Configuration options*

The folder `../install` (specified as the installation folder at fig. 1.3) will have the following structure

```
$ ls ../install
bin  include  lib  libexec  share
```

### 1.3.3 The LLVM debugger, its build and usage

The LLVM debugger, LLDB, has been created with a look at GDB (GNU debugger). Some of its commands repeat the counterparts from GDB. You may ask the question: "Why do we need a new debugger if we have a good one?" The answer can be found in the different architecture solutions used by GCC and LLVM. LLVM uses a modular architecture, and different parts of the compiler can be reused. For example, the Clang frontend can be reused in the debugger, resulting in actual support for modern C/C++ features. For example, the `print` command in `lldb` can specify any valid language constructions and you can use some modern C++ features with the `lldb print` command.

In contrast, GCC uses a monolithic architecture, and it's hard to separate the C/C++ frontend from other parts. Therefore, GDB has to implement language features separately, which may take some time before modern language features implemented in GCC become available in GDB.

You may find some info about LLDB build and typical usage scenario in the following example. We are going to create a separate folder for release build

```
$ cd llvm-project
$ mkdir release
$ cd release
```

We configure our project in Release mode and specify the `lldb` and `clang` projects only



---

```
cmake -G Ninja -DCMAKE_BUILD_TYPE=Release -DCMAKE_INSTALL_PREFIX=./install  
→ -DLLVM_TARGETS_TO_BUILD="X86" -DLLVM_ENABLE_PROJECTS="lldb;clang" ../llvm
```

---

We are going to build both Clang and LLDB using no more than 4 concurrent processes

```
$ ninja clang lldb -j4
```

You can install the created executables with the following command

```
$ ninja install-clang install-lldb
```

The binary will be installed into the folder specified via `-DCMAKE_INSTALL_PREFIX` config command argument.

We will use the following simple C++ program for the example debugger session

---

```
1 int main() {  
2     return 0;  
3 }
```

---

The program can be compiled using the following command

```
llvm-project/install/bin/clang main.cpp -o main -g -O0
```

As you may notice we don't use optimization (`-O0` option) and store debug info at the binary (`-g` option).

Typical debug session for the created executable is shown in fig. 1.4. There are several actions done:

- Run the debug session with `llvm-project/install/bin/lldb src/main`, where `main` is the executable we want to debug, see fig. 1.4, line 1.
- We set breakpoint at the `main` function, see fig. 1.4, line 4.
- Run the session with `r` command, see fig. 1.4, line 6.
- We can see that the process is interrupted at the breakpoint, see fig. 1.4, line 8, 12.
- We finish the session with `q` command, see fig. 1.4, line 14.

We are going to use LLDB as one of our tools for the Clang internals investigation. We will use the same sequence of commands that is shown in fig. 1.4. You can also use another debugger, such as GDB, that has a similar set of commands as LLDB.

---

```
1 $ llvm-project/install/bin/lldb main
2 (lldb) target create "./main"
3 ...
4 (lldb) b main
5 Breakpoint 1: where = main`main + 11 at main.cpp:2:3,...
6 (lldb) r
7 Process 1443051 launched: ...
8 Process 1443051 stopped
9 * thread #1, name = 'main', stop reason = breakpoint 1.1
10   frame #0: 0x0000555555555513b main`main at main.cpp:2:3
11     1   int main() {
12 -> 2       return 0;
13     3   }
14 (lldb) q
```

---

*Figure 1.4: LLDB session example*

## 1.4 Test project: syntax check with clang tool

For our first test project we will create a simple clang tool that runs compiler and checks syntax for provided source file. We will create so called out-of-tree LLVM project i.e. the project that will use LLVM but will be located outside the main LLVM source tree.

There are several actions to be done to create the project:

- The required LLVM libraries and headers have to be built and installed
- We have to create a build configuration file for our test project
- The source code that uses LLVM has to be created

We will start with the first step and will install clang support libraries and headers. We will use the following configuration command for CMake:

---

```
cmake -G Ninja -DCMAKE_BUILD_TYPE=Debug -DCMAKE_INSTALL_PREFIX=../install
→ -DLLVM_TARGETS_TO_BUILD="X86" -DLLVM_ENABLE_PROJECTS="clang"
→ -DLLVM_USE_LINKER=gold -DLLVM_USE_SPLIT_DWARF=ON -DBUILD_SHARED_LIBS=ON
→ ../llvm
```

---

*Figure 1.5: LLVM CMake configuration for simple syntax check clang tool*

As you may notice, we enabled only one project: `clang`, all other options are standard for our debug build. The command has to be run from a created `build` folder inside LLVM source tree, as it was suggested at [Section 1.3.1, Configuration with CMake](#).

The required libraries and headers can be installed with the following command

```
$ ninja install
```

The installation will be done into `install` folder as it was specified at `CMAKE_INSTALL_PREFIX` option.

We have to create two files for our project:

- `CMakeLists.txt` the build configuration file
- `TestProject.cpp` the project source code

The configuration file will accept a path to the LLVM install folder via `LLVM_HOME` environment variable. The configuration file is the following:

---

```
1 cmake_minimum_required(VERSION 3.16)
2 project("syntax-check")
3
4 if ( NOT DEFINED ENV{LLVM_HOME})
5     message(FATAL_ERROR "$LLVM_HOME is not defined")
6 else()
7     message(STATUS "$LLVM_HOME found: $ENV{LLVM_HOME}")
8     set(LLVM_HOME $ENV{LLVM_HOME} CACHE PATH "Root of LLVM installation")
9     set(LLVM_LIB ${LLVM_HOME}/lib)
10    set(LLVM_DIR ${LLVM_LIB}/cmake/llvm)
11    find_package(LLVM REQUIRED CONFIG)
12    include_directories(${LLVM_INCLUDE_DIRS})
13    link_directories(${LLVM_LIBRARY_DIRS})
14    set(SOURCE_FILES SyntaxCheck.cpp)
15    add_executable(syntax-check ${SOURCE_FILES})
16    set_target_properties(syntax-check PROPERTIES COMPILE_FLAGS "-fno-rtti")
17    target_link_libraries(syntax-check
18        LLVMSupport
19        clangBasic
20        clangFrontend
21        clangSerialization
22        clangTooling
23    )
24 endif()
```

---

*Figure 1.6: CMake file for simple syntax check clang tool*

The most important parts of the file are

- Line 2: We specify the project name (syntax-check). That also be the name for our executable
- Lines 4-7: Test for LLVM\_HOME environment variable
- Line 10: We set a path to LLVM CMake helpers
- Line 11: We load LLVM CMake package from the paths specified at line 10
- Line 14: We specify our source file that should be compiled

- Line 16: We setup an additional flag for compilation: `-fno-rtti`. The flag is required as soon as LLVM is built without RTTI. This is done in an effort to reduce code and executable size [Community, 2023b].
- Line 18-22 We specify the required libraries to be linked to our program

The source code for our tool is the following

---

```

1 #include "clang/Frontend/FrontendActions.h" // clang::SyntaxOnlyAction
2 #include "clang/Tooling/CommonOptionsParser.h"
3 #include "clang/Tooling/Tooling.h"
4 #include "llvm/Support/CommandLine.h" // llvm::cl::extrahelp
5
6 namespace {
7   llvm::cl::OptionCategory TestCategory("Test project");
8   llvm::cl::extrahelp
9     CommonHelp(clang::tooling::CommonOptionsParser::HelpMessage);
10 } // namespace
11
12 int main(int argc, const char **argv) {
13   llvm::Expected<clang::tooling::CommonOptionsParser> OptionsParser =
14     clang::tooling::CommonOptionsParser::create(argc, argv, TestCategory);
15   if (!OptionsParser) {
16     llvm::errs() << OptionsParser.takeError();
17     return 1;
18   }
19   clang::tooling::ClangTool Tool(OptionsParser->getCompilations(),
20                                   OptionsParser->getSourcePathList());
21   return Tool.run(
22     clang::tooling::newFrontendActionFactory<clang::SyntaxOnlyAction>()
23       .get());
24 }

```

---

The most important part of the file are

- Line 7-9: The majority of compiler tools have the same set of command line arguments. LLVM command line library [Community, 2023c] provides some API to process compiler command options. We setup the library at line 7. We also setup additional help message at lines 8-9.

- Line 13-18: We parse command line arguments
- Line 19-23: We create and run our clang tool
- Line 22: We use `clang::SyntaxOnlyAction` frontend action that will run syntax and semantic check on the input file. You can get more info about frontend actions later at [Section 2.3.1, Frontend action](#).

We have to specify a path to LLVM `install` folder to build our tool. As it was mentioned above, the path has to be specified via `LLVM_HOME` environment variable. Our configuration command (see fig. 1.5) specifies the path as `install` folder inside LLVM project source tree. Thus we can build our tool as follows

```
export LLVM_HOME=<...>/llvm-project/install
mkdir build
cd build
cmake -G Ninja ..
ninja
```

We can run the tool as follows

```
$ ./build/syntax-check --help
USAGE: syntax-check [options] <source0> [... <sourceN>]
...
```

The program will successively terminate if we run it on a valid C++ source file, but it will produce an error message if it's run on a broken C++ file:

```
$ ./syntax-check mainbroken.cpp
...
Running without flags.
mainbroken.cpp:2:11: error: expected ';' after return statement
    return 0
           ^
           ;
1 error generated.
mainbroken.cpp.
```

We can also run our tool under LLDB debugger.

---

```

1 $ llvm-project/install/bin/lldb ./syntax-check -- main.cpp
2 ...
3 (lldb) b clang::ParseAST
4 ...
5 (lldb) r
6 ...
7 Running without flags.
8 Process 608249 stopped
9 * thread #1, name = 'syntax-check', stop reason = breakpoint 1.1
10   frame #0: ... clang::ParseAST(...) at ParseAST.cpp:116:3
11     113
12     114         void clang::ParseAST(...) {
13     115             // Collect global stats on Decls/Stmts ...
14 -> 116             if (PrintStats) {
15     117                 Decl::EnableStatistics();
16     118                 Stmt::EnableStatistics();
17     119             }
18 (lldb) c
19 Process 608249 resuming
20 Process 608249 exited with status = 0 (0x00000000)
21 (lldb)

```

---

*Figure 1.7: LLDB session for clang tool test project*

We run `syntax-check` as the primary binary and set `main.cpp` source file as an argument for the tool (section 1.4, line 1). We also set a breakpoint at `clang::ParseAST` function (section 1.4, line 3). The function is the primary entry point for source code parsing. We run program at line 5 and continue the execution after breakpoint at line 18.

We will use the same debugging techniques later in the book when we will investigate clang source code.

## 1.5 Summary

In the chapter, we covered the history of the LLVM project, obtained the source code for LLVM, and explored its internal structure. We learned about the tools used to build LLVM, such as CMake and Ninja. We studied the various configuration options for building LLVM and how they can be used to optimize resources, including disk space. We built Clang and LLDB in release mode and used the resulting tools

to compile a basic program and run it under the debugger. We also created a simple clang tool and run it under LLDB debugger.

The next chapter will introduce you to compiler design architecture and how it appears in the context of Clang . We will primarily focus on the Clang frontend, but we will also cover the important concept of the Clang driver - the backbone that manages all stages of the compilation process, from parsing to linking.

## 1.6 Further reading

- Getting Started with the LLVM System, <https://llvm.org/docs/GettingStarted.html>
- Building LLVM with CMake, <https://llvm.org/docs/CMake.html>
- Clang Compiler User's Manual, <https://clang.llvm.org/docs/UsersManual.html>



# Clang architecture

In this chapter, we will examine the internal architecture of Clang and its relationship with other LLVM components. We will begin with an overview of the overall compiler architecture, with a specific focus on the clang driver. As the backbone of the compiler, the driver runs all compilation phases and controls their execution. Finally, we will concentrate on the frontend portion of the Clang compiler, which includes lexical and semantic analysis, and produces an Abstract Syntax Tree (AST) as its primary output. The AST forms the foundation for most clang tools, and we will examine it more closely in the next chapters.

The following topics will be covered in this chapter:

- Compiler overview
- Clang driver overview, including an explanation of the compilation phases and their execution
- Clang frontend overview, which covers frontend actions, preprocessor, parser, and sema

## 2.1 Getting started with compilers

Despite the fact that compilers are used to translate programs from one form to another, they can also be considered as large software systems that use various algorithms and data structures. The knowledge obtained from studying compilers is not specific to compilers and can be applied to other software as well. On the other hand, compilers are also a subject of active scientific research, and there are many unexplored areas and topics to investigate.

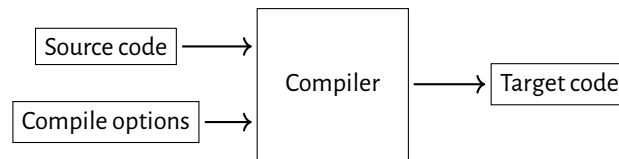


Figure 2.1: Compiler takes source code and compile options and transform them into a code on the target platform

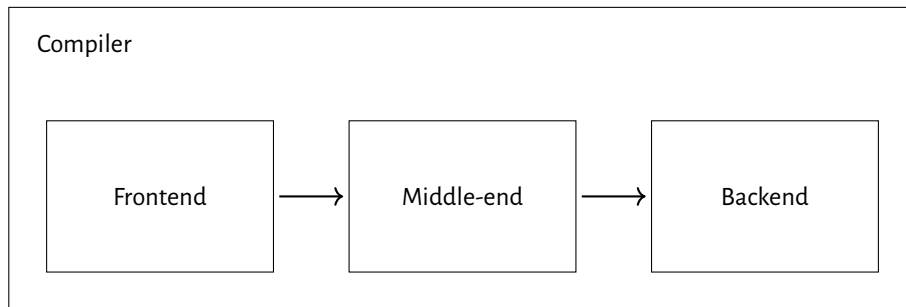


Figure 2.2: Typical compiler workflow: source program is passed via different stages: frontend, middle-end and backend

You can find some basic information about the internal structure of a compiler here. We will keep it as basic as possible so the information is applicable to any compiler, not just Clang. We will briefly cover all phases of compilation, which will help to understand Clang's position in the overall compiler architecture.

## 2.1.1 Explore compiler workflow

The primary function of a compiler is to convert a program written in a specific programming language (such as C/C++ or FORTRAN) into a format that can be executed on a target platform. This process involves the use of a compiler, which takes the source file and any compilation flags, and produces a build artifact, such as an executable or object file, as shown in fig. 2.1. The term "target platform" can have a broad meaning. It can refer to machine code that is executed on the same host, as is typically the case. But it can also refer to cross-compilation, where the compiler generates code for a different computer architecture than the host. For example, code for a mobile application or embedded application running on ARM can be generated using an Intel machine as the host. Additionally, the target platform is not limited to machine code only. For example, some early C++ compilers (such as "cc") would produce pure C code as output. This was done because, at the time, C was the most widely used and well-established programming language, and the C compiler was the most reliable way to generate machine code. This approach allowed early C++ programs to be run on a wide range of platforms since most systems already had a C compiler available. The produced C code could then be compiled into binary by another compiler.

We are going to focus on compilers that produce binary code, and a typical compiler workflow for such a compiler is shown in fig. 2.2. The stages of compilation can be described as follows:

- Frontend: The Frontend does lexical analysis and parsing, which includes both syntax analysis and semantic analysis. The syntax analysis assumes that your program is well-organized according to the language grammar rules. The semantic analysis performs checks on the program's meaning and rejects invalid programs, such as those that use wrong types.
- Middle-end: The Middle-end performs various optimizations on the intermediate representation (IR) code (LLVM-IR for clang).
- Backend: The Backend of a compiler takes the optimized or transformed IR and generates machine code or assembly code that can be executed by the target platform.

The source program is transformed into different forms as it passes through the various stages. For example, the Frontend produces IR code, which is then optimized by the Middle-end and finally converted into native code by the Backend (see fig. 2.3). Lets look into compiler frontend as the first component at the compiler's workflow.

## 2.1.2 Frontend

The primary goal for frontend is the source code conversion to IR. It's worth mentioning that frontend also transforms the source code into various forms before it produces the IR. Frontend will be our primary focus in the book, so we will examine its components. The first component of the frontend is the Lexer (see fig. 2.4). It converts the source code into a set of tokens, which are used to create a special data structure called the abstract syntax tree (AST). The final component, code generator (Codegen), traverses the AST and generates the IR from it.

We will use a simple C/C++ program that calculates the maximum of two numbers to demonstrate the workings of the Frontend. The code for the program is as follows:

---

```
1 int max(int a, int b) {  
2     if (a > b)  
3         return a;  
4     return b;  
5 }
```

---

*Figure 2.5: Test program for compiler frontend investigations*

The first component of the frontend is Lexer, lets examine it.

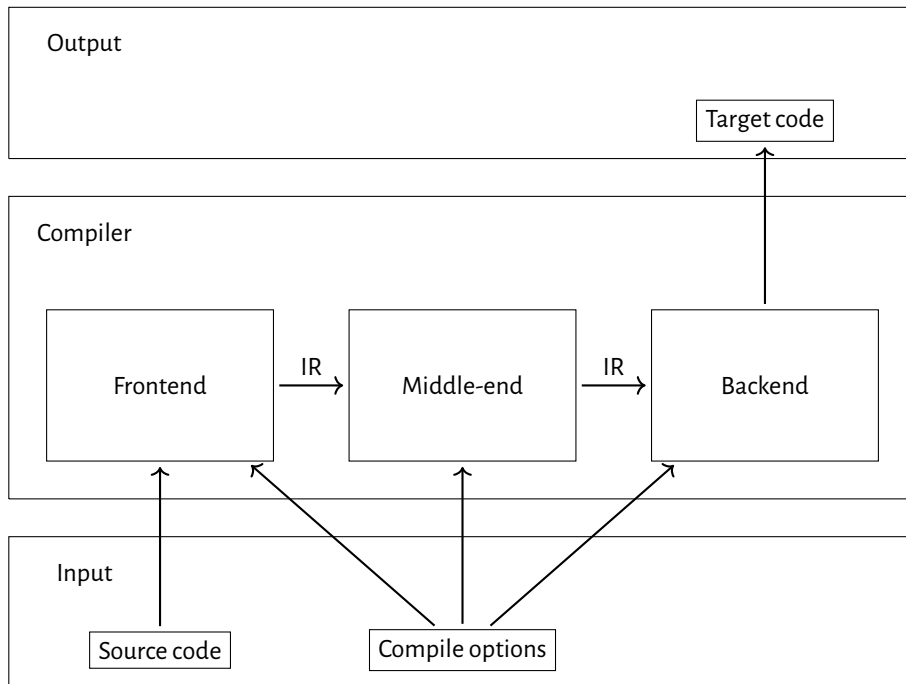


Figure 2.3: Source code transformation by compiler: Input data consists of Source code and Compile options. The source code is transformed by Frontend into IR (Intermediate representation). Middle-end does different optimizations on IR and passes the final (optimized) result to Backend. Backend generates the Target code. Frontend, Middle-end and Backend use Compile options as setting for the code transformations

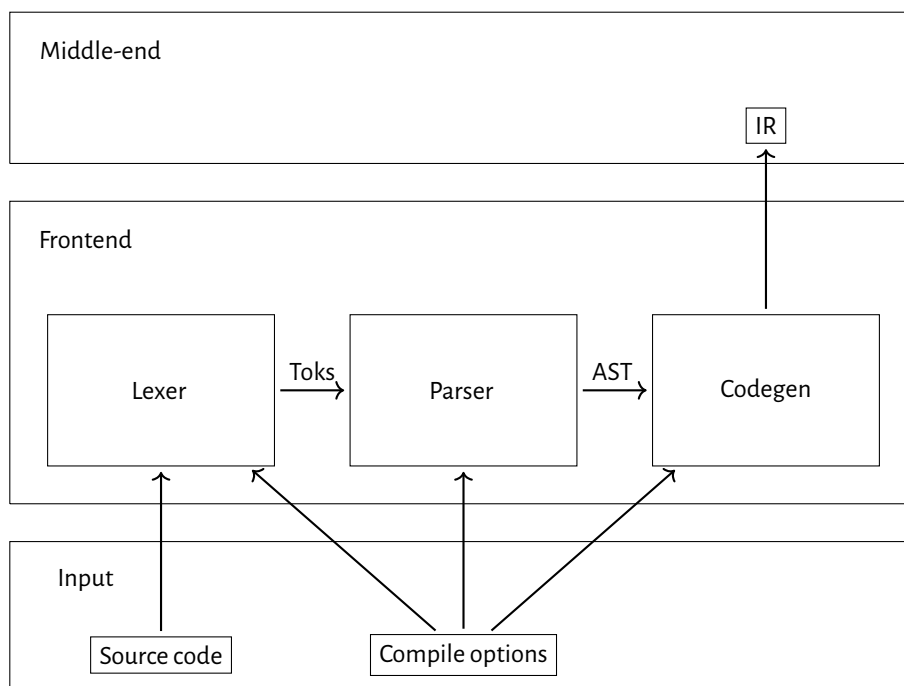


Figure 2.4: Compiler frontend: source code is transformed into a set of tokens (Toks) by Lexer. Parser takes the tokens and creates Abstract syntax tree (AST). Codegen generates IR from AST

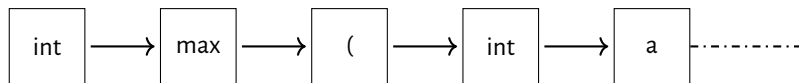


Figure 2.6: *Lexer* : the program source is converted into a stream of tokens

### 2.1.2.1 Lexer

The Frontend process starts with the *Lexer* , which converts the input source into a stream of tokens. In our example program (see fig. 2.5), the first token is the keyword `int` , which represents the integer type. This is followed by the identifier `max` for the function name. The next token is the left parenthesis `(` , and so on (see fig. 2.6).

### 2.1.2.2 Parser

The *Parser* is the next component following the *Lexer* . The primary output produced by the *Parser* is called as abstract syntax tree (AST). This tree represents the abstract syntactic structure of the source code written in a programming language. The *Parser* generates the AST by taking the stream of tokens produced by the *Lexer* as input and organizing them into a tree-like structure. Each node in the tree represents a construct in the source code, such as a statement or expression, and the edges between nodes represent the relationships between these constructs.

The AST for our example program is shown in fig. 2.7. As you can see, our function (`max` ) has two parameters (`a` and `b` ) and a body. The body is marked as a compound statement in fig. 2.7. The compound statement consists of other statements, such as `return` and `if` . The variables `a` and `b` are used in the bodies of these statements. You may also be interested in the real AST generated by Clang for the compound statement, the result of which is shown in fig. 2.8.

The *Parser* performs two activities:

1. Syntax analysis: the *Parser* constructs the AST by analyzing the syntax of the program.
2. Semantic analysis: the *Parser* analyzes the program semantically.

One of the primary goals of analysis is error detection, and the *Parser* produces an error message if it fails in syntax or semantic analysis. We can get a sense of this by considering what types of errors are detected by syntax analysis and which ones are detected by semantic analysis.

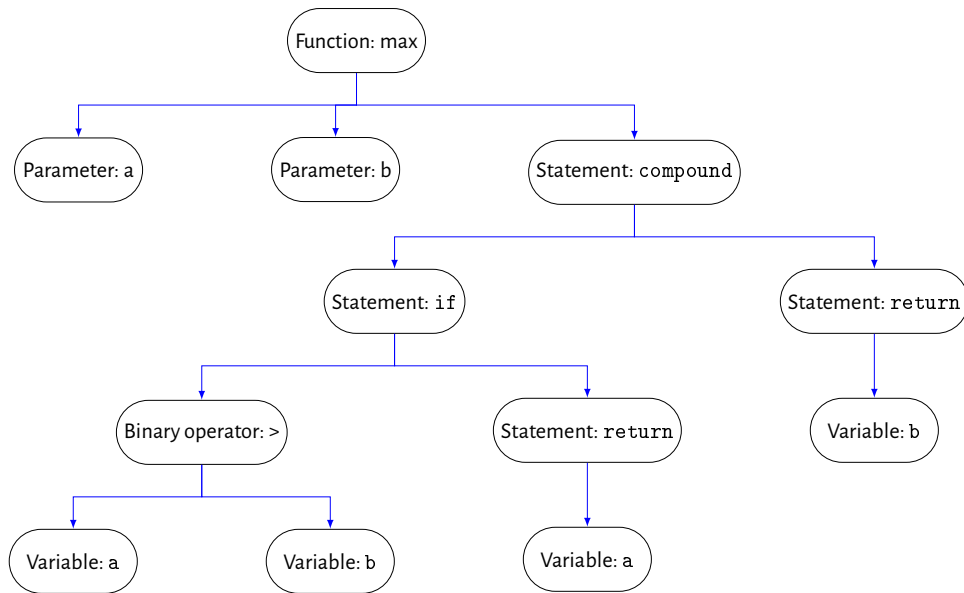


Figure 2.7: The AST for our example program that calculates maximum of 2 numbers

Syntax analysis assumes that the program should be correct in terms of the grammar specified for the language. For example, the following program is invalid in terms of syntax because a semicolon is missing at the last return statement:

---

```

1 int max(int a, int b) {
2     if (a > b)
3         return a;
4     return b // missing ;
5 }
  
```

---

Clang produces the following output for the program:

```

max_invalid_syntax.cpp:4:11: error: expected ';' after return statement
    return b // missing ;
           ^
           ;
  
```

On the other hand, a program can be syntactically correct but have no sense. The Parser should detect a semantic error in such cases. For instance, the following program has a semantic error related to the wrongly used type for the return value:

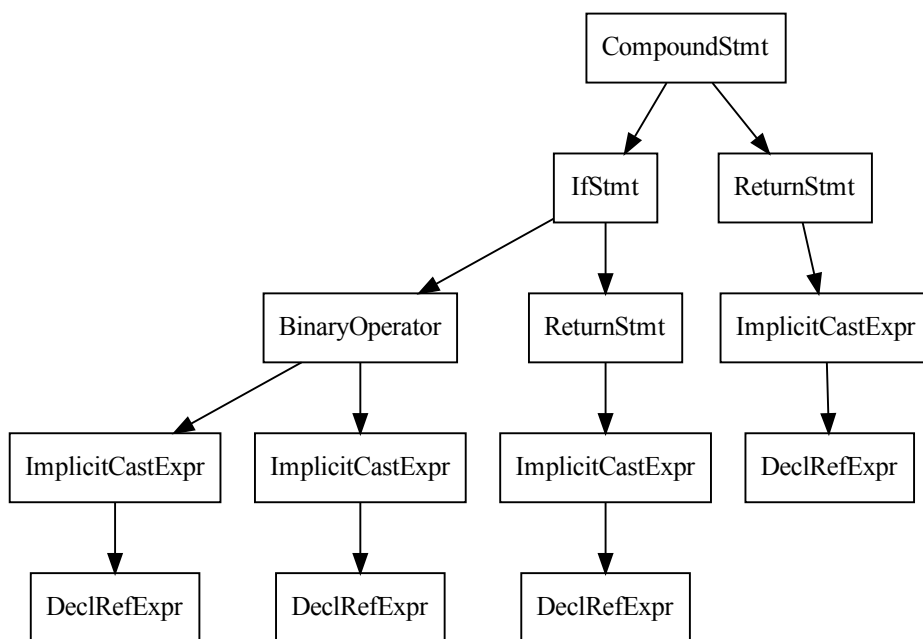


Figure 2.8: The AST for compound statement generated by Clang. The tree generated by `clang -cc1 -ast-view <...>` command



---

```
1 int max(int a, int b) {  
2     if (a > b)  
3         return a;  
4     return &b; // invalid return type  
5 }
```

---

Clang produces the following output for the program:

```
max_invalid_sema.cpp:4:10: error: cannot initialize return object of type \  
'int' with an rvalue of type 'int *'  
    return &b; // invalid return type  
           ^~
```

AST is mainly constructed as a result of syntax analysis, but for certain languages, such as C++, semantic analysis is also crucial for constructing the AST, particularly for C++ template instantiation.

During syntax analysis, the compiler verifies that the template declaration adheres to the language's grammar and syntax rules, including the proper use of keywords such as "template" and "typename", as well as the formation of the template parameters and body.

Semantic analysis, on the other hand, involves the compiler performing template instantiation, which generates the AST for specific instances of the template. It's worth noting that the semantic analysis of templates can be quite complex, as the compiler must perform tasks such as type checking, name resolution, and more for each template instantiation. Additionally, the instantiation process can be recursive and lead to a significant amount of code duplication, known as code bloat. To combat this, C++ compilers employ techniques such as template instantiation caching to minimize the amount of redundant code generated.

### 2.1.2.3 Codegen

The Codegen<sup>1</sup> or code generator, which is the final component of the compiler's frontend, has the primary goal of generating the Intermediate Representation (IR). For this purpose, the compiler traverses the AST generated by the parser and converts it into another source code that is called the Intermediate Representation or IR. The IR is a language-independent representation, allowing the same middle-end component to be used for different frontends (FORTRAN vs C++).

The use of Intermediate Representations (IRs) in compilers is a concept that has been around for several decades. The idea of using an intermediate representation to represent the source code of a program during compilation has evolved over time, and the exact date when IR was first introduced in compilers is not clear.

However, it is known that the first compilers in the 1950s and 1960s did not use IRs and instead translated

---

<sup>1</sup>It's worth mentioning that we also have another Codegen component as a part of Backend that generate the target code.

source code directly into machine code. By the 1960s and 1970s, researchers had begun experimenting with using IRs in compilers to improve the efficiency and flexibility of the compilation process.

One of the first widely used IRs was the three-address code, which was used in the mid-1960s in IBM/360's FORTRAN compiler. Other early examples of IRs include the register transfer language (RTL) and the static single assignment (SSA) form, which were introduced in the 1970s and 1980s respectively.

Today, the use of IRs in compilers is a standard practice, and many compilers use multiple IRs throughout the compilation process. This allows for more powerful optimization and code generation techniques to be applied.

## 2.2 Clang driver overview

When discussing compilers, we typically refer to a command-line utility that initiates and manages the compilation process. For example, to use the GNU Compiler Collection, one must call `gcc` to start the compilation process. Similarly, to compile a C++ program using Clang, one must call `clang` as the compiler. The program that controls the compilation process is known as the driver. The driver coordinates different stages of compilation and connects them together. In the book, we will be focusing on LLVM and using `clang` as the driver for the compilation process.

It may be confusing for readers that the same word, "clang," is used to refer to both the compiler front-end and the compilation driver. In contrast, with GCC, the driver and C++ compiler are separate executables<sup>2</sup>. However, "clang" is a single executable that functions as both the driver and the compiler front-end. To use Clang as the compiler front-end only, the special option `-cc1` must be passed to it.

### 2.2.1 Example program

We will use the simple "Hello world!" example program for our experiments with clang driver. The main source file is called `hello.cpp`. The file implements a trivial C++ program that prints "Hello world!" to the standard output:

---

```
1 #include <iostream>
2
3 int main() {
4     std::cout << "Hello world!" << std::endl;
5     return 0;
6 }
```

---

You can compile the source with

---

<sup>2</sup>The C/C++ compiler in GCC is a separate executable called "cc"

---

```
1 $ clang hello.cpp -o /tmp/hello -lstdc++
```

---

*Figure 2.9: Compilation for `hello.cpp`: We use `clang` executable as the compiler. We also use standard C++ library i.e. we link the executable with `-lstdc++`. The result is stored at `/tmp/hello`*

As you may see, we specified `-lstdc++` library option because we used `<iostream>` header from the standard C++ library. We also specified the output for executable (`/tmp/hello`) with `-o` option.

## 2.2.2 Compilation phases

We used 2 inputs for our example program. The first one is our source code, the second one is a shared library for standard C++ library. The clang driver should combine the inputs together, pass them via different phases of compilation process and finally provide the executable file on the target platform.

Clang uses the same typical compiler workflow as shown in fig. 2.2. You can ask Clang to show the phases using `-ccc-print-phases` additional argument

```
$ clang hello.cpp -o /tmp/hello -lstdc++ -ccc-print-phases
```

The output for the command is the following

```
+ 0: input, "hello.cpp", c++
+ 1: preprocessor, {0}, c++-cpp-output
+ 2: compiler, {1}, ir
+ 3: backend, {2}, assembler
+ 4: assembler, {3}, object
|- 5: input, "1%dM", object
6: linker, {4, 5}, image
```

We can visualize the output as shown in fig. 2.10. As we can see in fig. 2.10, the driver receives an input file `hello.cpp`, which is a C++ file. The file is processed by the preprocessor, and we obtain the preprocessor output (marked as `c++-cpp-output`). The result is compiled into IR form by the compiler, and then the backend converts it into assembly form. This form is later transformed into an object file. The final object file is combined with another object (`libstdc++`) to produce the final binary (`image`).

## 2.2.3 Tools execution

The phases are combined into several tool executions. The Clang driver invokes different programs to produce the final executable. Specifically, for our example, it calls the `clang` compiler and the `ld` linker. Both programs require additional arguments that are set up by the driver.

For instance, our example program (`hello.cpp`) includes the following header:

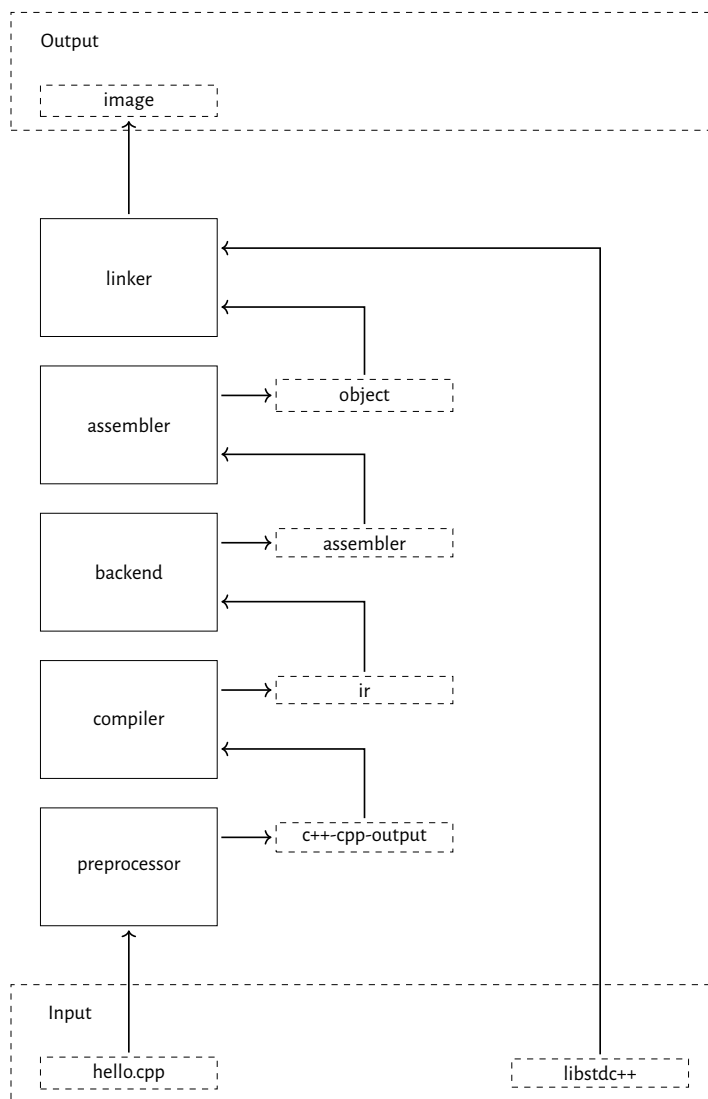


Figure 2.10: Clang driver phases

---

```

1 #include <iostream>
2 ...

```

---

We have not specified any additional arguments (such as search paths, for example `-I`) when we invoked the compilation. However, different architectures and operating systems might have different paths for locating headers.

On Fedora 37, the header is located in the `/usr/include/c++/12/iostream` folder. We can examine a detailed description of the process executed by the driver and the arguments used with the `###` option:

```
$ clang hello.cpp -o /tmp/hello -lstdc++ -###
```

The output for this command is quite extensive, and certain parts have been omitted here. Please refer to fig. 2.11.

---

```

1 clang version 16.0.0 (https://github.com/llvm/llvm-project.git ...)
2 "<...>/llvm-project/build/bin/clang-16"
3     "-cc1" ... \
4     "-internal-isystem" \
5     "/usr/include/c++/12" ... \
6     "-internal-isystem" \
7     "/usr/include/c++/12/x86_64-redhat-linux" ... \
8     "-internal-isystem" ... \
9     "<...>/llvm-project/build/lib/clang/16/include" ... \
10    "-internal-externc-isystem" \
11    "/usr/include" ... \
12    "-o" "/tmp/hello-XXX.o" "-x" "c++" "hello.cpp"
13    ".../bin/ld" ... \
14    "-o" "/tmp/hello" ... \
15    "/tmp/hello-XXX.o" \
16    "-lstdc++" ...

```

---

*Figure 2.11: Clang driver, tools execution, the host system is Fedora 37.*

As we can see in fig. 2.11, the driver initiates two processes: `clang-16` with the `-cc1` flag (see lines 2-12) and the linker `ld` (see lines 13-16). The clang compiler implicitly receives several search paths, as seen in lines 5, 7, 9 and 11. These paths are necessary for the inclusion of the `iostream` header in the test program.

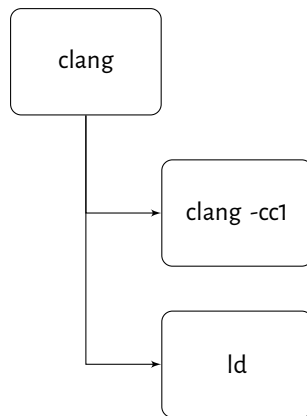


Figure 2.12: Clang driver, tools execution: clang driver runs 2 executables: The clang executable with `-cc1` flag and linker - `ld` executable

The output of the first executable (`/tmp/hello-XXX.o`) serves as input for the second one (see lines 12 and 15). The arguments `-lstdc++` and `-o /tmp/hello` are set for the linker, while the first argument (`hello.cpp`) is provided for the compiler invocation (first executable).

The process can be visualized as shown in fig. 2.12, where we can see that two executables are executed as part of the compilation process. The first one is `bin/clang-16` with a special flag (`-cc1`). The second one is the linker: `bin/ld`.

## 2.2.4 Combine all together

We can summarize the knowledge we have acquired so far using fig. 2.13. The figure illustrates two different processes started by the clang driver. The first one is `clang -cc1` (compiler), and the second one is `ld` (linker). The compiler process is the same executable as the clang driver (`clang`), but it is run with a special argument: `-cc1`. The compiler produces an object file that is then processed by the linker (`ld`) to generate the final binary.

In fig. 2.13, we can observe similar components of the compiler mentioned earlier (see [Section 2.1, Getting started with compilers](#)). However, the main difference is that the **preprocessor** (part of the lexer) is shown separately, while the frontend and middle-end are combined into the **compiler**. Additionally, the figure depicts an **assembler** that is executed by the driver to generate the object code. It is important to note that the assembler can be integrated, as shown in fig. 2.13, or it may require a separate process to be executed.

*Here is an example of specifying an external assembler using the `-c` (compile only) and `-o` (output file) options, along with the appropriate flags for your platform:*

```
$<...>/llvm-project/build/bin/clang -c hello.cpp -o /tmp/hello.o
```

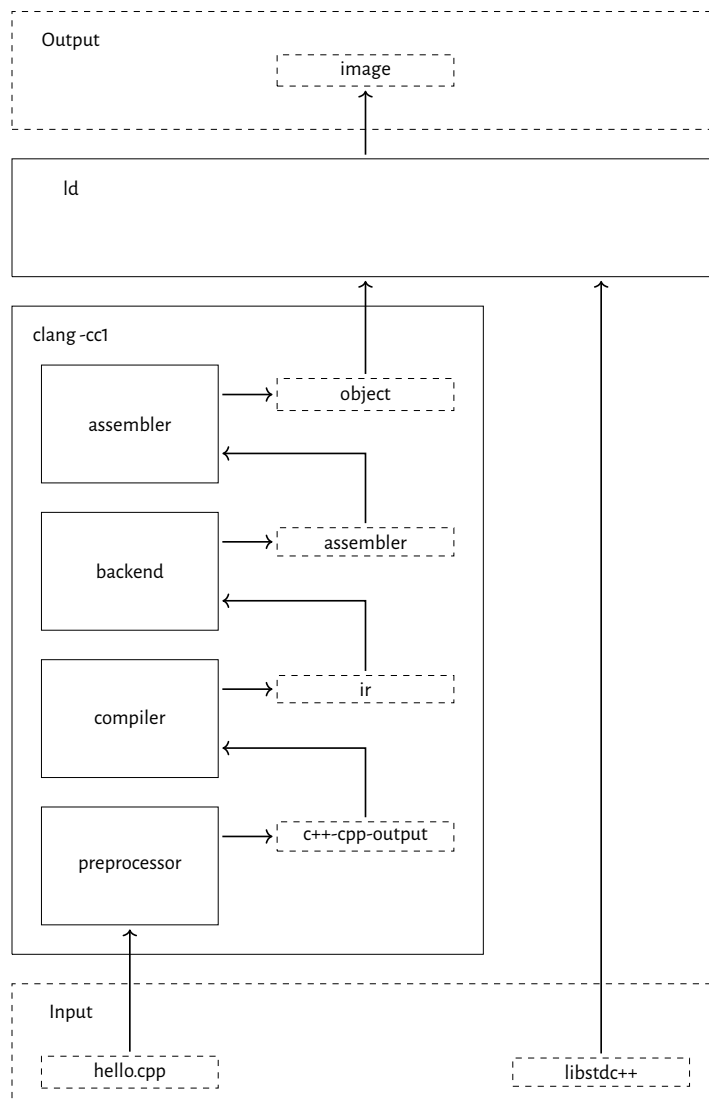


Figure 2.13: Clang driver: The driver got input file `hello.cpp` that is a C++ file. It starts 2 processes: `clang` and `ld`. The first one does real compilation and starts the integrated assembler. The last one is the linker (`ld`) that produces the final binary (`image`) from the result got from compiler and external library (`libstdc++`)

```
as -o /tmp/hello.o /tmp/hello.s
```

## 2.2.5 Debugging clang

We're going to step through a debugging session for our compilation process, illustrated in fig. 2.9. Our chosen point of interest, or breakpoint, is the `clang::ParseAST` function. In a typical debug session, which resembles the one outlined in fig. 1.4, you would feed command-line arguments following the `-` symbol. The command should look like this:

```
$lldb <...>/llvm-project/build/bin/clang -- hello.cpp -o /tmp/hello -lstdc++
```

In this case, `<...>` represents the directory path used to clone the LLVM project.

Unfortunately, this approach doesn't work with the Clang compiler:

---

```
1 $ lldb <...>/llvm-project/build/bin/clang -- src/part1/ch2_arch/hello.cpp -o
   → /tmp/hello.o -lstdc++
2 ...
3 (lldb) b clang::ParseAST
4 ...
5 (lldb) r
6 ...
7 2 locations added to breakpoint 1
8 ...
9 Process 247135 stopped and restarted: thread 1 received signal: SIGCHLD
10 Process 247135 stopped and restarted: thread 1 received signal: SIGCHLD
11 Process 247135 exited with status = 0 (0x00000000)
12 (lldb)
```

---

As we can see from line 7, the breakpoint was set but the process finished successfully (line 11) without any interruptions. In other words, our breakpoint didn't trigger in this instance.

Understanding the internals of clang-driver can help us identify the problem at hand. As mentioned earlier, the clang executable acts as a driver in this context, running two separate processes (refer to fig. 2.12). Therefore, if we wish to debug the compiler, we need to run it using the `-cc1` option.

*It's worth mentioning a certain optimization implemented in clang in 2019 [Ganea, "2019"]. When using the `-c` option, the clang driver doesn't spawn a new process for the compiler:*

```
$<...>/llvm-project/build/bin/clang -c hello.cpp -o /tmp/hello.o -###
clang version 16.0.0 ...
```



```

InstalledDir: <...>/llvm-project/build/bin
(in-process)
"<...>/llvm-project/build/bin/clang-16" "-cc1" ... "hello.cpp"
...

```

*As shown above, the clang driver does not spawn a new process and instead calls the "cc1" tool within the same process. This feature not only improves the compiler's performance but can also be leveraged for clang debugging.*

Upon using `-cc1` option and excluding the `-lstdc++` option (which is specific to the second process, the ld linker), the debugger will generate the following output:

---

```

1 $ lldb <...>/llvm-project-16/build/bin/clang -- -cc1 hello.cpp -o
   ↪ /tmp/hello.o
2 ...
3 (lldb) b clang::ParseAST
4 ...
5 (lldb) r
6 ...
7 2 locations added to breakpoint 1
8 Process 249890 stopped
9 * thread #1, name = 'clang', stop reason = breakpoint 1.1
10   frame #0: 0x00007fffe803eae0 ... at ParseAST.cpp:116:3
11   113
12   114 void clang::ParseAST(...) {
13   115     // Collect global stats on Decl/Stmts ...
14 -> 116     if (PrintStats) {
15   117         Decl::EnableStatistics();
16   118         Stmt::EnableStatistics();
17   119     }
18 (lldb) c
19 Process 249890 resuming
20 hello.cpp:1:10: fatal error: 'iostream' file not found
21 #include <iostream>
22     ~~~~~
23 1 error generated.
24 Process 249890 exited with status = 1 (0x00000001)
25 (lldb)

```

---

Thus we can see that we were able to successfully set the breakpoint, but the process ended with an error (see lines 20-24). This error arose because we omitted certain search paths, which are typically appended implicitly by the clang driver, necessary to find all the includes required for successful compilation.

We can successfully execute the process if we explicitly include all necessary arguments in the compiler invocation. Here's how to do that:

---

```
1 $ lladb <...>/llvm-project/build/bin/clang -- -cc1 -internal-isystem
   → /usr/include/c++/12 -internal-isystem
   → /usr/include/c++/12/x86_64-redhat-linux -internal-isystem
   → <...>/llvm-project/build/lib/clang/16/include -internal-externc-isystem
   → /usr/include hello.cpp -o /tmp/hello.o
2 ...
3 (lldb) b clang::ParseAST
4 ...
5 (lldb) r
6 ...
7 2 locations added to breakpoint 1
8 Process 251736 stopped
9 * thread #1, name = 'clang', stop reason = breakpoint 1.1
10   frame #0: 0x00007fffe803eae0 ... at ParseAST.cpp:116:3
11   113
12   114 void clang::ParseAST(...) {
13   115     // Collect global stats on Decl/Stmts ...
14 -> 116     if (PrintStats) {
15   117         Decl::EnableStatistics();
16   118         Stmt::EnableStatistics();
17   119     }
18 (lldb) c
19 Process 251736 resuming
20 Process 251736 exited with status = 0 (0x00000000)
21 (lldb)
```

---

In conclusion, we have successfully demonstrated the debugging of a clang compiler invocation. The techniques presented can be effectively employed for exploring the internals of a compiler and addressing compiler-related bugs.

## 2.3 Clang frontend overview

It's evident that the clang compiler toolchain conforms to the pattern widely described in various compiler books [Cooper and Torczon, 2012]. However, the clang's frontend part diverges significantly from a typical compiler frontend. The primary reason for this distinction is the complexity of the C++ language. Some features, like macros, can modify the source code itself, while others, like typedef, can influence the kind of token. Clang can also generate output in a variety of formats. For instance, the following command generates an aesthetically pleasing HTML view of the program shown in fig. 2.5:

```
$ clang -cc1 -emit-html max.cpp
```

Take note that we pass the argument to emit the HTML form of the source program to the clang frontend, specified with the `-cc1` option. Alternatively, you can pass an option to the frontend via the `-Xclang` option, which requires an additional argument representing the option itself. For example:

```
$ clang -fsyntax-only -Xclang -emit-html max.cpp
```

You may notice that in the command above, we utilized the `-fsyntax-only` option, instructing Clang to only execute the preprocessor, parser, and semantic analysis stages.

Accordingly, we can instruct the Clang frontend to perform different actions and produce varying types of output based on the provided compilation options. The base class for these actions is termed `FrontendAction`.

### 2.3.1 Frontend action

The Clang frontend is capable of executing only one frontend action at a time. A frontend action is a specific task or process that the frontend performs based on the provided compiler option. Below is a list of some possible frontend actions (the table only includes a subset of the available frontend actions):

| FrontendAction   | Compiler option                  | Description                            |
|------------------|----------------------------------|--|
| EmitObjAction    | <code>-emit-obj</code> (default) | Compile to an object file              |
| EmitBCAction     | <code>-emit-llvm-bc</code>       | Compile to LLVM bytecode               |
| EmitLLVMAction   | <code>-emit-llvm</code>          | Compile to LLVM readable form          |
| ASTPrintAction   | <code>-ast-print</code>          | Build ASTs and then pretty-print them. |
| HTMLPrintAction  | <code>-emit-html</code>          | Prints the program source in HTML form |
| DumpTokensAction | <code>-dump-tokens</code>        | Prints preprocessor tokens             |

Table 2.1: Frontend actions

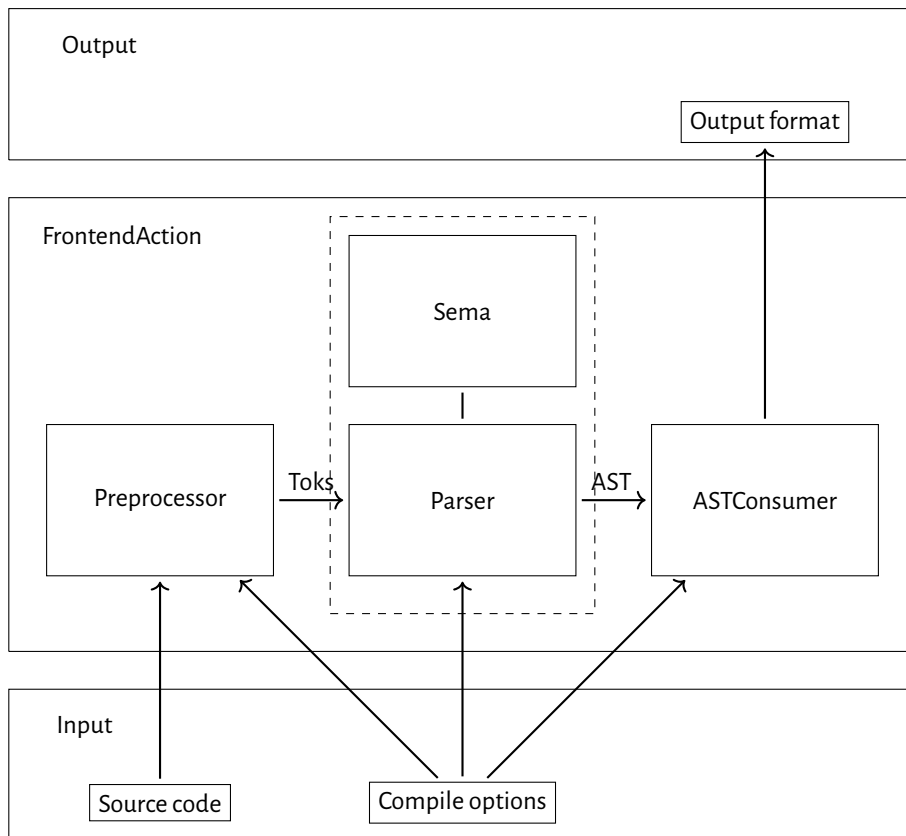


Figure 2.14: Clang frontend components

The diagram shown in fig. 2.14 illustrates the basic frontend architecture, which is similar to the architecture shown in fig. 2.4. However, there are notable differences specific to Clang.

One significant change is the naming of the lexer. In Clang, the lexer is referred to as the **Preprocessor**. This naming convention reflects the fact that the lexer implementation is encapsulated within the **Preprocessor** class. This alteration was inspired by the unique aspects of the C/C++ language, which includes special types of tokens (macros) that require specialized preprocessing.

Another noteworthy deviation is found in the parser component. While conventional compilers typically perform both syntax and semantic analysis within the parser, Clang distributes these tasks across different components. The **Parser** component focuses solely on syntax analysis, while the **Sema** component handles semantic analysis.

Furthermore, Clang offers the ability to produce output in different forms or formats. For example, the **CodeGenAction** class serves as the base class for various code generation actions, such as **EmitObjAction**.

or `EmitLLVMAction`.

We will use the code for `max` function from [fig. 2.5](#) for our future exploration of the Clang frontend's internals:

---

```

1 int max(int a, int b) {
2     if (a > b)
3         return a;
4     return b;
5 }
```

---

By utilizing the `-cc1` option, we can directly invoke the Clang frontend, bypassing the driver. This approach allows us to examine and analyze the inner workings of the Clang frontend in greater detail.

## 2.3.2 Preprocessor

The first part is the Lexer that is called as Preprocessor in Clang. Its primary goal is to convert the input program into a stream of tokens. You can print the token stream using the `-dump-tokens` options as follows

```
$ clang -cc1 -dump-tokens max.cpp
```

The output of the command is as shown below:

```

int 'int'      [StartOfLine]  Loc=<max.cpp:1:1>
identifier 'max' [LeadingSpace] Loc=<max.cpp:1:5>
l_paren '('    Loc=<max.cpp:1:8>
int 'int'      Loc=<max.cpp:1:9>
identifier 'a' [LeadingSpace] Loc=<max.cpp:1:13>
comma ','      Loc=<max.cpp:1:14>
int 'int'      [LeadingSpace] Loc=<max.cpp:1:16>
identifier 'b' [LeadingSpace] Loc=<max.cpp:1:20>
r_paren ')'    Loc=<max.cpp:1:21>
l_brace '{'    [LeadingSpace] Loc=<max.cpp:1:23>
if 'if' [StartOfLine] [LeadingSpace] Loc=<max.cpp:2:3>
l_paren '('    [LeadingSpace] Loc=<max.cpp:2:6>
identifier 'a' Loc=<max.cpp:2:7>
greater '>'    [LeadingSpace] Loc=<max.cpp:2:9>
identifier 'b' [LeadingSpace] Loc=<max.cpp:2:11>
r_paren ')'    Loc=<max.cpp:2:12>
return 'return' [StartOfLine] [LeadingSpace] Loc=<max.cpp:3:5>
identifier 'a' [LeadingSpace] Loc=<max.cpp:3:12>
```

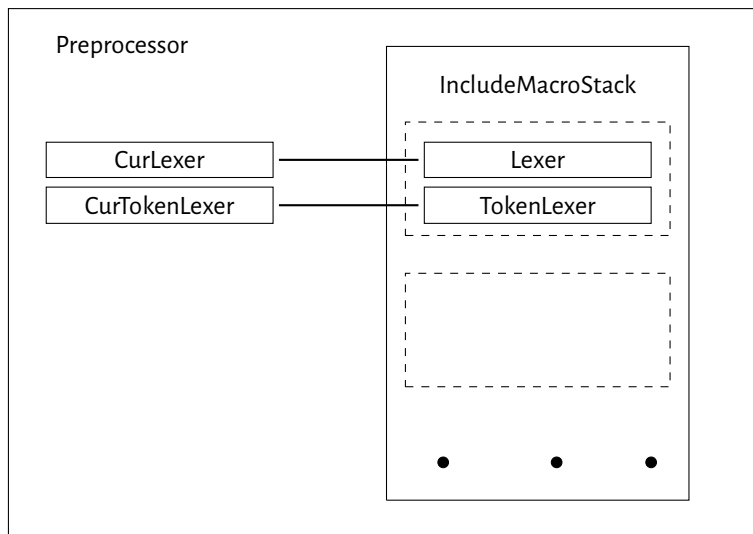


Figure 2.15: Preprocessor (clang lexer) class internals

```

semi ';'          Loc=<max.cpp:3:13>
return 'return'   [StartOfLine] [LeadingSpace]   Loc=<max.cpp:4:3>
identifier 'b'    [LeadingSpace] Loc=<max.cpp:4:10>
semi ';'          Loc=<max.cpp:4:11>
r_brace '}'       [StartOfLine] Loc=<max.cpp:5:1>
eof ''           Loc=<max.cpp:5:2>
  
```

As we can see, there are different types of tokens, such as language keywords (e.g., `int`, `return`), identifiers (e.g., `max`, `a`, `b`, etc.), and special symbols (e.g., semicolon, comma, etc.). The tokens for our small program are called **normal tokens**, which are returned by the lexer.

In addition to normal tokens, Clang has an additional type of token called **annotation tokens**. The primary difference is that these tokens also store additional semantic information. For instance, a sequence of normal tokens can be replaced by the parser with a single annotation token that contains information about the type or C++ scope. The primary reason for using such tokens is performance, as it allows for the prevention of reparsing when the parser needs to backtrack.

C/C++ language has some specifics that influence the internal implementation of the Preprocessor class. The first one is about macros. The Preprocessor class has two different helper classes to retrieve tokens:

- The `Lexer` class is used to convert a text buffer into a stream of tokens.
- The `TokenLexer` class is used to retrieve tokens from macro expansions.

It should be noted that only one of these helpers can be active at a time.

Another specific aspect of C/C++ is the `#include` directive<sup>3</sup>. In this case, we need to maintain a stack of includes, where each include can have its own `TokenLexer` or `Lexer`, depending on whether there is a macro expansion within it. As a result, the `Preprocessor` class keeps a stack of lexers (`IncludeMacroStack` class) for each `#include` directive, as shown in fig. 2.15.

### 2.3.3 Parser and Sema

The Parser and Sema are crucial components of the Clang compiler frontend. They handle the syntax and semantic analysis of the source code, producing an AST as output. This tree can be visualized for our test program using the command:

```
$ clang -cc1 -ast-dump max.cpp
```

The output of this command is shown below

```
TranslationUnitDecl 0xc4b578 <<invalid sloc>> <invalid sloc>
|-TypedefDecl 0xc4bde0 <<invalid sloc>> <invalid sloc> implicit __int128_t '__int128'
| `~BuiltinType 0xc4bb40 '__int128'
...
~-FunctionDecl 0xc91580 <max.cpp:1:1, line:5:1> line:1:5 max 'int (int, int)'
  |-ParmVarDecl 0xc91428 <col:9, col:13> col:13 used a 'int'
  |-ParmVarDecl 0xc914a8 <col:16, col:20> col:20 used b 'int'
  ~-CompoundStmt 0xc917b8 <col:23, line:5:1>
    |-IfStmt 0xc91750 <line:2:3, line:3:12>
      | |-BinaryOperator 0xc916e8 <line:2:7, col:11> 'bool' '>'
      | | |-ImplicitCastExpr 0xc916b8 <col:7> 'int' <LValueToRValue>
      | | | `~DeclRefExpr 0xc91678 <col:7> 'int' lvalue ParmVar 0xc91428 'a' 'int'
      | | `~ImplicitCastExpr 0xc916d0 <col:11> 'int' <LValueToRValue>
      | |   `~DeclRefExpr 0xc91698 <col:11> 'int' lvalue ParmVar 0xc914a8 'b' 'int'
      | `~ReturnStmt 0xc91740 <line:3:5, col:12>
      |   `~ImplicitCastExpr 0xc91728 <col:12> 'int' <LValueToRValue>
      |     `~DeclRefExpr 0xc91708 <col:12> 'int' lvalue ParmVar 0xc91428 'a' 'int'
    `~ReturnStmt 0xc917a8 <line:4:3, col:10>
      `~ImplicitCastExpr 0xc91790 <col:10> 'int' <LValueToRValue>
        `~DeclRefExpr 0xc91770 <col:10> 'int' lvalue ParmVar 0xc914a8 'b' 'int'
```

Clang utilizes a hand-written recursive-descent parser [Community, 2023a]. This parser can be considered simple, and this simplicity was one key reason for its selection. Additionally, the complex rules specified for the C/C++ languages necessitated an ad-hoc parser with easily adaptable rules.

Let's explore how this works with our example. Parsing begins with a top-level declaration known as a `TranslationUnitDecl`, representing a single translation unit. The C++ standard defines a translation unit as follows [for Standardization, 2020, lex.separate]:

<sup>3</sup>which is also applicable to the import directive

*A source file together with all the headers (16.5.1.2) and source files included (15.3) via the preprocessing directive `#include`, less any source lines skipped by any of the conditional inclusion (15.2) preprocessing directives, is called a translation unit.*

The parser first recognizes that the initial tokens from the source code correspond to a function definition, as defined in the C++ standard [for Standardization, 2020, dcl.fct.def.general]:

```
function-definition :
    ... declarator ... function-body
    ...
```

The corresponding code is below

---

```
1 int max(...) {
2     ...
3 }
```

---

The function definition necessitates a declarator and function body. We'll start with the declarator, defined in the C++ standard as [for Standardization, 2020, dcl.decl.general]:

```
declarator:
    ...
    ... parameters-and-qualifiers ...
    ...
parameters-and-qualifiers:
    ( parameter-declaration-clause ) ...
    ...
parameter-declaration-clause:
    parameter-declaration-list ...
parameter-declaration-list:
    parameter-declaration
    parameter-declaration-list , parameter-declaration
```

In other words, the declarator specifies a list of parameter declarations within brackets. The corresponding piece of code from the source is as follows:

---

```
1 ... (int a, int b)
2     ...
```

---

The function definition, as stated above, also requires a function body. The C++ standard specifies the function body as follows: [for Standardization, 2020, dcl.fct.def.general]



```
function-body:
    ... compound-statement
    ...
```

Thus the function body consists of a compound statement, which is defined as follows in the C++ standard [[for Standardization, 2020](#), stmt.block]

```
compound-statement:
    { statement-seq ... }
statement-seq:
    statement
    statement-seq statement
```

Therefore, it describes a sequence of statements enclosed within { . . . } brackets.

Our program has two types of statements: the conditional (if) statement and the return statement. These are represented in the C++ grammar definition as follows [[for Standardization, 2020](#), stmt.pre]:

```
statement:
    ...
    selection-statement
    ...
    jump-statement
    ...
```

In this context, the selection statement corresponds to the `if` condition in our program, while the jump statement corresponds to the `return` operator.

Let's examine the jump statement in more detail [[for Standardization, 2020](#), stmt.jump.general]:

```
jump-statement:
    ...
    return expr-or-braced-init-list;
    ...
```

where `expr-or-braced-init-list` is defined as [[for Standardization, 2020](#), dcl.init.general]:

```
expr-or-braced-init-list:
    expression
    ...
```

In this context, the `return` keyword is followed by an expression and a semicolon. In our case, there's an implicit cast expression that automatically converts the variable into the required type (`int`).

It can be enlightening to examine the parser's operation through the LLDB debugger:

---

```

1 $ lldb <...>/llvm-project/build/bin/clang -- -cc1 max.cpp
2 ...
3 (lldb) b clang::Parser::ParseReturnStatement
4 (lldb) r
5 ...
6 (lldb) c
7 ...
8 * thread #1, name = 'clang', stop reason = breakpoint 1.1
9   frame #0: ... clang::Parser::ParseReturnStatement(...) at
   ↳ ParseStmt.cpp:2358:3
10   2355 ///           'co_return' expression[opt] ';'
11   2356 ///           'co_return' braced-init-list ';'
12   2357 StmtResult Parser::ParseReturnStatement() {
13 -> 2358   assert((Tok.is(tok::kw_return) || Tok.is(tok::kw_co_return)) &&
14   2359           "Not a return stmt!");
15   2360   bool IsCoreturn = Tok.is(tok::kw_co_return);
16   2361   SourceLocation ReturnLoc = ConsumeToken(); // eat the 'return'.
17 (lldb) bt
18 * frame #0: ... clang::Parser::ParseReturnStatement( ...
19   ...
20   frame #2: ... clang::Parser::ParseStatementOrDeclaration( ...
21   frame #3: ... clang::Parser::ParseCompoundStatementBody( ...
22   frame #4: ... clang::Parser::ParseFunctionStatementBody( ...
23   frame #5: ... clang::Parser::ParseFunctionDefinition( ...
24 ...

```

---

*Figure 2.16: Second return statement parsing at max.cpp example program*

As you can see in fig. 2.16, line 3, we've set a breakpoint for the parsing of return statements<sup>4</sup>. Our program has two return statements. We bypass the first call (line 6) and halt at the second method invocation (line 13). The backtrace (from the 'bt' command at line 17) displays the call stack for the parsing process. This stack mirrors the parsing blocks we described earlier, adhering to the C++ grammar detailed in [for Standardization, 2020, lex.separate].

The parsing results in the generation of AST. We can also inspect the process of AST creation using the debugger. To do this, we need to set a corresponding breakpoint at the `clang::ReturnStmt::Create`

---

<sup>4</sup>Specifically, at the `clang::Parser::ParseReturnStatement` method

method:

---

```

1 $ lldb <...>/llvm-project/build/bin/clang -- -cc1 max.cpp
2 ...
3 (lldb) b clang::ReturnStmt::Create
4 (lldb) r
5 ...
6 (lldb) c
7 ...
8 * thread #1, name = 'clang', stop reason = breakpoint 1.1
9   frame #0: ... clang::ReturnStmt::Create(...) at Stmt.cpp:1204:8
10   1201
11   1202 ReturnStmt *ReturnStmt::Create( ... ) {
12 -> 1204   bool HasNRVOCandidate = NRVOCandidate != nullptr;
13   1205   ...
14   1206   ...
15   1207   return new (Mem) ReturnStmt(RL, E, NRVOCandidate);
16 (lldb) bt
17 * thread #1, name = 'clang', stop reason = breakpoint 1.1
18 * frame #0: ... clang::ReturnStmt::Create( ...
19   frame #1: ... clang::Sema::BuildReturnStmt( ...
20   frame #2: ... clang::Sema::ActOnReturnStmt( ...
21   frame #3: ... clang::Parser::ParseReturnStatement( ...
22   frame #4: ... clang::Parser::ParseStatementOrDeclarationAfterAttributes(
   →   ...
23   ...

```

---

As can be seen, the AST node for the return statement is created by the Sema component.

The beginning of the return statement parser can be located in frame 4: As we can observe, there is a reference to the C99 standard [[International Organization for Standardization \(ISO\), 1999](#)] for the corresponding statement. The standard [[International Organization for Standardization \(ISO\), 1999](#)] provides a detailed description of the statement and the process for handling it.

The code assumes that the current token is of type `tok::kw_return`, and in this case, the parser invokes the relevant `clang::Parser::ParseReturnStatement` method.

While the process of AST node creation can vary across different C++ constructs, it generally follows the pattern displayed in [fig. 2.17](#).

---

```

1 (lldb) f 4
2 frame #4: ... clang::Parser::ParseStatementOrDeclarationAfterAttributes( ...
3   325     break;
4   326     case tok::kw_return:                // C99 6.8.6.4: return-statement
5 -> 327     Res = ParseReturnStatement();
6   328     SemiError = "return";
7   329     break;
8   330     case tok::kw_co_return:              // C++ Coroutines: ...
9 (lldb)

```

---

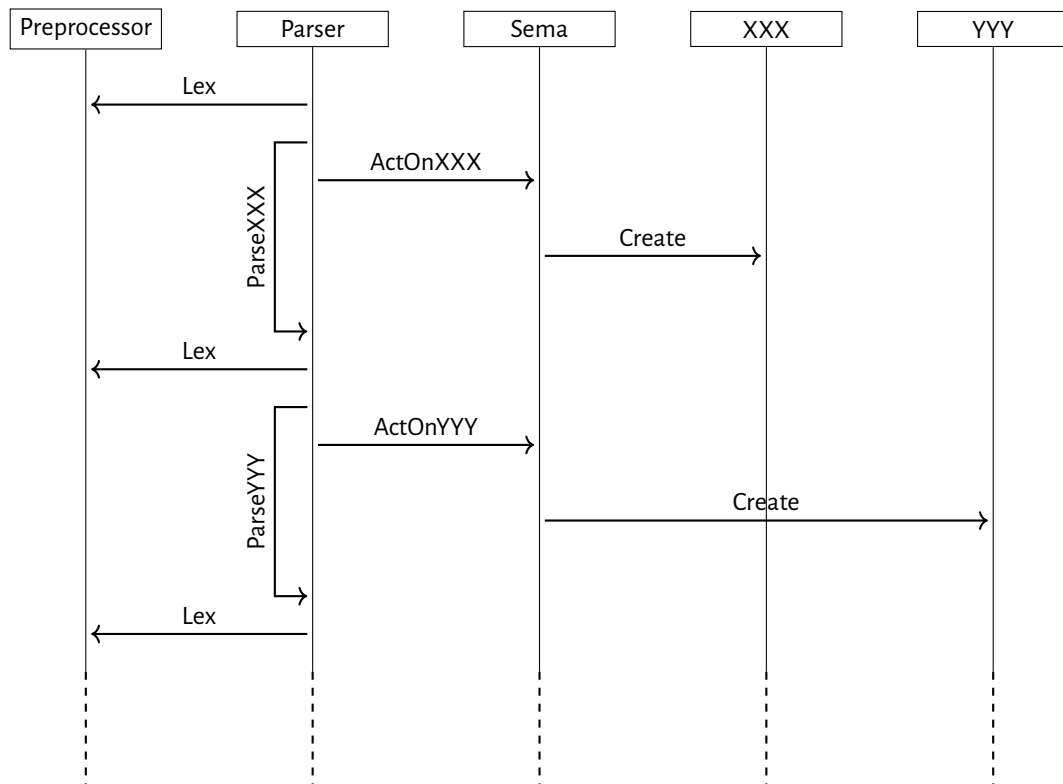


Figure 2.17: C++ parsing in Clang frontend

As can be seen, the Parser invokes the `Preprocessor::Lex` method to retrieve a token from the lexer.

It then calls a method corresponding to the token, for example, `Parser::ParseXXX` for the token `XXX`. This method then calls `Sema::ActOnXXX`, which creates the corresponding object using `XXX::Create`. The process is then repeated with a new token.

With this, we have now fully explored how the typical compiler frontend flow is implemented in Clang. We can see how the lexer component (the preprocessor) works in tandem with the parser (which comprises the parser and sema components) to produce the primary data structure for future code generation: the Abstract Syntax Tree (AST). The AST is not only essential for code generation but also for code analysis and modification. Clang provides easy access to the AST, thereby enabling the development of a diverse range of compiler tools.

## 2.4 Summary

In this chapter, we have acquired a basic understanding of compiler architecture and delved into the various stages of the compilation process, with a focus on the Clang driver. We have explored the internals of the Clang frontend, studying the preprocessor that transforms a program into a set of tokens, and the parser, which interacts with a component called 'Sema'. Together, these elements perform syntax and semantic analysis.

The upcoming chapter will center on the Clang Abstract Syntax Tree (AST)—the primary data structure employed in various Clang tools. We will discuss its construction and the methods for traversing it.

## 2.5 Further reading

- Working Draft, Standard for Programming Language C++: <https://eel.is/c++draft/>
- “Clang” CFE Internals Manual: <https://clang.llvm.org/docs/InternalsManual.html>
- Keith Cooper and Linda Torczon: Engineering A Compiler, 2012 [[Cooper and Torczon, 2012](#)]



# Clang AST

The Abstract Syntax Tree (AST) is a fundamental algorithmic structure used to represent the results of parsing. The AST serves as the framework for the Clang frontend and is the primary tool for various Clang utilities, including linters. Clang offers sophisticated tools for searching (or matching) various AST nodes. These tools are implemented using a DSL (domain-specific language). It's crucial to understand its implementation to use it effectively.

We will start with the basic data structures and the class hierarchy that Clang uses to construct the AST. Additionally, we will explore the methods used for AST traversal and highlight some helper classes that facilitate node matching during this traversal.

## 3.1 AST

The AST is usually depicted as a tree, with its leaf nodes corresponding to various objects, such as function declarations and loop bodies. Typically, the AST represents the result of syntax analysis, i.e., parsing. Clang's AST nodes were designed to be immutable. This design requires that the Clang AST store results from semantic analysis, meaning the Clang AST represents the outcomes of both syntax and semantic analyses.

*Although Clang also employs an AST, it's worth noting that the Clang AST is not a true tree. The presence of backward edges makes "graph" a more appropriate term for describing Clang's AST.*

Typical tree structure implemented in C++ has all nodes derived from a base class. Clang uses a different

approach. It splits different C++ constructions into separate groups with basic classes for each of them.

- Statements. `clang::Stmt` is the basic class for all statements. That includes ordinary statements such as if-statement (`clang::IfStmt` class) as well as expressions and other C++ constructions
- Declarations. `clang::Decl` is the basic class for declarations. This includes a variable, typedef, function, struct, etc. There is also a separate basic class for declarations with context i.e. declarations that might contain another declarations: `clang::DeclContext`. Translation units (`clang::TranslationUnit` class) and namespaces (`clang::NamespaceDecl` class) are typical examples for declarations with context.
- Types. C++ has rich type system. It includes basic types such as `int` for integers as well as custom defined types and types redefinition via `typedef` or `using`. Types in C++ can have qualifiers such as `const` and can represent different memory addressing modes aka pointers, references etc. Clang uses `clang::Type` as the basic class for types representations in AST.

It's worth noting that there are additional relations between the groups. For example, the `clang::DeclStmt` class, that inherits from `clang::Stmt`, has methods to retrieve corresponding declarations. Additionally, expressions (represented by the `clang::Expr` class) which inherit from `clang::Stmt` have methods to work with types. Let's look at all the groups in detail.

### 3.1.1 Statements

`Stmt` is the basic class for all statements. The statements can be combined into two sets (see fig. 3.1). The first one contains statements with values and an opposite group is for statements without values.

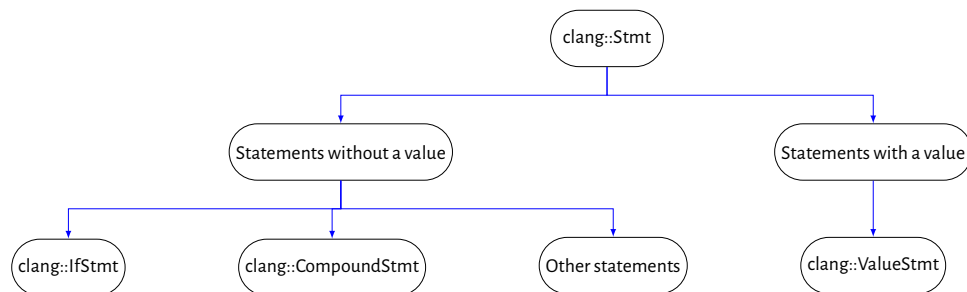


Figure 3.1: Clang AST: Statements

The group of statements without a value consist of different C++ constructions such as if-statement (`clang::IfStmt` class) or compound statement (`CompoundStmt` class). The majority of all statements falls into the group.

The group of statement with a value consists of one base class `clang::ValueStmt` that has several children such as `clang::LabelStmt` (for labels representation) or `clang::ExprStmt` (for expressions representation), see fig. 3.2.



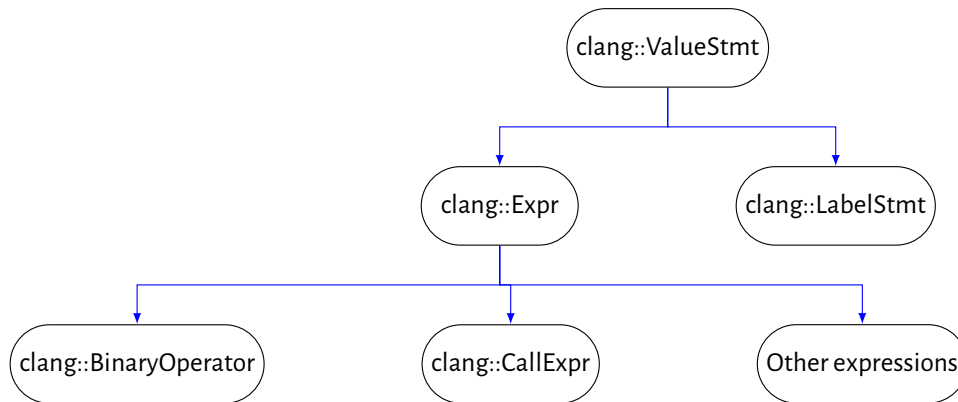


Figure 3.2: Clang AST: Statements with a value

### 3.1.2 Declarations

Declarations can also be combined into 2 primary groups: Declarations with context and without. Declarations with context can be considered as placeholders for other declarations. For example C++ namespace as well as translation unit or function declaration might contain other declarations. Declaration of a friend entity (`clang::DeclFriend`) can be considered as an example of a declaration without context.

It has to be noted that classes that are inherited from `DeclContext` have also `clang::Decl` as their top parent.

Some declarations can be redeclared as in the following example

---

```

1 extern int a;
2 int a = 1;
  
```

---

Such declarations have an additional parent that is implemented via `clang::Redeclarable<...>` template.

### 3.1.3 Types

C++ is a statically typed language, which means that the types of variables must be declared at compile-time. The types allow compiler to make a reasonable conclusion about the program meaning i.e. types are important part of semantic analysis. `clang::Type` is the basic class for types in Clang.

Types in C/C++ might have qualifiers that are called as CV-qualifiers at standard [for Standardization, 2020, `basic.type.qualifier`]. CV here is for 2 keywords `const` and `volatile` that can be used as the qualifier for a type.

C99 standard has an additional type qualifier *restrict* that is also supported by clang [International Organization for Standardization (ISO), 1999, 6.7.3]

Clang has a special class to support a type with qualifier: `clang::QualType` that is a pair of a pointer to `clang::Type` and a bit mask with information about the type qualifier. The class has a method to retrieve a pointer to the `clang::Type` and check different qualifiers. The code below shows how we can check a type for a const qualifier (the code is from LLVM release 16.0, `clang/lib/AST/ExprConstant.cpp`, line 3855)

---

```

1  bool checkConst(QualType QT) {
2      // Assigning to a const object has undefined behavior.
3      if (QT.isConstQualified()) {
4          Info.FFDiag(E, diag::note_constexpr_modify_const_type) << QT;
5          return false;
6      }
7      return true;
8  }
```

---

It's worth mentioning that `clang::QualType` has `operator->()` and `operator*()` implemented, i.e. it can be considered as a smart pointer for underlying `clang::Type` class.

In addition to the qualifiers type can have additional information that represents different memory address models, for instance there can be a pointer to an object or reference. `clang::Type` has the following helper methods to check different address models:

- `clang::Type::isPointerType()` for pointer type check
- `clang::Type::isReferenceType()` for reference type check

Types in C/C++ can also use aliases that are introduced by using `typedef` or `using` keywords. The following code defines `foo` and `bar` as aliases for `int` type

---

```

1  using foo = int;
2  typedef int bar;
```

---

Original types, `int` at our case, are called as canonical. You can test if the type is canonical or not using `clang::QualType::isCanonical()` method. `clang::QualType` also provides a method to retrieve the canonical type from an alias: `clang::QualType::getCanonicalType()`.

## 3.2 AST traversal

The compiler requires traversal of the AST to generate IR code. Thus, having a well-structured data structure for tree traversal is paramount for AST design. To put it another way, the design of the AST should prioritize facilitating easy tree traversal. A standard approach in many systems is to have a common base class for all AST nodes. This class typically provides a method to retrieve the node's children, allowing for tree traversal using popular algorithms like Breadth-First Search (BFS) [Cormen et al. \[2009\]](#). Clang, however, takes a different approach: its AST nodes don't share a common ancestor. This poses the question: how is tree traversal organized in Clang?

Clang employs three unique techniques:

- The Curiously Recurring Template Pattern (CRTP) for visitor class definition,
- Ad-hoc methods tailored specifically for different nodes,
- Macros, which can be perceived as the connecting layer between the ad-hoc methods and CRTP.

We will explore these techniques through a simple program designed to identify function definitions and display the function names together with their parameters.

### 3.2.1 DeclVisitor test tool

Our test tool will build upon the `clang::DeclVisitor` class, which is defined as a straightforward visitor class aiding in the creation of visitors for C/C++ declarations.

We will use the same CMake file as it was created for our first clang tool (see [fig. 1.6](#)). The sole addition for the new tool is the `clangAST` library. The resultant `CMakeLists.txt` is shown in [fig. 3.3](#):

---

```
2 project("declvisitor")
3
4 if ( NOT DEFINED ENV{LLVM_HOME})
5     message(FATAL_ERROR "$LLVM_HOME is not defined")
6 else()
7     message(STATUS "$LLVM_HOME found: $ENV{LLVM_HOME}")
8     set(LLVM_HOME $ENV{LLVM_HOME} CACHE PATH "Root of LLVM installation")
9     set(LLVM_LIB ${LLVM_HOME}/lib)
10    set(LLVM_DIR ${LLVM_LIB}/cmake/llvm)
11    find_package(LLVM REQUIRED CONFIG)
12    include_directories(${LLVM_INCLUDE_DIRS})
13    link_directories(${LLVM_LIBRARY_DIRS})
14    set(SOURCE_FILE DeclVisitor.cpp)
15    add_executable(declvisitor ${SOURCE_FILE})
16    set_target_properties(declvisitor PROPERTIES COMPILE_FLAGS "-fno-rtti")
17    target_link_libraries(declvisitor
18        LLVMSupport
19        clangAST
20        clangBasic
21        clangFrontend
22        clangSerialization
23        clangTooling
24    )
```

---

*Figure 3.3: CMakeLists.txt file for DeclVisitor test tool*

The main function of our tool is presented below:

---

```

1 #include "clang/Tooling/CommonOptionsParser.h"
2 #include "clang/Tooling/Tooling.h"
3 #include "llvm/Support/CommandLine.h" // llvm::cl::extrahelp
4
5 #include "FrontendAction.hpp"
6
7 namespace {
8   llvm::cl::OptionCategory TestCategory("Test project");
9   llvm::cl::extrahelp
10     CommonHelp(clang::tooling::CommonOptionsParser::HelpMessage);
11 } // namespace
12
13 int main(int argc, const char **argv) {
14   llvm::Expected<clang::tooling::CommonOptionsParser> OptionsParser =
15     clang::tooling::CommonOptionsParser::create(argc, argv, TestCategory);
16   if (!OptionsParser) {
17     llvm::errs() << OptionsParser.takeError();
18     return 1;
19   }
20   clang::tooling::ClangTool Tool(OptionsParser->getCompilations(),
21                                 OptionsParser->getSourcePathList());
22   return Tool.run(clang::tooling::newFrontendActionFactory<
23     clangbook::declvisitor::FrontendAction>()
24     .get());
25 }

```

---

Figure 3.4: The main function for DeclVisitor test tool

From lines 5 and 23, it's evident that we employ a custom frontend action specific to our project: `clangbook::declvisitor`. The code for this class is provided below

---

```
1 #include "Consumer.hpp"
2 #include "clang/Frontend/FrontendActions.h"
3
4 namespace clangbook {
5 namespace declvisitor {
6 class FrontendAction : public clang::ASTFrontendAction {
7 public:
8     virtual std::unique_ptr<clang::ASTConsumer>
9         CreateASTConsumer(clang::CompilerInstance &CI,
10                           clang::StringRef File) override {
11         return std::make_unique<Consumer>();
12     }
13 };
14 } // namespace declvisitor
15 } // namespace clangbook
```

---

*Figure 3.5: Custom FrontendAction class for DeclVisitor test tool*

You'll notice that we have overridden the `clang::ASTFrontendAction::CreateASTConsumer` function to instantiate our custom `clangbook::declvisitor::Consumer` (as highlighted in lines 9-12). The implementation for our tailored AST consumer is as follows:

---

```
1 #include "Visitor.hpp"
2 #include "clang/Frontend/ASTConsumers.h"
3
4 namespace clangbook {
5 namespace declvisitor {
6 class Consumer : public clang::ASTConsumer {
7 public:
8     Consumer() : V(std::make_unique<Visitor>()) {}
9
10    virtual void HandleTranslationUnit(clang::ASTContext &Context) override {
11        V->Visit(Context.getTranslationUnitDecl());
12    }
13
14 private:
15     std::unique_ptr<Visitor> V;
16 };
17 } // namespace declvisitor
18 } // namespace clangbook
```

---

*Figure 3.6: Consumer class for DeclVisitor test tool*

Here, we can see that we've created a sample visitor and invoked it using the `clang::ASTConsumer::HandleTranslationUnit` (see fig. 3.6, line 11). However, the most intriguing portion is the code for the visitor:

---

```
1 #include "clang/AST/DeclVisitor.h"
2
3 namespace clangbook {
4 namespace declvisitor {
5 class Visitor : public clang::DeclVisitor<Visitor> {
6 public:
7     void VisitFunctionDecl(const clang::FunctionDecl *FD) {
8         llvm::outs() << "Function: '" << FD->getName() << "'\n";
9         for (auto Param : FD->parameters()) {
10             Visit(Param);
11         }
12     }
13     void VisitParmVarDecl(const clang::ParmVarDecl *PVD) {
14         llvm::outs() << "\tParameter: '" << PVD->getName() << "'\n";
15     }
16     void VisitTranslationUnitDecl(const clang::TranslationUnitDecl *TU) {
17         for (auto Decl : TU->decls()) {
18             Visit(Decl);
19         }
20     }
21 };
22 } // namespace declvisitor
23 } // namespace clangbook
```

---

Figure 3.7: Visitor class implementation

We will explore the code in more depth later. For now, we observe that it prints the function name at line 8 and the parameter name at line 14.

We can compile our program using the same sequence of commands as we did for our test project, as detailed in [Section 1.4, Test project: syntax check with clang tool](#).



```
export LLVM_HOME=<...>/llvm-project/install
mkdir build
cd build
cmake -G Ninja -DCMAKE_BUILD_TYPE=Debug ..
ninja
```

*Figure 3.8: Configure and build commands for DeclVisitor test tool*

As you may notice we used `-DCMAKE_BUILD_TYPE=Debug` option for CMake. That is because we might want to investigate the result program under debugger.

We will use the same program we referenced in our previous investigations (see fig. 2.5) to also study AST traversal:

---

```
1 int max(int a, int b) {
2     if (a > b)
3         return a;
4     return b;
5 }
```

---

*Figure 3.9: Test program max.cpp*

This program consists of a single function, `max`, which takes two parameters: `a` and `b`, and returns the maximum of the two.

We can run our program as follows:

```
$ ./declvisitor max.cpp
...
Function: 'max'
    Parameter: 'a'
    Parameter: 'b'
```

Let's investigate the Visitor class implementation in detail.

## 3.2.2 Visitor implementation

Let's delve into the Visitor code (see fig. 3.7). Firstly, you'll notice an unusual construct where our class is derived from a base class parameterized by our own class:

---

```
5 class Visitor : public clang::DeclVisitor<Visitor> {
```

---

This construct is known as the Curiously Recurring Template Pattern, or CRTP.

The Visitor class has several callbacks that are triggered when a corresponding AST node is encountered. The first callback targets the AST node representing a function declaration

---

```
7  void VisitFunctionDecl(const clang::FunctionDecl *FD) {  
8      llvm::outs() << "Function: '" << FD->getName() << "'\n";  
9      for (auto Param : FD->parameters()) {  
10         Visit(Param);  
11     }  
12 }
```

---

*Figure 3.10: FunctionDecl callback*

As shown in fig. 3.10, the function name is printed on line 8. Our subsequent step involves printing the names of the parameters. To retrieve the function parameters, we can utilize the `parameters()` method from the `clang::FunctionDecl` class. This method was previously mentioned as an ad-hoc approach for AST traversal. Each AST node provides its own methods to access child nodes. Since we have an AST node of a specific type (i.e., `clang::FunctionDecl*`) as an argument, we can employ these methods.

The function parameter is passed to the `Visit(...)` method of the base class `clang::DeclVisitor<>`. This call is subsequently transformed into another callback, specifically for the `clang::ParmVarDecl` AST node.

---

```
13 void VisitParmVarDecl(const clang::ParmVarDecl *PVD) {  
14     llvm::outs() << "\tParameter: '" << PVD->getName() << "'\n";  
15 }
```

---

*Figure 3.11: ParmVarDecl callback*

You might be wondering how this conversion is achieved. The answer lies in a combination of the CRTP and C/C++ macros. To understand this, we need to dive into the `Visit()` method implementation of the `clang::DeclVisitor<>` class. This implementation heavily relies on C/C++ macros, so to get a glimpse of the actual code, we must expand these macros. This can be done using the `-E` compiler option. Let's make some modifications to our `CMakeLists.txt` and introduce a new custom target.

---

```

25  add_custom_command(
26      OUTPUT ${SOURCE_FILE}.preprocessed
27      COMMAND ${CMAKE_CXX_COMPILER} -E -I ${LLVM_HOME}/include
        ↳ ${CMAKE_CURRENT_SOURCE_DIR}/${SOURCE_FILE} >
        ↳ ${SOURCE_FILE}.preprocessed
28      DEPENDS ${SOURCE_FILE}
29      COMMENT "Preprocessing ${SOURCE_FILE}"
30  )
31  add_custom_target(preprocess ALL DEPENDS ${SOURCE_FILE}.preprocessed)

```

---

*Figure 3.12: Custom target to expand macros*

We can run the target as follows

```
ninja preprocess
```

The resulting file can be located in the build folder specified earlier, named

`DeclVisitor.cpp.preprocessed`. This is the folder we pointed to when executing the `cmake` command (see fig. 3.8). Within this file, the generated code for the `Visit()` method appears as follows:

---

```

1  RetTy Visit(typename Ptr<Decl>::type D) {
2      switch (D->getKind()) {
3          ...
4          case Decl::ParmVar: return
            ↳ static_cast<ImplClass*>(this)->VisitParmVarDecl(static_cast<typename
            ↳ Ptr<ParmVarDecl>::type>(D));
5          ...
6      }
7  }

```

---

This code showcases the use of the CRTP in Clang. In this context, CRTP is employed to redirect back to our `Visitor` class, which is referenced as `ImplClass`. CRTP allows the base class to call a method from an inherited class. This pattern can serve as an alternative to virtual functions and offers several advantages, the most notable being performance-related. Specifically, the method call is resolved at compile-time, eliminating the need for a vtable lookup associated with virtual method calls.

The code was generated using C/C++ macros, as demonstrated below. This particular code was sourced from the `clang/include/clang/AST/DeclVisitor.h` header.

---

```
1  #define DISPATCH(NAME, CLASS) \  
2  return  
→ static_cast<ImplClass*>(this)->Visit##NAME(static_cast<PTR(CLASS)>(D))
```

---

The NAME is replaced with the node name; in our case, it's `ParmVarDecl`.

The `DeclVisitor` is used to traverse C++ declarations. Clang also has `StmtVisitor` and `TypeVisitor` to traverse statements and types, respectively. These are built on the same principles as we considered in our example with the declaration visitor. However, these visitors come with several issues. They can only be used with specific groups of AST nodes. For instance, the `DeclVisitor` can only be used with descendants of the `Decl` class. Another limitation is that we are required to implement recursion. For example, we set up recursion to traverse the function declaration in lines 9-11 (fig. 3.7). The same recursion was employed to traverse declarations within the translation unit (see fig. 3.7, lines 17-19). This brings up another concern: it's possible to miss some parts of the recursion. For instance, our code will not function correctly if the `max` function declaration is specified inside a namespace. To address such scenarios, we would need to implement an additional visit method specifically for namespace declarations.

These challenges are addressed by the recursive visitor, which we will discuss shortly.

### 3.3 Recursive AST Visitor

Recursive AST visitors address the limitations observed with specialized visitors. We will create the same program, which searches for and prints function declarations along with their parameters, but we'll use a recursive visitor this time.

The `CMakeLists.txt` for recursive visitor test tool will be similar to used previously. The only project name and source file name was changed

---

```
1 cmake_minimum_required(VERSION 3.16)
2 project("recursivevisitor")
3
4 if ( NOT DEFINED ENV{LLVM_HOME})
5     message(FATAL_ERROR "$LLVM_HOME is not defined")
6 else()
7     message(STATUS "$LLVM_HOME found: $ENV{LLVM_HOME}")
8     set(LLVM_HOME $ENV{LLVM_HOME} CACHE PATH "Root of LLVM installation")
9     set(LLVM_LIB ${LLVM_HOME}/lib)
10    set(LLVM_DIR ${LLVM_LIB}/cmake/llvm)
11    find_package(LLVM REQUIRED CONFIG)
12    include_directories(${LLVM_INCLUDE_DIRS})
13    link_directories(${LLVM_LIBRARY_DIRS})
14    set(SOURCE_FILE RecursiveVisitor.cpp)
15    add_executable(recursivevisitor ${SOURCE_FILE})
16    set_target_properties(recursivevisitor PROPERTIES COMPILE_FLAGS
17        ↪ "-fno-rtti")
17    target_link_libraries(recursivevisitor
18        LLVMSupport
19        clangAST
20        clangBasic
21        clangFrontend
22        clangSerialization
23        clangTooling
24    )
25 endif()
```

---

Figure 3.13: CMakeLists.txt file for RecursiveVisitor test tool

The main function for our tool is similar to the ‘DeclVisitor’ one defined in fig. 3.4.

---

```
1 #include "clang/Tooling/CommonOptionsParser.h"
2 #include "clang/Tooling/Tooling.h"
3 #include "llvm/Support/CommandLine.h" // llvm::cl::extrahelp
4
5 #include "FrontendAction.hpp"
6
7 namespace {
8   llvm::cl::OptionCategory TestCategory("Test project");
9   llvm::cl::extrahelp
10     CommonHelp(clang::tooling::CommonOptionsParser::HelpMessage);
11 } // namespace
12
13 int main(int argc, const char **argv) {
14   llvm::Expected<clang::tooling::CommonOptionsParser> OptionsParser =
15     clang::tooling::CommonOptionsParser::create(argc, argv, TestCategory);
16   if (!OptionsParser) {
17     llvm::errs() << OptionsParser.takeError();
18     return 1;
19   }
20   clang::tooling::ClangTool Tool(OptionsParser->getCompilations(),
21                                 OptionsParser->getSourcePathList());
22   return Tool.run(clang::tooling::newFrontendActionFactory<
23     clangbook::recursivevisitor::FrontendAction>()
24     .get());
25 }
```

---

Figure 3.14: The main function for the RecursiveVisitor test tool.

As you can see, we changed only the namespace name for our custom frontend action at line 23.

The code for the frontend action and consumer is the same as in fig. 3.5 and fig. 3.6, with the only difference being the namespace change from `declvisitor` to `recursivevisitor`. The most interesting part of the program is the Visitor class implementation.

---

```

1 #include "clang/AST/RecursiveASTVisitor.h"
2
3 namespace clangbook {
4 namespace recursivevisitor {
5 class Visitor : public clang::RecursiveASTVisitor<Visitor> {
6 public:
7     bool VisitFunctionDecl(const clang::FunctionDecl *FD) {
8         llvm::outs() << "Function: '" << FD->getName() << "'\n";
9         return true;
10    }
11    bool VisitParmVarDecl(const clang::ParmVarDecl *PVD) {
12        llvm::outs() << "\tParameter: '" << PVD->getName() << "'\n";
13        return true;
14    }
15 };
16 } // namespace recursivevisitor
17 } // namespace clangbook

```

---

*Figure 3.15: Visitor class implementation*

There are several changes compared to the code for ‘DeclVisitor’ (see fig. 3.7). The first is that recursion isn’t implemented. We’ve only implemented the callbacks for nodes of interest to us. A reasonable question arises: how is the recursion controlled? The answer lies in another change: our callbacks now return a boolean result. The value `false` indicates that the recursion should stop, while `true` signals the visitor to continue the traversal.

The program can be compiled using the same sequence of commands as we used previously, see fig. 3.8.

We can run our program as follows:

```

$ ./recursivevisitor max.cpp
...
Function: 'max'
    Parameter: 'a'
    Parameter: 'b'

```

As we can see, it produces the same result as we obtained with the DeclVisitor implementation. The AST traversal techniques considered so far are not the only ways for AST traversal. Most of the tools that we will consider later will use a different approach based on AST matchers.

## 3.4 AST matchers

AST matchers provide another approach for locating specific AST nodes. They can be particularly useful in linters when searching for improper pattern usage or in refactoring tools when identifying AST nodes for modification.

We will create a simple program to test AST matchers. The program will identify a function definition with the name `max`.

We will use a slightly modified `CMakeLists.txt` file from the previous examples.



---

```
1 cmake_minimum_required(VERSION 3.16)
2 project("matchvisitor")
3
4 if ( NOT DEFINED ENV{LLVM_HOME})
5     message(FATAL_ERROR "$LLVM_HOME is not defined")
6 else()
7     message(STATUS "$LLVM_HOME found: $ENV{LLVM_HOME}")
8     set(LLVM_HOME $ENV{LLVM_HOME} CACHE PATH "Root of LLVM installation")
9     set(LLVM_LIB ${LLVM_HOME}/lib)
10    set(LLVM_DIR ${LLVM_LIB}/cmake/llvm)
11    find_package(LLVM REQUIRED CONFIG)
12    include_directories(${LLVM_INCLUDE_DIRS})
13    link_directories(${LLVM_LIBRARY_DIRS})
14    set(SOURCE_FILE MatchVisitor.cpp)
15    add_executable(matchvisitor ${SOURCE_FILE})
16    set_target_properties(matchvisitor PROPERTIES COMPILE_FLAGS "-fno-rtti")
17    target_link_libraries(matchvisitor
18        LLVMFrontendOpenMP
19        LLVMSupport
20        clangAST
21        clangASTMatchers
22        clangBasic
23        clangFrontend
24        clangSerialization
25        clangTooling
26    )
27 endif()
```

---

*Figure 3.16: CMakeLists.txt for AST matchers test tool*

There are two additional libraries were added: LLVMFrontendOpenMP and clangASTMatchers (see lines 18 and 21 in fig. 3.16).

The main function for our tool looks like

---

```

1  #include "clang/Tooling/CommonOptionsParser.h"
2  #include "clang/Tooling/Tooling.h"
3  #include "llvm/Support/CommandLine.h" // llvm::cl::extrahelp
4
5  #include "MatchCallback.hpp"
6
7  namespace {
8  llvm::cl::OptionCategory TestCategory("Test project");
9  llvm::cl::extrahelp
10     CommonHelp(clang::tooling::CommonOptionsParser::HelpMessage);
11 } // namespace
12
13 int main(int argc, const char **argv) {
14     llvm::Expected<clang::tooling::CommonOptionsParser> OptionsParser =
15         clang::tooling::CommonOptionsParser::create(argc, argv, TestCategory);
16     if (!OptionsParser) {
17         llvm::errs() << OptionsParser.takeError();
18         return 1;
19     }
20     clang::tooling::ClangTool Tool(OptionsParser->getCompilations(),
21                                   OptionsParser->getSourcePathList());
22     clangbook::matchvisitor::MatchCallback MC;
23     clang::ast_matchers::MatchFinder Finder;
24     Finder.addMatcher(clangbook::matchvisitor::M, &MC);
25
26     return Tool.run(clang::tooling::newFrontendActionFactory(&Finder).get());
27 }

```

---

Figure 3.17: The main function for AST matchers test tool

As you can observe (lines 22-24), we employ the `MatchFinder` class and define a custom callback (included via the header in line 5) that outlines the specific AST node we intend to match.

The callback is implemented as follows:

---

```

1 #include "clang/ASTMatchers/ASTMatchFinder.h"
2 #include "clang/ASTMatchers/ASTMatchers.h"
3
4 namespace clangbook {
5 namespace matchvisitor {
6 using namespace clang::ast_matchers;
7 static const char *MatchID = "match-id";
8 clang::ast_matchers::DeclarationMatcher M =
9     functionDecl(decl().bind(MatchID), matchesName("max"));
10
11 class MatchCallback : public clang::ast_matchers::MatchFinder::MatchCallback
12     ↪ {
13 public:
14     virtual void
15     run(const clang::ast_matchers::MatchFinder::MatchResult &Result) final {
16         if (const auto *FD =
17             ↪ Result.Nodes.getNodeAs<clang::FunctionDecl>(MatchID)) {
18             const auto &SM = *Result.SourceManager;
19             const auto &Loc = FD->getLocation();
20             llvm::outs() << "Found 'max' function at " << SM.getFilename(Loc) <<
21                 ↪ " "
22                 << SM.getSpellingLineNumber(Loc) << " "
23                 << SM.getSpellingColumnNumber(Loc) << "\n";
24         }
25     }
26 };
27 } // namespace matchvisitor
28 } // namespace clangbook

```

---

Figure 3.18: The match callback for AST matchers test tool

The most crucial section of the code is located at lines 7-9. Each matcher is identified by an ID, which, in our case, is 'match-id'. The matcher itself is defined in lines 8-9:

---

```
8 clang::ast_matchers::DeclarationMatcher M =  
9     functionDecl(decl().bind(MatchID), matchesName("max"));
```

---

This matcher seeks a function declaration using `functionDecl()` that has a specific name, as seen in `matchesName()`. We utilized a specialized Domain-Specific Language (DSL) to specify the matcher. The DSL is implemented using C++ macros. It's worth noting that the recursive AST visitor serves as the backbone for AST traversal inside the matchers implementation.

The program can be compiled using the same sequence of commands as we used previously, see fig. 3.8.

We will utilize a slightly modified version of the example shown in fig. 2.5, with an additional function added.

---

```
1 int max(int a, int b) {  
2     if (a > b) return a;  
3     return b;  
4 }  
5  
6 int min(int a, int b) {  
7     if (a > b) return b;  
8     return a;  
9 }
```

---

*Figure 3.19: Test program minmax.cpp for AST matchers*

When we run our test tool on the example, we will obtain the following output:

```
./matchvisitor minmax.cpp  
...  
Found 'max' function at minmax.cpp:1:5
```

As we can see, it has located only one function declaration with the name specified for the matcher.

The AST matchers reference page (<https://clang.llvm.org/docs/LibASTMatchersReference.html>) provides extensive information about various matchers and their usage.

The DSL for matchers is typically employed in custom Clang tools, such as clang-tidy (as discussed in [Chapter 5, clang-tidy linter framework](#)), but it can also be used as a standalone tool. A specialized program called `clang-query` enables the execution of different match queries.

## 3.5 Explore clang AST with clang-query

AST matchers are incredibly useful, and there's a utility that facilitates checking various matchers and analyzing the AST of your source code. This utility is known as `clang-query`. You can build and install this utility using the following command.

```
ninja install-clang-query
```

You can run the tool as follows

```
llvm-project/install/bin/clang-query minmax.cpp
```

We can use `match` command as follows

```
clang-query> match functionDecl(decl().bind("match-id"), matchesName("max"))
```

Match #1:

```
minmax.cpp:1:1: note: "match-id" binds here
```

```
int max(int a, int b) {
~~~~~
```

```
minmax.cpp:1:1: note: "root" binds here
```

```
int max(int a, int b) {
~~~~~
```

```
1 match.
```

```
clang-query>
```

The default output, referred to as `diag`, is available. Among several potential outputs, the most relevant one for us is `dump`. When the output is set to `dump`, `clang-query` will display the located AST node. For example, the following demonstrates how to match a function parameter named `a`:

```
clang-query> set output dump
```

```
clang-query> match parmVarDecl(hasName("a"))
```

Match #1:

```
Binding for "root":
```

```
ParmVarDecl 0x6775e48 <minmax.cpp:1:9, col:13> col:13 used a 'int'
```

Match #2:

```
Binding for "root":
```

```
ParmVarDecl 0x6776218 <minmax.cpp:6:9, col:13> col:13 used a 'int'
```

```
2 matches.
```

```
clang-query>
```

This tool proves useful when you wish to test a particular matcher or investigate a portion of the AST tree. We will utilize this tool to explore how Clang handles compilation errors.

## 3.6 Processing AST in the case of errors

One of the most interesting aspects of Clang pertains to error processing. Error processing encompasses error detection, the display of corresponding error messages, and potential error recovery. The latter is particularly intriguing in terms of the Clang AST. Error recovery occurs when Clang doesn't halt upon encountering a compilation error but continues to compile in order to detect additional issues.

Such behavior is beneficial for various reasons. The most evident one is user convenience. When programmers compile a program, they typically prefer to be informed about as many errors as possible in a single compilation run. If the compiler were to stop at the first error, the programmer would have to correct that error, recompile, then address the subsequent error, and recompile again, and so forth. This iterative process can be tedious and frustrating, especially with larger codebases or intricate errors.

Another compelling reason centers on IDE integration, which will be discussed in more detail in [Chapter 8, IDE support and code navigation](#). IDEs offer navigation support coupled with an integrated compiler. We will explore `clangd` as one of such tools. Editing code in IDEs often leads to compilation errors. Most errors are confined to specific sections of the code, and it might be suboptimal to cease navigation in such cases.

Clang employs various techniques for error recovery. For the syntax stage of parsing, it utilizes heuristics; for instance, if a user forgets to insert a semicolon, Clang may attempt to add it as part of the recovery process. The semantic stage of parsing is more intricate, and Clang implements unique techniques to manage recovery during this phase.

Consider a program (`maxerr.cpp`) that is syntactically correct but has a semantic error. For example, it might use an undeclared variable. In this program, refer to line 3 where the undeclared variable `ab` is used

---

```
1 int max(int a, int b) {  
2     if (a > b) {  
3         return ab;  
4     }  
5     return b;  
6 }
```

---

*Figure 3.20: The `maxerr.cpp` test program with a semantic error, undeclared variable*

We are interested in the AST result produced by Clang, and we will use `clang-query` to examine it, which can be run as follows:

```

llvm-project/install/bin/clang-query maxerr.cpp
...
maxerr.cpp:3:12: error: use of undeclared identifier 'ab'
    return ab;
           ^

```

From the output, we can see that clang-query displayed a compilation error detected by the compiler. It's worth noting that despite this, an AST was produced for the program, and we can examine it. We are particularly interested in the return statements and can use the corresponding matcher to highlight the relevant parts of the AST.

We will also set up the output to produce the AST and search for return statements that are of interest to us

```

clang-query> set output dump
clang-query> match returnStmt()

```

The resulting output identifies two return statements in our program: the first match on line 5 and the second match on line 3

Match #1:

```

Binding for "root":
ReturnStmt 0x6b63230 <maxerr.cpp:5:3, col:10>
~-ImplicitCastExpr 0x6b63218 <col:10> 'int' <LValueToRValue>
  ~-DeclRefExpr 0x6b631f8 <col:10> 'int' lvalue ParmVar 0x6b62ec8 'b' 'int'

```

Match #2:

```

Binding for "root":
ReturnStmt 0x6b631b0 <maxerr.cpp:3:5, col:12>
~-RecoveryExpr 0x6b63190 <col:12> '<dependent type>' contains-errors lvalue

```

2 matches.

*Figure 3.21: ReturnStmt node matches at maxerr.cpp test program*

As we can see, the first match corresponds to semantically correct code on line 5 and contains a reference to the parameter `a`. The second match is for line 3, which has a compilation error. Notably, Clang has inserted a special type of AST node: `RecoveryExpr`. It's worth noting that in certain situations, Clang might produce an incomplete AST. This can cause issues with Clang tools, such as lint checks. In instances of compilation errors, lint checks might yield unexpected results because Clang couldn't recover accurately from the compilation errors. We will revisit the problem when exploring the clang-tidy lint check framework in [Chapter 5, clang-tidy linter framework](#).

## 3.7 Summary

We explored the Clang AST, a major instrument for creating various Clang tools. We learned about the architectural design principles chosen for the implementation of Clang AST and investigated different methods for AST traversal. We delved into specialized traversal techniques, such as those for C/C++ declarations, and also looked into more universal techniques that employ recursive visitors and Clang AST matchers. Our exploration concluded with the `clang-query` tool and how it can be used for Clang AST exploration. Specifically, we used it to understand how Clang processes compilation errors.

The next chapter will discuss the basic libraries used in Clang and LLVM development. We will explore the LLVM code style and foundational Clang/LLVM classes, such as `SourceManager` and `SourceLocation`. We will also cover the TableGen library, which is used for code generation, and the LIT (LLVM Integration Test) framework.

## 3.8 Further reading

- How to write RecursiveASTVisitor: <https://clang.llvm.org/docs/RAVFrontendAction.html>
- AST Matcher Reference: <https://clang.llvm.org/docs/LibASTMatchersReference.html>



# Basic libraries and tools

LLVM has been written in C++ language and currently (since July 2022) uses the c++17 version of the C++ standard [Community, 2022b]. On the other side, it has a lot of internal implementations for fundamental containers with the primary goal of performance. Therefore, being familiar with the extensions is crucial if you want to work with LLVM and clang. In addition, LLVM also introduced additional development tools such as TableGen - DSL (domain-specific language) for structural data processing and LIT (LLVM test framework). You can find some info about the tools here.

We are planning to use a simple example project to demonstrate the tools. The example project will be Clang plugin that can estimate complexity of the source code being compiled and print a warning if the number of functions/methods exceeds the limit specified as a parameter.

## 4.1 LLVM coding style

TBD

## 4.2 Basic libs

TBD

## 4.3 TableGen

TBD

## 4.4 LLVM test framework

TBD

## 4.5 Summary

TBD

## 4.6 Further reading

- LLVM Coding Standards: <https://llvm.org/docs/CodingStandards.html> [Community, 2023b]

# Part 2:

# Clang tools

You can find some info about different clang tools [here](#). We will start with linters that are based on clang-tidy, continue with some advanced code analysis techniques (CFG and live time analysis). The next chapter will be about different refactoring tools such as clang-format. The last chapter will be about IDE support. We are going to investigate how VSCode can be extended with language server provided by LLVM (clangd).



# clang-tidy linter framework

There is an introduction to clang-tidy - the clang-based linter framework that uses AST to find anti-patterns in the C/C++/Objective-C code. First, we will start with a clang-tidy description, what kind of checks it has and how we can use them. Later we will investigate the clang-tidy architecture and how we can create our custom lint check.

## 5.1 clang-tidy overview and usage examples

TBD

## 5.2 clang-tidy internal design

TBD

## 5.3 Custom clang-tidy check

TBD

## 5.4 Results in the case of compilation errors

As we noticed in [Section 3.6, Processing AST in the case of errors](#), compilation errors might produce strange results for lint checks and there is an example.

# Advanced code analysis

he clang-tidy checks from the previous chapter can be considered based on an advanced matching provided by AST. However, it might not be enough when you want to detect some complex problems, such as lifetime issues. We will introduce advanced code analysis tools based on CFG (control flow graph). The clang static analyser is an excellent example of such devices, but we also have some CFG integration into clang-tidy. The chapter will start with typical usage examples and continue with implementation details. We will finish with our custom check that uses advanced techniques.

## 6.1 Usage cases

TBD

## 6.2 CFG and life time analysis

TBD

## 6.3 Custom CFG check

TBD





# Refactoring tools

The chapter is about code refactoring tools. For example, suppose you want to modify a set of files according to a new code style or want to detect and fix a specific problem in your project. The code modification tools can help here. We will start with a typical usage scenario and finish with our custom code modification tool.

## 7.1 Code modification and clang-format

TBD

## 7.2 Custom code modification tool

Idea: we are going to replace public vars in a class with custom get/set methods.

TBD



# IDE support and code navigation

The chapter is about Language Server Protocol (LSP) and how you can use it to extend your IDE. As the primary IDE, we will use VSCode. LLVM has its implementation of LSP as clangd. We will start with a typical usage scenario, continue with implementation details and finish with our custom extension for LSP that we must implement on both client (VSCode extension) and server (clangd) sides.

## 8.1 VSCode and LSP

TBD

## 8.2 clangd internals

TBD

## 8.3 Custom extension for LSP

TBD

# Part 3:

## Clang frontend features

You can find some info about some clang features [here](#). The features have the primary goal the compiler performance. Some of them (clang modules) are part of C++ standard but some (HeaderMaps) are clang specific feature.



# Features

You can find some info about different clang features in the chapter. Some of them are used with clang tools such as precompiled headers that allow fast editing at different IDEs. Others are specific clang features that are not well known but can provide some benefits, for instance, improving overall compilation performance. We will start with a typical usage scenario for each component and finish with its implementation details.

## 9.1 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.

### 9.1.1 User guide

Generate your pch file is simple [\[Community, 2022a\]](#). Suppose you have a header file with name **header.h**:

---

```
1 #pragma once
2
3 void foo() {
```

---

```
4 }
```

---

then you can generate a pch for it with

```
1 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

---

```
1 // test pchs
2
3 #include "header.h"
4
5 int main() {
6     foo();
7     return 0;
8 }
```

---

As you may see, the header is included as follows

---

```
1 ...
2 #include "header.h"
3 ...
```

---

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

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

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

```
1 $ lldb ~/local/llvm-project/build/bin/clang -- -cc1 -include-pch
  ↳ header.pch main.cpp -fsyntax-only
2 ...
```



```

3 | (lldb) b clang::ASTReader::ReadAST
4 | ...
5 | (lldb) r
6 | ...
7 | 4231 llvm::SaveAndRestore<SourceLocation>
8 | -> 4232 SetCurImportLocRAII(CurrentImportLoc, ImportLoc);
9 | 4233 llvm::SaveAndRestore<Optional<ModuleKind>>
   | ↪ SetCurModuleKindRAII(
10 | 4234 CurrentDeserializingModuleKind, Type);
11 | 4235
12 | (lldb) p FileName
13 | (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.

## 9.2 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<sup>1</sup>

### 9.2.1 User guide

The C++20 standard [for Standardization, 2020] introduced 2 concepts related to modules. The first one is ordinary modules described at section 10 of [for Standardization, 2020]. 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 [for Standardization, 2020]) will be call **explicit modules**. We will start with the implicit modules first.

---

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

## 9.2.2 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

```
1 module header1 {  
2     header "header1.h"  
3     export *  
4 }
```

The header paths in the modulemap file has to be either 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.

### 9.2.2.1 Explicit modules

TBD

### 9.2.2.2 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`:

---

```
1 #pragma once  
2  
3 int h1 = 1;
```

---

`header2.h`:

---

```
1 #pragma once
```

---

```
2
3 int h2 = 2;
```

---

The header1.h is included into the main.cpp

---

```
1 #include <iostream>
2
3 #include "header1.h"
4
5 int main() {
6     std::cout << "Header1 value: " << h1 << std::endl;
7     std::cout << "Header2 value: " << h2 << std::endl;
8 }
```

---

The code will not compile

```
1 clang++ -std=c++20 -fconstexpr-depth=1271242 main.cpp -o main
   ↪ -lstdc++
2 main.cpp:7:37: error: use of undeclared identifier 'h2'
3     std::cout << "Header2 value: " << h2 << std::endl;
4                                     ^
5 1 error generated.
```

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

---

```
1 module h1 {
2     header "header1.h"
3     export *
4     module h2 {
5         header "header2.h"
6         export *
7     }
8 }
```

---

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

### 9.2.3 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 caalls `clang::CompilerInst`

## 9.3 Header-Map files

TBD



[www.packt.com](http://www.packt.com)

Subscribe to our online digital library for full access to over 7,000 books and videos, as well as industry leading tools to help you plan your personal development and advance your career. For more information, please visit our website.

## Why subscribe?

- Spend less time learning and more time coding with practical eBooks and Videos from over 4,000 industry professionals
- Improve your learning with Skill Plans built especially for you
- Get a free eBook or video every month
- Fully searchable for easy access to vital information
- Copy and paste, print, and bookmark content

Did you know that Packt offers eBook versions of every book published, with PDF and ePub files available? You can upgrade to the eBook version at [packt.com](http://packt.com) and as a print book customer, you are entitled to a discount on the eBook copy. Get in touch with us at [customercare@packtpub.com](mailto:customercare@packtpub.com) for more details.

At [www.packt.com](http://www.packt.com), you can also read a collection of free technical articles, sign up for a range of free

## Other Books You Might Enjoy

If you enjoyed this book, you may be interested in these other books by Packt:



### Modern Computer Architecture and Organization

Jim Ledin

ISBN: 978-1-83898-439-7

- Get to grips with transistor technology and digital circuit principles
- Discover the functional elements of computer processors
- Understand pipelining and superscalar execution
- Work with floating-point data formats
- Understand the purpose and operation of the supervisor mode
- Implement a complete RISC-V processor in a low-cost FPGA
- Explore the techniques used in virtual machine implementation
- Write a quantum computing program and run it on a quantum computer



### Software Architecture with C# 9 and .NET 5 - Second Edition

Gabriel Baptista, Francesco Abbruzzese

ISBN: 978-1-80056-604-0

- Use different techniques to overcome real-world architectural challenges and solve design consideration issues
- Apply architectural approaches such as layered architecture, service-oriented architecture (SOA), and microservices

- Leverage tools such as containers, Docker, Kubernetes, and Blazor to manage microservices effectively
- Get up to speed with Azure tools and features for delivering global solutions
- Program and maintain Azure Functions using C# 9 and its latest features
- Understand when it is best to use test-driven development (TDD) as an approach for software development
- Write automated functional test cases
- Get the best of DevOps principles to enable CI/CD environments

## Packt is searching for authors like you

If you're interested in becoming an author for Packt, please visit [authors.packtpub.com](https://authors.packtpub.com) and apply today. We have worked with thousands of developers and tech professionals, just like you, to help them share their insight with the global tech community. You can make a general application, apply for a specific hot topic that we are recruiting an author for, or submit your own idea.

## Share Your Thoughts

Now you've finished **NameOfTheProduct**, we'd love to hear your thoughts! Scan the QR code below to go straight to the Amazon review page for this book and share your feedback.



<https://packt.link/r/<ISBN10P>>

Your review is important to us and the tech community and will help us make sure we're delivering excellent quality content.

Replace **NameOfTheProduct** with your title name.

Replace the dummy QR Code image with your product QR code image.

Replace *ISBN10P* with the 10-digit *ISBN-10P* from EPIC.

This page should be added after the Front Matter.





# Index

|       |  |  |             |                               |
|-------|--|--|-------------|-------------------------------|
| AST   | <b>A</b>                                       |  | header unit | <b>H</b>                      |
|       | 21, 33, 35, 37, 38, 52, 55, 103                |  |             |                               |
| Clang | <b>C</b>                                       |  | Lexer       | <b>L</b>                      |
|       | 9, 15, 17–19, 21–25, 31, 32, 38–43, 54, 55, 89 |  |             |                               |
| GDB   | <b>G</b>                                       |  | LIT         | 35, 37–39, 53                 |
|       |  |  | LLDB        | 20                            |
|       | 24, 25   |  |             | 15, 18, 20, 21, 24–26, 31, 57 |



# Bibliography

LLVM Community. Moving llvm projects to github, 2019. URL <https://llvm.org/docs/Proposals/GitHubMove.html>.

LLVM Community. Clang compiler user's manual, 2022a. URL <https://clang.llvm.org/docs/UsersManual.html>.

LLVM Community. [llvm] update c++ standard to 17, 2022b. URL <https://reviews.llvm.org/D130689>.

LLVM Community. Llvm releases, 2022c. URL <https://releases.llvm.org/>.

LLVM Community. Building llvm with cmake, 2023. URL <https://llvm.org/docs/CMake.html>.

"LLVM Community". "msvc compatibility", "2023". URL "<https://clang.llvm.org/docs/MSVCCompatibility.html>".

LLVM Community. Clang features, 2023a. URL <https://clang.llvm.org/features.html>.

LLVM Community. Llvm coding standards, 2023b. URL <https://llvm.org/docs/CodingStandards.html>.

LLVM Community. Commandline 2.0 library manual, 2023c. URL <https://llvm.org/docs/CommandLine.html>.

LLVM Community. Getting started with the llvm system, 2023d. URL <https://llvm.org/docs/GettingStarted.html>.

Keith Cooper and Linda Torczon. *Engineering A Compiler*. Elsevier Inc., 2nd edition, 2012. ISBN 978-0-12-088478-0.

Thomas H. Cormen, Charles E. Leiserson, Ronald L. Rivest, and Clifford Stein. *Introduction to Algorithms*. MIT press, 3rd edition, 2009.

International Organization for Standardization. International standard iso/iec 14882:2017(e) – programming languages – c++, 2017. URL <https://www.iso.org/standard/69466.html>.

International Organization for Standardization. International standard iso/iec14882:2020(e) – programming languages – c++, 2020. URL <https://www.iso.org/standard/73560.html>.

"Alexandre Ganea". [clang][driver] re-use the calling process instead of creating a new process for the cc1 invocation", "2019". URL "<https://reviews.llvm.org/D69825>".

International Organization for Standardization (ISO). ISO/IEC 9899:1999 - Programming languages - C, 1999. URL <https://www.iso.org/standard/23482.html>.

Chris Lattner and Vikram Adve. LLVM: A Compilation Framework for Lifelong Program Analysis & Transformation. *Proceedings of the 2004 International Symposium on Code Generation and Optimization (CGO'04)*, Mar 2004.

Wikipedia contributors. Google test — Wikipedia, the free encyclopedia, 2022a. URL [https://en.wikipedia.org/w/index.php?title=Google\\_Test&oldid=1123707063](https://en.wikipedia.org/w/index.php?title=Google_Test&oldid=1123707063). [Online; accessed 23-December-2022].

Wikipedia contributors. Opencv — Wikipedia, the free encyclopedia, 2022b. URL <https://en.wikipedia.org/w/index.php?title=OpenCV&oldid=1113479408>. [Online; accessed 23-December-2022].

Wikipedia contributors. Qt (software) — Wikipedia, the free encyclopedia, 2022c. URL [https://en.wikipedia.org/w/index.php?title=Qt\\_\(software\)&oldid=1125020271](https://en.wikipedia.org/w/index.php?title=Qt_(software)&oldid=1125020271). [Online; accessed 23-December-2022].