# SIMD, stdx::simd, vir-simd, and merged GR4 blocks

maximize throughput, minimize latency, express parallelism via data-parallel types

Dr. Matthias Kretz

**GSI** Helmholtz Center for Heavy Ion Research

European GNU Radio Days '24

 @mkretz@floss.social

 github.com/mattkretz

## Goals and non-goals for this talk

- Introduce SIMD
- Introduce data-parallel types: `std::experimental::simd`[1]
- Introduce extensions: vir-simd
- Walk through a simple signal-processing example
- This talk wants to be a starting point for you to use `stdx::simd` in a GR4 block

---

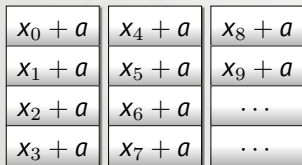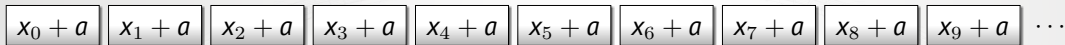[1]From now on shortened to: `stdx::simd` — via `namespace stdx = std::experimental;`

Motivation

Motivation
○●○○○○○○○○
stdx::simd Overview
○○○○○○○○○○○○○
Example: Signal Processing
○○○○○○○○○
SIMD and Merged Blocks
○○○○○○○○○
Outlook
○○○○
Summary
○

GSI

# SIMD – Single Instruction Multiple Data

in other words...

- multiple operations in one instruction, or
- execute the same work in less time



$$x_0 + a \quad x_1 + a \quad x_2 + a \quad x_3 + a \quad x_4 + a \quad x_5 + a \quad x_6 + a \quad x_7 + a \quad x_8 + a \quad x_9 + a \quad \cdots$$

Time

$$\begin{array}{ccc}
x_0 + a & x_4 + a & x_8 + a \\
x_1 + a & x_5 + a & x_9 + a \\
x_2 + a & x_6 + a & \cdots \\
x_3 + a & x_7 + a & \cdots
\end{array}$$

Time

Motivation | stdx::simd Overview | Example: Signal Processing | SIMD and Merged Blocks | Outlook | Summary
○○●○○○○○○○ | ○○○○○○○○○○○○ | ○○○○○○○○○ | ○○○○○○○○○ | ○○○○ | ○

GSI

## ILP — Instruction Level Paralellism

in other words...

- multiple instructions in one CPU cycle, or
- execute the same work in even less time



$$
\begin{array}{|c|c|c|}
\hline
x_0 + a & x_8 + a & x_{16} + a \\
x_1 + a & x_9 + a & x_{17} + a \\
x_2 + a & x_{10} + a & x_{18} + a \\
x_3 + a & x_{11} + a & x_{19} + a \\
\hline
x_4 + a & x_{12} + a & x_{20} + a \\
x_5 + a & x_{13} + a & x_{21} + a \\
x_6 + a & x_{14} + a & x_{22} + a \\
x_7 + a & x_{15} + a & \cdots \\
\hline
\end{array}
$$

$\longrightarrow$
Time

Take-Away #1

### `stdx::simd` expresses data-parallelism — SIMD and ILP

again, `stdx` is short for `std::experimental`

Example: Reaching Peak FLOP

Compiler Input:   sequential scalar code
Compiler Output:   no SIMD, no/little ILP

Motivation
○○○○○○●○○○

stdx::simd Overview
○○○○○○○○○○○○

Example: Signal Processing
○○○○○○○○○○

SIMD and Merged Blocks
○○○○○○○○○

Outlook
○○○○

Summary
○

GSI

## single-precision multiply-add
Linux, GCC 13, Intel Xeon W-2145 (2 AVX-512 FMA ports), C++26 `std::simd` prototype

```
1  void peak(benchmark::State& state) {
2    float x = 0.f;
3    fake_modify(x);
4    for (auto _ : state) {
5      x = x * 3.f + 1.f;
6    }
7    fake_read(x);
8  }
```

```
1  void peak(benchmark::State& state) {
2    simd<float, simd<float>::size() * 8> x = 0.f;
3    fake_modify(x);
4    for (auto _ : state) {
5      x = x * 3.f + 1.f;
6    }
7    fake_read(x);
8  }
```

|        | g++ -O3 -DNDEBUG | g++ -O3 -DNDEBUG -march=native |
|--------|------------------|-------------------------------|
| scalar | 0.25 FLOP/cycle  | 0.5 FLOP/cycle                |

Motivation
○○○○○○●○○○
stdx::simd Overview
○○○○○○○○○○○○○
Example: Signal Processing
○○○○○○○○○○
SIMD and Merged Blocks
○○○○○○○○○
Outlook
○○○○
Summary
○
GSI

## single-precision multiply-add
Linux, GCC 13, Intel Xeon W-2145 (2 AVX-512 FMA ports), C++26 `std::simd` prototype

```
1  void peak(benchmark::State& state) {
2    float x = 0.f;
3    fake_modify(x);
4    for (auto _ : state) {
5      x = x * 3.f + 1.f;
6    }
7    fake_read(x);
8  }
```

```
1  void peak(benchmark::State& state) {
2    simd<float, simd<float>::size() * 8> x = 0.f;
3    fake_modify(x);
4    for (auto _ : state) {
5      x = x * 3.f + 1.f;
6    }
7    fake_read(x);
8  }
```

|                | g++ -O3 -DNDEBUG | g++ -O3 -DNDEBUG -march=native |
|---------------:|:----------------:|:-----------------------------:|
| scalar         | 0.25 FLOP/cycle  | 0.5 FLOP/cycle                |
| data-parallel  | 8 FLOP/cycle     | 64 FLOP/cycle                 |

$\frac{64}{0.25} = 256$ times faster

A data-parallel type wider than the default can increase ILP!

## single-precision multiply-add
Linux, GCC 13, Intel Xeon W-2145 (2 AVX-512 FMA ports), C++26 `std::simd` prototype

```
1  void peak(benchmark::State& state) {
2    float x = 0.f;
3    fake_modify(x);
4    for (auto _ : state) {
5      x = x * 3.f + 1.f;
6    }
7    fake_read(x);
8  }
```

```
1  void peak(benchmark::State& state) {
2    simd<float, simd<float>::size() * 8> x = 0.f;
3    fake_modify(x);
4    for (auto _ : state) {
5      x = x * 3.f + 1.f;
6    }
7    fake_read(x);
8  }
```
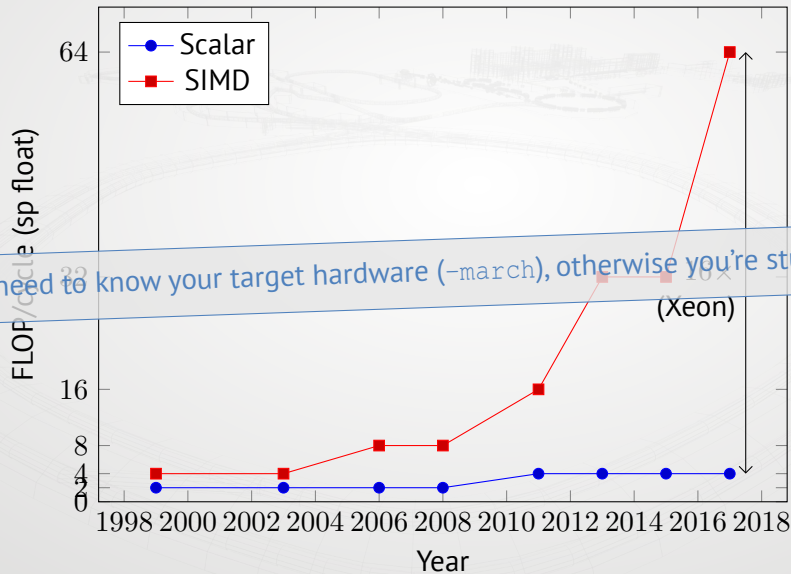
The naive benchmark was slow because the code didn't contain any parallelism: one long dependency chain

|               | g++ -O3 -DNDEBUG | g++ -O3 -DNDEBUG -march=native |
|---------------|------------------|--------------------------------|
| scalar        | 0.25 FLOP/cycle  | 0.5 FLOP/cycle                 |
| data-parallel | 8 FLOP/cycle     | 64 FLOP/cycle                  |

$\frac{64}{0.25} = 256$ times faster

A data-parallel type wider than the default can increase ILP!

x86

Compilers need to know your target hardware (−march), otherwise you're stuck in 2006!

## Our non-parallel reality

We write and use code as impossible to optimize as this naive benchmark all the time!

- Example: `float std::cos(float)`
- This interface is bad for performance, so compilers replace `std::cos` calls with `__builtin_cosf`.
  - allows compile-time evaluation for constant inputs
  - enables vectorization if the caller does multiple calls
- Compilers will not be modified to replace your library functions with compiler builtins.

🤔

Calling a function (over library boundaries) with a single input/output value inhibits parallelization (SIMD & ILP).

# An alternative

```
std::simd<float> std::cos(std::simd<float>)
```

Just sayin'

Consider using a "T or simd<T>" concept for more generality!

# An alternative

```
std::simd<float> std::cos(std::simd<float>)
```

Just sayin'

Consider using a "T or simd<T>" concept for more generality!

# stdx::simd Overview

Motivation    stdx::simd **Overview**    Example: Signal Processing    SIMD and Merged Blocks    Outlook    Summary
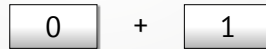○○○○○○○○○○    ○●○○○○○○○○○○○    ○○○○○○○○○○    ○○○○○○○○○    ○○○○    ○

GSI

## Data-Parallel Types

One variable stores $\mathcal{W}_T$ values.
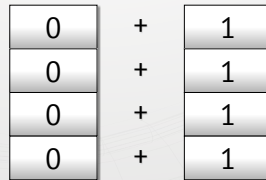One operator signifies $\mathcal{W}_T$ operations (element-wise).

$\mathcal{W}$ for "width"; depends on type $T$

```
int x = 0;
x += 1;
```

| 0 | + | 1 |

vs.

```
simd<int> x = 0;
x += 1;
```

| 0 | + | 1 |
| 0 | + | 1 |
| 0 | + | 1 |
| 0 | + | 1 |

Motivation
○○○○○○○○○○

stdx::simd **Overview**
○○●○○○○○○○○○

Example: Signal Processing
○○○○○○○○○○

SIMD and Merged Blocks
○○○○○○○○○

Outlook
○○○○

Summary
○

GSI

# stdx::simd::size()

- simd<T>::size() is a constant expression, denoting the number of elements.
- A default size() is chosen by the implementation, depending on the target.
- For most targets
  simd<T>::size() * sizeof(T) == simd<U>::size() * sizeof(U) holds:

| float | float | float | float | float | float | float | float |
|-------|-------|-------|-------|-------|-------|-------|-------|

| double | double | double | double |
|--------|--------|--------|--------|

If you really want to rely on it, add a static_assert

Motivation
○○○○○○○○○○
stdx::simd **Overview**
○○○●○○○○○○○○○
Example: Signal Processing
○○○○○○○○○○
SIMD and Merged Blocks
○○○○○○○○○
Outlook
○○○○
Summary
○
GSI

## Synopsis

```cpp
template <typename T, typename Abi = ...>
  class simd;

template <typename T>
  using native_simd = simd<T, ...>;
```

```cpp
template <typename T, typename Abi>
  class simd_mask;

template <typename T>
  using native_simd_mask = simd_mask<T, ...>;
```

- simd<T> behaves just like T (element-wise)

- T must be a *vectorizable type* — all arithmetic types except bool

- simd_mask<T> behaves like bool (element-wise)
  In contrast to bool, there are many different mask types:
  - storage: bit-masks vs. element-sized masks (array of bool is not unheard of),
  - SIMD width simd::size

- Abi determines width and ABI (i.e. how parameters are passed to functions)
  - Example: native_simd<int> on x86 can be one xmm register (sizeof == 16), one ymm
    register (sizeof == 32), or one zmm register (sizeof == 64).

Motivation
○○○○○○○○○○
stdx::simd **Overview**
○○○●○○○○○○○○○
Example: Signal Processing
○○○○○○○○○○
SIMD and Merged Blocks
○○○○○○○○○
Outlook
○○○○
Summary
○
GSI

## Synopsis

```cpp
template <typename T, typename Abi = ...>          template <typename T, typename Abi>
  class simd;                                         class simd_mask;

template <typename T>                              template <typename T>
  using native_simd = simd<T, ...>;                  using native_simd_mask = simd_mask<T, ...>;
```

- simd<T> behaves just like T (element-wise)

- T must be a *vectorizable type* — all arithmetic types except bool

- simd_mask<T> behaves like bool (element-wise)
  In contrast to bool, there are many different mask types:
  - storage: bit-masks vs. element-sized masks (array of bool is not unheard of),
  - SIMD width simd::size

- Abi determines width and ABI (i.e. how parameters are passed to functions)
  - Example: native_simd<int> on x86 can be one xmm register (sizeof == 16), one ymm register (sizeof == 32), or one zmm register (sizeof == 64).

## Synopsis

```
template <typename T, typename Abi = ...>          template <typename T, typename Abi>
  class simd;                                         class simd_mask;

template <typename T>                              template <typename T>
  using native_simd = simd<T, ...>;                  using native_simd_mask = simd_mask<T, ...>;
```

- simd<T> behaves just like T (element-wise)

- T must be a *vectorizable type* — all arithmetic types except bool

- simd_mask<T> behaves like bool (element-wise)
  In contrast to bool, there are many different mask types:
    - storage: bit-masks vs. element-sized masks (array of bool is not unheard of),
    - SIMD width simd::size

- Abi determines width and ABI (i.e. how parameters are passed to functions)
    - Example: native_simd<int> on x86 can be one xmm register (sizeof == 16), one ymm register (sizeof == 32), or one zmm register (sizeof == 64).

Motivation
0000000000
stdx::simd **Overview**
00000●00000000
Example: Signal Processing
0000000000
SIMD and Merged Blocks
000000000
Outlook
0000
Summary
0

GSI

## Constructor examples

```
stdx::simd<int> x0;
```
| ? | ? | ? | ... |

```
stdx::simd<int> x1{};
```
| 0 | 0 | 0 | ... |

```
stdx::simd<int> x2 = 1;
```
| 1 | 1 | 1 | ... |

```
stdx::simd<int> x3(p, stdx::element_aligned);
```
| p[0] | p[1] | p[2] | ... |

```
stdx::simd<int> iota([](int i) { return i; });
```
| 0 | 1 | 2 | ... |

Motivation
○○○○○○○○○○

stdx::simd **Overview**
○○○○○●○○○○○○○

Example: Signal Processing
○○○○○○○○○○

SIMD and Merged Blocks
○○○○○○○○○

Outlook
○○○○

Summary
○

GSI

## Constructor examples

```
stdx::simd<int> x0;
```

| ? | ? | ? | ... |

```
stdx::simd<int> x1{};
```

| 0 | 0 | 0 | ... |

```
stdx::simd<int> x2 = 1;
```

| 1 | 1 | 1 | ... |

```
stdx::simd<int> x3(p, stdx::element_aligned);
```

| p[0] | p[1] | p[2] | ... |

```
stdx::simd<int> iota([](int i) { return i; });
```

| 0 | 1 | 2 | ... |

## Loads & stores

```
1  class simd {
2    simd(const auto* ptr, auto flags);
3
4    void
5    copy_from(const auto* ptr, auto flags);
6
7    void
8    copy_to(auto* ptr, auto flags);
9  };
```
(simplified)

- Loads copy simd<T>::size() elements from a contiguous array chunk into the simd<T> elements.
- Stores do the reverse.
- Consider loads & stores equivalent to dereferencing an iterator in the scalar case.

## Arithmetic & math

```
1  void f(stdx::simd<float> x,
2          stdx::simd<float> y) {
3    x += y;      // x.size() additions
4    x = sqrt(x); // x.size() square roots
5    ...
   // etc. all operators and <cmath>
6  }
```

- Operations act element-wise
- Speed-up is often a factor of `simd<T>::size()`, but may be less, depending on hardware details.

| x0 | += | y0 |
|----|----|----|
| x1 | += | y1 |
| x2 | += | y2 |
| x3 | += | y3 |

## Same for compares (element-wise)

```
1  void f(stdx::simd<float> x, stdx::simd<float> y) {
2    if (x < y) { x = y; } // nonono, you don't write 'if (truefalsetruetrue)' either
3    where(x < y, x) = y; // x = y but only for the elements where x < y
4    if (all_of(x < y)) {...}      // this makes sense, yes
5  }
```

- Comparisons return `simd_mask`.
- `simd_mask` is not convertible to `bool`.
- `simd_mask` can be *reduced* to `bool` via `all_of`, `any_of`, or `none_of`.

```
if ( x0  <  y0  )   x0  =  y0  ;
if ( x1  <  y1  )   x1  =  y1  ;
if ( x2  <  y2  )   x2  =  y2  ;
if ( x3  <  y3  )   x3  =  y3  ;
```

Motivation
○○○○○○○○○○

stdx::simd **Overview**
○○○○○○○●○○○○

Example: Signal Processing
○○○○○○○○○○

SIMD and Merged Blocks
○○○○○○○○○

Outlook
○○○○

Summary
○

GSI

## Same for compares (element-wise)

```
1  void f(stdx::simd<float> x, stdx::simd<float> y) {
2    if (x < y) { x = y; } // nonono, you don't write 'if (truefalsetruetrue)' either
3    where(x < y, x) = y; // x = y but only for the elements where x < y
4    if (all_of(x < y)) {...}      // this makes sense, yes
5  }
```

- Comparisons return `simd_mask`.

- `simd_mask` is not convertible to `bool`.

- `simd_mask` can be *reduced* to `bool` via `all_of`, `any_of`, or `none_of`.

|       |    |   |    |   |    |   |    |   |
|-------|----|---|----|---|----|---|----|---|
| if (  | x0 | < | y0 | ) | x0 | = | y0 | ; |
| if (  | x1 | < | y1 | ) | x1 | = | y1 | ; |
| if (  | x2 | < | y2 | ) | x2 | = | y2 | ; |
| if (  | x3 | < | y3 | ) | x3 | = | y3 | ; |

## Permutations (vir-simd)

Permutations enable

- reductions,
- vectorization of loops with dependent iterations, and
- more efficient permutations of larger arrays.

### Example

```
1  stdx::native_simd<float> permute_even_odd(stdx::native_simd<float> x) {
2    return vir::simd_permute(x, [](auto idx) { return idx ^ 1; });
3  }
```

with AVX, compiles to:
vpermilps ymm0, ymm0, 177

## Reductions

- simd_mask reductions:
  all_of, any_of, none_of, popcount,
  find_first_set, find_last_set

- simd reductions:
  reduce, hmin, hmax

SIMD tree reduction:



### Example

```
1  void f(stdx::simd<float> x) {
2    float sum = stdx::reduce(x);
3    float product = stdx::reduce(x, std::multiplies());
4    float sum_of_pos_x = stdx::reduce(where(x > 0, x));
5    float minimum = stdx::hmin(x);
6    int min_idx = stdx::find_first_set(x == minimum);
7  }
```

Automatic type vectorization (vir-simd)

```
1  struct Point {
2    float x, y, z;
3  };
4  using PointV = vir::simdize<Point>;
5
6  // equivalent to
7  struct PointV {
8    simd<float> x, y, z;
9
10   // stdx::simd-like interface (loads & stores, broadcast, size, subscripting, ...)
11 };
```

A much better solution should be possible after *reflection* lands in the C++ standard.
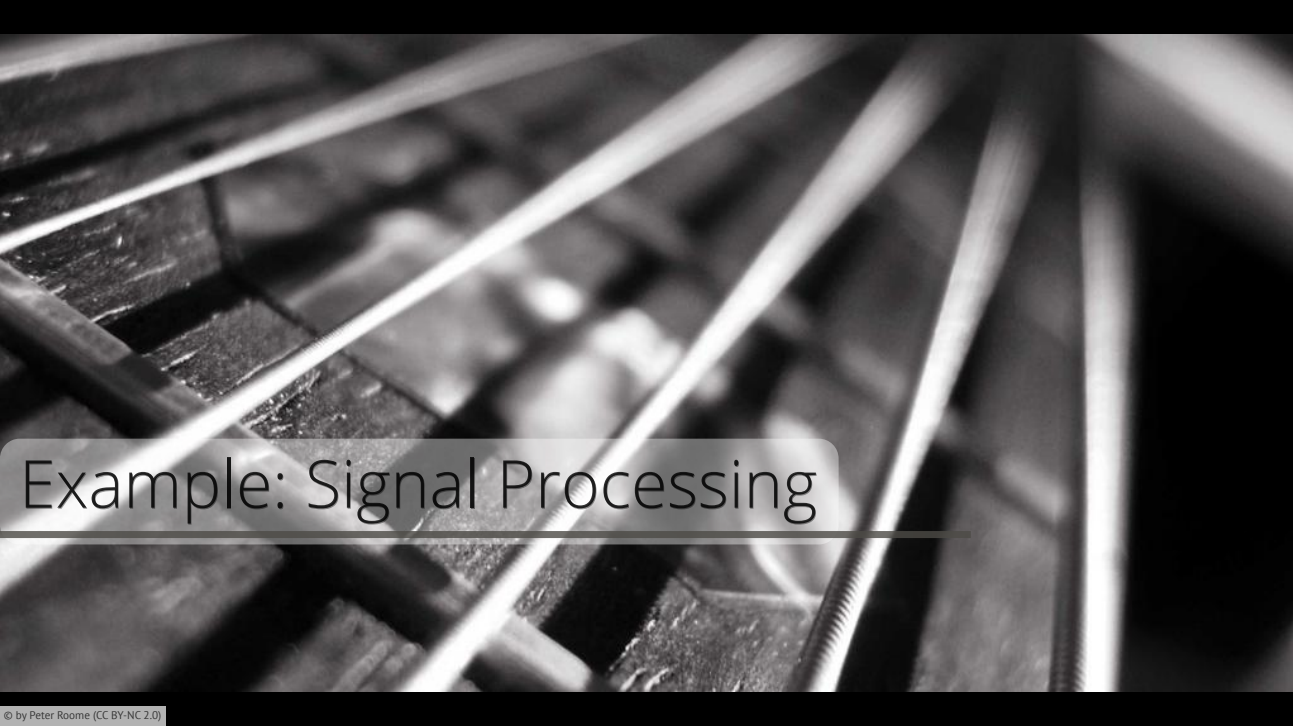
## Documentation

std::experimental::simd

https://en.cppreference.com/w/cpp/experimental/simd

vir-simd

https://github.com/mattkretz/vir-simd/blob/master/README.md
**or just** https://github.com/mattkretz/vir-simd

# Example: Signal Processing

Motivation
0000000000

stdx::simd Overview
000000000000

Example: Signal Processing
0●00000000

SIMD and Merged Blocks
000000000

Outlook
0000

Summary
0

GSI

# Example: signal processing

Let's start really simple, with a gain function:

```
1  void apply_gain(std::span<const float> in, std::span<float> out, float gain) {
2    assert(in.size() == out.size());
3    // implement here
4  }
```

Implementation using a for-loop:

```
     for (std::size_t i = 0; i < in.size(); ++i) {
       out[i] = in[i] * gain;
     }
```

Implementation using a standard algorithm:

```
     std::transform(in.begin(), in.end(), out.begin(), [&](float sample) {
       return sample * gain;
     });
```

## Example: signal processing

Let's start really simple, with a gain function:

```
void apply_gain(std::span<const float> in, std::span<float> out, float gain) {
  assert(in.size() == out.size());
  // implement here
}
```
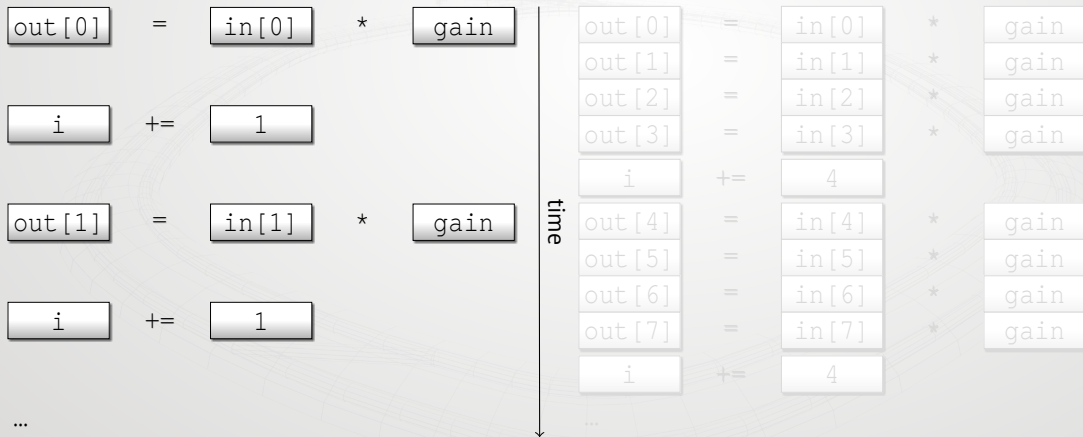
Implementation using a for-loop:

```
  for (std::size_t i = 0; i < in.size(); ++i) {
    out[i] = in[i] * gain;
  }
```

Implementation using a standard algorithm:

```
std::transform(in.begin(), in.end(), out.begin(), [&](float sample) {
  return sample * gain;
});
```

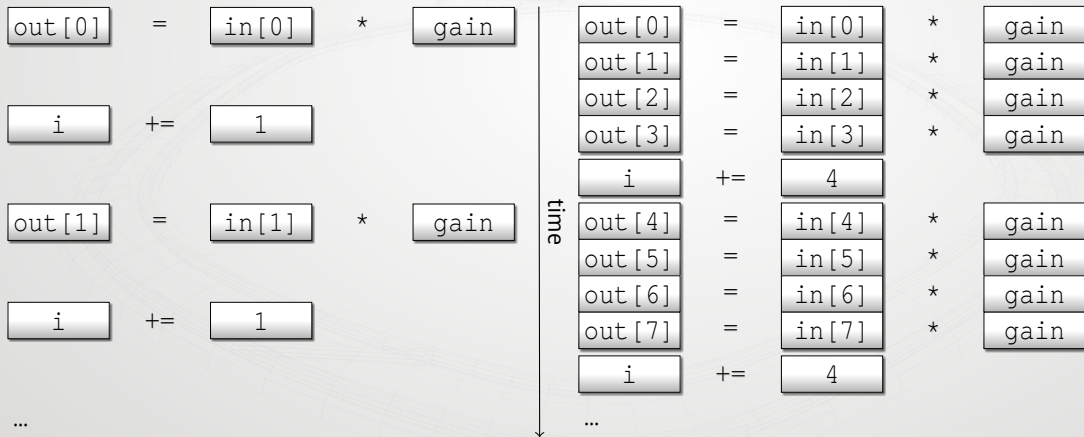Motivation        stdx::simd Overview        **Example: Signal Processing**        SIMD and Merged Blocks        Outlook        Summary
0000000000        000000000000               0●00000000                            000000000              0000            0

GSI

# Example: signal processing

Let's start really simple, with a gain function:

```
1  void apply_gain(std::span<const float> in, std::span<float> out, float gain) {
2    assert(in.size() == out.size());
3    // implement here
4  }
```

Implementation using a for-loop:

```
1  for (std::size_t i = 0; i < in.size(); ++i) {
2    out[i] = in[i] * gain;
3  }
```

Implementation using a standard algorithm:

```
1  std::transform(in.begin(), in.end(), out.begin(), [&](float sample) {
2    return sample * gain;
3  });
```

Motivation
0000000000

stdx::simd Overview
000000000000

Example: Signal Processing
00●0000000

SIMD and Merged Blocks
000000000

Outlook
0000

Summary
0

How can we use `stdx::simd` here?

Look at loops for data-parallelism.

| out[0] | = | in[0] | * | gain |
|--------|---|-------|---|------|

| i | += | 1 |
|---|----|---|

| out[1] | = | in[1] | * | gain |
|--------|---|-------|---|------|

| i | += | 1 |
|---|----|---|

...

time →

| out[0] | = | in[0] | * | gain |
| out[1] | = | in[1] | * | gain |
| out[2] | = | in[2] | * | gain |
| out[3] | = | in[3] | * | gain |
| i | += | 4 |

| out[4] | = | in[4] | * | gain |
| out[5] | = | in[5] | * | gain |
| out[6] | = | in[6] | * | gain |
| out[7] | = | in[7] | * | gain |
| i | += | 4 |

...

How can we use `stdx::simd` here?

Look at loops for data-parallelism.

| out[0] | = | in[0] | * | gain |

| i | += | 1 |

| out[1] | = | in[1] | * | gain |

| i | += | 1 |

...

time →

| out[0] | = | in[0] | * | gain |
| out[1] | = | in[1] | * | gain |
| out[2] | = | in[2] | * | gain |
| out[3] | = | in[3] | * | gain |
| i | += | 4 |
| out[4] | = | in[4] | * | gain |
| out[5] | = | in[5] | * | gain |
| out[6] | = | in[6] | * | gain |
| out[7] | = | in[7] | * | gain |
| i | += | 4 |

...

Motivation
○○○○○○○○○○

stdx::simd Overview
○○○○○○○○○○○○○

Example: Signal Processing
○○●○○○○○○○

SIMD and Merged Blocks
○○○○○○○○○

Outlook
○○○○

Summary
○

GSI

# How can we use `stdx::simd` here?

Look at loops for data-parallelism.

Motivation
○○○○○○○○○○

stdx::simd Overview
○○○○○○○○○○○○

Example: Signal Processing
○○○●○○○○○○

SIMD and Merged Blocks
○○○○○○○○○

Outlook
○○○○

Summary
○

GSI

# Apply gain using `stdx::simd`

```
1   for (std::size_t i = 0; i < in.size(); ++i) {
2     out[i] = in[i] * gain;
3
4
5   }
```

**1** Where's the `float`?

2 Let's turn the temporary into an lvalue.

3 We start with changing the `float` into a `simd<float>`. Not correct yet.

4 Introduced load and store. It compiles. What is missing?

5 Now we can go faster. But it's possible we invoke UB.
Why?

6 This condition avoids the out-of-bounds. But we might skip samples at the end.

Motivation
○○○○○○○○○○
stdx::simd Overview
○○○○○○○○○○○○
Example: Signal Processing
○○○●○○○○○○
SIMD and Merged Blocks
○○○○○○○○○
Outlook
○○○○
Summary
○
GSII

## Apply gain using `stdx::simd`

```
1    for (std::size_t i = 0; i < in.size(); ++i) {
2      float sample = in[i];
3      sample *= gain;
4      out[i] = sample;
5    }
```

1   Where's the `float`?

2   **Let's turn the temporary into an lvalue.**

3   We start with changing the `float` into a `simd<float>`. Not correct yet.

4   Introduced load and store. It compiles. What is missing?

5   Now we can go faster. But it's possible we invoke UB.
    Why?

6   This condition avoids the out-of-bounds. But we might skip samples at the end.

Motivation
○○○○○○○○○○

stdx::simd Overview
○○○○○○○○○○○○○

Example: Signal Processing
○○○●○○○○○○

SIMD and Merged Blocks
○○○○○○○○○

Outlook
○○○○

Summary
○

GSI

## Apply gain using `stdx::simd`

```
1    for (std::size_t i = 0; i < in.size(); ++i) {
2      simd<float> sample = in[i]; // wrong: copies in[i] to all elements of the simd
3      sample *= gain;  // OK
4      out[i] = sample; // Error: cannot convert simd<float> to float
5    }
```

1. Where's the `float`?

2. Let's turn the temporary into an lvalue.

3. **We start with changing the `float` into a `simd<float>`. Not correct yet.**

4. Introduced load and store. It compiles. What is missing?

5. Now we can go faster. But it's possible we invoke UB.
   Why?

6. This condition avoids the out-of-bounds. But we might skip samples at the end.

# Apply gain using `stdx::simd`

```
1   for (std::size_t i = 0; i < in.size(); ++i) {
2     auto sample = simd<float>(&in[i], stdx::element_aligned); // load
3     sample *= gain;
4     sample.copy_to(&out[i], stdx::element_aligned);            // store
5   }
```

1  Where's the `float`?

2  Let's turn the temporary into an lvalue.

3  We start with changing the `float` into a `simd<float>`. Not correct yet.

4  Introduced load and store. It compiles. What is missing?

5  Now we can go faster. But it's possible we invoke UB.
   Why?

6  This condition avoids the out-of-bounds. But we might skip samples at the end.

Motivation
○○○○○○○○○○

stdx::simd Overview
○○○○○○○○○○○○○

Example: Signal Processing
○○○●○○○○○○

SIMD and Merged Blocks
○○○○○○○○○

Outlook
○○○○

Summary
○

GSI

## Apply gain using `stdx::simd`

```
1   for (std::size_t i = 0; i < in.size(); ++i) {
2     auto sample = simd<float>(&in[i], stdx::element_aligned); // load
3     sample *= gain;
4     sample.copy_to(&out[i], stdx::element_aligned);           // store
5   }
```

1 Where's the float?

2 Let's turn the temporary into an lvalue.

3 We start with changing the float into a simd<float>. Not correct yet.

4 Introduced load and store. It compiles. Every sample is processed multiple times.

5 Now we can go faster. But it's possible we invoke UB.
  Why?

6 This condition avoids the out-of-bounds. But we might skip samples at the end.

Motivation  stdx::simd Overview  **Example: Signal Processing**  SIMD and Merged Blocks  Outlook  Summary
○○○○○○○○○○  ○○○○○○○○○○○○  ○○○●○○○○○○  ○○○○○○○○○  ○○○○  ○

GSI

## Apply gain using `stdx::simd`

```cpp
1   for (std::size_t i = 0; i < in.size(); i += simd<float>::size()) {
2     auto sample = simd<float>(&in[i], stdx::element_aligned);
3     sample *= gain;
4     sample.copy_to(&out[i], stdx::element_aligned);
5   }
```

1  Where's the float?

2  Let's turn the temporary into an lvalue.

3  We start with changing the float into a simd<float>. Not correct yet.

4  Introduced load and store. It compiles. Every sample is processed multiple times.

5  Now we can go faster. But it's possible we invoke UB.
   Why?

6  This condition avoids the out-of-bounds. But we might skip samples at the end.

# Apply gain using stdx::simd

```
1  for (std::size_t i = 0; i < in.size(); i += simd<float>::size()) {
2    auto sample = simd<float>(&in[i], stdx::element_aligned);
3    sample *= gain;
4    sample.copy_to(&out[i], stdx::element_aligned);
5  }
```

1 Where's the float?

2 Let's turn the temporary into an lvalue.

3 We start with changing the float into a simd<float>. Not correct yet.

4 Introduced load and store. It compiles. Every sample is processed multiple times.

5 Now we can go faster. But it's possible we invoke UB.

   4 < 7 is true, but accessing in[4 + simd::size - 1] is out of bounds.

6 This condition avoids the out-of-bounds. But we might skip samples at the end.

Motivation
○○○○○○○○○○

stdx::simd Overview
○○○○○○○○○○○○

Example: Signal Processing
○○○●○○○○○○

SIMD and Merged Blocks
○○○○○○○○○

Outlook
○○○○

Summary
○

GSI

## Apply gain using `stdx::simd`

```
1   for (std::size_t i = 0; i + simd<float>::size() <= in.size(); i += simd<float>::size()
2     auto sample = simd<float>(&in[i], stdx::element_aligned);
3     sample *= gain;
4     sample.copy_to(&out[i], stdx::element_aligned);
5   }
```

1. Where's the float?

2. Let's turn the temporary into an lvalue.

3. We start with changing the float into a simd<float>. Not correct yet.

4. Introduced load and store. It compiles. Every sample is processed multiple times.

5. Now we can go faster. But it's possible we invoke UB.
   `4 < 7` is true, but accessing `in[4 + simd::size - 1]` is out of bounds.

6. This condition avoids the out-of-bounds. But we might skip samples at the end.

Motivation | stdx::simd Overview | Example: Signal Processing | SIMD and Merged Blocks | Outlook | Summary
○○○○○○○○○○ | ○○○○○○○○○○○○ | ○○○●○○○○○○ | ○○○○○○○○○ | ○○○○ | ○

GSI

# Apply gain using `stdx::simd`

```
1  for (std::size_t i = 0; i + simd<float>::size() <= in.size(); i += simd<float>::size())
2    auto sample = simd<float>(&in[i], stdx::element_aligned);
3    sample *= gain;
4    sample.copy_to(&out[i], stdx::element_aligned);
5  }
```

*Sigh ... and we're running out of slide space*

1. Where's the f...
2. Let's turn the temporary into an tr...
3. We start with changing the float into a simd<float>. Not cor... yet.
4. Introduced load and store. It compiles. Every sample is processed multiple times.
5. Now we can go faster. But it's possible we invoke UB.
   4 < 7 is true, but accessing in[4 + simd::size - 1] is out of bounds.
6. This condition avoids the out-of-bounds. But we might skip samples at the end.

Motivation
○○○○○○○○○○
stdx::simd Overview
○○○○○○○○○○○○
Example: Signal Processing
○○○○●○○○○○
SIMD and Merged Blocks
○○○○○○○○○
Outlook
○○○○
Summary
○

GSII

Apply gain using `stdx::simd` with epilogue

```
1    using floatv = stdx::simd<float>;
2    std::size_t i = 0;
3    for (; i + floatv::size() <= in.size(); i += floatv::size()) {
4      auto sample = simd<float>(&in[i], stdx::element_aligned);
5      sample *= gain;
6      sample.copy_to(&out[i], stdx::element_aligned);
7    }
8    //
9    // TODO: process samples in range [i, in.size())
10   //
```

❶ Introduce an alias for our `simd` type. Recommended practice anyway.

❷ In SIMD programming, this "handle the remainder" is called an *epilogue.*

## Apply gain using `stdx::simd` with epilogue

```
1    using floatv = stdx::simd<float>;
2    std::size_t i = 0;
3    for (; i + floatv::size() <= in.size(); i += floatv::size()) {
4      auto sample = simd<float>(&in[i], stdx::element_aligned);
5      sample *= gain;
6      sample.copy_to(&out[i], stdx::element_aligned);
7    }
8    for (; i < in.size(); ++i) {
9      out[i] = in[i] * gain;
10   }
```

❶ Introduce an alias for our simd type. Recommended practice anyway.

❷ In SIMD programming, this "handle the remainder" is called an *epilogue*.

# Apply gain using `stdx::simd` with epilogue

```
1   using floatv = stdx::simd<float>;
2   std::size_t i = 0;
3   for (; i + floatv::size() <= in.size(); i += floatv::size()) {
4     auto sample = simd<float>(&in[i], stdx::element_aligned);
5     sample *= gain;
6     sample.copy_to(&out[i], stdx::element_aligned);
7   }
8   for (; i < in.size(); ++i) {
9     out[i] = in[i] * gain;
10  }
```

*Phew. Now that you've seen this... I have good news...*

① Introduce an alias for our `simd` type. Recommended practice anyway.

② In SIMD programming, this "handle the remainder" is called an *epilogue.*

### Good news 👍

# GR4 provides the above for you,
# simply implement processOne and you're done!

### Example

```
1  template <gr::meta::t_or_simd<float> V>
2  V processOne(V samples) {
3    return samples * gain;
4  }
```

You can also overload on simd:
    processOne(float x) { return x * gain; }

template <typename  >
simd<float,  > processOne(simd<float,  > x) { return x * gain; }

Motivation
0000000000

stdx::simd Overview
000000000000

Example: Signal Processing
0000000000

SIMD and Merged Blocks
000000000

Outlook
0000

Summary
0

GSI

### Good news 👍

## GR4 provides the above for you,
## simply implement processOne and you're done!

### Example

```
template <gr::meta::t_or_simd<float> V>
V processOne(V samples) {
  return samples * gain;
}
```

You can also overload on simd:

```
float processOne(float x) { return x * gain; }

template <typename Abi>
simd<float, Abi> processOne(simd<float, Abi> x) { return x * gain; }
```

### More good news 👍👍

vir-simd provides algorithms that implement the above for you!

### Example

Start with:

```
1  void apply_gain(std::span<const float> in, std::span<float> out, float gain) {
2    std::transform(std::execution::seq,
3                   in.begin(), in.end(), out.begin(), [&](float x) {
4      return x * gain;
5    });
6  }
```

Motivation
○○○○○○○○○○

stdx::simd Overview
○○○○○○○○○○○○○

Example: Signal Processing
○○○○○○○●○○○

SIMD and Merged Blocks
○○○○○○○○○

Outlook
○○○○

Summary
○

GSI

### More good news 👍👍

vir-simd provides algorithms that implement the above for you!

### Example

A minor tweak with huge effect:

```
1  void apply_gain(std::span<const float> in, std::span<float> out, float gain) {
2    std::transform(vir::execution::simd,
3                   in.begin(), in.end(), out.begin(), [&](auto x) {
4      return x * gain;
5    });
6  }
```

A Benchmark

# In Code

## trivial data-parallelism

given `data` of type `std::vector<int>`

```
1 std::for_each(std::execution::unseq,
2   data.begin(), data.end(),
3   [](int& x) {
4     x += 1;
5   });
```

```
1 std::for_each(vir::execution::simd,
2   data.begin(), data.end(),
3   [](auto& x) {
4     x += 1;
5   });
```

vir::execution::simd is part of the vir-simd library.

Motivation
○○○○○○○○○○

stdx::simd Overview
○○○○○○○○○○○○

Example: Signal Processing
○○○○○○○●○

SIMD and Merged Blocks
○○○○○○○○○

Outlook
○○○○

Summary
○

GSI

# In Code

### trivial data-parallelism

given `data` of type `std::vector<int>`

```
1  std::for_each(std::execution::unseq,
2    data.begin(), data.end(),
3    [](int& x) {
4      x += 1;
5    });
```

```
1  std::for_each(vir::execution::simd,
2    data.begin(), data.end(),
3    [](auto& x) {
4      x += 1;
5    });
```

`vir::execution::simd` is part of the vir-simd library.

# Results (GCC 13, Intel Skylake i7-8550U)

```
1  for (int& x : data) {
2    x += 1;
3  }
```

# Results (GCC 13, Intel Skylake i7-8550U)

```
1  for (int& x : data) {
2    x += 1;
3  }
```

```
1  std::for_each(std::execution::unseq,
2    data.begin(), data.end(),
3    [](int& x) { x += 1; }
4  );
```

# Results (GCC 13, Intel Skylake i7-8550U)
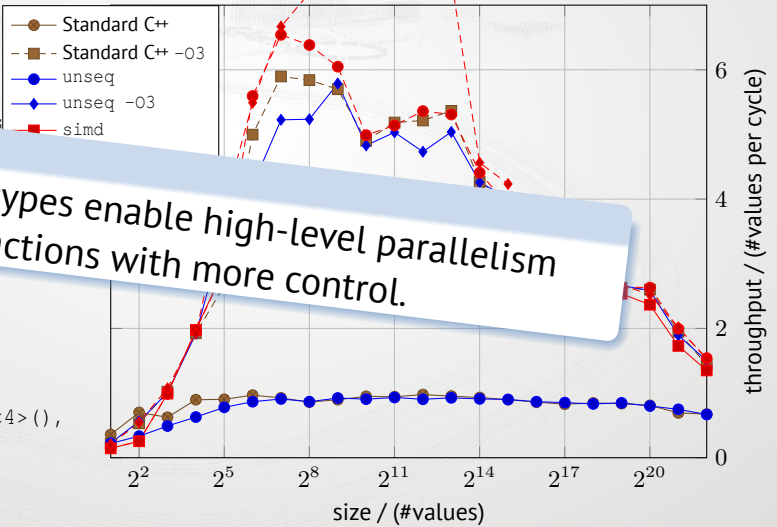
```
1  for (int& x : data) {
2    x += 1;
3  }
```

```
1  std::for_each(std::execution::unseq,
2    data.begin(), data.end(),
3    [](int& x) { x += 1; }
4  );
```

```
1  std::for_each(vir::execution::simd,
2    data.begin(), data.end(),
3    [](auto& x) { x += 1; }
4  );
```

# Results (GCC 13, Intel Skylake i7-8550U)

```
1  for (int& x : data) {
2    x += 1;
3  }
```

```
1  std::for_each(std::execution::unseq,
2    data.begin(), data.end(),
3    [](int& x) { x += 1; }
4  );
```

```
1  std::for_each(vir::execution::simd,
2    data.begin(), data.end(),
3    [](auto& x) { x += 1; }
4  );
```

# Results (GCC 13, Intel Skylake i7-8550U)

```
1  for (int& x : data) {
2    x += 1;
3  }

1  std::for_each(std::execution::uns
2    data.begin(), data.end(),
3    [](int& x) { x += 1; }
4  );

1  std::for_each(vir::execution::simd,
2    data.begin(), data.end(),
3    [](auto& x) { x += 1; }
4  );
```

# Results (GCC 13, Intel Skylake i7-8550U)

```
1  for (int& x : data) {
2      x += 1;
3  }
```

```
1  std::for_each(std::execution::uns
2      data.begin(), data.end(),
3      [](int& x) { x += 1; }
4  );
```

```
1  std::for_each(vir::execution::simd,
2      data.begin(), data.end(),
3      [](auto& x) { x += 1; }
4  );
```

GSI

# Results (GCC 13, Intel Skylake i7-8550U)

```
1  for (int& x : data) {
2    x += 1;
3  }

1  std::for_each(std::execution::uns
2    data.begin(), data.end(),
3    [](int& x) { x += 1; }
4  );

1  std::for_each(vir::execution::simd,
2    data.begin(), data.end(),
3    [](auto& x) { x += 1; }
4  );

1  std::for_each(
2    vir::execution::simd.unroll_by<4>(),
3    data.begin(), data.end(),
4    [](auto& x) { x += 1; }
5  );
```



max. store throughput

Legend:
- Standard C++
- Standard C++ -O3
- unseq
- unseq -O3
- simd
- simd -O3
- simd unroll 4

throughput / (#values per cycle)

size / (#values)

# Results (GCC 13, Intel Skylake i7-8550U)

```
1  for (int& x : data) {
2    x += 1;
3  }
```

```
1  std::for_each(...
2    data.be...
3    [](int...
4  );
```

```
1  std::for_each(vir::execution::...
2    data.begin(), data.end(),
3    [](auto& x) { x += 1; }
4  );
```

```
1  std::for_each(
2    vir::execution::simd.unroll_by<4>(),
3    data.begin(), data.end(),
4    [](auto& x) { x += 1; }
5  );
```

max. store throughput



Legend:
- Standard C++
- Standard C++ −O3
- unseq
- unseq −O3
- simd

x-axis: size / (#values) — $2^2$, $2^5$, $2^8$, $2^{11}$, $2^{14}$, $2^{17}$, $2^{20}$

y-axis: throughput / (#values per cycle)

**Take-Away #2**
Data-parallel types enable high-level parallelism abstractions with more control.

# Results (GCC 13, Intel Skylake i7-8550U)



max. store throughput

```
1  for (int& x : data) {
2    x += 1;
3  }
```

```
1  std::for_each(... unse...
2    data.be...
3    [](int...
4  );
```

```
1  std::for_each(
2    da...
3    [...
4  );
```

```
1  std::for_each(
2    vir::execution::simd.unroll_by<4>(),
3    data.begin(), data.end(),
4    [](auto& x) { x += 1; }
5  );
```

**Take-Away #2**

Data-parallel types enable high-level abst...

**Take-Away #3**

Data-parallel types compose, allowing separation of concern.

separation of range algorithm from element operations, in this case...

Legend:
- Standard C++
- Standard C++ −O3
- unseq
- unseq −O3
- simd

throughput / (#values per cycle)

size / (#values)

$2^2$  $2^5$  $2^8$  $2^{11}$  $2^{14}$  $2^{17}$  $2^{20}$

# SIMD and Merged Blocks

# Amdahl's law

> "the overall performance improvement gained by optimizing a single part of a system is limited by the fraction of time that the improved part is actually used."
>
> [Reddy, Martin (2011). API Design for C++. doi:10.1016/2010-0-65832-9]

In other words:
The overall speedup from using SIMD of a flow graph may be frustratingly low, even though SIMD improves throughput & latency per block.

Because ...

# Amdahl's law

> "the overall performance improvement gained by optimizing a single part of a system is limited by the fraction of time that the improved part is actually used."
>
> [Reddy, Martin (2011). API Design for C++. doi:10.1016/2010-0-65832-9]

In other words:
The overall speedup from using SIMD of a flow graph may be frustratingly low, even though SIMD improves throughput & latency per block.

Because ...

Motivation
○○○○○○○○○○

stdx::simd Overview
○○○○○○○○○○○○○

Example: Signal Processing
○○○○○○○○○

SIMD and Merged Blocks
○○●○○○○○○

Outlook
○○○○

Summary
○

GSI

time

| invoke block a | a: process samples | invoke block b | b: process samples |

Run-time configurable communication between blocks has a cost:

- block output written to memory,
- block scheduling
- buffer management (avoid over- and underruns),
- block input read from memory
- (additional cache/memory latencies)

So, we need to write our application as one loop? 🤔 😱

Alternative:

- let 'b' process the results from 'a' immediately.
- This requires connection at compile-time but does not require implementing a new block that merges their functionality.

time

| invoke block a | a: process samples | invoke block b | b: process samples |

Run-time configurable communication between blocks has a cost:

- block output written to memory,
- block scheduling
- buffer management (avoid over- and underruns),
- block input read from memory
- (additional cache/memory latencies)

So, we need to write our application as one loop? 🤔 😱
Alternative:

- let 'b' process the results from 'a' immediately.
- This requires connection at compile-time but does not require implementing a new block that merges their functionality.

Motivation
○○○○○○○○○○
stdx::simd Overview
○○○○○○○○○○○○
Example: Signal Processing
○○○○○○○○○○
SIMD and Merged Blocks
○○○●○○○○○
Outlook
○○○○
Summary
○

GSI

time

| invoke block a | a: process samples | invoke block b | b: process samples |

after compile-time connection:

| invoke block a+b | a: process samples | b: process samples |

with a good chance of interleaving computation:

| invoke block a+b | a: process samples |
|  | b: process samples |

Motivation
stdx::simd Overview
Example: Signal Processing
SIMD and Merged Blocks
Outlook
Summary

- a mergable block implements `processOne` instead of `processBulk`
- "One" means "one object", *not* "one sample"
- "One" can be a `simd<T>` object (multiple samples) or a `T` object (one sample)
- thus, a better picture looks like:

time

| invoke block a | a.work() loops over | invoke block b | b.work() loops over |
| | a.processOne(...) | | b.processOne(...) |

and turns into:

- a mergable block implements `processOne` instead of `processBulk`
- "`One`" means "one object", *not* "one sample"
- "`One`" can be a `simd<T>` object (multiple samples) or a `T` object (one sample)
- thus, a better picture looks like:

time

| invoke block a | a.work() loops over a.processOne(...) | invoke block b | b.work() loops over b.processOne(...) |

and turns into:

| invoke merged block a+b | ab.work() loops over b.processOne( a.processOne(...) ) |

🏆 Performance

A merged block is an efficient loop implementation

. . .

🏆 Performance

## A merged block is an efficient loop implementation

🏆 Abstraction

## ... via composition of existing components!

# In Pictures: not merged, no SIMD

Motivation
○○○○○○○○○○

stdx::simd Overview
○○○○○○○○○○○○○

Example: Signal Processing
○○○○○○○○○○

SIMD and Merged Blocks
○○○○○○○●○○

Outlook
○○○○

Summary
○

# In Pictures: not merged, no SIMD

# In Pictures: not merged, no SIMD



a.processOne

Motivation
○○○○○○○○○○

stdx::simd Overview
○○○○○○○○○○○○○

Example: Signal Processing
○○○○○○○○○○

SIMD and Merged Blocks
○○○○○○○●○○

Outlook
○○○○

Summary
○

GSI

# In Pictures: not merged, no SIMD

Motivation
○○○○○○○○○○

stdx::simd Overview
○○○○○○○○○○○○

Example: Signal Processing
○○○○○○○○○○

SIMD and Merged Blocks
○○○○○○○●○○

Outlook
○○○○

Summary
○

GSI

## In Pictures: not merged, no SIMD

# In Pictures: not merged, no SIMD

a.processOne

a.processOne

a.processOne

a.processOne

# In Pictures: not merged, no SIMD

# In Pictures: not merged, no SIMD

# In Pictures: not merged, no SIMD

# In Pictures: not merged, no SIMD

Motivation
○○○○○○○○○○

stdx::simd Overview
○○○○○○○○○○○○○

Example: Signal Processing
○○○○○○○○○○

SIMD and Merged Blocks
○○○○○○○●○○

Outlook
○○○○

Summary
○

GSI

# In Pictures: not merged, no SIMD

Motivation
○○○○○○○○○○

stdx::simd Overview
○○○○○○○○○○○○

Example: Signal Processing
○○○○○○○○○○

SIMD and Merged Blocks
○○○○○○○●○

Outlook
○○○○

Summary
○

GSI

# In Pictures: not merged, with 4-wide SIMD

# In Pictures: not merged, with 4-wide SIMD



a.processOne
a.processOne
a.processOne
a.processOne

b.processOne
b.processOne
b.processOne
b.processOne

Motivation
0000000000

stdx::simd Overview
000000000000

Example: Signal Processing
0000000000

SIMD and Merged Blocks
000000000●0

Outlook
0000

Summary
0

GSI

# In Pictures: not merged, with 4-wide SIMD

→ a.processOne → b.processOne
→ a.processOne → b.processOne
→ a.processOne → b.processOne
→ a.processOne → b.processOne

Bandwidth increased by a factor of 4! 🏆

# In Pictures: not merged, with 4-wide SIMD

# In Pictures: merged, with 4-wide SIMD

# In Pictures: merged, with 4-wide SIMD

b.processOne(a.processOne)

b.processOne(a.processOne)

b.processOne(a.processOne)

b.processOne(a.processOne)

*Hard to go any faster between two given buffers!*

Outlook

Motivation
○○○○○○○○○○

stdx::simd Overview
○○○○○○○○○○○○○

Example: Signal Processing
○○○○○○○○○○

SIMD and Merged Blocks
○○○○○○○○○

Outlook
○●○○

Summary
○

GSI

## vectorizable types

- TS ( C++17 )
  - arithmetic types other than `bool`
- P1928 ( C++26, in library wording )
  - arithmetic types other than `bool` and `long double`
- P2663 ( C++26, in library wording, waiting on the above )
  - adds `std::complex` to vectorizable types (`simd<complex>` *not* `complex<simd>`)
- P2964 ( C++26 is unlikely, scope & design review )
  - add scoped and unscoped enums to vectorizable types (e.g. `std::byte`)
  - add user-defined numeric types to vectorizable types (fixed-point, saturating integers, etc.)

Motivation
○○○○○○○○○○
stdx::simd Overview
○○○○○○○○○○○○
Example: Signal Processing
○○○○○○○○○○
SIMD and Merged Blocks
○○○○○○○○○
Outlook
○○●○
Summary
○
GSI

# We want more interoperability

- Conversion to/from `std::array`, `std::span`, or generally any contiguous range.

- Mask conversion to/from `std::bitset`.

- Initializer list constructor for `simd`.

Tricky because of portability concerns.

- Make `simd` and `simd_mask` (read-only) ranges.
  - ⇒ Formatting (`std::format`).
  - ⇒ Makes it easy to flatten a `vector<simd<T>>` into a range of `T`.

Motivation
○○○○○○○○○○

stdx::simd Overview
○○○○○○○○○○○○○

Example: Signal Processing
○○○○○○○○○○

SIMD and Merged Blocks
○○○○○○○○○

**Outlook**
○○○●

Summary
○

# Better gather & scatter integration

### Example

```
1  void f(std::vector<float>& data, t_or_simd<int> auto indexes) {
2    data[indexes] = std::sin(data[indexes]);
3  }
```

1. `stdx::simd` is for you — because you care about latency & throughput!
2. `stdx::simd` expresses data-parallelism — compiled to SIMD (and ILP[2])!
3. Focus on data-parallelism not SIMD instructions/registers!
4. `stdx::simd` enables separation of
   - serial execution,
   - synchronously parallel execution (SIMD and ILP), and
   - asynchronously parallel execution (threads).
5. `stdx::simd` guides you to design scalable and portable parallelization.
6. vir-simd implements high-level standard algorithms to simplify SIMD access to ranges of scalars.
7. vir-simd implements high-level & specialized parallelism abstractions with more control where loop-vectorization fails.
8. `stdx::simd` provides an API & ABI for vectorization across translation units (including library boundaries).

---

[2]much better with C++26

## High-throughput computing without overhead

```
1  mov      eax, DWORD PTR [rdi]
2  imul     eax, eax
3  mov      DWORD PTR [rdi], eax
```

- src/dst: array of integers
- throughput: 0.5/1/1 cycles (Intel)
- integer multiplications: 1

```
1  vmovdqu32    zmm0, ZMMWORD PTR [rdi]
2  vpmulld      zmm0, zmm0, zmm0
3  vmovdqu32    ZMMWORD PTR [rdi], zmm0
```

- src/dst: array of integers
- throughput: 0.5/1/1 cycles (Intel)
- integer multiplications: 16

### Take-Away #4

SIMD is relevant for low-latency, not only high-throughput

# Data Structures for SIMD Processing

Consider this class template definition:

```cpp
template <typename T> struct Point {
  T x, y, z;

  Point normalized() const {
    using std::sqrt;
    const auto scale = 1 / sqrt(x * x + y * y + z * z);
    return {x * scale, y * scale, z * scale};
  }
};
```

Appendix

More Motivation
○

Data Structures
○●○○

std::simd
○○○○○○○○

GSI



**A**rray **o**f **S**truct

`vector<Point<float>>`

## Array of Struct

`vector<Point<float>>`

Array of Struct

`vector<Point<float>>`

Struct of Array

`Point<valarray<float>>`

Array of Struct

vector<Point<float>>

Struct of Array

Point<valarray<float>>

Array of vectorized Struct

vector<Point<simd<float>>>

$\mathcal{W}_{\text{float}} = 1$

$\mathcal{W}_{\text{float}} = 2$

## Example: normalize *N* 3-D points

```
1  template <typename T> struct Point {
2    T x, y, z;
3
4    Point normalized() const {
5      using std::sqrt;
6      const auto scale
7        = 1 / sqrt(x*x + y*y + z*z);
8      return {x * scale,
9              y * scale,
10             z * scale};
11   }
12 };
```

```
1  void aos(const std::vector<Point<float>>& points) {
2    for (auto&p : points) {
3      p = p.normalized();
4    }
5  }
```

```
1  void soa(const Point<std::valarray<float>>& points) {
2    points = points.normalized();
3  }
```

```
1  void aovs(const std::vector<Point<  sim<float>>>&
2            points) {
3    for (auto&p : points) {
4      p = p.normalized();
5    }
6  }
```

## Example: normalize *N* 3-D points

```
1  template <typename T> struct Point {
2    T x, y, z;
3
4    Point normalized() const {
5      using std::sqrt;
6      const auto scale
7        = 1 / sqrt(x*x + y*y + z*z);
8      return {x * scale,
9              y * scale,
10             z * scale};
11   }
12 };
```

```
1  void aos(const std::vector<Point<float>>& points) {
2    for (auto&p : points) {
3      p = p.normalized();
4    }
5  }
```

```
1  void soa(const Point<std::valarray<float>>& points) {
2    points = points.normalized();
3  }
```

```
1  void aovs(const std::vector<Point<      simd<float>>>&
2            points) {
3    for (auto&p : points) {
4      p = p.normalized();
5    }
6  }
```

## Example: normalize *N* 3-D points

```
1  template <typename T> struct Point {
2    T x, y, z;
3
4    Point normalized() const {
5      using std::sqrt;
6      const auto scale
7        = 1 / sqrt(x*x + y*y + z*z);
8      return {x * scale,
9              y * scale,
10             z * scale};
11   }
12 };
```

```
1  void aos(const std::vector<Point<float>>& points) {
2    for (auto&p : points) {
3      p = p.normalized();
4    }
5  }
```

```
1  void soa(const Point<std::valarray<float>>& points) {
2    points = points.normalized();
3  }
```
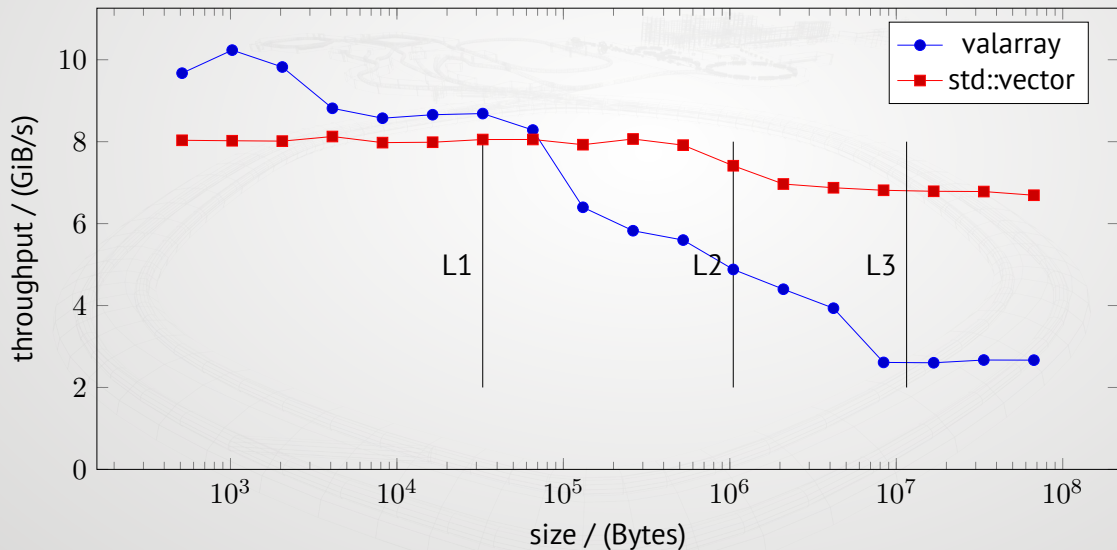
```
1  void aovs(const std::vector<Point<std::simd<float>>>&
2            points) {
3    for (auto&p : points) {
4      p = p.normalized();
5    }
6  }
```
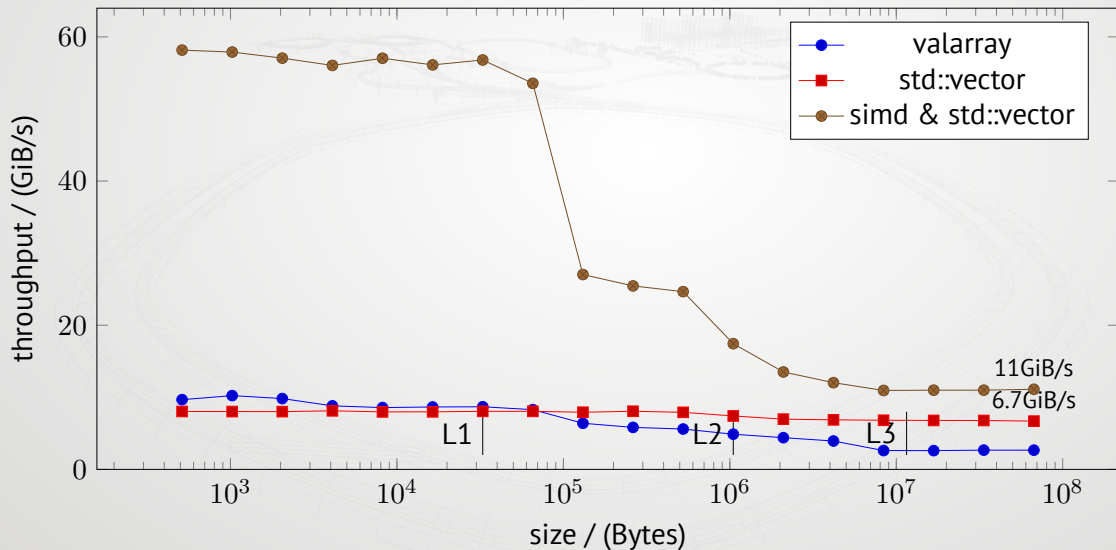
# Normalization benchmark
Linux, Intel Xeon W-2145 (2 AVX-512 FMA ports)

# Normalization benchmark
Linux, Intel Xeon W-2145 (2 AVX-512 FMA ports)

### Example

One multiplication:

```
float f(float x) {
  return x * 2.f;
}
```
https://godbolt.org/z/1TY9jbqqj

Several multiplications in parallel:

```
simd<float> f(simd<float> x) {
  return x * 2.f;
}
```

```
; Intel, AVX-512:
f(float):
  vaddss xmm0, xmm0, xmm0
  ret

; aarch64:
f(float):
  fadd s0, s0, s0
  ret
```

```
; Intel, AVX-512:
f(std::simd<float, std::__detail::_VecBltnBtmsk<64> >):
  vaddps zmm0, zmm0, zmm0
  ret

; aarch64:
f(std::simd<float, std::__detail::_VecBuiltin<16> >):
  fadd v0.4s, v0.4s, v0.4s
  ret
```

## Pass by value or const-ref?

- You pass `int` and `float` by value not const-ref... 🤔
- ...so you pass `simd<int>` and `simd<float>` by value! (exceptions apply)

### by reference 😱

```
1 void f(simd<float>& result, const simd<float>& x) {
2   result = x * 2.f;
3 }
```

```
1 vmovaps zmm0, ZMMWORD PTR [rsi]
2 vaddps zmm0, zmm0, zmm0
3 vmovaps ZMMWORD PTR [rdi], zmm0
4 vzeroupper
5 ret
```

### by value 🏆

```
1 simd<float> f(simd<float> x) {
2   return x * 2.f;
3 }
```

```
1 vaddps zmm0, zmm0, zmm0
2 ret
```

## Data-Parallel Conditionals

### Example

One compare and 0 or 1 assignments:
```
float f(float x) {
  if (x > 0.f) { x *= 2.f; }
  return x;
}
```

$\mathcal{W}_{float}$ compares and $0 - \mathcal{W}_{float}$ assignments in parallel:
```
simd<float> f(simd<float> x) {
  return simd_select(x > 0.f, x * 2.f, x);
}
```

return x > 0.f ? x * 2.f : x; – anyone?
The TS uses where-expressions instead.

- Compares yield $\mathcal{W}_T$ boolean answers
- Return type of compares: std::simd_mask<T, N>
- Reduction functions: all_of, any_of, none_of
- simd code typically uses no/few branches, relying on masked assignment instead

## Data-Parallel Conditionals

### Example

One compare and 0 or 1 assignments:
```
float f(float x) {
  return x > 0.f ? x * 2.f : x;
}
```

$\mathcal{W}_{\texttt{float}}$ compares and $0-\mathcal{W}_{\texttt{float}}$ assignments in parallel:
```
simd<float> f(simd<float> x) {
  return simd_select(x > 0.f, x * 2.f, x);
}
```

return x > 0.f ? x * 2.f : x; — anyone?
The TS uses where-expressions instead.

- Compares yield $\mathcal{W}_{\texttt{T}}$ boolean answers
- Return type of compares: `std::simd_mask<T, N>`
- Reduction functions: `all_of`, `any_of`, `none_of`
- `simd` code typically uses no/few branches, relying on masked assignment instead

## Data-Parallel Conditionals

### Example

One compare and 0 or 1 assignments:
```
float f(float x) {
  return x > 0.f ? x * 2.f : x;
}
```

$\mathcal{W}_{\text{float}}$ compares and $0 - \mathcal{W}_{\text{float}}$ assignments in parallel:
```
simd<float> f(simd<float> x) {
  return simd_select(x > 0.f, x * 2.f, x);
}
```

```
return x > 0.f ? x * 2.f : x; — anyone?
```
The TS uses `where`-expressions instead.

- Compares yield $\mathcal{W}_{\text{T}}$ boolean answers
- Return type of compares: `std::simd_mask<T, N>`
- Reduction functions: `all_of`, `any_of`, `none_of`
- `simd` code typically uses no/few branches, relying on masked assignment instead

## ABI tag / width default

- Compiler flags determine the default ABI tag / SIMD width.

- `simd<T>` sets the ABI tag to the widest efficient $\mathcal{W}_T$ for your `-march=` setting. It also influences the representation of `simd_mask` (i.e. `sizeof(mask)` may be very different).
  The TS uses the wrong default for the ABI tag. The TS gives you the lowest common denominator for all possible implementations of the target architecture. Use `native_simd<T>` with the TS!

- The ABI tag enables support for future ISA extensions without breaking existing code.
  The dreaded ABI break becomes an ABI addition...

### Consequence

The `std::simd` and `simd_mask` ABI depends on `-m` flags!

## Constructors (simplified)

```
1  template <typename T, typename Abi = ...>
2  class basic_simd {
3    basic_simd() = default;
4    basic_simd(T);
5    basic_simd(std::contiguous_iterator auto, Flags = ...);
6    basic_simd(Generator);
7  }
```

- The defaulted *default* constructor allows uninitialized and zero-initialized objects.

- The *broadcast* constructor initializes all elements with the given value.

  requires a value-preserving conversion

- The *load* constructor reads $\mathcal{W}_T$ elements starting from the given iterator.

  Flags can hint about alignment and opt in to non-value-preserving conversions

- The *generator* constructor initializes each element via the given generator function.
  The generator function is called with `std::integral_constant<std::size_t, i>`, where i is the
  index of the element to be initialized.

## Loads & stores (epilogue)

SIMD code typically needs an *epilogue*:

```
1  void f(std::vector<float>& data) {
2    using floatv = std::simd<float>;
3    auto it = data.begin();
4    for (; it <= data.end() - floatv::size();
5         it += floatv::size()) {
6      std::sin(floatv(it)).copy_to(it);
7    }
8    for (; it < data.end(); ++it) {
9      *it = std::sin(*it);
10   }
11 }
```

- Having to write the epilogue every time is *error prone*.
- P1928 provides the low-level primitives, enabling *library-based high-level abstractions*. E.g.
  - SIMD iterator adaptor,
  - SIMD execution policy,
  - your ideas...

## P0350 Outlook – SIMD execution policy

```
1  void f(std::vector<float>& data) {
2    std::for_each(std::execution::simd, data.begin(), data.end(), [](auto& v) {
3      v = std::sin(v);
4    });
5  }
```

- Lambda called with `stdx::native_simd<float>`.

- Epilogue: called with `stdx::simd<float, Abi>` with different `Abi` so that the remainder of `data` is processed with minimal calls to the lambda.

## Subscripting

Loads & stores are great, but sometimes you just want to access it like an array.

```
1  void f(std::simd<float> x) {
2    for (int i = 0; i < x.size(); ++i) {
3      x[i] *= 2.f;
4      x[i] = foo(x[i]);
5      auto ref = x[i];
6      ref = foo(x[i]); // ERROR: no assignment
7    }
8  }
```

- non-const subscripting returns a
  basic_simd::reference
- implements all non-const operators, i.e.
  (compound) assignment, increment and
  decrement, and also swap.

- all of the above functions are rvalue-ref
  qualified, i.e. are *only allowed on*
  *temporaries*

- For std::vector<float> x the type of
  val in auto val = x is float, not
  float&. I.e. we expect a *decay* of the
  reference (proxy) to the element type.
  Another paper I still have to write and
  defend in the committee. 😉

## Subscripting

Loads & stores are great, but sometimes you just want to access it like an array.

```
1  void f(std::simd<float> x) {
2    for (int i = 0; i < x.size(); ++i) {
3      x[i] *= 2.f;
4      x[i] = foo(x[i]);
5      auto ref = x[i];
6      ref = foo(x[i]); // ERROR: no assignment
7    }
8  }
```

- non-`const` subscripting returns a
  `basic_simd::reference`
- implements all non-const operators, i.e.
  (compound) assignment, increment and
  decrement, and also `swap`.

- all of the above functions are rvalue-ref
  qualified, i.e. are *only allowed on
  temporaries*

- For `std::vector<float> x` the type of
  `val` in `auto val = x` is `float`, not
  `float&`. I.e. we expect a *decay* of the
  reference (proxy) to the element type.
  Another paper I still have to write and
  defend in the committee 😉

## Subscripting

Loads & stores are great, but sometimes you just want to access it like an array.

```
1  void f(std::simd<float> x) {
2    for (int i = 0; i < x.size(); ++i) {
3      x[i] *= 2.f;
4      x[i] = foo(x[i]);
5      auto ref = x[i];
6      ref = foo(x[i]); // ERROR: no assignment
7    }
8  }
```

- non-const subscripting returns a basic_simd::reference
- implements all non-const operators, i.e. (compound) assignment, increment and decrement, and also swap.

- all of the above functions are rvalue-ref qualified, i.e. are *only allowed on temporaries*
- For std::vector<float> x the type of val in auto val = x is float, not float&. I.e. we expect a *decay* of the reference (proxy) to the element type. Another paper I still have to write and defend in the committee. 😉