

Refactoring Open Source Project OCC

A Journey of Build Time Optimization

Athanasiadis Panagiotis

Department of Computer Science and Engineering, University of Ioannina

April 22, 2026

Open
Comp
Chem

Introduction to OCC

What is OCC?

A next-generation quantum chemistry and crystallography program designed for modern computational workflows, implemented in C++17 with Python bindings.

Quantum Chemistry Features:

- ▶ Hartree-Fock (RHF, UHF)
- ▶ Density Functional Theory
- ▶ Density fitting (RI-JK)
- ▶ Property calculations
- ▶ Implicit solvation (SMD)

Crystal Structure Analysis:

- ▶ CIF file processing
- ▶ Periodic bond detection
- ▶ Energy calculations
- ▶ Hirshfeld surfaces
- ▶ Dimer identification

OCC:Project Structure

Code Composition:

- ▶ **90.8%** C++ (core library)
- ▶ **10,000** lines of C++
- ▶ **3.7%** JavaScript (web tools)
- ▶ **2.8%** HTML (documentation)
- ▶ **1.6%** Python (bindings)
- ▶ **1.0%** CMake (build system)

Development Activity:

- ▶ 1,044 commits
- ▶ 45 releases
- ▶ 5 active contributors

Community Engagement

21 stars • 875/month downloads • Active development since 2020

Project Organization:

- ▶ `include/` - Header files
- ▶ `src/` - Implementation
- ▶ `tests/` - Test suite
- ▶ `docs/` - Documentation
- ▶ `examples/` - Usage examples
- ▶ `3rdparty/` - Dependencies

Dependencies: 11+ modern C++ libraries including Eigen3, libxc, libcint, gemmi

Project Goal: Improving Build Times

Current Challenge

With 10,000 lines of C++ code across multiple files, build times can significantly impact development productivity.

Build Time Impact:

- ▶ Slows down development cycles
- ▶ Increases testing feedback time
- ▶ Affects CI/CD pipeline efficiency
- ▶ Impacts developer experience

Contributing Factors:

- ▶ Large C++ codebase (90.8%)
- ▶ Template-heavy code (Eigen3)
- ▶ Multiple dependencies (11+)
- ▶ Complex file structure

Project Objective

Goal: Analyze and optimize the OCC build system to reduce compilation times while maintaining code quality and functionality.

This involves identifying compilation bottlenecks and studying modern C++ build optimization techniques.

C++ Compilation Process

Understanding the Build Pipeline

Compilation Stages:

1. Preprocessing

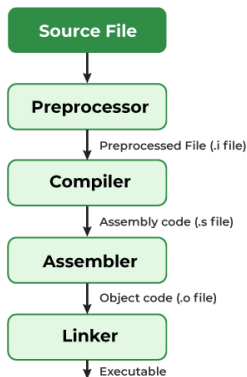
- ▶ Header file inclusion
- ▶ Macro expansion
- ▶ Conditional compilation

2. Compilation

- ▶ Parsing source code
- ▶ Semantic analysis
- ▶ Code generation (object files)

3. Linking

- ▶ Combine object files
- ▶ Resolve symbols
- ▶ Generate executable



Chem

Methodology: Analysis Tools & Techniques

Identifying Compilation Bottlenecks (Part 1)

1. Build Time Profiling

Ninja Log Analysis

- ▶ Parse `.ninja_log` files generated during build
- ▶ Extract per-file compilation times
- ▶ Identify slowest translation units
- ▶ Track build time trends across iterations

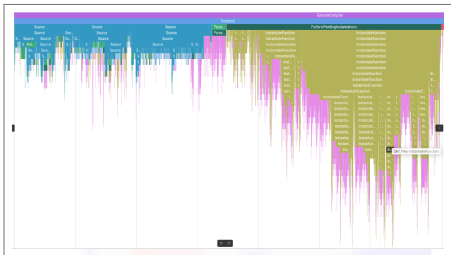
Clang Build Analyzer

- ▶ Generate time-trace profiles using `-ftime-trace`
- ▶ Detailed breakdown of compilation stages
- ▶ Template instantiation time analysis
- ▶ Header file inclusion impact measurement
- ▶ Function-level compilation time tracking

Methodology: Visualization & Code Review

Identifying Compilation Bottlenecks (Part 2)

Flame Graph Visualization



Hierarchical view showing which files and templates consume the most build time

Manual Code Review

▶ Header Structure

- ▶ Analyze include hierarchies
- ▶ Identify circular dependencies
- ▶ Check for unnecessary includes

▶ Template Usage

- ▶ Heavy template instantiations

Methodology: Static Analysis Tools

CppDepend

Issue	# Issues	Added	Fixed	Elements	Group
⚠️ Avoid methods too big, too complex	CppDepend.Core.CodeModelEng1.NumberOfMatch			methods	Project Rules \ Code Smells
⚠️ Avoid methods potentially poorly commented	CppDepend.Core.CodeModelEng1.NumberOfMatch			methods	Project Rules \ Code Smells
⚠️ From now, all methods added should respect basic quality principles	CppDepend.Core.CodeModelEng1.NumberOfMatch			methods	Project Rules \ Code Smells Regression
⚠️ Assignment to Variable without Use (Unused Variable)	CppDepend.Core.CodeModelEng1.NumberOfMatch			fields	Project Rules \ Coding Standards \ C# Coding Standard
⚠️ Don't assign static fields from instance methods	CppDepend.Core.CodeModelEng1.NumberOfMatch			fields	Project Rules \ Object Oriented Design
⚠️ Projects with poor cohesion (RelationalCohesion)	CppDepend.Core.CodeModelEng1.NumberOfMatch			projects	Project Rules \ Object Oriented Design
⚠️ Redundant code: Found a statement that begins with type constant.	CppDepend.Core.CodeModelEng1.NumberOfMatch			fields	Project Rules \ Misc
⚠️ Unassigned variable name can be negative as it is unnecessary to test it.	CppDepend.Core.CodeModelEng1.NumberOfMatch			fields	Project Rules \ Misc
⚠️ Potentially dead Methods	CppDepend.Core.CodeModelEng1.NumberOfMatch			methods	Project Rules \ Dead Code
⚠️ Potentially dead Fields	CppDepend.Core.CodeModelEng1.NumberOfMatch			fields	Project Rules \ Dead Code
⚠️ Instance fields should be prefixed with 'I,'.	CppDepend.Core.CodeModelEng1.NumberOfMatch			fields	Project Rules \ Naming Conventions
⚠️ Static fields should be prefixed with 'S,'.	CppDepend.Core.CodeModelEng1.NumberOfMatch			fields	Project Rules \ Naming Conventions
⚠️ Types name should begin with an upper character	CppDepend.Core.CodeModelEng1.NumberOfMatch			types	Project Rules \ Naming Conventions
⚠️ Avoid methods with name too long	CppDepend.Core.CodeModelEng1.NumberOfMatch			methods	Project Rules \ Naming Conventions

Showing 1 to 14 of 16 entries

CodeChecker

Checker name	Severity	Number of report
bugprone-forwarding-reference-overload	3	0
bugprone-incorrect-math	1	0
bugprone-integer-division	6	0
bugprone-implicit-widening-const	3	0
bugprone-integer-literal-rotate	4	0
bugprone-sizeof-expression	1	0
bugprone-use-after-move	1	0
core-dcl35-cpp	1	0
core-dcl39-cpp	3	0
clang-diagnostic-bitwise-instead-of-logical	1	0
clang-diagnostic-deprecated-cpp	18	0
clang-diagnostic-deprecated-cpp-with-user-provided-cpp	12	0
clang-diagnostic-error	4	0
clang-diagnostic-ignored-qualifiers	14	0
clang-diagnostic-implicit-cast	1	0
clang-diagnostic-raising-field-initializer	23	0
clang-diagnostic-range-loop-construct	1	0
clang-diagnostic-reorder-cast	6	0
clang-diagnostic-sign-compare	488	0
clang-diagnostic-unused-but-set-variable	2	0
clang-diagnostic-unused-function	2	0
clang-diagnostic-unused-lambda-capture	4	0
clang-diagnostic-unused-local-symbol	1	0
clang-diagnostic-unused-parameter	88	0
clang-diagnostic-unused-private-field	6	0
clang-diagnostic-unused-variable	188	0
core-CallAndThrow	2	0
core-DeadStores	1	0
core-NullDereference	2	0
core-unsafeArithmetic	7	0

Methodology: Quantifying Build Performance

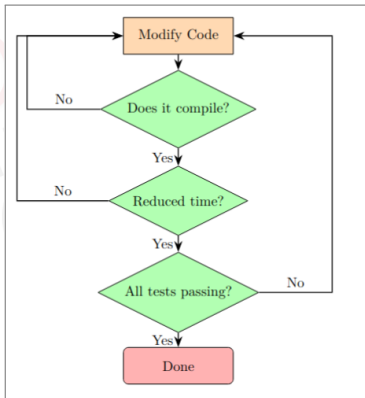
Objective: Establish a reliable baseline to measure the impact of refactoring on occ compilation speed.

Measurement Tool: Linux time

- ▶ Utilized the GNU time utility.
- ▶ **Key Metrics Tracked:**
 - ▶ **Real (Wall) Time:** The actual elapsed time perceived by the developer.

** Final results averaged over 10 clean builds to mitigate caching variance.*

Benchmarking Protocol:



Headers and Includes Audit

The Problem:

- ▶ Unnecessary includes in headers and source files
- ▶ STL headers are thousands of lines long
- ▶ Compiler copies every included file recursively
- ▶ No automatic duplicate detection across compilation units

Methodology:

- ▶ Bottom-up approach: source files first, then headers
- ▶ Used dependency graphs of core classes as guide
- ▶ Manual removal with compilation verification

Results: Removed approximately **330** unused includes from **167** files

Ccache and LLD Integration

Ccache (Compiler Cache):

- ▶ Speeds up recompilation by caching previous builds
- ▶ Hashes source code and compiler options
- ▶ Returns cached object files when matches are found
- ▶ Critical for rapid testing workflow
- ▶ Avoids long recompilation times during development

LLD (LLVM Linker):

- ▶ Drop-in replacement for system linkers
- ▶ Significantly faster than traditional linkers

Implementation: Configured in CMakeLists.txt to automatically use both tools when available on user's system

Forward Declarations

What are Forward Declarations?

- ▶ Technique to break unnecessary header dependencies
- ▶ Informs compiler of type's existence without internal details
- ▶ Works for functions, classes, and other types

Key Benefit: Reducing Recompilation Cascade

- ▶ Modified headers only trigger recompilation of files that explicitly include them
- ▶ Decouples source files from implementation details
- ▶ Small changes in low-level headers don't trigger project-wide rebuilds

Impact: Critical for rapid testing strategy and incremental development

Compiler Differences

Key Insight:

- ▶ Different compilers produce different machine code for same source
- ▶ Build times can vary significantly between compilers

Benchmark Results:

- ▶ GCC vs Clang comparison conducted
- ▶ **GCC was faster by 1 minute average**

Compilers Not Tested:

- ▶ Platform-specific: MinGW-w64 (Windows)
- ▶ Hardware-specific: Intel ICC
- ▶ Proprietary: Microsoft MSVC
- ▶ Outdated compilers

Suggestion: Use GCC for optimal build performance

Unity Builds

What are Unity Builds?

- ▶ Compilation technique merging multiple source files into single compilation units
- ▶ Reduces redundant header parsing overhead

Benefits:

- ▶ Compiler sees larger codebase portion → better optimizations
- ▶ Fewer object files → reduced linker workload
- ▶ Significant build time reduction

Implementation:

- ▶ Configured through CMake
- ▶ Batch size of 8 files provided best performance

Challenges:

- ▶ Fixed multiple One Definition Rule (ODR) violations
- ▶ Some files required exclusion from unity builds

Results: Substantial Build Time Reduction

Test System Configuration:

- ▶ Intel Core i7-3740QM processor
- ▶ 16GB DDR3 RAM

Full Build Time Results:

Before: 22 minutes (average)

↓ **70.5% reduction** ↓

After: 6.5 minutes (average)

More Information:

- ▶ Pull Request #45 (merged) on official GitHub
- ▶ Detailed code changes and implementation notes

Pull Request #45: Code Changes

Improving compile time #45

Edit <> Code

Merged **peterspackman** merged 22 commits into **peterspackman:main** from **KonstantinosVard:main** 3 days ago

Conversation 3 Commits 22 Checks 0 Files changed 162

+207 -417



PanagiotisAthanasidis commented on Dec 15, 2025

Contributor

Audited the includes of source/header files for unused headers, implemented checking for ccache and lld in the Cmakefile, implemented Unity builds and created some forward declarations.



PanagiotisAthanasidis and others added 20 commits 2 months ago

- testing commit
- Test commit
- test
- Merge branch 'main' of https://github.com/KonstantinosVard/occ_lambda_B
- test on branch
- added ccache and lld checks
- moved PES::solve_linear_system to .cpp to compile once
- Removed unnecessary includes
- Removed unnecessary libraries from integral_engine_df
- small changes to cint interface.h
- Removed unnecessary includes from smd_salvation.cpp, occ_cube.cpp, mo.cpp
- Removed unnecessary libraries
- Removed unnecessary libraries mostly linear_algebra double includes
- Removed unnecessary libraries
- Removed unnecessary libraries from source files
- Removed unnecessary libraries from source files (final)

cc40ba4

fd813f1

f74b4c3

61a05e0

1710bee

095d1af

357a0a4

7c6e569

48b8525

fed5cc0

@bfdecb

a6bfd19

6b9fb28

bba5b2b

84338ea

e6c1810

Reviewers

No reviews

Assignees

No one assigned

Labels

None yet

Projects

None yet

Milestone

No milestone

Development

Successfully merging this pull request may close these issues.

None yet

Notifications

Customize

Unsubscribe

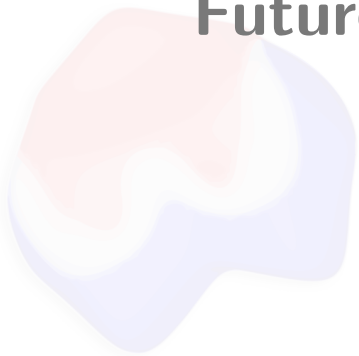
You're receiving notifications because you authored the thread.

3 participants



Allow edits by maintainers

Future Work



Open

Comp

Chem

Future Work: C++20 Modules

C++20 Modules:

- ▶ Modern replacement for traditional header files
- ▶ Promising for reducing build times

Current Limitations:

- ▶ Not ready for mainstream adoption
- ▶ Most third-party libraries still use header files
- ▶ Standard library module support incomplete
- ▶ Significant refactoring effort required

Status: Monitoring ecosystem maturity before adoption

Open
Comp
Chem

Future Work: PIMPL Pattern

Pointer to Implementation (PIMPL) Idiom:

- ▶ Separates class interface from implementation
- ▶ Hides private members in separate struct/class
- ▶ Accessed via opaque pointer

Benefits:

- ▶ Creates "compilation firewall"
- ▶ Implementation changes don't trigger client code recompilation
- ▶ Stable Binary Interface (ABI)
- ▶ Improved build times in large projects

Open
Comp
Chem

Future Work: Linear Algebra Library Replacement

Current Challenge: Eigen Library

- ▶ Header-only library → slow build times
- ▶ Extensive template instantiation overhead

Potential Solutions:

- ▶ Evaluate compiled linear algebra alternatives
- ▶ Consider libraries with binary distributions
- ▶ Balance performance vs. compilation speed

Goal: Find alternative without sacrificing numerical performance

Open
=
Comp
Chem

Questions?

Thank you for your attention!

Contact: pathanasiadis@cs.uoi.gr

Project Repository:

<https://github.com/peterspackman/occ>