# Benchmarking Modern C++ Abstraction Penalty

February 25th, 2019

**Marcin Zalewski, Andrew Lumsdaine**

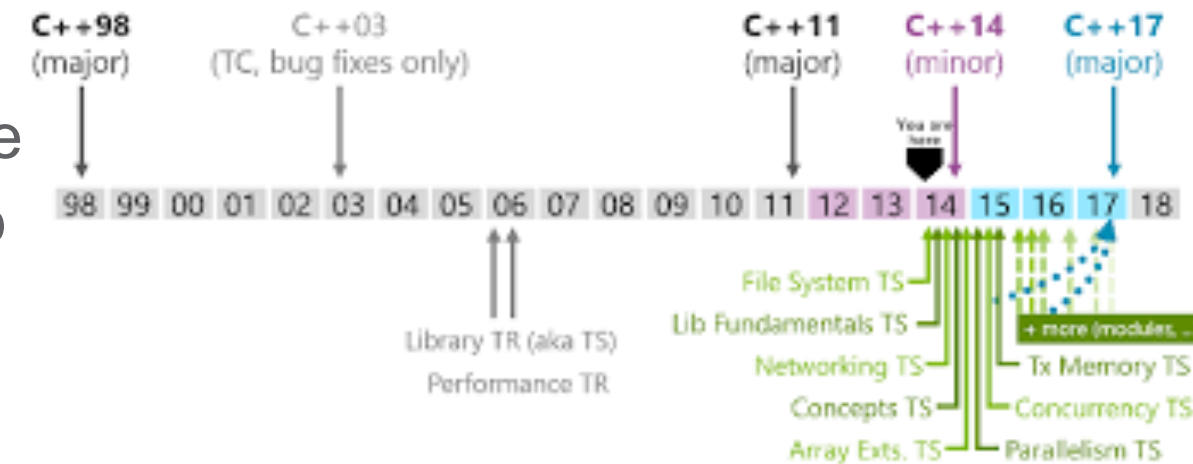# The Good, The Bad, and The Ugly

- Modern C++ is here!
  - C++98/C++03 was better than most
  - We waited for something new for a very long time
  - Now the compilers and the community caught up
  - Abstraction is (or should be) piling on

- However…
  - Have you ever heard how C must be faster because it is lower level?
  - More abstraction must mean that your code will be slower?
  - Of course, we C++ programmers know how wrong all of the above is*

- Oh no
  - Compilers vary widely
  - Who knows how much does abstraction cost?
  - How often do programmers think they are smarter then the compiler?

* It may take hours to compile

# So What Is That Abstraction You Are Talking About?

- Abstraction
  - Low-Level language constructs
    - ✓ Move semantics
    - ✓ Type deduction
    - ✓ Lambdas
    - ✓ Structured bindings
    - ✓ Range for loops
    - ✓ Variadic templates
    - ✓ Template aliases
    - ✓ Class enums
    - ✓ Static assert
  - Higher-level abstraction
    - ✓ Tuple types
    - ✓ Hash tables
    - ✓ Smart pointers
    - ✓ Threading
    - ✓ Polymorphic function object wrappers
    - ✓ Type traits for metaprogramming
  - DSLs
    - ✓ Exploit C++ features to create an embedded DSL
    - ✓ Library interfaces that closely correspond to the problem domain
    - ✓ Modern feel: use modern C++ features

# What Do We Want to Find Out?

- **Should we avoid abstraction?**
  - Ideal abstraction
  - May be pretty far from hand-optimized code
  - Often there is a temptation to pre-optimize
  - Lack of confidence in the compiler's ability
  - Compilation time
  - Compile-time errors

- How to measure abstraction penalty?
  - In general, this topic receives little attention
  - There are few benchmarks available
  - Compilers differ widely (even between versions)

- Let's look at some benchmarks

# Example: Abstract Sparse Matvec

- Let's start with a simple familiar loop nest:

```cpp
1  zeroize(y);
2
3  std::vector<size_t> indices(N*N+1);
4  std::vector<size_t> indexed(0);
5  std::vector<double> values(0);
6  piscetize(indices, indexed, values, N, N);
7
8  for (size_t i = 0; i < N*N; ++i) {
9    for (size_t j = indices[i]; j < indices[i+1]; ++j) {
10     y[i] += values[j] * x[indexed[j]];
11   }
12 }
```

# Matvec: What's The Intent?

- y=A*x: For every element A(i, j), multiply it with x[j] and accumulate to y[i]

```
1  edge_list<directed, double> A(0);
2  piscetize(A, N, N);
3  adjacency<0, double> C(A);
4
5  for (auto&& [i, j, v] : spmv_range(C)) {
6    y[i] += v * x[j];
7  }
```

- Can this possibly work well?

# How About This?

- Use std::for_each

```
1 edge_list<directed, double> A(0);
2 piscetize(A, N, N);
3 adjacency<0, double> C(A);
4
5 std::for_each(D.begin(), D.end(),
6   [&](std::tuple<size_t, size_t, double>&& a) {
7     y[std::get<0>(a)] += std::get<2>(a) * x[std::get<1>(a)];
8   }
9 );
```

# What is spmv_range?

```
1  v_range_iterator& operator++() {
2    ++u_begin;
3    if (u_begin == u_end) {
4      if (++first != last) {
5        u_begin = (*first).begin();
6        u_end = (*first).end();
7      }
8    }
9    return *this;
10 }
11
12 auto operator*() {
13   return std::tuple{first - the_range_.the_graph_.begin(),
14                     std::get<0>(*u_begin),
15                     std::get<1>(*u_begin)};
16 }
17
18 auto operator==(const v_range_iterator& b) const
19   { return first == b.first; }
```
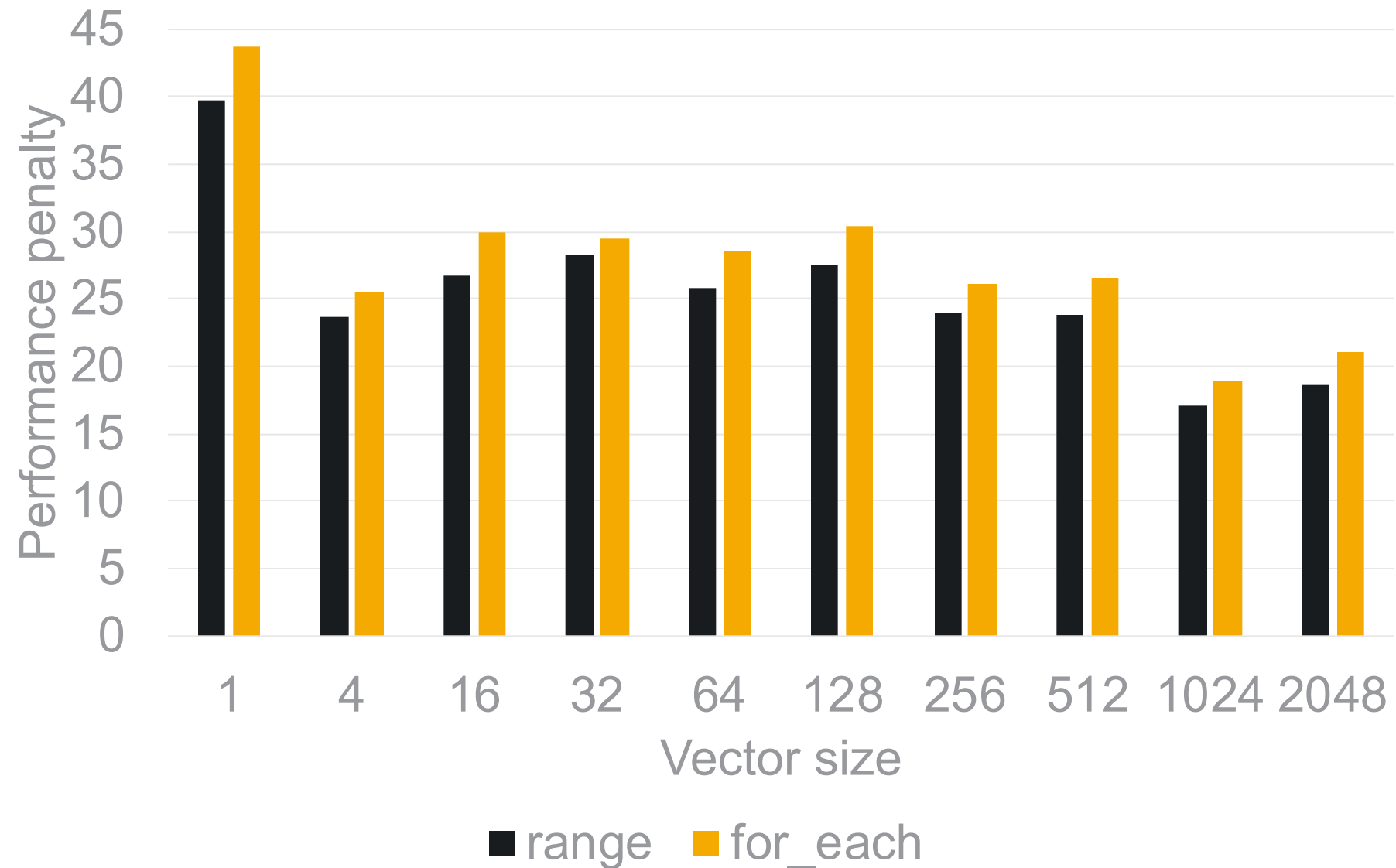
```
1  template<typename Graph>
2  class spmv_range {
3  public:
4    spmv_range(Graph& g) : the_graph_(g) {}
5    spmv_range(const spmv_range& b) : the_graph_(b.the_graph_) {}
6    auto& operator=(const spmv_range& b) {
7      the_graph_ = b.the_graph_;
8      return *this;
9    }
10
11   class v_range_iterator {
12   private:
13     spmv_range<Graph>&           the_range_;
14     typename Graph::outer_iterator first, last;
15     typename Graph::inner_iterator u_begin, u_end;
16
17   public:
18     v_range_iterator(spmv_range<Graph>& range, size_t offset = 0)
19       : the_range_{range}, first{the_range_.the_graph_.begin() + offset}, last{the_range_.the_graph_.end()} {
20       if (first != last) {
21         u_begin = (*first).begin();
22         u_end   = (*first).end();
23       }
24     }
25
26     v_range_iterator(const v_range_iterator& b)
27       : the_range_(b.the_range_), first(b.first), last(b.last), u_begin(b.u_begin), u_end(b.u_end) {}
28
29     v_range_iterator& operator++() {
30       ++u_begin;
31       if (u_begin == u_end) {
32         if (++first != last) {
33           u_begin =
34             (*first).begin();    // FIXME:  This was commented out -- why?  Might not need it for CSR -- but in general, we do
35           u_end = (*first).end();
36         }
37       }
38
39       return *this;
40     }
41
42     // auto operator*() { return std::tuple<vertex_index_t, vertex_index_t&&>(last - first, std::get<0>(*u_begin)); }
43
44     auto operator*() {
45       return std::tuple<typename std::iterator_traits<typename Graph::outer_iterator>::difference_type, vertex_index_t, double&>(
46         first - the_range_.the_graph_.begin(), std::get<0>(*u_begin), std::get<1>(*u_begin));
47     }
48
49     auto operator==(const v_range_iterator& b) const { return first == b.first; }
50     auto operator!=(const v_range_iterator& b) const { return first != b.first; }
51
52     auto operator<(const v_range_iterator& b) const { return first < b.first; }
53     auto operator-(const v_range_iterator& b) const { return first - b.first; }
54     auto operator+(unsigned long step) const { return v_range_iterator(the_range_, step); }
55
56     auto& operator=(const v_range_iterator& b) {
57       the_range_ = b.the_range_;
58       first      = b.first;
59       last       = b.last;
60       u_begin    = b.u_begin;
61       u_end      = b.u_end;
62       return *this;
63     }
64
65     //    typedef std::random_access_iterator_tag        iterator_category;
66     typedef std::forward_iterator_tag              iterator_category;
67     typedef size_t                                 difference_type;
68     typedef std::tuple<vertex_index_t, vertex_index_t> value_type;
69     typedef value_type*                            pointer;
70     typedef value_type&                            reference;
71   };
72
73   typedef v_range_iterator iterator;
74
75   auto begin() { return v_range_iterator(*this); }
76   auto end() { return v_range_iterator(*this, the_graph_.size()); }
77
78 private:
79   Graph& the_graph_;
80 };
```
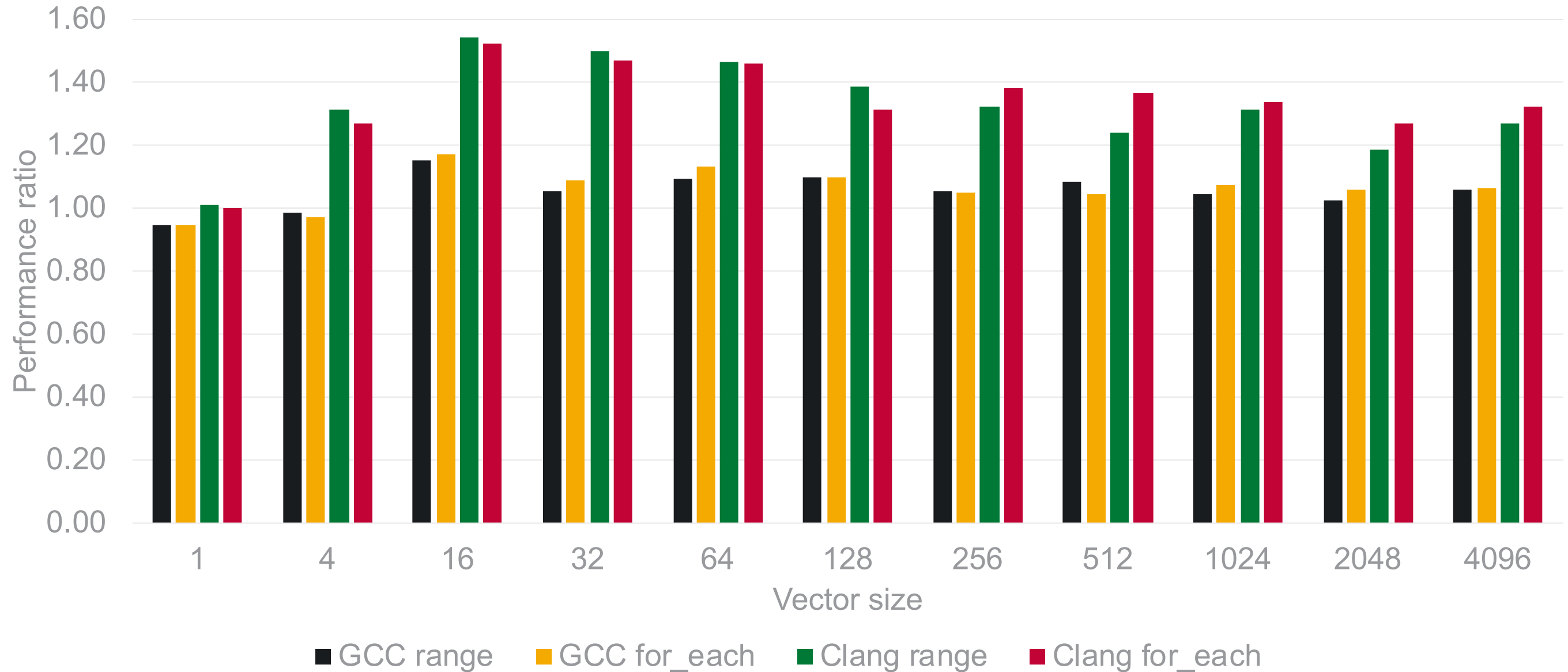
none
8

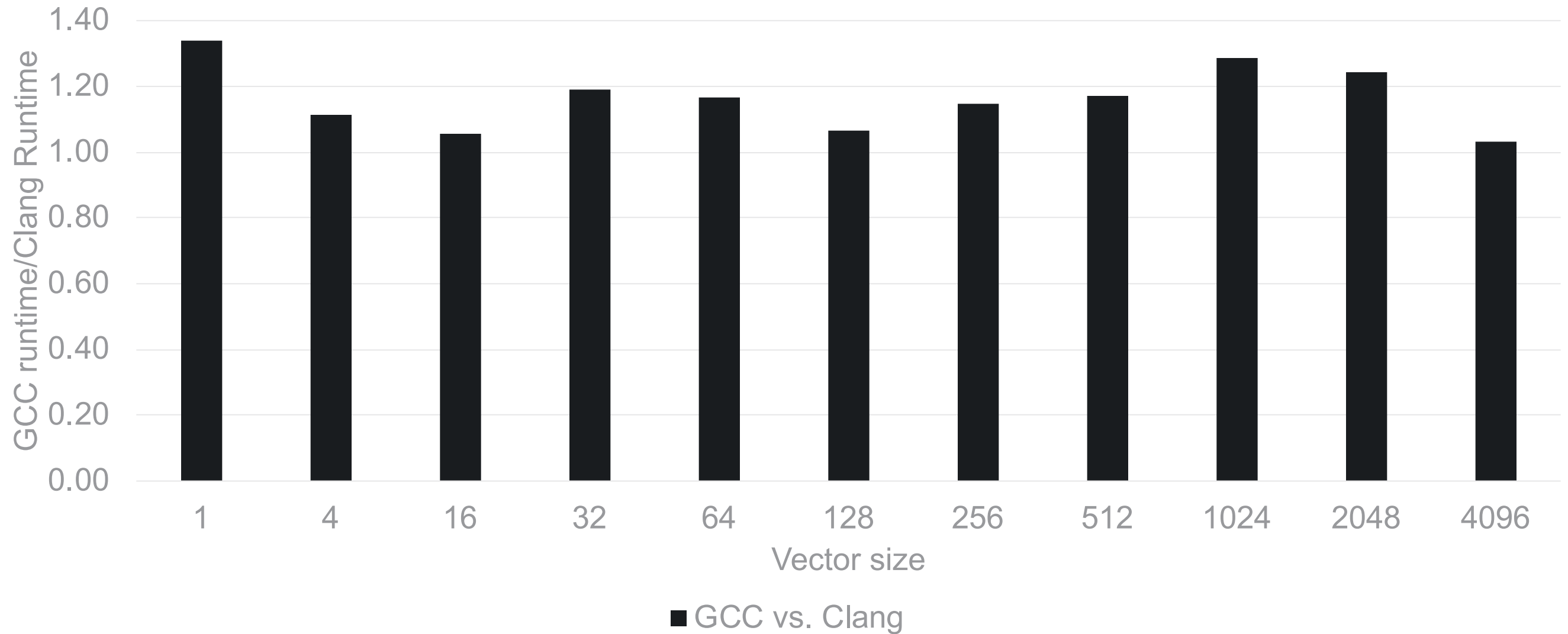# So How Does It Fare?



Unoptimized code penalty

# Compilers Can Do Amazing Things

## Performance Penalty

Legend: GCC range, GCC for_each, Clang range, Clang for_each

X-axis: Vector size (1, 4, 16, 32, 64, 128, 256, 512, 1024, 2048, 4096)
Y-axis: Performance ratio (0.00 to 1.60)

# Clang vs. GCC



GCC vs. Clang

# Matvec: Abstraction For The Win

- While there is a small performance loss, the abstraction is very close

- Writing loops by hand is much more error prone

- Saving time in productiveness may be worth a small performance hit
  - Future compilers should do better or at least not worse
  - Future readers of the application code will be grateful
  - The library on the other hand is a bit more complicated

- Optimization takes time: O3 is 3-4x slower to compile than O0, but let's not worry about it
  - Improvements in the language (e.g., modules)
  - Improvements in compiler technology

# Another Example: Tensor Template Library (TTL)

- Joint work with Luke Dalessandro and Alexander Winter

```
1  Tensor<4,3,double> A;
2  A[1][2][0][1] = 1;
3  double a[3][3][3][3];
4  a[1][2][0][1] = 1;
```

```
1  constexpr const Index<'i'> i;
2  constexpr const Index<'j'> j;
3  constexpr const Index<'k'> k;
4  constexpr const Index<'l'> l;
5
6  Tensor<2,2,double> A{}, B{};
7  Tensor<4,2,double> C;
8
9  C(i,j,k,l) = 1.9 * A(i,j) * B(k,l);
```

```
1  Tensor<2,3,double> A{}, B{};
2  Tensor<2,3,double> C;
3
4  auto AxB = A(i,j) * B(j,k);
5  auto BxA = B(i,j) * A(j,k);
6
7  C(i,k) = 0.5 * (AxB + BxA);
8
9  for (int i = 0; i < 3; ++i) {
10    for (int k = 0; k < 3; k++) {
11      C(i,k) = 0.0;
12      for (int j = 0; i < 3; ++j) {
13        C(i,k) += 0.5 * (A(i,j) * B(j,k) +
14                         B(i,j) * A(j,k));
15      }
16    }
17  }
```

# Benchmarking TTL

- Inner products

```
1 cArr[ii]() = A(i) * B(i);
```

```
1 for(size_t aa = 0; aa < iter; aa++){
2   for (size_t ii = 0; ii < D; ii++) {
3     arrFL[aa] += arrA[ii] * arrB[ii];
4   }
5 }
```

```
1 for(size_t aa = 0; aa < iter; aa++){
2   for (size_t ii = 0; ii < D; ii++){
3     for(size_t jj = 0; jj < D; jj++){
4       for(size_t kk = 0; kk < D; kk++){
5         for(size_t ll = 0; ll < D; ll++){
6           arrFL[aa] += arrA[ii][jj][kk][ll] *
7                        arrB[ii][jj][kk][ll];
8         }
9       }
10      }
11    }
12 }
```

```
1 cArr[ii]() = A(i,j,k,l) * B(i,j,k,l);
```

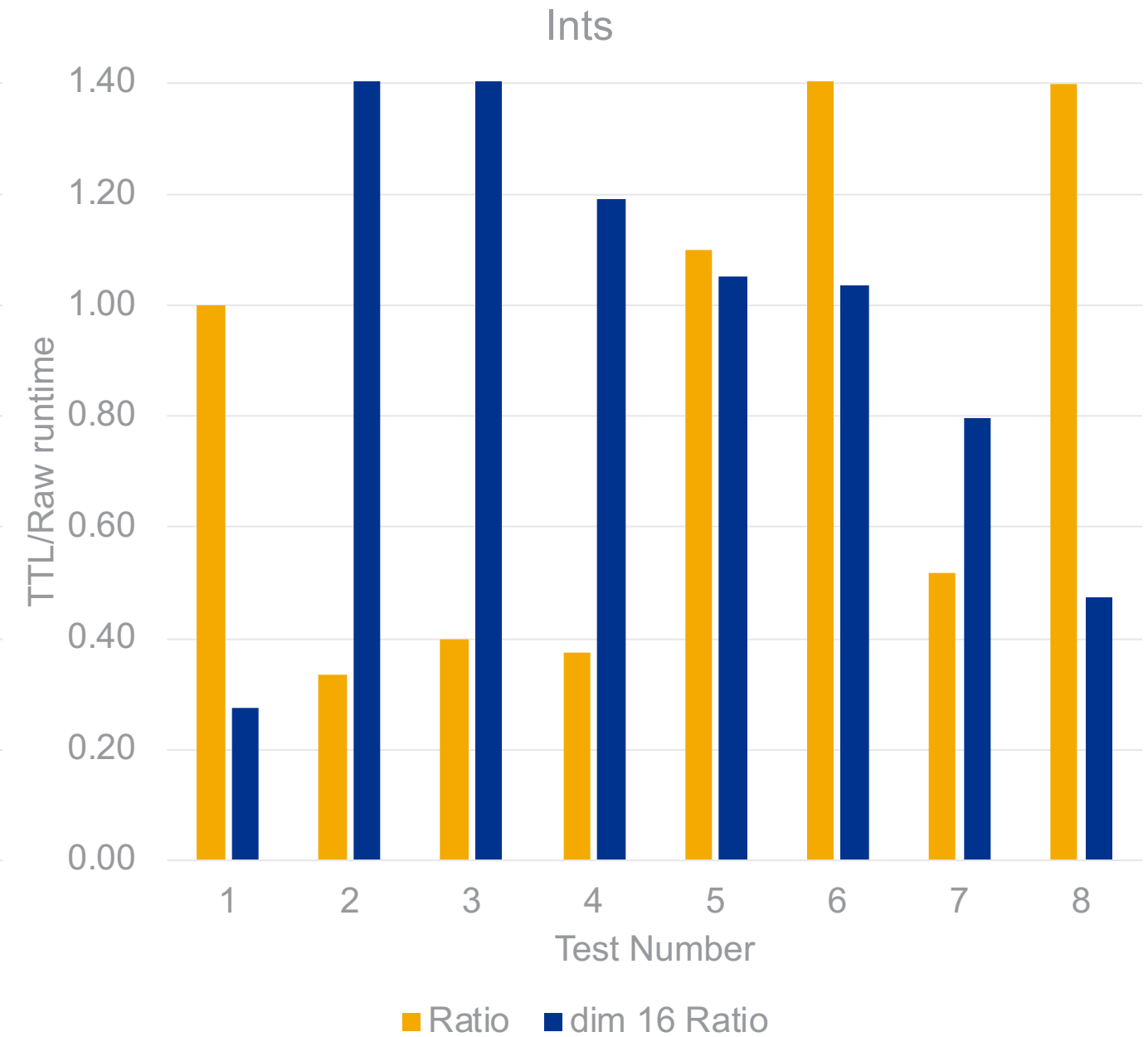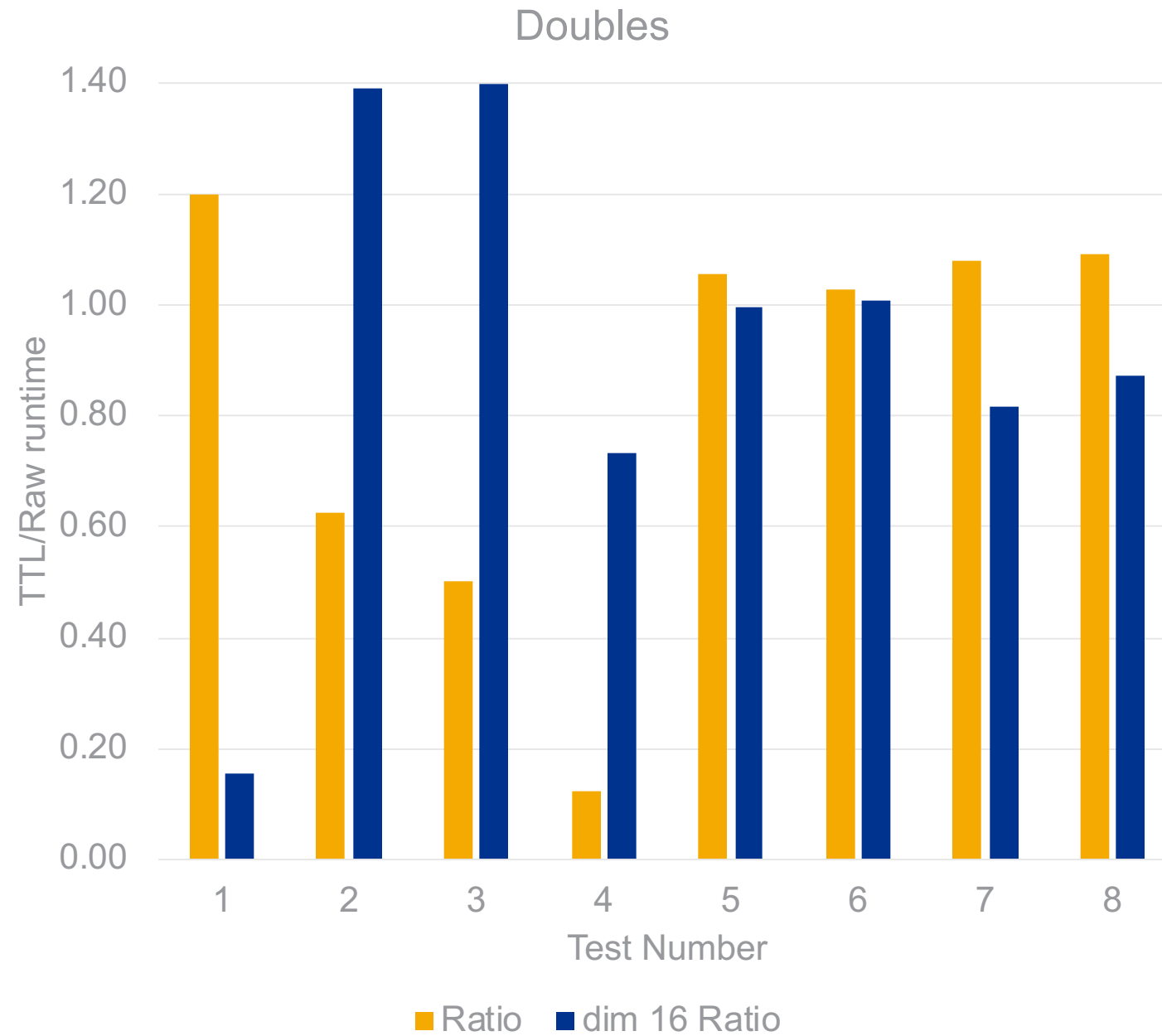# Benchmarking TTL

- Outer products

```
1  cArr[ii](i,j,k,l) = A(i,j) * B(k,l);
```

```
1  for(size_t aa = 0; aa < iter; aa++){
2    for(size_t ii = 0; ii < D; ii++){
3      for(size_t jj = 0; jj < D; jj++){
4        for(size_t kk = 0; kk < D; kk++){
5          for(size_t ll = 0; ll < D; ll++)
6            arrFL[aa][ii][jj][kk][ll] =
7              arrA[ii][jj] * arrB[kk][ll];
8          }
9        }
10     }
11 }
```

# TTL Benchmark Results

```
1  Operation         Rank      Format   Runtime (s)
2  Dimension: 2
3  Inner Product     1         TTL      6e-06
4  Inner Product     1         Raw      5e-06
5  Inner Product     2         TTL      5e-06
6  Inner Product     2         raw      8e-06
7  Inner Product     3         TTL      5e-06
8  Inner Product     3         raw      1e-05
9  Inner Product     4         TTL      4e-06
10 Inner Product     4         raw      3.2e-05
11 Outer Product     1->2      TTL      1.9e-05
12 Outer Product     1->2      Raw      1.8e-05
13 Outer Product     1,2->3    TTL      3.9e-05
14 Outer Product     1,2->3    Raw      3.8e-05
15 Outer Product     1s->3     TTL      4.1e-05
16 Outer Product     1s->3     Raw      3.8e-05
17 Outer Product     2->4      TTL      8.2e-05
18 Outer Product     2->4      raw      7.5e-05
```
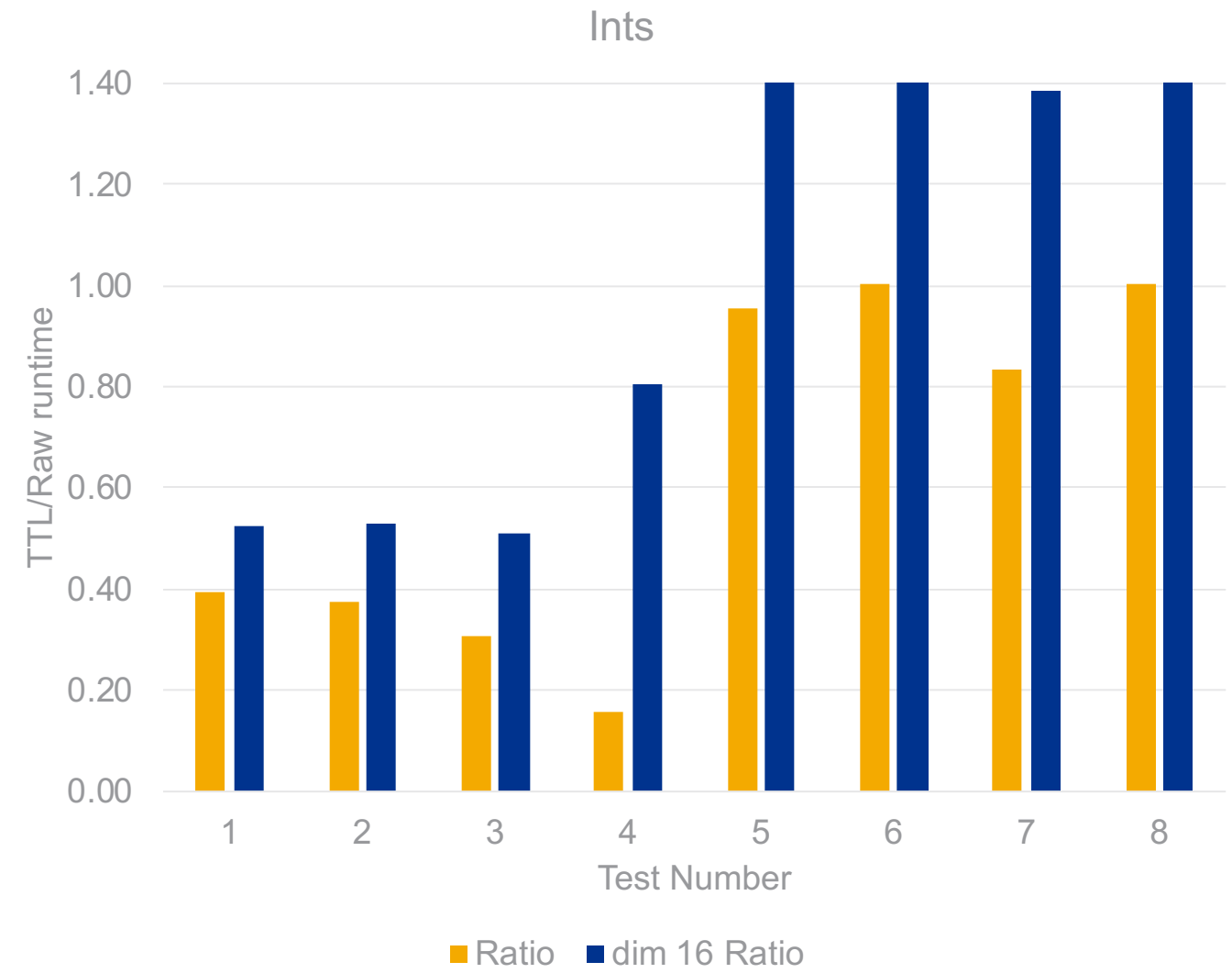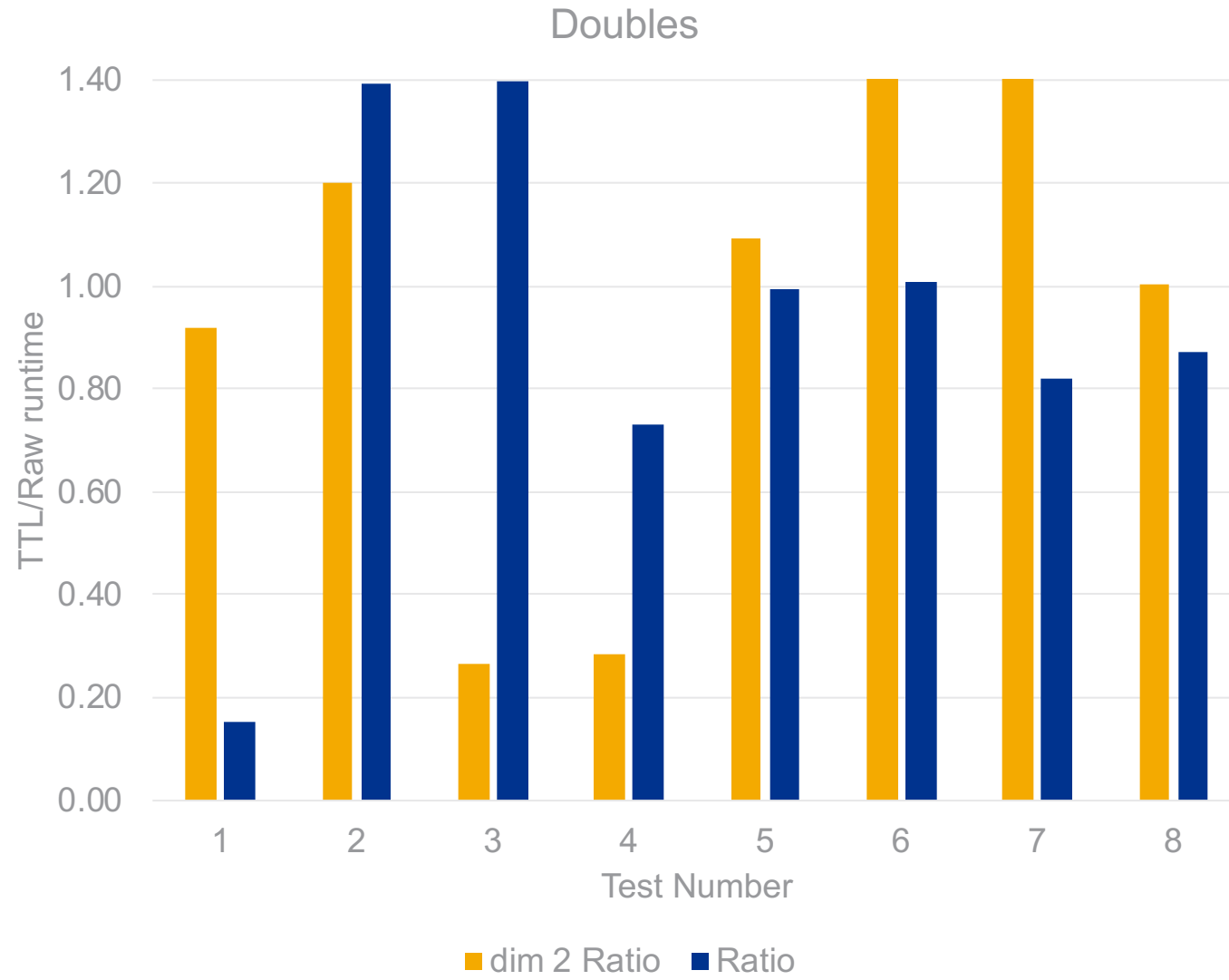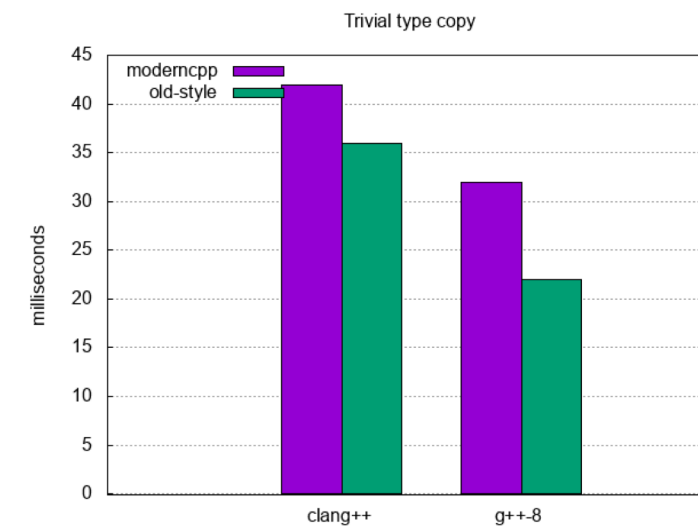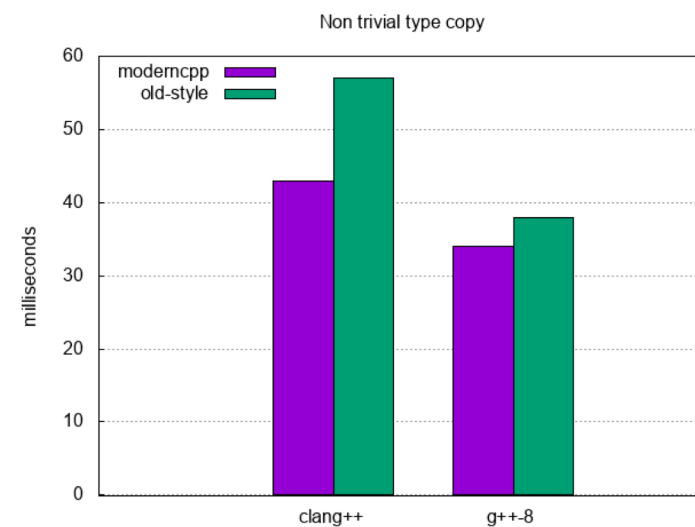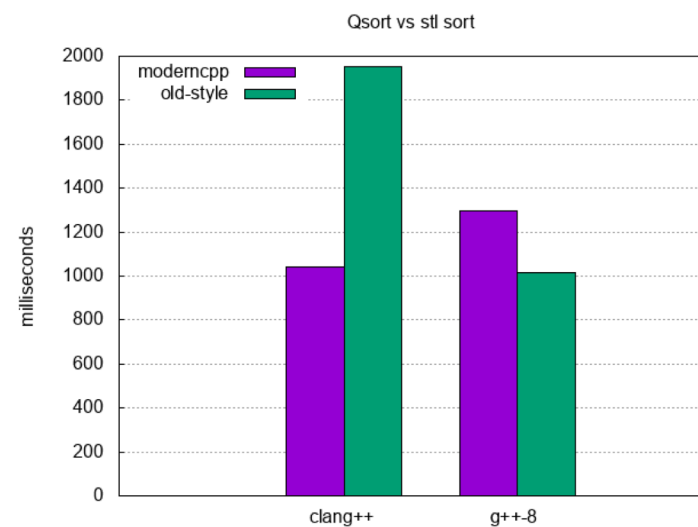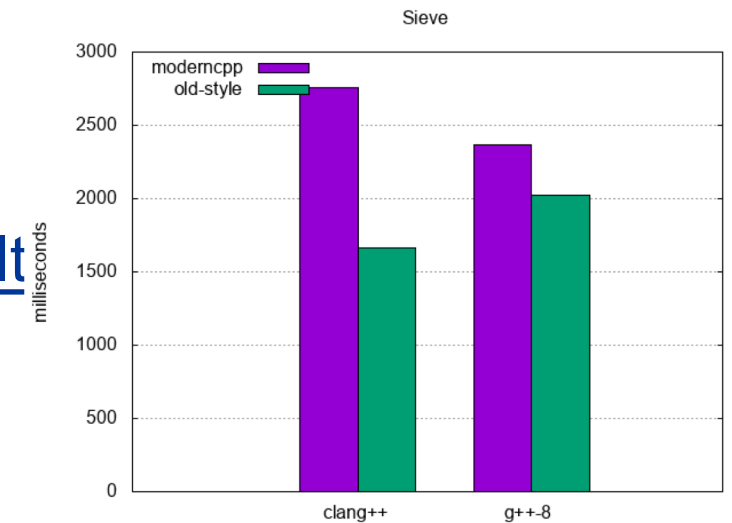
# TTL Benchmark Results



GCC

# TTL Benchmark Results



Doubles

Ints

Clang

# Other Benchmarks


Formatted read


Sieve

- Chris Cox's CppPerformanceBenchmarks
  - https://gitlab.com/chriscox/CppPerformanceBenchmarks
  - Lots of good tests for low-level language functionality
- German Diago Gomez's the-cpp-abstraction-penalty
  - https://github.com/germandiagogomez/the-cpp-abstraction-penalt
  - Fewer benchmarks, but tests more algorithms


Qsort vs stl sort


Non trivial type copy


Trivial type copy

# Don't Give Up: Use Abstraction Unless Proven Wrong

- Abstraction penalty is tough to quantify

- In general, trust the compiler and code the "ideal version"
  - Code readability is worth it

- Identify problems based on need
  - Profile
  - Benchmark abstraction vs. hand written code
  - Unfortunately, maybe look at assembly
  - Maybe the abstraction can be implemented better if it is lagging

- There are very few concentrated efforts to keep current information

- Submit your benchmarks or start your own collection

- Big question: GPU

# Thank you

# Bonus

COMPILER EXPLORER

Add... ▾  More ▾

Support diversity in C++ with #include <C++> ✕

C++ source #1 ✕                                    x86-64 gcc 8.2 (Editor #1, Compiler #1) C++ ✕

A ▾  Save/Load  + Add new... ▾  CppInsights        C++ ▾     x86-64 gcc 8.2 ▾  ✓  -O3 -g -std=c++17

A ▾  □ 11010  ☑ .LX0:  □ lib.f:  ☑ .text  ☑ //  □ \s+  □ Intel  ☑ Demangle

□ Libraries ▾  + Add new... ▾  ⚙ Add tool... ▾

```cpp
1   #include <optional>
2
3   template<class InputIt, class OutputIt>
4   OutputIt copy0(InputIt first, InputIt last, OutputIt out)
5   {
6       while (first != last)
7       {
8           *out++ = *first++;
9       }
10      return out;
11  }
12
13  template<class InputIt, class Sink>
14  Sink copy1(InputIt first, InputIt last, Sink sink)
15  {
16      while (first != last)
17      {
18          sink(*first++);
19      }
20      return sink;
21  }
22
23  template<class Source, class Sink>
24  Sink copy2(Source source, Sink sink)
25  {
26      while (auto val = source())
27      {
28          sink(*val);
29      }
30      return sink;
31  }
32
33  extern char * y;
34
35  int main()
36  {
37      char x[10000] = {0};
38
39      copy0(&x[0],&x[10000], y);
40      copy1(&x[0],&x[10000], [p=y](int v) mutable {*p++ = v;});
41      copy2([p=&x[0], i=0]()mutable {return i < 10000 ? std::optional<char>(p[i++]) :
42              [p=y](int v) mutable {*p++ = v;});
43  }
```

```asm
1   main:
2       subq    $10008, %rsp
3       movl    $10000, %edx
4       xorl    %esi, %esi
5       movq    %rsp, %rdi
6       call    memset
7       movq    y(%rip), %rdi
8       movq    %rsp, %rsi
9       movl    $10000, %edx
10      call    memcpy
11      movq    y(%rip), %rdi
12      movq    %rsp, %rsi
13      movl    $10000, %edx
14      call    memcpy
15      movq    y(%rip), %rdi
16      movq    %rsp, %rsi
17      movl    $10000, %edx
18      call    memcpy
19      xorl    %eax, %eax
20      addq    $10008, %rsp
21      ret
```