Portable C++ for Modern, Heterogenous Mixed-Precision Methods

Neil Lindquist

SIAM CSE23

March 1, 2023







Modern Numerical Methods

Precision

- Double
- Single
- Half
- Bfloat16
- Integer?
- Posits?

Hardware

- CPU
- GPUs
 - NVIDIA
 - AMD
 - Intel

Algorithm variants

- GMRES vs FGMRES
- Real vs Complex



Mixed Precision Experiments

- Many combinations of precisions to test
 - My GMRES work: 14 combinations in paper
 - Only single, double
 - GMRES-IR5¹: 112 "meaningful" combinations
 - Bfloat16, half, single, double, quad



Implementing

- Templates allow generic implementation
 - Portability layers, e.g., Kokkos
- Ideally: retain performance of separate codes
 - E.g., avoid "conversions" in uniform-precision
- Readable and portable
- → C++ features



- Sometimes need array in high and low prec.
 - In uniform-precision, just use reference
- Easy for types with reference semantics
 - E.g., Kokkos View
 - Just need different constructor



```
using hi t = ...; using lo t = ...;
Kokkos::View<hi_t, ...> w_hi = ...;
Kokkos::View<lo t, ...> w_lo;
if (std::is same v<hi t, lo t>) {
  w_lo = w_hi; 		 Error when hi t≠lo t
} else {
  w lo = Kokkos::View<lo t, ...>(w hi.extent(0));
```



Compile time if statement

- if constexpr (condition)
 - Since C++17
- Condition must be constexpr
- No template substitution for untaken branches
 - Can use type-specific functions



```
using hi t = ...; using lo t = ...;
Kokkos::View<hi_t, ...> w_hi = ...;
Kokkos::View<lo t, ...> w_lo;
if constexpr (std::is same v<hi t, lo t>) {
 w lo = w hi;
} else {
  w lo = Kokkos::View<lo t, ...>(w hi .extent(0));
```



Compile time if statement

```
if constexpr (is_complex<T>) {
  phase = val / std::abs(val);
} else {
  phase = val >= T(0) ? 1 : -1;
```



- Sometimes need array in high and low prec.
 - In uniform-precision, just use reference
- Harder for types with value semantics
 - Want to avoid copy operator
 - → conditionally make it a reference



```
using hi_t = ...; using lo_t = ...;
constexpr bool same_type_p = std::is_same_v<hi_t, lo_t>;
using vec lo t = std::conditional t<same type p,
                               Vect<hi t>&, Vect<lo t>>;
Vect<hi t> w hi = ...;
if constexpr (same_type_p){
  w lo = w hi;
} else {

    Move to helper function

  w_lo = Vect<lo_t>(w_hi.n());
```



Conditional types for real vs complex

- In rank revealing factorizations
 - SVD gives real-valued singular values
 - Pivoted QR gives complex-valued estimates
 - std::conditional_t allows same postprocessing code



Optional arguments

- Optional keyword arguments
 - Easily add new functionality w/out new routines
 - Users can focus on relevant args
 - Name-value pairs help clarity
- C++ only supports this via structs



Optional Arguments in SLATE

using namespace slate;

```
Options opts;
opts[Option::Target] = Target::Devices;
opts[Option::Lookahead] = 2;
opts[Option::PivotThreshold] = 0.1;
lu_solve(A, b, opts);
```



Optional Arguments in SLATE

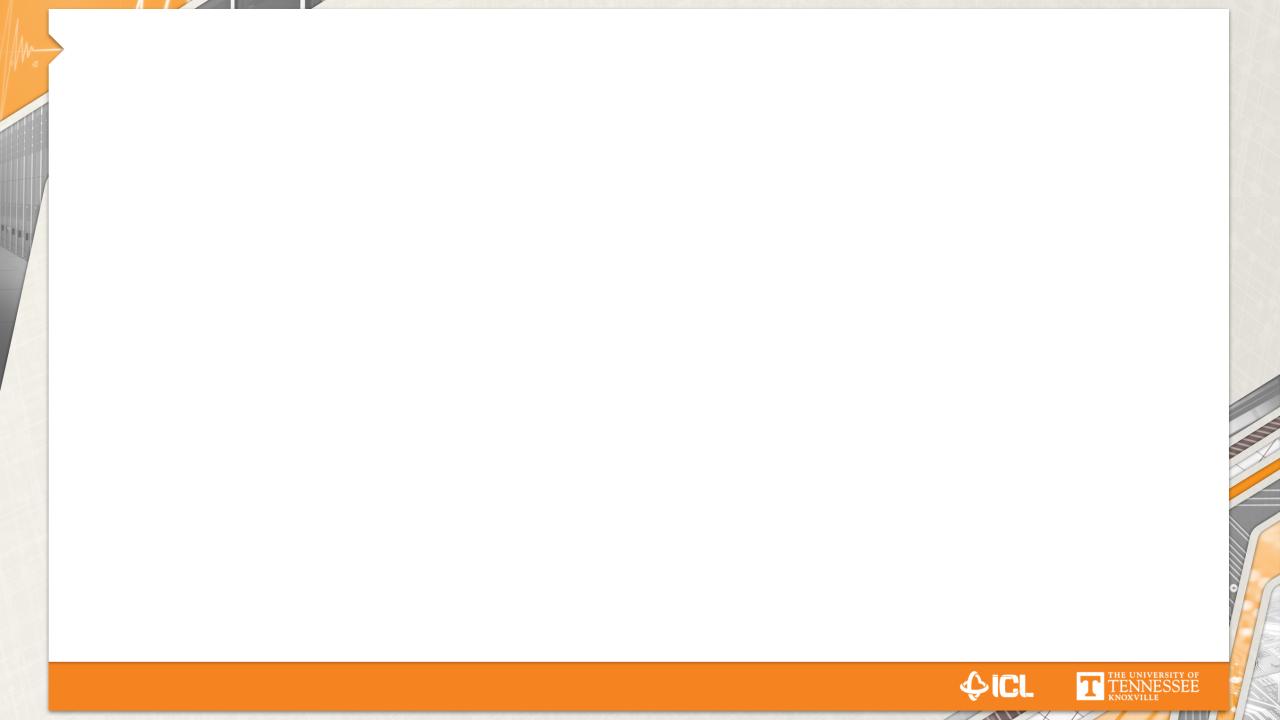
• C++ list initialization give better API using namespace slate;



Conclusions

- Combinatorial explosions of implementations
 - Precisions, accelerators, algorithmic variants
- Modern C++ features help
 - Reduce code duplication
 - Retain performance





```
template<class T, bool useFirst, class U, class V>
T make from selector(U arg1, V arg2) {
   // U, V are std::tuple<...>
    if constexpr (useFirst) {
        return std::make from tuple<T>(arg1);
    } else {
        return std::make from tuple<T>(arg2);
```

```
using hi t = ...; using lo t = ...;
constexpr bool same type p = std::is same v<hi t, lo t>;
using vec lo t = std::conditional t<same type p,
                               Vect<hi t>&, Vect<lo t>>;
Vect<hi t> w hi = ...;
vec lo t w lo = make_from_selector
                         <vec_lo_t, same_type_p>
                         (std::tuple(w), std::tuple(n));
```

Caveats

- Too many templates
 - Increases compile times
 - Increases mental burden
 - Results in poor compiler error messages
- Can't easily instantiate many combinations of template parameters
 - Often results in larger headers/repeated compilation

