

Techniques avancées de génération de code pour le parallélisme

Advanced techniques for parallel code generation

Thèse de doctorat de l'université Paris-Saclay

École doctorale n° d'accréditation, dénomination et sigle

Spécialité de doctorat: voir annexe

Graduate School : voir annexe. Référent : voir annexe

Thèse préparée dans la (ou les) unité(s) de recherche **STIC** (voir annexe), sous la direction
de **Joel FALCOU**, titre du directeur ou de la directrice de thèse

Thèse soutenue à Paris-Saclay, le JJ mois AAAA, par

Jules PÉNUCHOT

Composition du jury

Membres du jury avec voix délibérative

Prénom NOM

Titre, Affiliation

Prénom NOM

Titre, Affiliation

Prénom NOM

Titre, Affiliation

Prénom NOM

Titre, Affiliation

Prénom NOM

Titre, Affiliation

Président ou Présidente

Rapporteur & Examineur / trice

Rapporteur & Examineur / trice

Examineur ou Examinatrice

Examineur ou Examinatrice

Titre: Techniques avancées de génération de code pour le parallélisme

Mots clés: Métaprogrammation, compilation, C++

Résumé: Mettre le sum ici

Title: Advanced techniques for parallel code generation

Keywords: Metaprogramming, compilation, C++

Abstract: Mettre l'abstract ici

Contents

I	Current state of metaprogramming for high performance computing	5
1	Introduction	7
2	State of the art of Metaprogramming	9
2.1	Metaprogramming styles and languages	9
2.1.1	A short history of metaprogramming	9
2.1.2	Metaprogramming in contemporary languages	9
2.2	Metarpogramming and High Performance Computing (HPC) libraries	10
2.2.1	Core and application-specific libraries	10
2.2.2	High performance computing libraries	11
2.3	Conclusion	15
2.4	Language constructs of C++ metaprogramming	15
2.4.1	Input and output of metaprograms	15
2.4.2	Compile time logic	18
3	Code generation at low level	19
3.1	Introduction	19
3.2	High-performance generative programming	20
3.2.1	Metaprogramming as code generation principles	20
3.2.2	Application to High Performance Computing code generation	23
3.3	The matrix-vector multiplication routine: gemv	27
3.4	Performance results of the generated gemv codes	29
3.4.1	On X86 Intel processor	29
3.4.2	On ARM processor	30
3.5	Conclusion	31
4	Code generation at expression level	33
4.1	Expression level	33
4.1.1	Expression templates: a type-based representation for math formulas	33
4.1.2	Optimized code generation	33
II	C++ metaprogramming beyond templates	35
5	Establishing a methodology for compile time benchmarking	37
5.1	Compile time benchmarking: state of the art	37
5.1.1	Metabench	38

5.1.2	Templight	38
5.1.3	Clang's built-in profiler	38
5.1.4	Conclusion	39
5.2	ctbench design and features	40
5.2.1	CMake API for benchmark and graph target declarations	40
5.2.2	ctbench internals	41
5.2.3	Software quality	42
5.2.4	Conclusion	43
5.3	A ctbench use case	43
5.4	Related projects	46
5.5	Acknowledgements	46
6	Constexpr parsing for high performance computing	49
6.1	Introduction	49
6.2	A technical background of C++ DSELs	49
6.2.1	Constexpr programming	49
6.2.2	C++ Domain Specific Embedded Languages	52
6.3	Code generation from constexpr allocated structures	53
6.3.1	Code generation from pointer tree data structures	53
6.3.2	Using algorithms with serialized outputs	62
6.3.3	Preliminary observations	65
6.4	Brainfuck parsing and code generation	66
6.4.1	Constexpr Brainfuck parser and AST	66
6.4.2	A variety of techniques to generate code from a dynamic AST	67
6.5	Math parsing and high performance code generation	80
6.5.1	The Shunting-Yard algorithm	80
6.5.2	Code generation from a postfix math notation	80
6.5.3	Using Blaze for high performance code generation	80
6.5.4	Conclusion: a complete toolchain for High Performance Computing code generation from math formulas	80
6.6	Conclusion	80
7	Appendix	83
.1	Poacher	83
.1.1	Brainfuck AST definition header	83
.1.2	Brainfuck parser implementation header	87

Part I

Current state of metaprogramming for high performance computing

1 - Introduction

This thesis is about metaprogramming techniques for parallel code generation. It aims to study the boundary between compile-time parsing of Domain Specific Embedded Languages (DSELs) and high performance code generation in C++.

The main motivation is to provide tools, libraries, and guidelines for embedding mathematical languages in C++ with the hope that it can be useful to build a cohesive set of tools to make generic HPC more accessible to non-expert audiences. This goal is even more important as new computing architectures emerge over time. Developing high performance programs requires tuned implementations that can be achieved by either specializing implementations for the target platforms, or using libraries that provide specialized abstractions at various levels and for various domains.

2 - State of the art of Metaprogramming

In this state of the art I will first give an overview of metaprogramming in historic and contemporary languages. Then I will focus on the state of the art of C++ metaprogramming, and notable high performance computing libraries as they are essential for the scope of my thesis.

2.1 . Metaprogramming styles and languages

Metaprogramming is not an old concept. It exists in many forms: both the C preprocessor and C++ templates can be used to take code as input, and generate code as well. However it is clear that these two mechanisms are completely different. The C preprocessor can only read and write tokens with very basic language, while C++ templates can manipulate types and complex C++ values with a much richer semantic.

Our goal in this section will be to cover different languages with their own metaprogramming models, as well as HPC-focused libraries that use metaprogramming to achieve notable levels of genericity.

2.1.1 .A short history of metaprogramming

2.1.2 .Metaprogramming in contemporary languages

Metaprogramming perpetuates itself in contemporary languages, with some being more widespread than others.

MetaOCaml MetaOCaml[[metaocaml](#)] implements quoting and splicing *i.e.* the ability to essentially copy and paste expressions, as well as staged compilation to evaluate statements at compile-time. This enables code generation to occur both at runtime and at compile-time.

DLang more or less extends the C++ Metaprogramming model. It leverages templates and compile time function evaluation just like its predecessor.

Compile-time evaluation is much more permissive and mixins enable to generate code in a more direct way than C++. Dlang mixins allow injecting code in functions and structures in two ways: using template mixins which are pre-parsed constructs that can be injected later, as well as string mixins that allow strings containing Dlang code to be compiled and inserted directly into programs.

Rust proposes metaprogramming through macros, generics, and traits.

Rust macros are more powerful than those proposed in C and C++. They have a

Braid [braid] is language that implements metaprogramming through multi-staged programming (like MetaOCaml) for heterogeneous real-time 3D graphics. It is however unmaintained.

Julia [julia] is a dynamic language that uses the LLVM

Terra Terra[terra] implements a very explicit metaprogramming model. The language is based on LUA, and exploits the dynamic nature of the language together with LLVM Just-In-Time (JIT) compilation to allow code generation to happen at runtime. It implements multi-staged compilation and splicing just like MetaOCaml.

Additionally, Terra can be embedded in other languages through its C API. Overall it is a very versatile and experimental take on metaprogramming, but the lack of interoperability with C++ templates makes it hard to justify its use for HPC applications.

As we will see in section ??, Graphical Processing Unit (GPU) computing libraries rely heavily on C++ metaprogramming to provide building blocks for portable high performance compute kernels.

2.2 . Metarpogramming and HPC libraries

2.2.1 . Core and application-specific libraries

As previously said C++ templates can be seen as a functional language. Over time a range of libraries emerged, aiming to provide functionalities similar to regular language such as containers and algorithms for use in template metaprograms. Notable examples of such libraries are MPL[mpl], Brigand[brigand], and mp11[mp11].

Libraries for more specific uses were also implemented, such as Spirit[spirit] for writing parsers (not for compile time parsing), Compile Time Regular Expression (CTRE) [ctre] for compiling regular expressions, and Compile Time Parser Generator (CTPG) [ctpg] for generating LR1 parsers (also not for compile time parsing).

The benefits of metaprogrammed libraries are:

- Performance: notably in the case of CTRE. Regular expressions are usually interpreted at runtime, which adds a measurable overhead to text processing. CTRE shows leading performance, on par with Rust's regex library which also works by compiling regular expressions.

- Language integration: since these are C++ libraries, their APIs can take advantage of C++ operator overloading and lambdas. In CTPG, these are used to provide a domain-specific language that is close to what parser generators like YACC or Bison provide, though it is still regular C++ code which can be put inside any function body. Using a C++ API makes these libraries easier to learn as the syntax is already familiar to their users.
- Streamlined toolchain: as they only require to be included as headers. This avoids complicating compilation toolchains by requiring additional programs to be installed and integrated to the build system.

2.2.2 .High performance computing libraries

There are many C++ HPC libraries that use metaprogramming to achieve high levels of performance while being generic and proposing high level APIs. This subsection aims to provide an overview of notable libraries that fit this description as well as recent efforts that are still under development.

- **Eigen** [eigen] is the first major C++ library to implement Expression templates for the generation of high performance math computing. Expression templates is a C++ design pattern that consists in representing math expressions with type template trees. We will discuss them later in 4.1.

```
#include <Eigen/Dense>
#include <iostream>

using Eigen::MatrixXd;
using Eigen::VectorXd;

int main() {
    MatrixXd m = MatrixXd::Random(3, 3);
    m = (m + MatrixXd::Constant(3, 3, 1.2)) * 50;
    std::cout << "m =" << std::endl << m << std::endl;
    VectorXd v(3);
    v << 1, 2, 3;
    std::cout << "m * v =" << std::endl
              << m * v << std::endl;
}
```

- **Blaze** [blazelib] is a successor of Eigen that implements so-called "Smart Expression Templates" which extends upon the concept of expression templates implemented by Eigen. It aims to provide a more performant

and extensible HPC library. However, Eigen is not set in stone and its designed has since been updated.

```
#include <blaze/Math.h>
#include <iostream>

using blaze::DynamicVector;
using blaze::StaticVector;

int main() {
    StaticVector<int, 3UL> a{4, -2, 5};
    DynamicVector<int> b(3UL);

    b[0] = 2;
    b[1] = 5;
    b[2] = -3;

    DynamicVector<int> c = a + b;

    std::cout << "c =\n" << c << "\n";
}
```

- **NT2 [nt2]** is a research project that aims to provide a complete numerical toolbox that leverages metaprogramming to develop portable HPC applications with a Matlab-like interface while still achieving state-of-the-art computing performance.

```
#include <nt2/include/functions/ones.hpp>
#include <nt2/table.hpp>

using namespace nt2;

int main() {
    table<double> x;
    table<double> y = ones(4, 4);

    x = 40.0 * y + 2.0;

    NT2_DISPLAY(x);

    return 0;
}
```

- **EVE [eve]** provides generic abstractions over SIMD instructions as well as SIMD-optimized generic algorithms for the development of high performance and portable SIMD code [**hpcs2018-matvec**].

```

#include <eve/module/core.hpp>
#include <eve/wide.hpp>
#include <iostream>

int main() {
    eve::wide<float> x(
        [] (auto i, auto) { return 1.f + i; });
    std::cout << "x      = " << x << "\n";
    std::cout << "2*x    = " << x + x << "\n";
    std::cout << "x^0.5 = " << eve::sqrt(x) << "\n";

    return 0;
}

```

- **HPX [hpx]** is a C++ parallel and distributed runtime library. It can execute small parallel tasks efficiently and distribute larger distributed tasks with a work following data execution model. Its parallel and distributed APIs as well as its parallel implementation of the standard library (based on its own parallel runtime) use metaprogramming for algorithmic genericity.

```

std::uint64_t fibonacci(std::uint64_t n) {
    if (n < 2)
        return n;

    hpx::future<std::uint64_t> n1 =
        hpx::async(fibonacci, n - 1);
    hpx::future<std::uint64_t> n2 =
        hpx::async(fibonacci, n - 2);

    // wait for the futures to return their values
    return n1.get() + n2.get();
}

```

- **Thrust [thrust]** implements GPU-accelerated equivalents of the Standard Library's algorithms, while CUB [**cub**] provides GPU-optimized algorithm skeletons for generic programming on NVIDIA GPUs. AMD and Intel implement their equivalents for their own platforms, respectively ROCm and OneAPI.

```

#include <algorithm>
#include <thrust/device_vector.h>
#include <thrust/functional.h>
#include <thrust/generate.h>

```

```

#include <thrust/host_vector.h>
#include <thrust/reduce.h>

int main(void) {
    // generate random data serially
    thrust::host_vector<int> h_vec(100);
    std::generate(h_vec.begin(), h_vec.end(), rand);

    // transfer to device and compute sum
    thrust::device_vector<int> d_vec = h_vec;
    int x = thrust::reduce(d_vec.begin(), d_vec.end(),
                           0, thrust::plus<int>());

    return 0;
}

```

- **faer** [faer]

```

use faer::{mat, Mat, prelude::*};

// empty 0x0 matrix
let m0: Mat<f64> = Mat::new();

// zeroed 4x3 matrix
let m1: Mat<f64> = Mat::zeros(4, 3);

// 3x3 identity matrix
let m2 = Mat::from_fn(
    3, 3,
    |i, j| if i == j { 1.0 } else { 0.0 });

// 4x2 matrix with custom data
let m3 = mat![
    [4.93, 2.41],
    [5.43, 4.33],
    [9.83, 1.59],
    [7.13, 5.02_f64],
];

// compute the qr decomposition of a matrix
let qr_decomposition = m3.qr();

```

2.3 . Conclusion

All these examples show that metaprogramming is not just a gimmick. There is enough interest in leveraging metaprogramming for generic HPC programming that many frameworks, libraries, or languages have been developed around this paradigm and maintained sometimes for more than a decade. Its adoption in languages like Braid, Julia, Terra, Rust, or Dlang shows that it is not limited to C++.

However, most of generic HPC libraries and frameworks that are developed and maintained today are based on C++, and C++ itself is likely to evolve to provide more ways to generate programs as shown by the recent standard proposal for Scalable Reflection in C++ [**scalable-reflection**].

- Metaprogramming isn't a new idea
- Some languages provide advanced metaprogramming capabilities
- C++ has solid metaprogramming constructs, and a complete HPC ecosystem (libraries, compilers, etc.)

2.4 . Language constructs of C++ metaprogramming

In this section I will cover the constructs of C++ metaprogramming. Metaprograms are defined as programs that take code as input and/or output, therefore I will categorize metaprogramming constructs as: those used for input and output, and those that enable compile time computing.

This distinction will be important to explain the work presented in this thesis, as it focuses a lot on methods to translate compile time compute results into C++ programs.

2.4.1 . Input and output of metaprograms

Template metaprogramming

C++ template metaprogramming [**abrahams:2004**] is a technique based on the use of the template type system of C++ to perform arbitrary computation at compile time. This property of C++ templates is due to the fact they define a Turing-complete sub-language manipulating types and integral constants at compile time [**unruh:1994**]. Due to the fact that template code generation is performed at compile time, uses constants and supports pattern-matching and recursion thanks to template partial specialization, it can also be looked at as a pure functional language [**haeri:2012**].

Templates are an interesting technique for generative programming. As they are Turing-complete, one can design a set of template metaprograms acting as a Domain Specific Language (DSL) compiler run at compile time and generating temporary C++ code fragment as an output. The resulting temporary source code is then merged with the rest of the source code and finally

processed by the classic compilation process. Through this technique, compile time constants, data structures and complete functions can be manipulated. The execution of metaprograms by the compiler enables the library to implement domain-specific optimizations that lead to a complete domain oriented code generation. Such a technique can be hosted by several languages featuring metaprogramming features (incidental or by design) like D [template:dlang], Haskell [sheard:2002] and OCaml [serot:2008].

Listing 2.1: Example of compile time computation using C++ templates

```
// Template type accepting an integer
// as a non-type template parameter
template <unsigned N> struct fibonacci_t;

// General definition
template <unsigned N> struct fibonacci_t {
    static constexpr unsigned value =
        fibonacci_t<N - 2>::value + fibonacci_t<N - 1>::
        value;
};

// Specialization for cases 0 and 1
template <> struct fibonacci_t<0> {
    static constexpr unsigned value = 0;
};

template <> struct fibonacci_t<1> {
    static constexpr unsigned value = 1;
};

std::array<int, fibonacci_t<5>::value> some_array;
```

Listing 2.1 shows basic principles of C++ template metaprogramming. The `fibonacci_t` type template accepts a Non-Type Template Parameter (NTTP) called N , and exposes the N^{th} element of the Fibonacci series as its value static member. The template has 3 definitions: a generic one to calculate elements for $N > 1$, and two specializations for elements of ranks 0 and 1.

Different kinds of templates C++ templates offer ways to output code for both data structures and functions.

- Type templates

Listing 2.2: Type template example

```
/// Generic anotated type
template <typename T> struct named_value_t {
```



```

    std::string name;
    T value;
};

```

Listing 2.3: Type alias template example

```

/// Alias template for nested named value
template <typename T>
using nested_named_value_t =
    named_value_t<named_value_t<T>>;

```

Templates can be used with type aliases too, as shown in listing 2.3.

- Function templates

Listing 2.4: Function template example

```

/// Returns a value annotated with
/// its string representation
template <typename T>
named_value_t<T> make_named_value(T value) {
    return named_value_t<T>(std::to_string(value),
                             std::move(value));
}

```

- Variable templates

Listing 2.5: Variable template example

```

/// Returns the default value of T annotated with
/// its string representation
template <typename T>
named_value_t<T> annotated_default_v =
    make_named_value(T{});

named_value_t<int> val = annotated_default_v<int>;

```

Different kinds of parameters

- Type parameters
- NTTP
- Parameter packs

Advanced C++ template mechanisms

- Substitution Failure Is Not An Error (SFINAE)
- Template specialization
- Parameter packs

Overloading

2.4.2 . Compile time logic

Template Metaprogramming and constant evaluation

- Computation using types
- Computation using values

constexpr functions

- `constexpr` functions
- `constexpr` memory allocation

The use of `constexpr` functions for compile time programming is preferable for many reasons: familiar to most C++ developers (and therefore more maintainable), it allows the use of types to enforce semantics properly as opposed to type-based metaprogramming, and it scales much better than Template Metaprogramming (TMP).

3 - Code generation at low level

BLAS-level functions are the cornerstone of a large subset of applications. If a large body of work surrounding efficient and large-scale implementation of some routines such as `gemv` exists, the interest for small-sized, highly-optimized versions of those routines emerged. In this section, we propose to show how a modern C++ approach based on generative programming techniques such as vectorization and loop unrolling in the framework of metaprogramming can deliver efficient automatically generated codes for such routines, that are competitive with existing, hand-tuned library kernels with a very low programming effort compared to writing assembly code. In particular, we analyze the performance of automatically generated small-sized `gemv` kernels for both Intel x86 and ARM processors. We show through a performance comparison with the OpenBLAS `gemv` kernel on small matrices of sizes ranging from 4 to 32 that our C++ kernels are very efficient and may have a performance that is up to 3 times better than that of OpenBLAS `gemv`.

3.1 . Introduction

The efforts of optimizing the performance of BLAS routines fall into two main directions. The first direction is about writing very specific assembly code. This is the case for almost all the vendor libraries including Intel MKL[hpcs1], AMD ACML[hpcs2] etc. To provide the users with efficient BLAS routines, the vendors usually implement their own routines for their own hardware using assembly code with specific optimizations which is a low level solution that gives the developers full control over both the instruction scheduling and the register allocation. This makes these routines highly architecture dependent and needing considerable efforts to maintain the performance portability on the new architecture generations. Moreover, the developed source codes are generally complex. The second direction is based on using modern generative programming techniques which have the advantage of being independent from the architecture specifications and as a consequence easy to maintain since it is the same source code which is used to automatically generate a specific code for a specific target architecture. With respect to the second direction, some solutions have been proposed in recent years. However, they only solve partially the trade-off between the abstraction level and the efficiency of the generated codes. This is for example the case of the approach followed by the Formal Linear Algebra Methods Environment (FLAME) with the Libflame library[hpcs3]. Thus, it offers a framework to develop dense linear solvers using algorithmic skeletons[hpcs4] and an API which is more user-

friendly than LAPACK, giving satisfactory performance results. A more generic approach is the one followed in recent years by C++ libraries built around expression templates[hpcs5] or other generative programming[hpcs6] principles. In this section, we will focus on such an approach. To show the interest of this approach, we consider as example the matrix-vector multiplication kernel (gemv) which is crucial for the performance of both linear solvers and eigen and singular value problems. Achieving performance running a matrix-vector multiplication kernel on small problems is challenging as we can see through the current state-of-the-art implementation results. Moreover, the different CPU architectures bring further challenges for its design and optimization.

In this section, we describe how we obtained optimized generated C++ codes for the gemv routine to make it reach its peak performance on different target CPU architectures (Intel x86 and ARM in this section). Our gemv C++ generated kernels achieve uniform performances that outperform in the most of the cases the peaks of the state-of-the-art OpenBLAS gemv kernel for small matrices of sizes ranging from 4 to 128. This chapter is organized as follows: in Section II, we describe various generative programming techniques including metaprogramming principles as well as the main strategies to get efficient small-scale BLAS functions. In Section III, we show how to apply these programming techniques and principles in order to develop an efficient small-scale gemv kernel that will be used to automatically generates efficient C++ codes for different target architectures. Then, we present in Section IV performance comparisons on two CPU architectures (Intel x86 and ARM) between our generated gemv kernels for small matrices and the OpenBLAS gemv. Finally, concluding remarks and perspectives are given in Section V.

3.2 . High-performance generative programming

The quality and performance of BLAS like code require the ability to write tight and highly-optimized code. If the raw assembly of low-level C has been the language of choice for decades, our position is that the proper use of the zero abstraction property of C++ can lead to the design of a single, generic yet efficient code base for many BLAS like functions. To do so, we will rely on two main elements: a proper C++ Single Instruction Multiple Data (SIMD) layer and a set of code generation techniques stemmed from Generic Programming.

3.2.1 . Metaprogramming as code generation principles

Metaprogramming[hpcs7] is about the design and the implementation of programs whose input and output are themselves programs. This term encompasses a large body of idioms, some language dependent, which can be used to define, manipulate or introspect arbitrary code fragment. One typical usage of metaprogramming is the design of libraries which expose

a user API with an arbitrary high abstraction level while being able to get compiled to a very “close to the metal” implementation, often rivaling with a handwritten expert code on a given machine[hpcs8]. If metaprogramming is used in languages as different as C++ [hpcs9], D[hpcs10], OCaml[hpcs11] or Haskell[hpcs12], a subset of basic notions emerges:

- **Code fragment generation:** Any metaprogrammable language has a way to build an object that represents a piece of code. The granularity of this fragment of code may vary –ranging from statement to a complete class definition–but the end results is the same: to provide an entry level entity to the metaprogramming system. In some languages, such as MetaOCaml for example, a special syntax can be provided to construct such fragment. In some others, code fragment are represented as a string containing the code itself.
- **Code processing:** Code fragments are meant to be combined, introspected or replicated in order to let the developer rearrange these fragments and as a consequence to provide a given service. Those processing steps can be done either via a special syntax construct, like the MetaOCaml code aggregation operator, or can use a similar syntax than a regular code.
- **Code expansion:** Once the initial code fragments have been processed, the last step is to turn them into an actual code. This is often done in an explicit manner by using a function or a syntax construct provided by the metaprogramming layer to trigger the code generation. Note that this code generation can either lead to a code ready to be compiled – like in Haskell or C++ – or a code that can be run – like in OCaml– if the generation phase is done at runtime. Metaprogramming also includes other code generation techniques such as domain-specific languages and compilation infrastructures based on source-to-source compilers, which are actually able to perform the same techniques proposed by this section. Such systems includes:
- **SYCL[hpcs13]:** a single-source abstraction over OpenCL for heterogeneous systems. By using OpenCL, one can effectively write SYCL code for a large selection of architectures including SIMD capable CPU.
- **BLIS[hpcs14]:** a framework for generating BLAS like operations in ISO C99 from a small subset of kernels that can be retargeted for different back-end. Its performances are on par with open-source solutions like OpenBLAS and ATLAS.
- **LGEN[hpcs15]:** a compiler that produces performance-optimized basic linear algebra computations on matrices of fixed sizes with or without

structure composed from matrix multiplication, matrix addition, matrix transposition, and scalar multiplication. Based on polyhedral analysis using CLoog, the generated code outperforms MKL and Eigen.

In this section, we will focus on language-based metaprogramming techniques so that the proposed method can be used in various compilers and OS settings as long as the compiler follows a given standard. Classical design of metaprograms in C++ usually relies on complex template types that forced the compiler to follow intricate path during type deduction in order to take advantage of the Turing completeness of the template definition. By using template partial specialization and recursive definition, one could implement arbitrary transform on types in order to converge towards a code ready to be generated. Code fragments were usually static class member function which encapsulated the basic code block to be replicated and generated. If the efficiency of the code generated was as expected, the maintenance cost of the generating code was usually high. Template metaprograms were complex to write and read as the logic of the code generation was buried behind heaps of non-trivial syntax. Some progress was made by some infrastructure library like MPL[hpcs16] or Fusion, but the learning gap was still high.

With the standard C++ revision in 2014 and 2017, this strategy was renewed with three new C++ features:

- Polymorphic, variadic anonymous functions: C++ 11 introduced the notion of local, anonymous functions (also known as lambda functions) in the language. Their primary goal was to simplify the use of standard algorithms by providing a compact syntax to define a function in a local scope, hence raising code locality. C++ 14 added the support for polymorphic lambdas, *i.e.* anonymous functions behaving like function templates by accepting arguments of arbitrary types, and variadic lambdas, *i.e.* anonymous functions accepting a list of arguments of arbitrary size. Listing 3.1 demonstrates this particular feature.

Listing 3.1: Sample polymorphic lambda definition

```
// Variadic function object building array
auto array_from =
    [](auto... values) {
        // sizeof... retrieves the
        // number of arguments
        return std::array<double,
                               sizeof...(values)>{
            values...};
    }
// Build an array of 4 double
auto data = array_from(1, 2, 3., 4.5f);
```

- Fold expressions: C++ 11 introduced the `...` operator which was able to enumerate a variadic list of functions or template arguments in a comma-separated code fragment. Its main use was to provide the syntactic support required to write a code with variadic template arguments. However, Niebler and Parent showed that this can be used to generate far more complex code when paired with other language constructs. Both code replication and a crude form of code unrolling were possible. However, it required the use of some counter-intuitive structure. C++ 17 extends this notation to work with an arbitrary binary operator. Listing 3.2 illustrates an example for this feature.

Listing 3.2: C++ 17 fold expressions

```
template<typename... Args>
auto reduce(Args&&... args) {
    // Automatically unroll the args into a sum
    return (args + ...);
}
```

Tuples Introduced by C++ 11, tuple usage in C++ was simplified by providing new tuple related functions in C++ 17 that make tuple a fully programmable struct-like entity. The transition between tuple and structure is then handled via the new structured binding syntax that allow the compile-time deconstruction of structures and tuples in a set of disjoint variables, thus making interface dealing with tuples easier to use. Listing 3.3 gives an example about tuples.

Listing 3.3: Tuple and structured bindings

```
// Build a tuple from values
auto data = std::make_tuple(3.f, 5, "test");

// Direct access to tuple data
std::get<0>(data) = 6.5f;

// Structured binding access
auto&[a,b,c] = data;

// Add 3 to the second tuple's element
b += 3;
```

3.2.2 . Application to High Performance Computing code generation

The main strategies to get efficient small-scale BLAS functions are on one hand the usage of the specific instructions set (mainly SIMD instructions set)

of the target architecture that is vectorization and on the other hand the controlled unrolling of the inner loop to ensure proper register and pipeline usage.

Vectorization Vectorization can be achieved either using the specific instructions set of each vendor or by relying on auto-vectorization. In our case, to ensure homogeneous performances across the different target architectures, we relied on the Boost.SIMD[hpcs17] package to generate SIMD code for all our architectures. Boost.SIMD relies on C++ metaprogramming to act as a zero-cost abstraction over SIMD operations in a large number of contexts. The SIMD code is then as easily written as a scalar version of the code and deliver 95% to 99% of the peak performances for the L1 cache hot data. The main advantage of the Boost.SIMD package lies in the fact that both scalar and SIMD code can be expressed with the same subset of functions. The vector nature of the operations will be triggered by the use of a dedicated type – pack – representing the best hardware register type for a given type on a given architecture that leads to optimized code generation.

Listing 3.4 demonstrates how a naive implementation of a vectorized dot product can simply be derived from using Boost.SIMD types and range adapters, polymorphic lambdas and standard algorithm.

Listing 3.4: Sample Boost.SIMD code

```
template <typename T>
auto simd_dot(T *in1, T *in2, std::size_t count) {
    // Adapt [in,in+count[ as a vectorizable range
    auto r1 = simd::segmented_range(in1, in1 + count);
    auto r2 = simd::segmented_range(in2, in2 + count);

    // Extract sub-ranges
    auto h1 = r1.head, h2 = r2.head;
    auto t1 = r1.tail, t2 = r2.tail;

    // sum and product polymorphic functions
    auto sum = [](auto a, auto b) { return a + b; };
    auto prod = [](auto r, auto v) { return r * v; };

    // Process vectorizable & scalar sub-ranges
    auto vr = std::transform_reduce(
        h1.begin(), h1.end(), h2.begin(), prod, sum,
        simd::pack<T>{});

    auto sr = std::transform_reduce(
        t1.begin(), t1.end(), t2.begin(), prod, sum,
```



```

        T{}});

    // Compute final dot product
    return sr + simd::sum(vr);
}

```

Note how the Boost.SIMD abstraction shields the end user to have to handle any architecture specific idioms and how it integrates with standard algorithms, hence simplifying the design of more complex algorithms. Another point is that, by relying on higher-level library instead of SIMD pragma, Boost.SIMD guarantees the quality of the vectorization across compilers and compiler versions. It also leads to a cleaner and easier to maintain codes, relying only on standard C++ constructs.

Loop unrolling The notion of unrolling requires a proper abstraction. Loop unrolling requires three elements: the code fragment to repeat, the code replication process and the iteration space declaration. Their mapping into C++ code is as follows:

- The code fragment in itself, which represents the original loop body, is stored inside a polymorphic lambda function. This lambda function will takes a polymorphic argument which will represent the current value of the iteration variable. This value is passed as an instance of `std::integral_constant` which allows to turn an arbitrary compile-time constant integer into a type. By doing so, we are able to propagate the constness of the iteration variable as far as possible inside the code fragment of the loop body.
- The unrolling process itself relies on the fold expression mechanism. By using the sequencing operator, also known as operator comma, the compiler can unroll arbitrary expressions separated by the comma operator. The comma operator will take care of the order of evaluation and behave as an unrollable statement.
- The iteration space need to be specified as an entity supporting expansion via `...` and containing the actual value of the iteration space. Standard C++ provides the `std::integral_sequence<N...>` class that acts as a variadic container of integral constant. It can be generated via one helper meta-function such as `std::make_integral_sequence<T,N>` and passed directly to a variadic function template. All these elements can then be combined into a `for_constexpr` function detailed in listing 3.5.

The function proceed to compute the proper integral constant sequence from the Start, End and D compile-time integral constant. As `std::integral_sequence`

<N...> enumerates values from 0 to N, we need to pass the start index and iteration 1 pragma are compiler-dependent and can be ignored increment as separate constants. The actual index is then computed at the unrolling site. To prevent unwanted copies and ensure inlining, all elements are passed to the function as a rvalue-reference or a universal reference.

Listing 3.5: Compile-time unroller

```
template <int Start, int D, typename Body,
         int... Step>
void for_constexpr(
    Body body, std::integer_sequence<int, Step...>,
    std::integral_constant<int, Start>,
    std::integral_constant<int, D>) {
    (body(std::integral_constant<int,
                                   Start + D * Step>{}),
     ...);
}

template <int Start, int End, int D = 1,
         typename Body>
void for_constexpr(Body body) {
    constexpr auto size = End - Start;
    for_constexpr(
        std::move(body),
        std::make_integer_sequence<int, size>{},
        std::integral_constant<int, Start>{},
        std::integral_constant<int, D>{});
}
```

A sample usage of the `for_constexpr` function is given in Listing 3.6 in a function printing every element from a `std::tuple`.

Listing 3.6: Tuple member walkthrough via compile-time unrolling

```
template<typename Tuple>
void print_tuple(Tuple const& t) {
    constexpr auto size = std::tuple_size<Tuple>::value;
    for_constexpr<0, size>([&](auto i) {
        std::cout << std::get<i>(t) << "\n";
    });
}
```

Note that this implementation exposes some interesting properties:

- As `for_constexpr` calls are simple function call, they can be nested in arbitrary manners.

- Relying on `std::integral_constant` to carry the iteration index gives access to its automatic conversion to integer. This means the iteration index can be used in both compile-time and runtime contexts.
- Code generation quality will still be optimized by the compiler, thus letting all other non-unrolling related optimizations to be applied.

One can argue about the advantage of such a method compared to relying on the compiler unrolling algorithm or using non-standard unrolling pragma. In both cases, our method ensure that the unrolling is always done at the fullest extend and does not rely on non-standard extensions.

3.3 . The matrix-vector multiplication routine: `gemv`

Level 2 BLAS routines such as `gemv` have a low computational intensity compared to Level 3 BLAS operations such as `gemm`. For that reason, in many dense linear algebra algorithms in particular for one sided factorizations such as Cholesky, LU, and QR decompositions some techniques are used to accumulate several Level 2 BLAS operations when possible in order to perform them as one Level 3 BLAS operation[hpcs18]. However, for the two-sided factorizations, and despite the use of similar techniques, the fraction of the Level 2 BLAS floating point operations is still important. For instance, for both the bidiagonal and tridiagonal reductions, this fraction is around 50%[hpcs19]. Thus, having optimized implementations for these routines on different architectures remains important to improve the performance of several algorithms and applications. Moreover, small-scale BLAS kernels are useful for some batched computations[hpcs20].

Here, we consider the matrix-vector multiplication routine for general dense matrices, `gemv`, which performs either $y := \alpha Ax + \beta y$ or $y := \alpha ATx + \beta y$, where A is an $m \times n$ matrix, α and β are scalars, and y and x are vectors. In this section, we focus on matrices of small sizes ranging from 4 to 512 as this range of sizes encompasses the size of most L1 cache memory, thus allowing a maximal throughput for SIMD computation units. The algorithm we present in Listing 3.7 is optimized for a column-major matrix. For space consideration, we will only focus on the core processing of the computation, *i.e.* the SIMD part, as the computation of the scalar tail on the last columns and rows can be trivially inferred from there.

Our optimized code relies on two types representing statically-sized matrix and vector, namely `mat<T,H,W>` and `vec<T,N>`. Those types carry their height and width as template parameters so that all size related values can be derived from them. The code shown in Listing 7 is made up of three main steps as detailed in Figure 1: Broadcast of each element of the vector in dif-

ferent registersSIMDSIMD Scalar Fig. 1. An example of matrix vector multiplication showing the SIMD/scalar computation boundaries. The matrix is 9×9 of simple precision floats so we can put 4 elements per SIMD register.

1. The computation of SIMD/scalar boundaries based on the static size of the matrix and the size of the current SIMD registers. Those computations are done in constexpr contexts to ensure their usability in the upcoming unrolling steps.
2. A first level of unrolling that takes care of iterating over all set of columns that are able to fit into SIMD registers. This unrolling is done so that both the corresponding columns of the matrix and the elements of the vector can respectively be loaded and broadcasted into SIMD registers.

Listing 3.7: Unrolled gemv kernel

```
template <typename T, std::size_t H,
          std::size_t W>
void gemv(mat<T, H, W> &mat, vec<T, W> &vec,
          vec<T, W> &r) {
    using pack_t = bs::pack<T>;
    constexpr auto sz = pack_t::static_size;

    // Separating simd/scalar parts
    constexpr auto c_simd = W - W % sz;
    constexpr auto r_simd = H - H % sz;

    for_constexpr<0, c_simd, sz> ( [](auto j) {
        pack_t pvec(&vec[j]);
        pack_t mulp_arr[sz];

        // Broadcasting vectors once and for all
        for_constexpr<0, sz> ( [&](auto idx) {
            mulp_arr[idx] = simd::broadcast<idx>(pvec);
        });

        // Walk through SIMD rows
        for_constexpr<0, r_simd> ( [&](auto I) {
            pack_t resp(&res[i + (I * sz)]);

            // Walk through columns
            for_constexpr<0, sz> ( [&](auto J) {
                pack_t matp(&mat(i + (I * sz), j + J));
                resp += matp * mulp_arr[J];
                simd::store(resp, &r[i + (I * sz)]);
            });
        });
    });
}
```

```

    }

    // Scalar code follows ...
}

```

3. A second level of unrolling that pass through all the available SIMD registers loadable from a given column. We rely on an overloaded operator `()` on the matrix to compute the proper position to load from. As usual with Boost.SIMD, the actual computation is run with scalar-like syntax using regular operators.

It is important to notice how close the actual unrolled code is to an equivalent code that would use regular for loops. This verisimilitude shows that modern metaprogramming matured enough so that the frontier between regular runtime programming and compile-time computation and code generation is very thin. The effort to fix bugs in such code or to upgrade it to new algorithms is comparable to the effort required by a regular code. The notion of code fragment detailed in Section II helps us to encapsulate those complex metaprogramming cases into an API based on function calls.

3.4 . Performance results of the generated gemv codes

To validate our approach, we consider two main targeted processor architectures: an x86 Intel processor i5-7200 and an ARM processor AMD A1100 with Cortex A57. We compare the performance of the generated gemv codes to that of the gemv kernel of the OpenBLAS library based on GotoBLAS2 1.13[hpcs21]. We use GCC 7.2[hpcs22] with maximal level of optimization.

In the following experiments, we only show results for simple precision floats with column major data, but we obtained similar results for the double precision case, as well as the row major data. All the results displayed below are obtained using one thread. All those results has been obtained using Google Benchmark micro-benchmark facilities. Every experiments have been repeated for a duration of 3s, using the median time as an approximation of the most frequent measured time.

3.4.1 . On X86 Intel processor

Fig. 2. GEMV performance on Intel i5-7200 processor using SIMD Extensions set (SSE-4.2)

In Figure 2, we compare the performance of our implementation using the SIMD Extensions set SSE 4.2 and a similarly configured OpenBLAS gemv kernel. The obtained results show that the performances of our automatically generated code is up to 2 times better for matrices of sizes ranging from 4×4 elements to 16×16 elements. However, for matrices of size 32×32 elements

and 64×64 elements, the OpenBLAS gemv kernel gives a better performances, especially for the 64×64 case. This is because the OpenBLAS library uses a dedicated gemv kernel with specific optimizations and prefetching strategies that our generic solution can not emulate. Beyond this size (64×64 elements), the $L1 \rightarrow L2$ cache misses cause a performance drop for both our generated code and the OpenBLAS gemv kernel. Nevertheless, our generated code sustains a better throughput for matrices of sizes 128×128 elements. For matrices of size 256×256 elements, the register usage starts to cause spill to the stack, showing that our solution can not be arbitrarily extended further to larger matrix sizes.

Fig. 3. GEMV performance on Intel i5-7200 processor using Advanced Vector Extensions (AVX)

In Figure 3, we compare the performance of our generated gemv code using Advanced Vector Extensions AVX2 to the performance of a similarly configured OpenBLAS gemv kernel. Again, the performances of our implementation are close to that of OpenBLAS and are even quite better for matrices of small sizes ranging from 4 to 16 elements. For example, for a matrix of size 8 elements, the automatically generated code has a performance that is 3 times better than the OpenBLAS gemv kernel (15.78 Gflop/s vs 5.06 Gflop/s). Two phenomenons appear however. The first one is that the increased number of the AVX registers compared to the SSE ones makes the effect of register spilling less prevalent. The second one is that the code generated for the special 64×64 elements case [hpcs23] in OpenBLAS has a little advantage compared to our automatically generated code. Finally, we note the fact that, for matrices of size above 512×512 elements, we stop being competitive due to the large amount of registers our fully unrolled approach would require.

In both cases, the register pressure is clearly identified as a limitation. One possible way to fix this issue will be to rely on partial loop unrolling and using compile-time informations about architecture to decide the level of unrolling to apply for a given size on a given architecture.

3.4.2 . On ARM processor

The comparison between our automatically generated code and the ARM OpenBLAS gemv kernel is given in Figure 4. Contrary to the x86 Intel processor, we sustain a comparable yet slightly better throughput than the OpenBLAS gemv kernel. The analysis of the generated assembly code shows that our method of guiding the compiler and letting it do fine grained optimizations generates a better code than the hand-written assembly approach of the OpenBLAS library.

Fig. 4. GEMV performance on ARM Cortex A57 processor

We exhibit performance drops similar to OpenBLAS due to $L1 \rightarrow L2$ misses. Register spilling also happens once we reach 512×512 elements.

The combination of our template based unrolling and Boost.SIMD shows that it is indeed possible to generate ARM NEON code from high-level C++ with zero abstraction cost.

3.5 . Conclusion

This chapter presented the details of generating an optimized level 2 BLAS routine `gemv`. As a key difference with respect to highly tuned OpenBLAS routine, our generated code is designed to give the best performance with a minimum programming effort for rather small matrices that fit into the L1 cache. Compared to the best open source BLAS library, OpenBLAS, the automatically generated `gemv` codes show competitive speedups for most of the matrix sizes tested in this section, that is for sizes ranging from 4 to 512. Therefore through this section and the example of the level 2 BLAS routine `gemv`, we showed that it is possible to employ modern generative programming techniques for the implementation of dense linear algebra routines that are highly optimized for small matrix sizes.

One of our next steps is to provide such metaprogrammed code generation process to tackle larger matrix sizes by precomputing an optimal tiling availability of prefetch and architecture, turning large-scale `gemv` into an optimized sequence of statically sized small-scale `gemvs`. Such a technique is also naturally applicable to more complex algorithms, such the matrix-matrix multiplication, `gemm`, where tiling is paramount, or LAPACK-level functions where the compile-time optimization may lead to an easier to parallelize solvers or decompositions.

Another objective will be to adapt such techniques to cooperate with tuning framework, hence providing the required level of performance with less input from the user.

4 - Code generation at expression level

4.1 . Expression level

In this chapter I will cover the use of TMP for higher levels of abstraction. Templates can be used to represent whole mathematical expressions at compile time by creating type-based arborescences. This type of representation is called an Expression Template (ET) [**veldhuizen:1995**].

Combined with compile-time mechanisms such as function overloading, specialization, and operator overloading, ETs can be used to implement expression level DSELs and convert complex mathematical expressions into high performance code.

There are two main libraries that are able to do just this: Eigen [**eigen**] and Blaze [**blazelib**], which were cited in 2.2. In this section, I will introduce the basics of C++ TMP,

4.1.1 . Expression templates: a type-based representation for math formulas

Expression templates are template trees that represent math formulas. They are generated using operator and function overloading from expressions that must be known at compile time.

```
template<typename Left, typename Right>
struct add_t{};
```

4.1.2 . Optimized code generation

Enables a whole range of optimizations:

Eliminating temporaries with lazy evaluation

Leveraging BLAS libraries

Parallelism

SIMD

Multithreading

Graphical Processing Unit support efforts

Lots of rewrite to do

Potentially a job for source rewriting tools

- Eventually: GPU code generation, although Blaze needs a significant rewrite for that. Source rewriting tools might be a good fit for that job.

Expression templates can provide expression level APIs for HPC libraries. Still two limitations:

- Slow compilation times
- C++ syntax only

Expression templates are aging (pretty well but still). Newer C++ standards provide metaprogramming features that can fundamentally change the way we write metaprograms.

The next part of my thesis will focus on how to leverage these features to implement DSEL of arbitrary syntax, and the study of their impact on compilation times.

Part II

C++ metaprogramming beyond templates

5 - Establishing a methodology for compile time benchmarking

With TMP libraries like Eigen[**eigen**], Blaze[**blazelib**], or CTRE[**ctre**] becoming larger and more common, we're seeing increasing computing needs at compile time. These needs might grow even larger as C++ embeds more features over time to support and extend this kind of practices, like compile time containers[**more-constexpr-containers**] or static reflection[**static-reflection**]. However, there is still no clear cut methodology to compare the performance impact of different metaprogramming strategies. But as new C++ features allow for new techniques with claimed better compile time performance, no proper methodology is provided to back up those claims.

In this chapter I introduce `ctbench`, which is a set of tools for compile time benchmarking and analysis in C++. It aims to provide developer-friendly tools to declare and run benchmarks, then aggregate, filter out, and plot the data to analyze it. As such, `ctbench` is meant to become the first layer of a proper scientific methodology for analyzing compile time program behavior.

We'll first have a look at current tools for compile time profiling and benchmarking and establish the limits of current tooling, then I'll explain what `ctbench` brings to overcome these limits.

`ctbench` was first presented at the CPPP 2021 conference[**ctbench-cppp21**] which is the main C++ technical conference in France. It is being used to benchmark examples from the `poacher`[**poacher**] project, which was briefly presented at the Meeting C++ 2022[**meetingcpp22**] technical conference.

5.1 . Compile time benchmarking: state of the art

C++ TMP raised interest for allowing computing libraries to offer great performance with a very high level of abstraction. Instead of building representations of calculations at runtime for interpretation, they are built at compile time to be transformed directly into programs.

As metaprogramming became easier with C++ 11 and C++ 17, it became more frequent. Consequently developers now have to bear with longer compilation times, often without being able to explain them. Therefore being able to measure compilation times becomes increasingly important and being able to profile and explain them as well.

This need turned into a variety of projects that aim to bring novel techniques to analyze the compile time performance of C++ metaprograms and metaprogramming techniques beyond black box A/B comparisons.

5.1.1 . Metabench

Metabench[**metabench**] instantiates variably sized benchmarks using embedded Ruby (ERB) templating and plots compiler execution time, allowing scaling analyses of metaprograms. Its output is a series of web-based interactive graphs.

The ERB templates are used to generate C++ programs for measuring their total compilation time. This approach is compiler-agnostic, it allows to compare not just metaprograms but compilers as well.

However Metabench is currently unmaintained, and its design choices make it difficult to modify and reuse it for the assessment of metaprogramming techniques performance in a scientific context:

- non-CMake scripts are embedded in CMake strings, making it difficult to modify and debug them,
- benchmarks are not standalone, they are embedded in the Metabench project itself,
- web-based visualizations cannot be embedded in \LaTeX documents.

5.1.2 . Templight

Templight[**templight**] is the first effort to instrument Clang with a profiler. It aims to provide tools to measure resource usage (CPU time, memory usage, and file I/O) of template instantiations.

5.1.3 . Clang's built-in profiler

Additionally, Clang has a built-in profiler[**time-trace**] that provides in-depth time measurements of various compilation steps, which can be enabled by passing the `-ftime-trace` flag. Its output contains data that can be directly linked to symbols in the source code, making it easier to study the impact of specific symbols on various stages of compilation. The output format is a JSON file meant to be compatible with Chrome's flame graph visualizer, that contains a series of time events with optional metadata like the mangled C++ symbol or the file related to an event. The profiling data can then be visualized using tools such as Google's Perfetto UI as shown in figure 5.1.

The JSON files have a rather simple structure. They contain a `traceEvents` field that holds an array of JSON objects, each one of them containing a `name`, a `dur` (duration), and `ts` (timestamp) field. It may also contain additional data in the field located at `/args/data`, as seen in listing 5.1. Duration and timestamps are expressed in microseconds.

Listing 5.1: Time trace event generated by Clang's internal profiler

```
{  
  "pid": 8696 ,
```

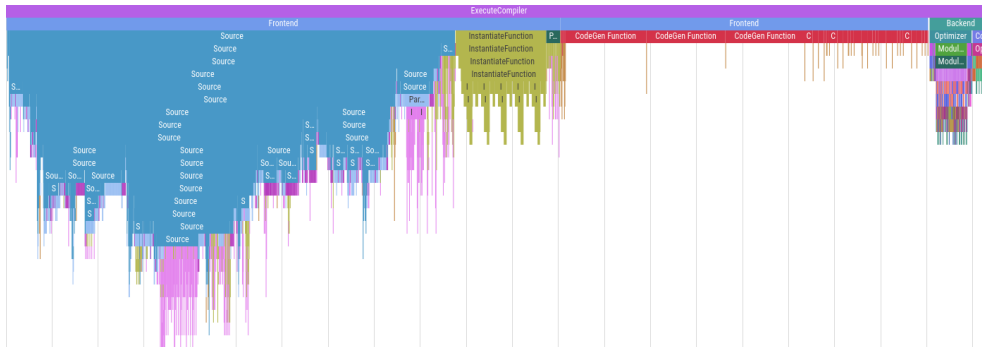


Figure 5.1: Clang time trace file in Perfetto UI

```
"tid": 8696,
"ph": "X",
"ts": 3238,
"dur": 7,
"name": "Source",
"args": {
  "detail": "/usr/include/bits/wordsize.h"
}
}
```

In this example the `/args/detail` field indicates the file that's being processed during this trace event. This field may also contain details related to symbols being processed for events like `InstantiateFunction`.

Clang's profiler data is very exhaustive and insightful, however there is no tooling to make sense of it in the context of variable size compile time benchmarks. `ctbench` tries to bridge the gap by providing a tool to analyze this valuable data. It also improves upon existing tools by providing a solution that's easy to integrate into existing CMake projects, and generates graphs in various formats that are trivially embeddable in documents like research papers, web pages, or documentations. Additionally, relying on persistent configuration, benchmark declaration and description files provides strong guarantees for benchmark reproducibility, as opposed to web tools or interactive profilers.

5.1.4 . Conclusion

While Metabench could have been a good candidate for assessing the scalability of metaprogramming techniques, its design choices make it simpler to write a new tool that better fits the needs of scientific assessment of metaprogramming techniques.

`ctbench`'s API and overall concept was largely inspired by this project for its simplicity and leverages Clang's profiling capabilities that are similar to those introduced by Templight.

5.2 . ctbench design and features

ctbench implements a new methodology for the analysis of compilation times: it allows users to define C++ sizable benchmarks to analyze the scaling performance of C++ metaprogramming techniques, and compare techniques against each other.

The project was designed to be a well maintained, easy to install, and easy to use alternative to current compile time benchmarking tools. It is developed with decent software quality in mind to offer a robust external API, and a well documented internal API to enable the reuse of its components.

All these goals are kept in mind to ensure that ctbench is not just a single use tool for this thesis, but a sustainable open-source project that enables researchers and developers to automatize the generation of reusable graphs and data artifacts.

5.2.1 . CMake API for benchmark and graph target declarations

The choice of using CMake as an entry point is to ensure that ctbench is not just a research tool made for a single set of benchmarks, but an easy-to-use solution for all C++ developers willing to analyze the compile time performance of their metaprograms.

The public CMake API of ctbench is meant to be as stable as reasonable. It is relatively simple, and all of it is implemented in a single CMake script file. It declares the following:

- Benchmark declaration functions, for the declaration and instantiation of user-defined sizable benchmark cases:
 - `ctbench_add_benchmark` takes a user-defined benchmark name, a C++ source file, as well as benchmark range parameters: iteration begin, end, and step parameters, as well as the number of samples per benchmark iteration. Benchmark iteration targets declared using this function are compiled with their sized passed as the `BENCHMARK_SIZE` define.
 - `ctbench_add_benchmark_for_range`, similar to `ctbench_add_benchmark` except for the fact that benchmark range parameters are taken as a single list. This interface only exists to provide a more compact function call.
 - `ctbench_add_benchmark_for_size_list`, similar to the previous ones, provides a way to define benchmarks for a given size list instead of a range.
 - And finally, `ctbench_add_custom_benchmark` inherits `ctbench_add_benchmark`'s interface with an addition: a callback function name must be

passed as a parameter. It will be called for each benchmark iteration target definition with the name and size of the target. This allows users to set compiler parameters other than `ctbench`'s pre-processor directive.

- `ctbench_add_graph` allows the declaration of a benchmark. It takes a user-defined name, a path to a JSON configuration file, an output destination, and a list of benchmark names defined using the functions mentioned above. `ctbench` expects JSON files to include information relative to the graphs themselves: the plotter being used, predicates to filter out time trace events, output format, and so on. Each plotter has its own set of parameters.
- `ctbench_set_compiler` and `ctbench_unset_compiler` are a pair of commands that using different compilers for benchmarks. CMake does not allow using more than one compiler within a CMake build, but `ctbench`'s compiler launcher can be used through these commands to work around that limitation and run benchmarks for compiler performance comparisons.
- The `ctbench-graph-all` target, which allows to build all the graphs declared with `ctbench` functions.

Now that we have a summary of `ctbench`'s high level API, it is time to take a look at its internals.

5.2.2 . `ctbench` internals

Besides the CMake scripts, `ctbench` is organized into several components:

- **grapher**: a C++ library and CLI tool that reads benchmark data and draws plots. It
- **compiler-launcher**: a simple program that handles compiler flags in order to move JSON profiling output to the right place for grapher to retrieve it, measures compiler execution time if required, and invokes the desired compiler as part of the `ctbench_set_compiler` CMake command implementation.

grapher: reading and plotting benchmark results

The grapher sub-project provides data structures for reading

- benchmark data structures and algorithms `<benchmark target name>/<iteration size>/<iteration number>.json`
- plotter engines

- reusable svg
- predicate interface
- CLI
- sciplot

compiler-launcher: working around CMake's limitations

CMake has its own limitations that would make the implementation of ct-bench impractical, or even impossible without a compiler launcher. Thankfully, CMake does support compiler launchers through the `CMAKE_CXX_COMPILER_LAUNCHER` target property, allowing us to work around those issues.

The main issue that compiler-launcher was meant to solve is the lack of a CMake interface to retrieve the path of the binaries it generates (which is needed to locate JSON trace files), and the lack of a Clang option to set the output path of time trace files until Clang 16.

The functionality of the wrapper was since extended to support more use cases.

- find out where json files are located with `-ftime-trace`
- CLI option parsing
- compiler execution time measurement

5.2.3 . Software quality

An emphasis was put on software quality since the beginning of the project. As mentioned, the goal was not just to develop a plotting tool for a single compilation time analysis, but to provide a methodology along with a robust implementation to improve compilation time analysis as a whole.

- The use of common C++ development tools for the project facilitates makes it easier for users to contribute to the project.
- The use of a GitHub CI for building and testing the project guarantees that project remains functional at all times on all supported platforms. As of writing, the project is continuously built and tested for Ubuntu 23.04 and Arch Linux. The test environment is fully reproducible as well thanks to the use of Docker for its setup.
- compiler checks: clang-tidy and clang-format
- Packages are provided through the AUR and the vcpkg repository
- Accepted in JOSS, which puts an emphasis on software quality. the review itself is meant to improve the quality of the software thanks to the reviewers' feedback.

5.2.4 . Conclusion

ctbench is a well

5.3 . A ctbench use case

This example covers a short yet practical example of ctbench usage. We want to calculate the sum of a series of integers known at compile-time, using a type template to store unsigned integer values at compile-time.

We will be comparing the compile-time performance of two implementations:

- one based on a recursive function template,
- and one based on C++ 11 parameter pack expansion.

First we need to include `utility` to instantiate our benchmark according to the size parameter using `std::make_index_sequence`, and define the compile-time container type for an unsigned integer:

```
#include <utility>

/// Compile-time std::size_t
template <std::size_t N> struct ct_uint_t {
    static constexpr std::size_t value = N;
};
```

The first version of the metaprogram is based on a recursive template function:

```
/// Recursive compile-time sum implementation
template<typename ... Ts> constexpr auto sum();

template <> constexpr auto sum() { return ct_uint_t
    <0>{}; }
template <typename T> constexpr auto sum(T const &) {
    return T{}; }

template <typename T, typename... Ts>
constexpr auto sum(T const &, Ts const &...tl) {
    return ct_uint_t<T::value + decltype(sum(tl...))::
        value>{};
}
```

And the other version relies on C++ 11 parameter pack expansion:

```
/// Expansion compile-time sum implementation
template<typename ... Ts> constexpr auto sum();
```

```

template <> constexpr auto sum() { return ct_uint_t
    <0>{}; }

template <typename... Ts> constexpr auto sum(Ts const
    &...) {
    return ct_uint_t<(Ts::value + ... + 0)>{};
}

```

Both versions share the same interface, and thus the same driver code as well. The driver code takes care of scaling the benchmark according to `BENCHMARK_SIZE`, which is defined by `ctbench` through the CMake API:

```

// Driver code

template <typename = void> constexpr auto foo() {
    return []<std::size_t... Is>(std::index_sequence<Is
        ...>) {
        return sum(ct_uint_t<Is>{}...);
    }
    (std::make_index_sequence<BENCHMARK_SIZE>{});
}

[[maybe_unused]] constexpr std::size_t result =
    decltype(foo())::value;

```

The CMake code needed to run the benchmarks is the following:

```

ctbench_add_benchmark(
    variadic_sum.expansion variadic_sum/expansion.cpp ${
        BENCHMARK_START}
    ${BENCHMARK_STOP} ${BENCHMARK_STEP} ${
        BENCHMARK_ITERATIONS})

ctbench_add_benchmark(
    variadic_sum.recursive variadic_sum/recursive.cpp ${
        BENCHMARK_START}
    ${BENCHMARK_STOP} ${BENCHMARK_STEP} ${
        BENCHMARK_ITERATIONS})

```

Then a graph target can be declared:

```

ctbench_add_graph(variadic_sum-compare-graph compare-
    all.json
                variadic_sum.expansion variadic_sum.
                recursive)

```

with `compare-all.json` containing the following:

```

{
  "plotter": "compare_by",
  "legend_title": "Timings",
  "x_label": "Benchmark size factor",
  "y_label": "Time (microsecond)",
  "draw_average": true,
  "demangle": false,
  "draw_points": false,
  "width": 800,
  "height": 400,
  "key_ptrs": ["/name", "/args/detail"],
  "value_ptr": "/dur",
  "plot_file_extensions": [".svg"]
}

```

This configuration file uses the `compare_by` plotter to generate one plot for each pair of elements designated by the JSON pointers in `key_ptrs`, namely `/name` and `/args/detail`. The first pointer designates the LLVM timer name, and the second *may* refer to metadata such a C++ symbol, or a C++ source filename. The `demangle` option may be used to demangle C++ symbols using LLVM.

The result is a series of graphs, each one designating a particular timer event, specific to a source or a symbol whenever it is possible (ie. whenever metadata is present in the `/args/detail` value of a timer event). Each graph compares the evolution of these timer events in function of the benchmark size.

The graphs following were generated on a Lenovo T470 with an Intel i5 6300U and 8GB of RAM. The compiler is Clang 14.0.6, and Pyperf[**pyperf**] was used to turn off CPU frequency scaling to improve measurement precision.

The first timer we want to look at is `ExecuteCompiler`, which is the total compilation time. Starting from there we can go down in the timer event hierarchy to take a look at frontend and backend execution times.

The backend is not being impacted here, supposedly because this is purely a compile-time program and the output program is empty. However this might not be the case for all metaprograms, and metaprograms might have different impacts on the backend as they may generate programs in different ways (ie. generate more symbols, larger symbols, more data structures, etc.).

The Total Instantiate function timer is an interesting one as it explicitly targets function instantiation time. Note that timers that are prefixed with "Total" measure the total time spent in a timer section, regardless of the specific symbol or source associated to its individual timer events.

Finally, we can take a look at `InstantiateFunction/foovoid.png` which measures the `InstantiateFunction` event time specifically for `foo<void>()`, which

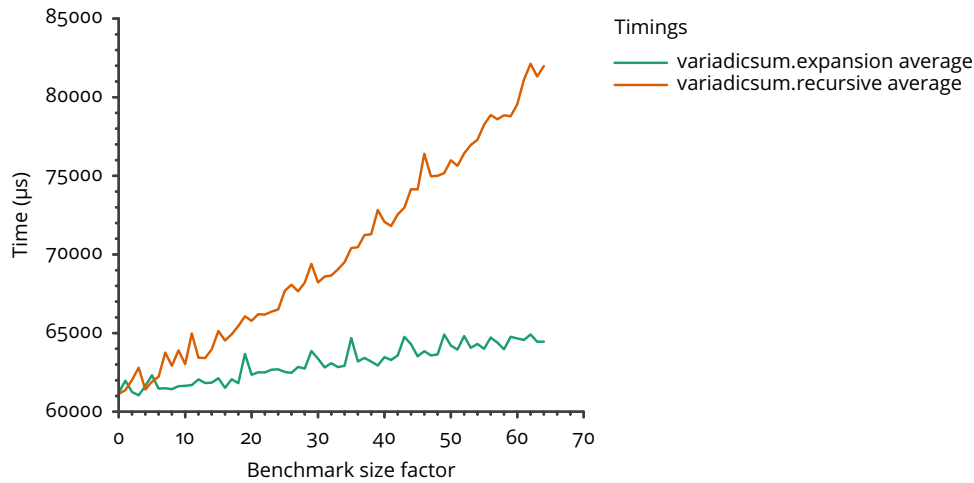


Figure 5.2: ExecuteCompiler

is our driver template function. Using Perfetto UI to look at the timer event hierarchy, we can validate that the timer event for this specific symbol includes the `InstantiateFunction` time for all the symbols that may be instantiated within this function.

This level of detail and granularity in the analysis of compile-time benchmarks was never reached before, and may help us set good practices to improve the compile-time performance of metaprograms.

5.4 .Related projects

- Poacher (<https://github.com/jpenuchot/poacher>): Experimental const-expr parsing and code generation for the integration of arbitrary syntax DSL in C++ 20
- Rule of Cheese (<https://github.com/jpenuchot/rule-of-cheese>): A collection of compile time microbenchmarks to help set better C++ metaprogramming guidelines to improve compile time performance

5.5 .Acknowledgements

We acknowledge contributions from Philippe Virouleau and Paul Keir for their insightful suggestions.

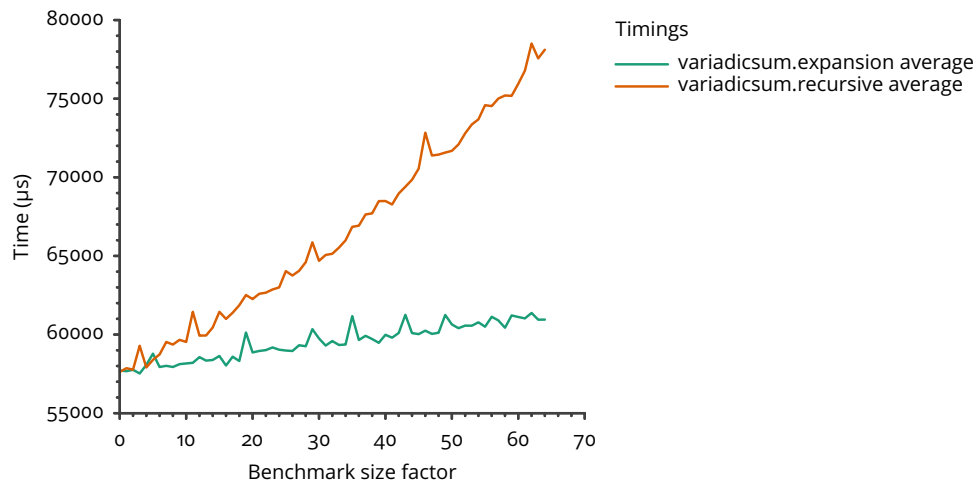


Figure 5.3: Total Frontend

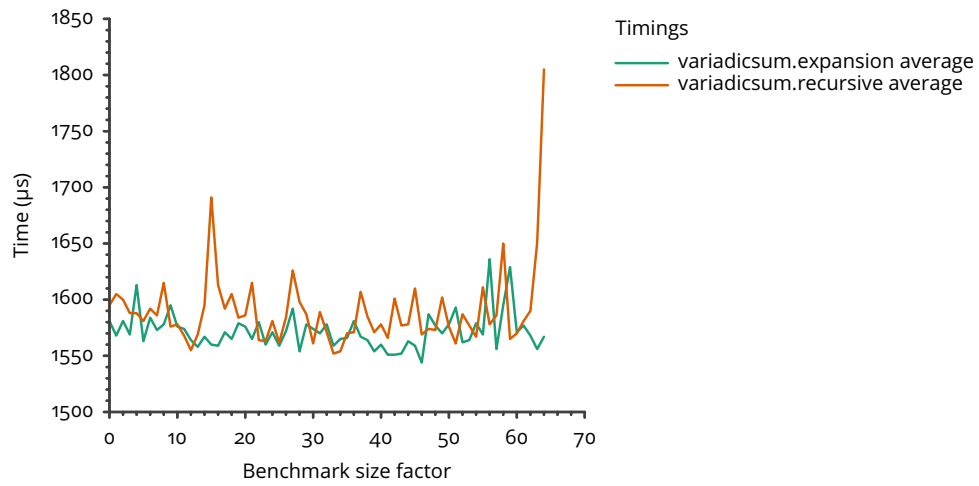


Figure 5.4: Total Backend

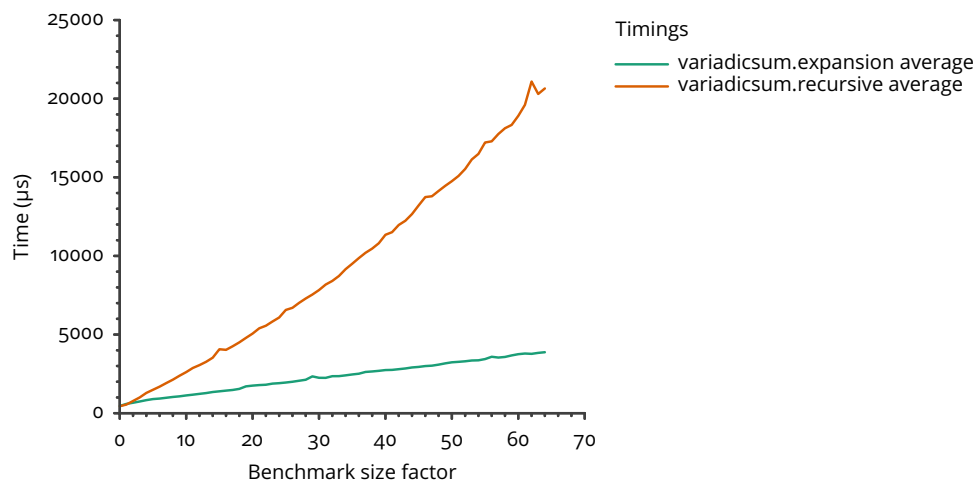


Figure 5.5: Total InstantiateFunction

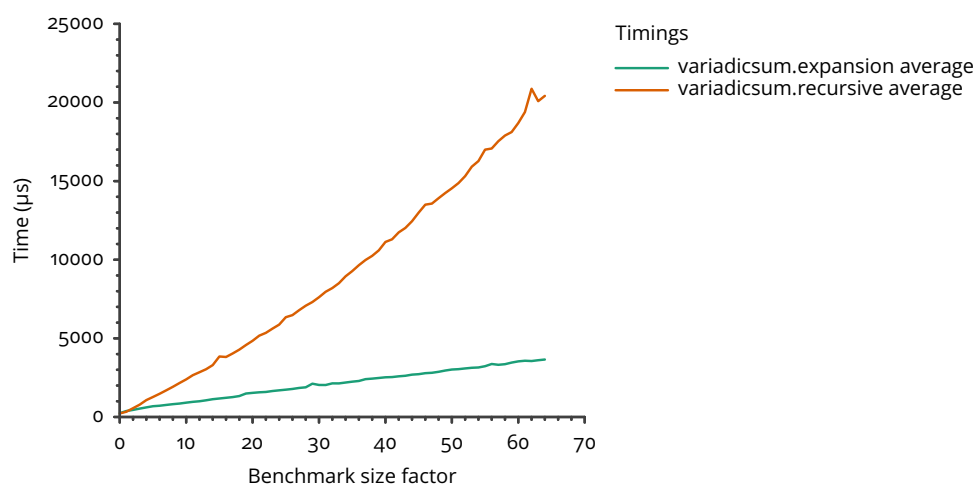


Figure 5.6: InstantiateFunction foovoid

6 - Constexpr parsing for high performance computing

6.1 .Introduction

C++ is often touted as a *Zero-Cost Abstraction* language due to some of its design philosophy and its ability to compile abstraction to a very efficient binary code. Some more radical techniques can be used to force the compiler to interpret C++ code as a DSEL. TMP is such a technique and it spawned a large corpus of related idioms from compile time function evaluation to lazy evaluation via ETs.

In the field of High Performance Computing, C++ users are often driven to use libraries built on top of those idioms like Eigen[**eigen**] or Blaze[**blazelib, iglberger2012_2**]. They all suffer from a major limitation: by being tied to the natural C++ syntax, they can't express nor embed arbitrary languages.

In this thesis, we try to demonstrate that the new features of C++ 23 related to compile time programming are able to help developers designing DSELS with arbitrary syntax by leveraging `constexpr` computations, compile time dynamic objects and lazy evaluation through lambda functions. After contextualizing our contribution in the general DSELS domain, this thesis will explain the core new techniques enabled by C++ 23 and how we can apply to build two different DSELS with their own non-C++ syntax. We'll also explore the performances of said DSELS in term of compile time to assess their usability in realistic code.

6.2 . A technical background of C++ DSELS

By definition, a DSL is a computer language specialized to a particular application domain, contrary to a general-purpose language, which is broadly applicable across domains, and lacks specialized features for a particular domain. DSELS are a subclass of DSL that rely on an existing general-purpose language to host it. DSELS then reuse the host language syntax and tool ecosystem to be compiled or interpreted.

In C++, the compile time process of generating new code is TMP to ensure performance and correctness.

6.2.1 . Constexpr programming

Expressions that generate results for use at compile time are called constant expressions. In C++ 11, the `constexpr` keyword was introduced to qualify functions or variables that may be used in constant expressions. Conse-

quently, regular functions and variables became usable for the evaluation of NTTPs.

However not all functions and variables can be qualified as `constexpr`, and not all constant expression results can be used as NTTPs.

Since the introduction of the `constexpr` specifier, the requirements on functions for being `constexpr` specifiable have constantly been relaxed as new C++ standards were adopted. In C++ 20, notable changes made dynamic memory allocations[**constexpr-memory**] and virtual `constexpr` function calls[**virtual-constexpr**] allowed in compile time `constexpr` function executions.

These two additions make dynamic memory and heritage-based polymorphism possible in `constexpr` functions. Therefore more regular C++ code can be executed at compile time, including parsers for arbitrary languages.

constexpr functions Functions that are `constexpr` qualified can still be used in other contexts than constant evaluation happening at compile time. In non-constant evaluation, `constexpr` functions can still call non-`constexpr` functions. But in constant evaluations, `constexpr` functions must only call other `constexpr` functions. This applies to methods as well. In order to make a C++ class or structure fully usable in constant evaluations, its methods — including the constructors and destructor— must be `constexpr`.

constexpr variables Variables that are `constexpr` qualified can be used in constant expressions. Note that they are different from non-`constexpr` variables used in `constexpr` functions. There are more requirements on `constexpr` variables. Their values must be literal, meaning that memory allocated in `constexpr` function bodies cannot be stored in `constexpr` variables.

constexpr memory allocation Starting from C++ 20, `std::allocate` and `std::deallocate` are `constexpr` functions, allowing memory allocations to happen in constant evaluations.

However `constexpr` allocated memory is not transient, *i.e.* memory allocated in constant expressions cannot be stored in `constexpr` variables, and NTTPs cannot hold pointers to `constexpr` allocated memory either.

Note that this restriction does not mean that data stored in `constexpr` memory cannot be passed through. There are techniques to use data in `constexpr` allocated memory

Listing 6.1: Illustration of constraints on `constexpr` allocated memory

```
// Cplusplus function generate returns a non-literal value
constexpr std::vector<int> generate() { return {1,2,3,4,5}; }
```

```

// Function template foo takes a polymorphic NTTP
template<auto bar> constexpr int foo() { return 1; }

// generate's return value cannot be stored in a
// constexpr variable
// or used as a NTTP, but it can be used to produce
// other literal

// constexpr auto a = generate();           // ERROR
constexpr auto b = generate().size();      // OK
// constexpr auto c = foo<generate>();     // ERROR
constexpr auto d = foo<&generate>();       // OK

```

Let's have a closer look at the four assignment cases:

- Case a: generate's return value is non-literal and therefore cannot be stored in a constexpr variable.
- Case b: generate's return value is used in a constant expression to produce a literal value. Therefore the expression's result can be stored in a constexpr variable.
- Case c: similarly to case a, generate's return value cannot be used as an NTTP because it is not a literal value.
- Case d: function references are allowed as NTTPs.

Notice how the last example works around restrictions of constexpr allocations by using a generator function instead of passing the non-empty `std::vector<int>` value directly. This technique along with the definition of lambdas can be used to explore more complex structures returned by constexpr functions such as pointer trees.

Moreover, constexpr allocated memory being non transient does not mean that its content cannot be transferred to NTTP compatible data structures.

Listing 6.2: Transformation from a dynamic vector to a constexpr static array

```

constexpr auto generate_as_array() {
    constexpr std::size_t array_size = generate().size();
    ;
    std::array<int, array_size> res{};
    std::ranges::copy(generate(), res.begin());
    return res;
}

```

```
constexpr auto e = generate_as_array();           // OK
constexpr auto f = foo<generate_as_array()>();    // OK
```

Listing 6.2 shows how `generate`'s result can be evaluated into a static array. Static arrays are literal as long as the values they hold are literal. Therefore the result of `generate_as_array` can be stored in a `constexpr` variable or used directly as an NTTP for code generation.

constexpr virtual functions This feature allows calls to virtual functions in constant expressions [**virtual-constexpr**]. This allows heritage-based polymorphism in `constexpr` programming when used with `constexpr` allocation of polymorphic types.

Cest: standard-like containers for `constexpr` programming

With dynamic memory and virtual functions making their way to `constexpr` programming, the next logical step is to make standard C++ containers available in constant evaluations. Modifying the C++ standard to accomodate that need takes time. However, implementing standard-like containers that are good enough for prototyping and experimenting with `constexpr` programs can be done fast enough.

The C'est [**cest**] library was created with this goal in mind, filling the gap in the C++ standard, and accommodating the need for `constexpr` compatible standard containers. As of today most of the containers implemented in C'est are available in up-to-date standard libraries, but it is still a useful library for compatibility with platforms that do not have C++ 23 compatible compilers and standard libraries such as old versions of Debian.

It was instrumental for this thesis as the research work I present here started a long time before C++ 23 was adopted and standard libraries and compilers started implementing it.

I contributed to its development myself by implementing a `constexpr` version of `std::unique_ptr`. This collaboration led Joel Falcou and myself to meet with its initiator, Paul Keir, in Glasgow and Paris. Our collaboration also led to speak together at Meeting C++ 2022 [**meetingcpp22**] where we talked about C'est and its use cases for `constexpr` programming research.

6.2.2 . C++ Domain Specific Embedded Languages

DSELs in C++ use TMP via the *Expression Template* idiom. **Expression Templates** [**veldhuizen:1995**, **vandevoorde:2002**] is a technique implementing a form of **delayed evaluation** in C++ [**spinellis:2001**]. Expression Templates are built around the *recursive type composition* idiom [**jarvi:1998**] that allows the construction, at compile time, of a type representing the Abstract Syntax Tree (AST) of an arbitrary statement. This is done by overloading functions

and operators on those types so they return a lightweight object. The object encodes the current operation in the AST being built in its type instead of performing any kind of computation. Once reconstructed, this AST can be transformed into arbitrary code fragments using Template Metaprograms.

As of today, most C++ EDSLs rely on *Expression Templates* and therefore are limited to the C++ syntax. New techniques are becoming more popular through the use of `constexpr` strings to embed arbitrary DSELs. One major example is the CTRE [ctre] that implements most of the Perl Compatible Regular Expression (PCRE) syntax. However, CTRE still relies on type-based TMP to parse regular expressions and transform them into regular expression evaluators.

6.3 . Code generation from constexpr allocated structures

In this section, we will present different solutions to generate code from non-literal data structures generated with `constexpr` functions.

Non-literal data structures allow more flexibility through the use of dynamic memory, and allow simpler code as they do not require contraptions to get around the lack of dynamic allocations.

However, keep in mind that there are still limitations. The solutions evaluated in this thesis are still workarounds. What they allow is the use of dynamically-sized structures for code generation.

The first subsection covers a case where we need to convert a pointer tree returned by a `constexpr` function into code.

The second one will cover a case where the tree is returned in a serialized representation.

6.3.1 . Code generation from pointer tree data structures

In this subsection, we introduce three techniques that will allow us to use a pointer tree generated from a `constexpr` function as a template parameter for code generation.

To illustrate them, we use a minimalistic use case. A generator function returns a pointer tree representation of addition and constant nodes, which we will use to generate functions that evaluate the tree itself.

Listing 6.3: `tree_t` type definition

```
/// Enum type to differentiate nodes
enum node_kind_t {
    constant_v,
    add_v,
};

/// Node base type
```

```

struct node_base_t {
    node_kind_t kind;

    constexpr node_base_t(node_kind_t kind_) : kind(
        kind_) {}
    constexpr virtual ~node_base_t() = default;
};

/// Tree pointer type
using tree_ptr_t = std::unique_ptr<node_base_t>;

/// Checks that an object is a tree generator
template <typename T>
concept tree_generator = requires(T fun) {
    { fun() } -> std::same_as<tree_ptr_t>;
};

/// Constant node type
struct constant_t : node_base_t {
    int value;
    constexpr constant_t(int value_)
        : node_base_t(constant_v), value(value_) {}
};

/// Addition node type
struct add_t : node_base_t {
    tree_ptr_t left;
    tree_ptr_t right;
    constexpr add_t(tree_ptr_t left_, tree_ptr_t right_)
        : node_base_t(add_v), left(std::move(left_)),
          right(std::move(right_)) {}
};

/// Generates an arbitrary tree
constexpr tree_ptr_t gen_tree() {
    return std::make_unique<add_t>(
        std::make_unique<add_t>(std::make_unique<
            constant_t>(1),
                                std::make_unique<
                                    constant_t>(2)),
        std::make_unique<constant_t>(3));
}

```

Listing 6.3 shows the type definitions, a concept to match tree generator functions, as well as the generator function itself. This is a common way to

represent trees in C++, but the limits mentioned in 6.2.1 make it impossible to use the result of `gen_tree` directly as a template parameter to generate code.

The techniques we will use to pass the result as a template parameter are:

- passing functions that return nodes as NTTPs,
- converting the tree into an expression template representation,
- serializing the tree into a dynamically allocated array, then converting the dynamic array into a static array that can be used as a NTTP,

The compilation performance measurements in ?? will rely on the same data passing techniques, but with more complex examples such as embedded compilation of Brainfuck programs, and of \LaTeX math formulae into high performance math computation kernels.

Pass-by-generator

One way to use dynamically allocated data structures as template parameters is to pass their generator functions instead of their values. You may have noted in listing 6.1 that despite `generate`'s return value being non-literal, the function itself can be passed as a NTTP.

Its result can be used to produce literal `constexpr` results, and the function itself can be used in generator lambda functions defined at compile-time.

Listing 6.4: Pass-by-generator

```
/// Accumulates the value from a tree returned by
/// TreeGenerator.
/// TreeGenerator() is expected to return a std::
/// unique_ptr<tree_t>.
template <tree_generator auto Fun> constexpr auto
codegen() {
    static_assert(Fun() != nullptr, "Ill-formed tree");

    constexpr node_kind_t Kind = Fun()->kind;

    if constexpr (Kind == add_v) {
        // Recursive codegen for left and right children:
        // for each child, a generator function is
        // generated
        // and passed to codegen.
        auto eval_left = codegen<[]>() {
            return static_cast<add_t &&>(*Fun()).left;
        }>();
```

```

    auto eval_right = codegen<[]>() {
        return static_cast<add_t &&>(*Fun()).right;
    }>();

    return [=]() { return eval_left() + eval_right();
    };
}

else if constexpr (Kind == constant_v) {
    constexpr auto Constant =
        static_cast<constant_t const &>(*Fun()).value;
    return [=]() { return Constant; };
}
}

```

In listing 6.4, we show how a `constexpr` pointer tree result can be visited recursively using generator lambdas to pass the subnodes' values, and used to generate code.

This technique is fairly simple to implement as it does not require any transformation into an ad hoc data structure to pass the tree as a type or NTTP.

The downside of using this value passing technique is that the number of calls of the generator function is proportional to the number of nodes. Experiments in ?? highlight the scaling issues induced by this code generation method. And while it is very quick to implement, there are still difficulties related to `constexpr` memory constraints and compiler or library support. GCC 13.2.1 is still unable to compile such code

Pass-by-generator + ET

Why through? Interoperability with type-based metaprogramming libraries.

Types:

```

/// Type representation of a constant
template <int Value> struct et_constant_t {};
/// Type representation of an addition
template <typename Left, typename Right> struct
    et_add_t {};

```

Codegen:

```

/// Accumulates the value from a tree returned by
    TreeGenerator.
/// TreeGenerator() is expected to return a std::
    unique_ptr<tree_t>.
template <tree_generator auto Fun>

```



```
constexpr auto to_expression_template() {
    static_assert(Fun() != nullptr, "Ill-formed tree");

    constexpr node_kind_t Kind = Fun()->kind;

    if constexpr (Kind == add_v) {
        // Recursive type generation using the pass-by-
        // generator technique
        using TypeLeft = decltype(to_expression_template
            <[]() {
                return static_cast<add_t &&>(*Fun()).left;
            }>());
        using TypeRight = decltype(to_expression_template
            <[]() {
                return static_cast<add_t &&>(*Fun()).right;
            }>());

        return et_add_t<TypeLeft, TypeRight>{};
    }

    else if constexpr (Kind == constant_v) {
        constexpr auto Value = static_cast<constant_t &>(*
            Fun()).value;
        return et_constant_t<Value>{};
    }
}

template <int Value>
constexpr auto codegen_impl(et_constant_t<Value> /*
    unused*/) {
    return []() { return Value; };
}

template <typename ExpressionLeft, typename
    ExpressionRight>
constexpr auto
codegen_impl(et_add_t<ExpressionLeft, ExpressionRight>
    /*unused*/) {
    auto eval_left = codegen_impl(ExpressionLeft{});
    auto eval_right = codegen_impl(ExpressionRight{});
    return [=]() { return eval_left() + eval_right(); };
}

template <tree_generator auto Fun> constexpr auto
codegen() {
```

```

using ExpressionTemplate = decltype(
    to_expression_template<Fun>());
return codegen_impl(ExpressionTemplate{});
}

```

FLAT

To overcome the performance issues of the previously introduced techniques, we can try a different approach. Instead of passing trees as they are, they can be transformed into static arrays which can be used as NTPs.

```

/// Serialized representation of a constant
struct flat_constant_t {
    int value;
};

/// Serialized representation of an addition
struct flat_add_t {
    std::size_t left;
    std::size_t right;
};

/// Serialized representation of a node
using flat_node_t = std::variant<flat_add_t,
    flat_constant_t>;

/// Defining max std::size_t value as an equivalent to
    std::nullptr.
constexpr std::size_t null_index = std::size_t(0) - 1;

```

This requires the tree to be serialized first. For the sake of demonstration, we will serialize the tree into an ad hoc data representation that is identical to the original one, except pointers are replaced with `std::size_t` indexes, and `nullptr` is replaced with an arbitrary value called `null_index` as shown in listing 6.3.1.

Heritage polymorphism is also replaced by `std::variant`

These nodes will be stored in `std::vector` containers, and the indexes will refer to the position of other nodes within the container. Note that our tree nodes are not polymorphic. If needed, `std::variant` could have been used to have polymorphic nodes in the serialized representation.

Listing 6.5: constexpr tree serialization implementation

```

/// Serializes the current subtree and returns
/// the top node's index to the caller.

```

```

constexpr std::size_t serialize_impl(tree_ptr_t const
    &top,
                                     std::vector<
                                     flat_node_t> &
                                     out) {
    // nullptr translates directly to null_index
    if (top == nullptr) {
        return null_index;
    }

    // Allocating space for the destination node
    std::size_t dst_index = out.size();
    out.emplace_back();

    if (top->kind == add_v) {
        auto const &typed_top = static_cast<add_t const
            &>(*top);

        // Serializing left and right subtrees,
        // initializing the new node
        out[dst_index] = {
            flat_add_t{.left = serialize_impl(typed_top.
                left, out),
                    .right = serialize_impl(typed_top.
                right, out)}};
    }

    if (top->kind == constant_v) {
        auto const &typed_top = static_cast<constant_t
            const &>(*top);
        out[dst_index] = {flat_constant_t{.value =
            typed_top.value}};
    }

    return dst_index;
}

/// Returns a serialized representation of a pointer
tree.
constexpr std::vector<flat_node_t> serialize(
    tree_ptr_t const &tree) {
    std::vector<flat_node_t> result;
    serialize_impl(tree, result);
    return result;
}

```

Listing 6.5 shows the implementation of the serialization step. The implementation itself has nothing particular, except for the functions being `constexpr`.

Once the data is serialized, it can be converted into a static array container. This can be done because the generator function is `constexpr`, therefore it can be used to produce `constexpr` values. The size of the resulting `std::vector` can be stored at compile time to set the size of a static array.

Listing 6.6: Definition of `serialize_as_array`

```
/// Evaluates a tree generator function into a
    serialized array.
template <tree_generator auto Fun>
constexpr auto serialize_as_array() {
    constexpr std::size_t Size = serialize(Fun()).size()
        ;
    std::array<flat_node_t, Size> result;
    std::ranges::copy(serialize(Fun()), result.begin());
    return result;
}
```

Listing 6.6 shows the implementation of a helper function that takes a `constexpr` tree generator function as an input, serializes the result, and returns it as a static array.

Static arrays are not dynamically allocated, therefore they can be used as NTTPs if their values do not hold pointers to `constexpr` allocated memory either. In our case, the elements only hold integers, so a `std::array` of `flat_node_t` elements can be used as a NTTPs.

Listing 6.7: Code generation implementation for the flat backend

```
/// Generates code from a tree serialized into a
    static array,
/// with CurrentIndex being the index of the starting
    node which
/// defaults to 0, ie. the top node of the whole tree.
template <auto const &Tree, std::size_t Index>
constexpr auto codegen_aux() {
    static_assert(Index != null_index, "Ill-formed tree
        (null index)");

    if constexpr (std::holds_alternative<flat_add_t>(
        Tree[Index])) {
        constexpr auto top = std::get<flat_add_t>(Tree[
            Index]);
    }
```

```

    // Recursive code generation for left and right
    children
    auto eval_left = codegen_aux<Tree, top.left>();
    auto eval_right = codegen_aux<Tree, top.right>();

    // Code generation for current node
    return [=]() { return eval_left() + eval_right();
    };
}

else if constexpr (std::holds_alternative<
    flat_constant_t>(
        Tree[Index])) {
    constexpr auto top = std::get<flat_constant_t>(
        Tree[Index]);
    constexpr int Value = top.value;
    return []() { return Value; };
}
}

/// Stores serialized representations of tree
generators' results.
template <tree_generator auto Fun>
static constexpr auto tree_as_array =
    serialize_as_array<Fun>();

/// Takes a tree generator function as non-type
template parameter
/// and generates the lambda associated to it.
template <tree_generator auto Fun> constexpr auto
codegen() {
    return codegen_aux<tree_as_array<Fun>, 0>();
}

```

To complete the implementation, we must implement a code generation function that accepts a serialized tree as an input as shown in listing 6.7. Note that this function is almost identical to the one shown in listing 6.4. The major difference is that `TreeGenerator` is called only twice regardless of the size of the tree. This allows much better scaling as we will see in ???. The downside is that it requires the implementation of an ad hoc data structure and a serialization function, which might be more or less complex depending on the complexity of the original tree structure.

6.3.2 . Using algorithms with serialized outputs

Current token	Action	Stack
2	Stack 2	2
3	Stack 3	2, 3
2	Stack 2	2, 3, 2
*	Multiply 2 and 3	2, 6
+	Add 2 and 6	8

Figure 6.1: formula reading example

Parsing algorithms may output serialized data. In this case, the serialization step described in 6.3.1 is not needed, and the result can be converted into a static array. This makes the code generation process rather straightforward as no complicated transformation is needed, while still scaling decently as we will see in ?? where we will be using a Shunting Yard parser [**shunting-yard**] to parse math formulae to a , which is its postfix notation.

Once converted into its postfix notation, a formula can be read using the following method:

- read symbols in order,
- put constants and variables on the top of a stack,
- when a function f or operator of arity N is being read, unstack N values and stack the result of f applied to the N operands.

Figure 6.1 shows a formula reading example with the formula $2 + 3 * 2$, or $2\ 3\ 2\ * \ +$ in reverse polish notation.

Starting from there, we will see how code can be generated using RPN representations of addition trees in C++.

Listing 6.8: RPN example base type and function definitions

```

/// Type for RPN representation of a constant
struct rpn_constant_t {
    int value;
};

/// Type for RPN representation of an addition
struct rpn_add_t {};

/// Type for RPN representation of an arbitrary symbol
using rpn_node_t = std::variant<rpn_constant_t,
    rpn_add_t>;

/// RPN equivalent of gen_tree

```

```
constexpr std::vector<rpn_node_t> gen_rpn_tree() {
    return {rpn_constant_t{1}, rpn_constant_t{2},
            rpn_add_t{},
            rpn_constant_t{3}, rpn_add_t{}};
}
```

In listing 6.8, we have the type definitions for an RPN representation of an addition tree as well as `gen_rpn_tree` which returns an RPN equivalent of `gen_tree`'s result.

Similar to the flat backend, an `eval_as_array` takes care of evaluating the `std::vector` result into a statically allocated array.

Listing 6.9: Codegen implementation for RPN formulae

```
/// Codegen implementation.
/// Reads tokens one by one, updates the stack
/// consequently
/// by consuming and/or stacking operands.
/// Operands are functions that evaluate parts of the
/// subtree.
template <auto const &RPNTree, std::size_t Index = 0>
constexpr auto codegen_impl(kumi::product_type auto
    stack) {
    // The end result is the last top stack operand
    if constexpr (Index == RPNTree.size()) {
        static_assert(stack.size() == 1, "Invalid tree");
        return kumi::back(stack);
    }

    else if constexpr (std::holds_alternative<
        rpn_constant_t>(
            RPNTree[Index])) {
        // Append the constant operand
        auto new_operand = [=]() {
            return get<rpn_constant_t>(RPNTree[Index]).value
                ;
        };
        return codegen_impl<RPNTree, Index + 1>(
            kumi::push_back(stack, new_operand));
    }

    else if constexpr (std::holds_alternative<rpn_add_t>
        >(
            RPNTree[Index])) {
        // Fetching 2 top elements and popping them off
        the stack
```

```

    auto left = kumi::get<stack.size() - 1>(stack);
    auto right = kumi::get<stack.size() - 2>(stack);
    auto stack_remainder = kumi::pop_back(kumi::
        pop_back(stack));

    // Append new operand and process next element
    auto new_operand = [=]() { return left() + right()
        ; };
    return codegen_impl<RPNTree, Index + 1>(
        kumi::push_back(stack_remainder, new_operand))
        ;
}
}

/// Stores static array representations of RPN
/// generators' results.
template <auto Fun>
static constexpr auto rpn_as_array = eval_as_array<Fun
    >();

/// Code generation function.
template <auto Fun> constexpr auto codegen() {
    return codegen_impl<rpn_as_array<Fun>, 0>(kumi::
        tuple{});
}

```

In listing 6.9, we have function definitions for the implementation of `codegen` which takes an RPN generator function, and generates a function that evaluates the tree.

The code generation process happens by updating an operand stack represented by a `kumi::tuple`. It is a standard-like tuple type with additional element access, extraction, and modification functions. We are using it to store functions that evaluate parts of the subtree.

Symbols are read in order, and the stack is updated depending on the symbol:

- When a terminal is read, a function that evaluates its value is stacked. In this case terminals are simply constants, but they can be anything a C++ lambda can return such as a variable or another function's result.
- When a function or operator of arity N is read, N operands will be consumed from the top of the stack and a new operand will be stacked. The new operand evaluates the consumed operands passed to the function corresponding to the symbol being read.

Once all symbols are read, there should be only one operand remaining on the stack: the function that evaluates the whole tree.

6.3.3 . Preliminary observations

The presented code snippets already allow us to draw a few observations.

- Code generation from `constexpr` dynamic structures is not a trivial process. It still requires advanced C++ knowledge to understand the limits of what is or isn't possible to do with NTTPs and `constexpr` functions.

While the examples show working examples for code generation from `constexpr` function results with `constexpr` allocated memory, they do not show the time it takes to comply with C++ `constexpr` restrictions.

For example, condition results in `if constexpr` statements must be stored in a `constexpr` variable if the result is evaluated from an temporary object that contains `constexpr` allocated memory.

- The easiest way to get around NTTP constraints, *i.e.* the pass-by-generator (PBG) technique, is very costly in terms of compilation times.

Section ?? covers that scaling issue in more detail. The overall assessment is that this technique can be used for small metaprograms, but it fails to scale properly as larger ones are being considered due to its quadratic compilation time complexity.

In our Brainfuck metacompilation examples, we were able to trigger Clang's timeout using this technique when compiling a Mandelbrot visualizer whereas the so-called "flat" backend was able to generate the program in less than a minute.

- Compiler support for `constexpr` programming and constant evaluation in general is still very inconsistent across compilers. GCC still has issues with `if constexpr` where templates are instantiated even when they are in a discarded statement.
- Library support is also limited. Most of the containers from the C++ standard library are not usable in `constexpr` functions simply because the C++ standard lacks `constexpr` qualifications for their methods (and `constexpr` qualification is not implicit).

As of today, the only exceptions are the main containers such as `std::vector`, `std::string` (since C++20), or `std::unique_ptr` (since C++23).

The C'est [**cest**] library however provides more standard-like containers that are usable in `constexpr` functions in C++20, which is enough to make the previous examples run on Clang in C++20.

So far we can conclude that while being doable, code generation from `constexpr` function results with dynamic memory is still not accessible to all C++ programmers.

An effort on the language itself would be needed to allow easier data passing from `constexpr` allocated memory to NTTP.

Alternatively, the language could allow code generation

6.4 . Brainfuck parsing and code generation

Now that we introduced the various techniques to generate programs from pointer trees generated by `constexpr` functions, we will use them in the context of compile time parsing and code generation for the Brainfuck language. Therefore use data structures and code generation techniques introduced in subsection 6.3.1.

6.4.1 . Constexpr Brainfuck parser and AST

AST

The Brainfuck AST is defined in the header shown in appendix .1.1. The header file also contains helper function definitions to handle AST nodes safely, such as `visit` which will be used in one of the code generation backends.

Here are the main data types:

- `node_interface_t`, which is a common base type for all AST nodes.
- `ast_token_t`, which represents a single AST token.
- `ast_block_t`, which represents an AST block, which simply is a `std::vector` of `std::unique_ptr<node_interface_t>`.
- `ast_while_t`, which represents a while conditional block. The instruction block itself is contained in an `ast_block_t` value.

The implementation of the AST are available in appendix .1.1 where you can observe that all the types are implemented as they would be for a regular Brainfuck parser, except all their methods are `constexpr`.

Listing 6.10: Definition of the AST visitor function

```
template <typename F>
constexpr auto visit(F f, ast_node_ptr_t const &p) {
    switch (p->get_kind()) {
        case ast_token_v:
            return f(static_cast<ast_token_t const &>(*p));
        case ast_block_v:
            return f(static_cast<ast_block_t const &>(*p));
```

```
case ast_while_v:
    return f(static_cast<ast_while_t const &>(*p));
}
}
```

The `visit` function implementation also looks like a regular C++ function as shown in listing 6.10. It is a higher order function that allows recursive operations on the AST to be carried in a type-safe manner.

Parser

The Brainfuck parser, again, looks like nothing special. For that reason I will not get into the implementation details. The function definition is available in appendix .1.2.

On the surface: the parser takes a pair of begin and end iterators as an input. It parses Brainfuck tokens until it reaches the end iterator or a while end token, and returns a pair containing an iterator pointing after the last parsed token and the parsing result.

When a while begin token is reached, it calls itself recursively and resumes parsing at the position of the iterator returned by the callee, which is right after the while block.

The main parsing function implementation (including the function prototype) is very condensed: it fits in 40 lines of code with a max line width set to 84. It is no different from a regular Brainfuck parsing function except for it being `constexpr`, and it can actually be used as a regular C++ program.

These make it much easier to debug as it can be ran through a C++ debugger like GDB or LLDB, and also more maintainable as it does not require any template metaprogramming experience to understand the implementation. Additionally, `constexpr` execution enforces checks on memory allocations and deallocations as well as memory bound checking. Therefore testing functions in `constexpr` contexts can help finding memory safety issues.

6.4.2 . A variety of techniques to generate code from a dynamic AST

Once the `constexpr` parser is implemented, the next step consists in figuring out how to transform its result, which contains dynamic memory, into C++ code.

As you may remember from subsection 6.2.1, there is no direct way to use values holding pointers to dynamic memory directly as NTTPs. Therefore it must be conveyed by other means or transformed into literal values to be used as template parameters for C++ code generation.

I implemented several of these workarounds to compare them. This will give us a clearer idea of their implementation difficulty, and they will enable

us to run compilation time benchmarks to compare their compilation time performance.

Pass-by-generator + ET

The first technique I implemented was passing AST nodes through lambdas to convert them into expression templates.

Pass-by-generator

actually very hard to implement because `std::unique_pointer` lifetime management gets more difficult when combined with `constexpr` constraints.

Serializing the AST into a literal value

The last technique I will discuss, which is by far the most efficient in regard to compilation time, consists in transforming the AST into a literal value that can be used as a NTTP.

As I've discussed in 6.2.1, a `constexpr std::vector` return value can be transformed into a static array with a trivial transformation. As long as a static array holds literal values, the array itself remains a literal value as well, and can therefore be used as a NTTP.

In order to do this, an intermediate serialized representation of the AST which must only contain literal values must be implemented. It is a fairly trivial task:

- Polymorphism through the inheritance of `node_interface_t` can be replaced with the use of `std::variant`, which is literal as long as the value it holds is literal as well.
- Pointers contained in AST nodes can be replaced by integer indexes pointing to internal values of the array (static or dynamic) in which serialized AST nodes are stored.

Listing 6.11: PBG intermediate representation type definitions

```
/// Represents a single instruction token
struct flat_token_t {
    token_t token;
};

/// Block descriptor at the beginning of every block
/// of adjacent instructions.
struct flat_block_descriptor_t {
    size_t size;
```

```

};

/// Represents a while instruction, pointing to
/// another instruction block
struct flat_while_t {
    size_t block_begin;
};

/// Polymorphic representation of a node
using flat_node_t =
    std::variant<flat_token_t, flat_block_descriptor_t
        ,
        flat_while_t>;

/// AST container type
using flat_ast_t = std::vector<flat_node_t>;

/// NTTP-compatible AST container type
template <size_t N>
using fixed_flat_ast_t = std::array<flat_node_t, N>;

```

In listing 6.11, I show how this was implemented in the backend that will be used for the benchmarks. Note that the serialized AST can be stored either as a `flat_ast_t` for convenience in constant evaluations, which can then be transformed into a `fixed_flat_ast_t<N>` to be passed as a NTTP.

To serialize the AST, we will rely on a function called `block_gen` which visits the AST recursively and generates serializes the AST progressively.

Listing 6.12: `block_gen_state_t` definition

```

/// Support structure for generate_blocks function
struct block_gen_state_t {
    /// Flat AST blocks, result of block_gen
    std::vector<flat_ast_t> blocks;

    /// Keeping track of which block is being generated
    size_t block_pos = 0;

    /// Keeping track of the size of the AST
    size_t total_size = 0;
};

```

In listing 6.12 we begin by defining a structure called `block_gen_state_t` that contains the intermediate result.

Listing 6.13: Generic `block_gen` implementation

```

/// Extracts an AST into a vector of blocks of
/// contiguous operations
constexpr void block_gen(ast_node_ptr_t const &,
                        block_gen_state_t &);

// Node-specific overloads...

constexpr void block_gen(ast_node_ptr_t const &p,
                        block_gen_state_t &s) {
    visit([&s](auto const &v) { block_gen(v, s); }, p);
}

```

Listing 6.13 shows the implementation of the generic version of `block_gen`. It consists in a forward declaration followed by node-specific implementations which call the generic version recursively. The generic version relies on the `visit` function, and for that reason overloads for derivatives of `node_interface_t` must be declared first. To make it simpler, I chose to define them as well to avoid unnecessary forward declarations.

Listing 6.14: `block_gen` implementation

```

/// block_gen for a single AST token. Basically just a
/// push_back.
constexpr void block_gen(ast_token_t const &tok,
                        block_gen_state_t &s) {
    s.blocks[s.block_pos].push_back(
        flat_token_t{tok.token});
}

/// block_gen for an AST block.
constexpr void block_gen(ast_block_t const &blo,
                        block_gen_state_t &s) {
    // Save & update block pos before adding a block
    size_t const previous_pos = s.block_pos;
    s.block_pos = s.blocks.size();
    s.blocks.emplace_back();

    // Preallocating
    s.blocks[s.block_pos].reserve(
        blo.content.size() + 1);
    s.total_size += blo.content.size() + 1;

    // Adding block descriptor as a prefix
    s.blocks[s.block_pos].push_back(
        flat_block_descriptor_t{
            blo.content.size()});
}

```

```

// Flattening instructions
for (ast_node_ptr_t const &node :
    blo.content) {
    block_gen(node, s);
}

// Restoring block pos after recursive block
// processing
s.block_pos = previous_pos;
}

/// block_gen for a while instruction.
constexpr void block_gen(ast_while_t const &whi,
                        block_gen_state_t &s) {
    s.blocks[s.block_pos].push_back(
        flat_while_t{s.blocks.size()});
    block_gen(whi.block, s);
}

```

Listing 6.14 shows the implementation of `block_gen` overloads for all types derived from `node_interface_t`.

The most important one is the implementation for `ast_block_t`. This is the overload that generates new blocks in the `blocks` vector of the `block_gen_state_t` structure.

It works by simply allocating a new block in the `block_gen_state_t` structure, and traverse each node by calling `block_gen` recursively. This ensures the elements contained in AST block nodes remain contiguous which is instrumental to make code generation easier later on.

The elements contained in these blocks have their indexes pointing to the elements of the `blocks` vector of the `block_gen_state_t` structure, but the expected result is a single `std::vector` containing all AST nodes.

The `flatten` function is responsible for transforming the contents of `block_gen_state_t` into such a representation.

Listing 6.15: `flatten` implementation

```

constexpr flat_ast_t
flatten(ast_node_ptr_t const &parser_input) {
    flat_ast_t serialized_ast;

    // Extracting as vector of blocks
    block_gen_state_t bg_result;
    block_gen(parser_input, bg_result);

    // Small optimization to avoid reallocations

```

```

    serialized_ast.reserve(bg_result.total_size);

    // block_map[i] gives the index of
    // bg_result.blocks[i] in the serialized
    // representation
    std::vector<size_t> block_map;
    block_map.reserve(bg_result.blocks.size());

    // Step 1: flattening
    for (flat_ast_t const &block : bg_result.blocks) {
        // Updating block_map
        block_map.push_back(serialized_ast.size());

        // Appending instructions
        std::ranges::copy(
            block, std::back_inserter(serialized_ast));
    }

    // Step 2: linking
    for (flat_node_t &node : serialized_ast) {
        if (std::holds_alternative<flat_while_t>(node)) {
            flat_while_t &w_ref =
                std::get<flat_while_t>(node);
            w_ref.block_begin =
                block_map[w_ref.block_begin];
        }
    }

    return serialized_ast;
}

```

Listing 6.15 shows the implementation of `flatten`, which simply calls `block_gen`, chains the blocks together, and translates block indexes to account for the new layout.

The postfix serialization layout (ie. the fact that blocks are laid out one after the other, and not one inside the other) ensures that traversing the AST remains trivial. Note that derializing the blocks into an infix representation would essentially get us back to the unparsed representation of Brainfuck programs.

From there, the only thing that needs to be done is to evaluate this dynamic array into a static one to make it usable as an NTTP.

Listing 6.16: `parse_to_fixed_flat_ast` implementation

```

/// Parses a BF program into a fixed_flat_ast_t value.
template <auto const &ProgramString>

```



```
constexpr auto parse_to_fixed_flat_ast() {
    // Getting AST vector size into a constexpr variable
    constexpr size_t AstArraySize =
        flatten(parser::parse_ast(ProgramString))
            .size();

    // Initializing static size array
    fixed_flat_ast_t<AstArraySize> arr;
    std::ranges::copy(
        flatten(parser::parse_ast(ProgramString)),
        arr.begin());

    return arr;
}
```

Listing 6.16 shows the implementation of the final function that takes a string as a program and transforms it into a `fixed_flat_ast_t`.

So far, this is the only function that takes a template parameter as an input, and thus where the distinction between `constexpr` programming and TMP begins.

All the remaining work consists in implementing a function that takes the serialized AST as a template parameter and generates a program from it.

Listing 6.17: `program_state_t` definition

```
struct program_state_t {
    constexpr program_state_t() : i(0) {
        for (auto &c : data) {
            c = 0;
        }
    }

    std::array<char, 30000> data;
    std::size_t i;
};
```

The functions generated by the codegen functions must take a `program_state_t` as a reference to execute Brainfuck code. The program state structure definition is shown in listing 6.17.

The implementation principle for the code generation functions is to visit nodes recursively to compose and generate lambdas. We will explore two ways to implement the AST traversal: one based on a monolithic implementation based on function overloading to differentiate each type of node, and another one based on `if constexpr`.

- The overloaded version consists in a code generation entry point func-

tion that selects instruction-specific code generation functions depending on the type of token stored in the current instruction.

Listing 6.18: Overloaded `codegen` entry point function

```
// Necessary forward declaration
template <auto const &Ast,
          size_t InstructionPos = 0>
constexpr auto codegen();

// Specialization implementations...

/// Generic code generation entrypoint
template <auto const &Ast, size_t InstructionPos>
constexpr auto codegen() {
    constexpr flat_node_t Instr =
        Ast[InstructionPos];

    // Calling specialized codegen versions
    // using function overloading
    return codegen<Ast, InstructionPos>(<
        decltype(get<Instr.index>()(Instr)){});
}
```

Listing 6.18 shows the implementation of the entry point function whose only purpose is to unwrap the type of the current element and use it to call the appropriate code generation function depending on the instruction type. The dispatch is done automatically through function overloading.

Listing 6.19: `codegen` specialization for `flat_token_t` elements

```
/// Code generation implementation
/// for a single instruction
template <auto const &Ast,
          size_t InstructionPos = 0>
constexpr auto codegen(flat_token_t) {
    // Extracting token value
    constexpr flat_token_t Token =
        get<flat_token_t>(Ast[InstructionPos]);

    // Returning code for a single Brainfuck
    // instruction

    // >
    if constexpr (Token.token ==
        pointer_increase_v) {
```

```

    return [](program_state_t &s) { ++s.i; };
}
// <
else if constexpr (Token.token ==
                    pointer_decrease_v) {
    return [](program_state_t &s) { --s.i; };
}

// More instructions...
}

```

In listing 6.19 we can see the implementation of a code generation function for simple instructions, *i.e.* any instruction represented by a token that isn't a while block delimiter.

The dispatch over the tokens is done using `if constexpr`, but it could have been done using a variable template with a series of specializations.

Note that the token passed as a regular parameter is only used for the overload selection. The token value used to evaluate conditions in the `if constexpr` statements is sourced from the AST passed as a NTTP, therefore it is still `constexpr`.

Listing 6.20: codegen specialization for `flat_block_descriptor_t` elements

```

/// Code generation implementation
/// for a code block
template <auto const &Ast,
          size_t InstructionPos = 0>
constexpr auto codegen(flat_block_descriptor_t) {
    return [](program_state_t &s) {
        // Generating an index sequence type
        // with a size equal to the code block size.
        // It will be passed to the template lambda
        // to expand its indexes.
        auto index_sequence =
            std::make_index_sequence<
                get<flat_block_descriptor_t>(
                    Ast[InstructionPos])
                    .size>{};

        // Static unrolling on the block's
        // instructions, made possible by the
        // contiguity of its elements
        [&]<size_t... Indexes>(

```

```

        std::index_sequence<Indexes...>) {
        // Expansion on the index to generate code
        // for each node and invoke it with the
        // program state
        (... , codegen<Ast, 1 + InstructionPos +
                Indexes>()(s));
    }(index_sequence);
};
}

```

The trickiest part of code generation is generating code blocks, as shown in listing 6.20. Doing so requires the use of a compile time unrolling technique based on C++ parameter packs. Iteration over the elements must be done that way to keep the index `constexpr` and use it as a NTTP to generate code.

The result of this metafunction is an anonymous function that evaluates all the Brainfuck code within a block. At this point the only function that remains is the `flat_while_t` overload.

Listing 6.21: `codegen` specialization for `flat_while_t` elements

```

/// Code generation implementation
/// for a while block
template <auto const &Ast,
        size_t InstructionPos = 0>
constexpr auto codegen(flat_while_t) {
    return [](program_state_t &s) {
        while (s.data[s.i]) {
            codegen<Ast, get<flat_while_t>(
                Ast[InstructionPos])
                .block_begin>()(s);
        }
    };
}

```

The `codegen` implementation for a while block 6.21 is trivial: it returns a function that runs a while loop as defined by the Brainfuck language specification, and the body itself is the code generation result for the block element it points to.

This implementation is a good way to show how code generation from a NTTP can be implemented for a serialized AST.

- A monolithic implementation of `codegen` can be implemented by replacing overloading with `if constexpr` and `std::holds_alternative`.

Listing 6.22: `if constexpr` based codegen implementation

```
template <auto const &Ast,
          size_t InstructionPos = 0>
constexpr auto codegen() {
    constexpr flat_node_t Instr =
        Ast[InstructionPos];

    if constexpr (holds_alternative<flat_token_t>(
        Instr)) {
        constexpr flat_token_t Token =
            get<flat_token_t>(Instr);

        // Single token code generation...
    }

    else if constexpr (holds_alternative<
        flat_block_descriptor_t>(
        Instr)) {
        constexpr flat_block_descriptor_t
            BlockDescriptor =
                get<flat_block_descriptor_t>(Instr);
        // Block code generation...
    }

    else if constexpr (holds_alternative<
        flat_while_t>(Instr)) {
        constexpr flat_while_t While =
            get<flat_while_t>(Instr);
        // While loop code generation...
    }
}
```

Listing 6.22 shows the `if constexpr` based code generation implementation. Note that bits of code were cut to make the listing shorter, but they are the same as the overloaded `codegen` overload function bodies. This `codegen` implementation remains functionally identical to the previous one.

Difficulties

- Embedding text from the original string
- `if constexpr` requiring the definition of a `constexpr` variable

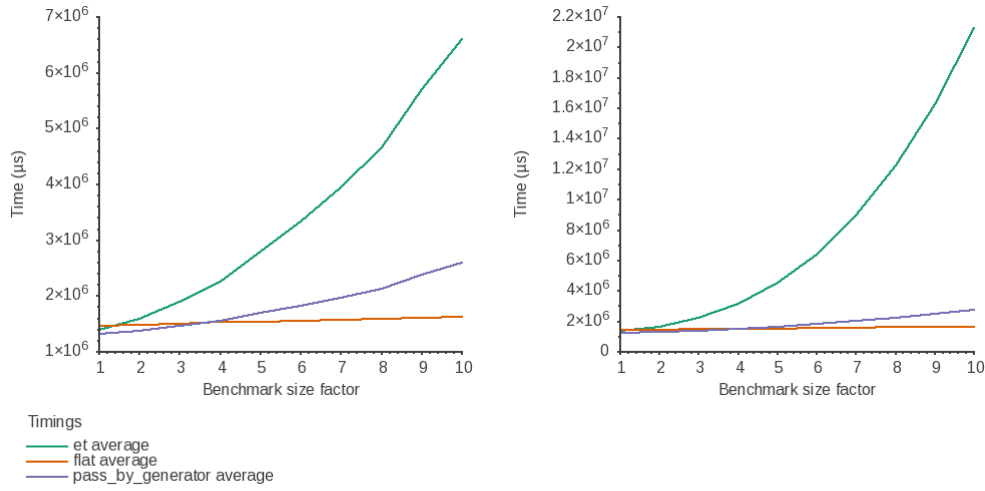


Figure 6.2: Compiler execution time measurements for consecutive loops (left) and nested loops (right)

Implementation complexity

Small synthetic variable-size benchmarks

We first begin by running two variable-sized benchmarks, consisting in measuring compiler execution time as the AST widens, and as the AST deepens.

The first variable-sized benchmark consists in generating a valid BF AST by concatenating strings to generate a succession of BF while loops in a `constexpr` string. This benchmark was instantiated with sizes going from 1 to 10 with a step of 1, with 10 timing iterations for each size.

The second benchmark generates a string with nested loops, making the AST deeper as the benchmark size increases instead of making it wider as in the previous case.

Both benchmarks generate programs of the same size so comparisons can be made properly.

Figure 6.2 both highlight considerably higher compiler execution times for the expression template based backend, high enough to suggest that the use of expression templates induces an overhead higher than parsing and generating Brainfuck programs using the PBG backend. However the PBG backend still has a compile time overhead much higher than the flat backend, which shows near constant compiler execution times on these small scale benchmarks.

Finally, AST deepening has a much higher impact on compile times than AST widening with the expression template backends, whereas the other backends seem to scale similarly as the AST grows wider or deeper.

Backend	Hello World	Hello World x2	Mandelbrot
Flat (Overloaded)	???	???	???
Flat (Monolithic)	???	???	???
Pass-by-generator	???	???	Failure (timeout)
Expression template	???	???	Failure (timeout)

Figure 6.3: BF compile time measurements in seconds

Large Brainfuck programs

The following benchmarks consist in measuring compiler execution times for compiling Brainfuck code examples. These example programs are also used to validate the metacompiler's backend implementations by compiling them and verifying their output.

- A Hello World program (106 tokens).
- The same Hello World program, ran twice (212 tokens).
- A Mandelbrot set fractal viewer (11672 tokens).

The measurements in figure 6.3 help us better understand how various metacompiling techniques behave at scale.

- **Flat backends:** they are not the easiest to implement due to the additional intermediate representation and serialization step that need to be implemented.

However, they present clear benefits when it comes to compilation times. So far they are the only ones that can be used at scale. The Mandelbrot example is supposed to illustrate an extreme case where DSEs are used to integrate large programs (approximately 11'000 AST nodes), and yet both "flat" implementations manage to keep compilation times well under a minute.

- **PBG:** The pass-by-generator backend has the shortest implementation. However, it might not be the least difficult to implement.

The code examples do not reflect the time spent debugging `constexpr` allocated memory errors. While the pass-by-generator may be a decent route for rapid prototyping to compile simple embedded expressions or programs, it is not fast enough for

- **ET:** This backend was originally meant to be a demonstrator for the interoperability of `constexpr` memory and ETs. It shows no particular advantage compared to the two other backends: it is slower than the PBG backend while requiring additional effort to implement the ET intermediate representation.

The "Flat" backend shows very good performance on all examples, including the Mandelbrot example that is about 100 times larger than the Hello World example. However the other cases highlight severe scaling issues and tend to confirm our previous hypothesis being that using generator functions to pass values makes the code generation quadratic. Finally, the ET backend performance highlights heavy performance impact when expression templates are being used, which is likely due to the complexity of the mechanisms expression templates involve like SFINAE and overload resolution.

Conclusions

6.5 . Math parsing and high performance code generation

6.5.1 . The Shunting-Yard algorithm

As seen in

6.5.2 . Code generation from a postfix math notation

6.5.3 . Using Blaze for high performance code generation

6.5.4 . Conclusion: a complete toolchain for High Performance Computing code generation from math formulas

6.6 .Conclusion

We wanted to demonstrate that using `constexpr` code to implement parsers for DSEL of arbitrary syntax in C++ 23 is possible despite limitations on `constexpr` memory allocation, and that doing so is possible with reasonable impact on compilation times.

We achieved that by implementing a `constexpr` parser for the Brainfuck language, with code generation backends implementing three different strategies to transform `constexpr` program representations into code using function generators, ETs, and non-type template parameters. We also demonstrated the interoperability of these `constexpr` parsers by implementing a parser for mathematical languages that can be used as a frontend for existing high performance C++ computation libraries.

Our benchmarks highlight compilation time scaling issues with pass-by-generator and ET code generation strategies for large programs, and excellent scaling capabilities for non-type template parameter based code generation strategies. These results can be used to decide which strategy to adopt for the implementation of future DSEL based on `constexpr` parsers based on considerations for compilation times or implementation complexity.

Going forward, `constexpr` parser generators could help reduce DSEL implementation time and help embed more languages into C++ 23. Further research has to be made to determine the impact of such generators on DSEL implementation complexity and compilation times.

7 -Appendix

.1 . Poacher

.1.1 . Brainfuck AST definition header

Listing 1: brainfuck/include/brainfuck/ast.hpp

```
#pragma once

#include <stdint>
#include <memory>
#include <vector>

namespace brainfuck {

//-----
// TOKEN TYPE

/// Represents a Brainfuck token.
enum token_t : char {
    pointer_increase_v = '>', // ++ptr;
    pointer_decrease_v = '<', // --ptr;
    pointee_increase_v = '+', // ++*ptr;
    pointee_decrease_v = '-', // --*ptr;
    put_v = '.', // putchar(*ptr);
    get_v = ',', // *ptr=getchar();
    while_begin_v = '[', // while (*ptr) {
    while_end_v = ']', // }
    nop_v, // nop
};

/// Converts a char into its corresponding
/// Brainfuck token_t value.
constexpr enum token_t to_token(char c) {
    switch (c) {
        case pointer_increase_v:
            return pointer_increase_v;
        case pointer_decrease_v:
            return pointer_decrease_v;
        case pointee_increase_v:
            return pointee_increase_v;
        case pointee_decrease_v:
            return pointee_decrease_v;
    }
}
```

```

    case put_v:
        return put_v;
    case get_v:
        return get_v;
    case while_begin_v:
        return while_begin_v;
    case while_end_v:
        return while_end_v;
    }
    return nop_v;
}

//-----
// AST

/// Holds the underlying node type
enum ast_node_kind_t : std::uint8_t {
    /// AST token node
    ast_token_v,
    /// AST block node
    ast_block_v,
    /// AST while node
    ast_while_v,
};

/// Parent class for any AST node type,
/// holds its type as an ast_node_kind_t
struct node_interface_t {
private:
    ast_node_kind_t kind_;

protected:
    constexpr node_interface_t(ast_node_kind_t kind)
        : kind_(kind){};

public:
    /// Returns the node kind tag.
    constexpr ast_node_kind_t get_kind() const {
        return kind_;
    }
    constexpr virtual ~node_interface_t() = default;
};

// Helpers

/// Token vector helper

```

```

using token_vec_t = std::vector<token_t>;

/// AST node pointer helper type
using ast_node_ptr_t =
    std::unique_ptr<node_interface_t>;

/// AST node vector helper type
using ast_node_vec_t = std::vector<ast_node_ptr_t>;

// !Helpers

/// AST node type for single Brainfuck tokens
struct ast_token_t : node_interface_t {
    token_t token;

    constexpr ast_token_t(token_t token_)
        : node_interface_t(ast_token_v), token(token_)
    {}
};

/// AST node type for Brainfuck code blocks
struct ast_block_t : node_interface_t {
    using node_ptr_t = ast_node_ptr_t;

    ast_node_vec_t content;

    constexpr ast_block_t(ast_node_vec_t &&content_)
        : node_interface_t(ast_block_v),
          content(std::move(content_)) {}

    constexpr ast_block_t(ast_block_t &&v) = default;
    constexpr ast_block_t &
    operator=(ast_block_t &&v) = default;

    constexpr ast_block_t(ast_block_t const &v) =
        delete;
    constexpr ast_block_t &
    operator=(ast_block_t const &v) = delete;
};

/// AST node type for Brainfuck while loop
struct ast_while_t : node_interface_t {
    ast_block_t block;

    constexpr ast_while_t(ast_block_t &&block_)

```

```

        : node_interface_t(ast_while_v),
          block(std::move(block_)) {}
};

template <typename F>
constexpr auto visit(F f, ast_node_ptr_t const &p) {
    switch (p->get_kind()) {
        case ast_token_v:
            return f(static_cast<ast_token_t const &>(*p));
        case ast_block_v:
            return f(static_cast<ast_block_t const &>(*p));
        case ast_while_v:
            return f(static_cast<ast_while_t const &>(*p));
    }
}

} // namespace brainfuck

```

.1.2 . Brainfuck parser implementation header

Listing 2: brainfuck/include/brainfuck/parser.hpp

```
#pragma once

#include <algorithm>

#include <iterator>
#include <memory>
#include <string>

#include <brainfuck/ast.hpp>

namespace brainfuck::parser {

    /// Parser implementation
    namespace impl {

        /// Converts a string into a list of BF tokens
        constexpr token_vec_t
        lex_tokens(std::string const& input) {
            token_vec_t result;
            result.reserve(input.size());
            std::transform(input.begin(), input.end(),
                           std::back_inserter(result),
                           [](auto current_character) {
                               return to_token(current_character);
                           });
            return result;
        }

        /// Parses BF code until the end of the block (or the
        /// end of the formula, ie. parse_end), then returns
        /// an iterator to the last parsed token or parse_end
        /// if the parser has parsed all the tokens.
        constexpr std::tuple<ast_block_t,
                             token_vec_t::const_iterator>
        parse_block(token_vec_t::const_iterator parse_begin,
                    token_vec_t::const_iterator parse_end) {
            using input_it_t = token_vec_t::const_iterator;

            ast_node_vec_t block_content;

            for (; parse_begin != parse_end; parse_begin++) {
                // While end bracket: return block content and
```

```

// while block end position
if (*parse_begin == while_end_v) {
    return {std::move(block_content), parse_begin};
}

// While begin bracket: recurse,
// then continue parsing from the end of the block
else if (*parse_begin == while_begin_v) {
    // Parse while body
    auto [while_block_content, while_block_end] =
        parse_block(parse_begin + 1, parse_end);

    block_content.push_back(
        std::make_unique<ast_while_t>(
            std::move(while_block_content)));

    parse_begin = while_block_end;
}

// Any other token that is not a nop instruction:
// add it to the AST
else if (*parse_begin != nop_v) {
    block_content.push_back(ast_node_ptr_t(
        std::make_unique<ast_token_t>(
            *parse_begin)));
}
}

return {ast_block_t(std::move(block_content)),
        parse_end};
}

} // namespace impl

/// Driver function for the token parser
constexpr ast_node_ptr_t
parse_ast(std::string const &input) {
    token_vec_t const tok = impl::lex_tokens(input);
    ast_block_t parse_result = get<ast_block_t>(
        impl::parse_block(tok.begin(), tok.end()));
    return std::make_unique<ast_block_t>(
        std::move(parse_result));
}

} // namespace brainfuck::parser

```