Introduction
○○○○○

Standardization
○○○○○○○○○○○○○○○

Design
○○○○○

EDSL
○○○○○○○○○○○○○○

Extents
○○○○○○○○○○○○○○

Going beyond
○○○○○

Conclusion
○○○

# The Shapes of Multi-Dimensional Arrays

Vincent Reverdy

September 17th, 2020

# Table of contents

## Introduction

1. **Introduction**

2. Standardization

3. Design

4. EDSL

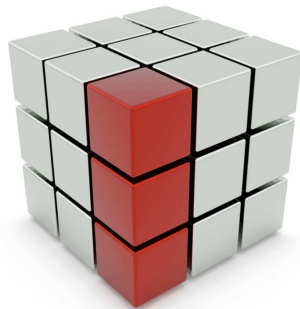5. Extents

6. Going beyond

7. Conclusion

Multidimensional arrays and linear algebra

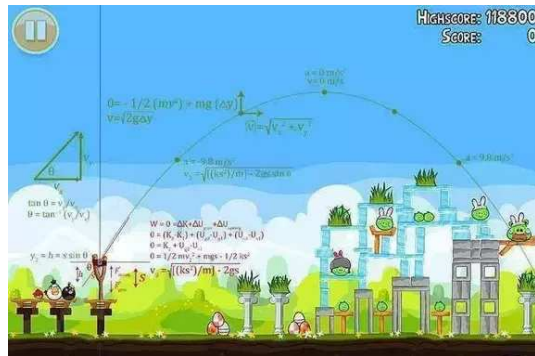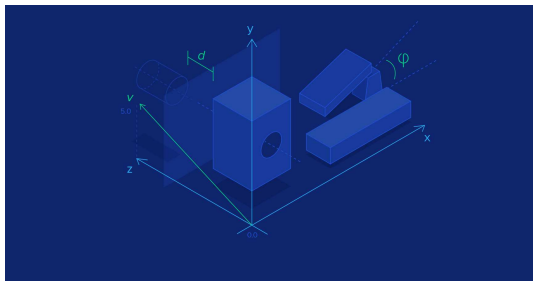$$\begin{pmatrix} 0.5 \\ 1.2 \\ 4.1 \end{pmatrix} \qquad \begin{pmatrix} 0.5 & 4.5 & -0.3 \\ 1.2 & -0.9 & 1.4 \\ 4.1 & -3.4 & 0.1 \end{pmatrix}$$

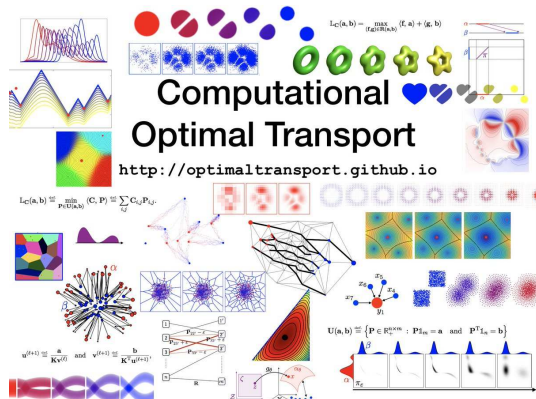<div align="center">1D          2D</div>



3D

Computer graphics and videogames: small fixed-size vectors and matrices

Personal use case with optimal transport: up to 6D dynamic arrays with billions of elements



© www.deus-consortium.org

Summary: multidimensional stuff are useful

Introduction | **Standardization** | Design | EDSL | Extents | Going beyond | Conclusion
00000 | ●000000000000 | 00000 | 0000000000000 | 0000000000000 | 00000 | 000

Standardization

1. Introduction

2. **Standardization**

3. Design

4. EDSL

5. Extents

6. Going beyond

7. Conclusion

Introduction
00000
Standardization
○●○○○○○○○○○○○○
Design
00000
EDSL
0000000000000
Extents
0000000000000
Going beyond
00000
Conclusion
000

## A standardization effort

### On mdspan

- P0332: Relaxed Incomplete Multidimensional Array Type Declaration
- P0009: mdspan: A Non-Owning Multidimensional Array Reference
- P1684: mdarray: An Owning Multidimensional Array Analog of mdspan

### On linear algebra in general

- P1417: Historical lessons for C++ linear algebra library standardization
- P1166: What do we need from a linear algebra library?
- P1385: A proposal to add linear algebra support to the C++ standard library
- P1635: A Design for an Inter-Operable and Customizable Linear Algebra Library

### On BLAS (Basic Linear Algebra Subprograms)

- P1673: A free function linear algebra interface based on the BLAS
- P1674: Evolving a Standard C++ Linear Algebra Library from the BLAS

### Miscellaneous

- P1416: Linear Algebra for Data Science and Machine Learning
- P1770: On vectors, tensors, matrices, and hypermatrices

## Naming

### Originally proposed in P0009 (with `array` instead of `span`)

- `sci_array`
- `numeric_array`
- `vla_array`
- `multidimensional_array`
- `multidim_array`
- `mdarray`/`md_array`
- `multiarray`/`multi_array` (`multi` means something different in `multi`-containers)

### Alternatives

- `ndarray`/`nd_array`: because the rank is finite and known at compile-time (used in mathematics and in python)
- `hyperarray`/`hyper_array`: because the mathematical generalization of a matrix is called an hypermatrix (same with hypercubes, hyperspheres, . . . )

## Layered approach

### Standardization by layers

- non-owning references to arrays
- basic algorithms/operations
- owning arrays
- overloaded operators
- geometric aspects
- . . .

Introduction
00000

**Standardization**
0000●000000000

Design
00000

EDSL
0000000000000

Extents
0000000000000

Going beyond
00000

Conclusion
000

The `mdspan` proposal: synopsis

<div align="center">Proposed <code>mdspan</code> synopsis as of P0009R10</div>

```
 1  namespace std {
 2    // [mdspan.extents], class template extents
 3    template<ptrdiff_t... Extents>
 4      class extents;
 5
 6    // [mdspan.layout], Layout mapping policies
 7    class layout_left;
 8    class layout_right;
 9    class layout_stride;
10
11    // [mdspan.accessor.basic]
12    template<class ElementType>
13      class accessor_basic;
14
15    // [mdspan.basic], class template mdspan
16    template<class ElementType, class Extents, class LayoutPolicy = layout_right,
17             class AccessorPolicy = accessor_basic<ElementType>>
18      class basic_mdspan;
19
20    template<class T, ptrdiff_t... Extents>
21      using mdspan = basic_mdspan<T, extents<Extents...>>;
22
23    // ...
24  }
```

## The `mdspan` proposal: `extents`

### Proposed `extents` synopsis as of P0009R10

```cpp
template<ptrdiff_t... Extents>
class extents {
public:
  using index_type = ptrdiff_t;

  // [mdspan.extents.cons], Constructors and assignment
  constexpr extents() noexcept = default;
  constexpr extents(const extents&) noexcept = default;
  constexpr extents& operator=(const extents&) noexcept = default;

  template<ptrdiff_t... OtherExtents>
    constexpr extents(const extents<OtherExtents...>&) noexcept;
  template<class... IndexType>
    constexpr extents(IndexType...) noexcept;
  template<class IndexType>
    constexpr extents(const array<IndexType, rank_dynamic()>&) noexcept;
  template<ptrdiff_t... OtherExtents>
    constexpr extents& operator=(const extents<OtherExtents...>&) noexcept;

  // [mdspan.extents.obs], Observers of the domain multidimensional index space
  static constexpr size_t rank() noexcept;
  static constexpr size_t rank_dynamic() noexcept;
  static constexpr index_type static_extent(size_t);
  constexpr index_type extent(size_t) const noexcept;

private:
  static constexpr size_t dynamic_index(size_t) noexcept; // exposition only
  array<index_type, rank_dynamic()> dynamic_extents_{}; // exposition only
};
```

Introduction
○○○○○

Standardization
○○○○○○○○○●○○○○○○○○

Design
○○○○○

EDSL
○○○○○○○○○○○○○○○○

Extents
○○○○○○○○○○○○○○○○

Going beyond
○○○○○

Conclusion
○○○

## The `mdspan` proposal: `layout`

Proposed `layout_left` synopsis as of P0009R10

```cpp
struct layout_left {
  template<class Extents>
  class mapping {
    using index_type = typename Extents::index_type; // exposition only
  public:
    constexpr mapping() noexcept = default;
    constexpr mapping(const mapping&) noexcept = default;
    constexpr mapping(const Extents&) noexcept;
    template<class OtherExtents> constexpr mapping(const mapping<OtherExtents>&) noexcept;

    constexpr mapping& operator=(const mapping&) noexcept = default;
    template<class OtherExtents> constexpr mapping& operator=(const mapping<OtherExtents>&) noexcept;

    constexpr Extents extents() const noexcept { return extents_; }
    constexpr index_type required_span_size() const noexcept;
    template<class... Indices> index_type operator()(Indices...) const noexcept;
    static constexpr bool is_always_unique() noexcept { return true; }
    static constexpr bool is_always_contiguous() noexcept { return true; }
    static constexpr bool is_always_strided() noexcept { return true; }
    constexpr bool is_unique() const noexcept { return true; }
    constexpr bool is_contiguous() const noexcept { return true; }
    constexpr bool is_strided() const noexcept { return true; }
    index_type stride(size_t) const noexcept;
    template<class OtherExtents> constexpr bool operator==(const mapping<OtherExtents>&) const noexcept;
    template<class OtherExtents> constexpr bool operator!=(const mapping<OtherExtents>& rhs) const noexcept

  private:
    Extents extents_{}; // exposition only
  };
};
```

Introduction
○○○○○

Standardization
○○○○○○○○●○○○○○○

Design
○○○○○

EDSL
○○○○○○○○○○○○○○○

Extents
○○○○○○○○○○○○○○○

Going beyond
○○○○○

Conclusion
○○○

## The `mdspan` proposal: `accessor`

### Proposed `accessor_basic` synopsis as of P0009R10

```
 1  template<class ElementType>
 2    struct accessor_basic {
 3      using offset_policy = accessor_basic;
 4      using element_type = ElementType;
 5      using reference = ElementType&;
 6      using pointer = ElementType*;
 7
 8      constexpr typename offset_policy::pointer
 9        offset(pointer p, ptrdiff_t i) const noexcept;
10
11      constexpr reference access(pointer p, ptrdiff_t i) const noexcept;
12
13      constexpr pointer decay(pointer p) const noexcept;
14    };
```

Introduction
○○○○○
Standardization
○○○○○○○○○●○○○○○
Design
○○○○○
EDSL
○○○○○○○○○○○○○○○
Extents
○○○○○○○○○○○○○○
Going beyond
○○○○○
Conclusion
○○○

# The `mdspan` proposal: `mdspan` (1/3)

## Proposed `mdspan` synopsis as of P0009R10

```
 1  template<class ElementType, class Extents, class LayoutPolicy, class AccessorPolicy>
 2  class basic_mdspan {
 3  public:
 4
 5    // Domain and codomain types
 6    using extents_type = Extents;
 7    using layout_type = LayoutPolicy;
 8    using accessor_type = AccessorPolicy;
 9    using mapping_type = typename layout_type::template mapping_type<extents_type>;
10    using element_type = typename accessor_type::element_type;
11    using value_type = remove_cv_t<element_type>;
12    using index_type = ptrdiff_t ;
13    using difference_type = ptrdiff_t;
14    using pointer = typename accessor_type::pointer;
15    using reference = typename accessor_type::reference;
16
17    // [mdspan.basic.cons], basic_mdspan constructors, assignment, and destructor
18    constexpr basic_mdspan() noexcept = default;
19    constexpr basic_mdspan(const basic_mdspan&) noexcept = default;
20    constexpr basic_mdspan(basic_mdspan&&) noexcept = default;
21    template<class... IndexType>
22      explicit constexpr basic_mdspan(pointer p, IndexType... dynamic_extents);
23    template<class IndexType, size_t N>
24      explicit constexpr basic_mdspan(pointer p, const array<IndexType, N>& dynamic_extents);
25    constexpr basic_mdspan(pointer p, const mapping_type& m);
26    constexpr basic_mdspan(pointer p, const mapping_type& m, const accessor_type& a);
27    template<class OtherElementType, class OtherExtents, class OtherLayoutPolicy, class OtherAccessorPolicy>
28      constexpr basic_mdspan(
29        const basic_mdspan<OtherElementType, OtherExtents, OtherLayoutPolicy, OtherAccessorPolicy>& other
30      );
31
32    // ...
33  };
```

## The `mdspan` proposal: `mdspan` (2/3)

### Proposed `mdspan` synopsis as of P0009R10

```cpp
1  template<class ElementType, class Extents, class LayoutPolicy, class AccessorPolicy>
2  class basic_mdspan {
3  public:
4    // ...
5    constexpr basic_mdspan& operator=(const basic_mdspan&) noexcept = default;
6    constexpr basic_mdspan& operator=(basic_mdspan&&) noexcept = default;
7    template<class OtherElementType, class OtherExtents, class OtherLayoutPolicy, class OtherAccessorPolicy>
8      constexpr basic_mdspan& operator=(
9        const basic_mdspan<OtherElementType, OtherExtents, OtherLayoutPolicy, OtherAccessorPolicy>& other
10     ) noexcept;
11
12   // [mdspan.basic.mapping], basic_mdspan mapping domain multidimensional index to access codomain element
13   constexpr reference operator[](index_type) const noexcept;
14   template<class... IndexType>
15     constexpr reference operator()(IndexType... indices) const noexcept;
16   template<class IndexType, size_t N>
17     constexpr reference operator()(const array<IndexType, N>& indices) const noexcept;
18   // ...
19 };
```

## The `mdspan` proposal: `mdspan` (3/3)

### Proposed `mdspan` synopsis as of P0009R10

```cpp
template<class ElementType, class Extents, class LayoutPolicy, class AccessorPolicy>
class basic_mdspan {
public:
  // ...
  accessor_type accessor() const;

  static constexpr int rank() noexcept;
  static constexpr int rank_dynamic() noexcept;
  static constexpr index_type static_extent(size_t r) noexcept;

  constexpr Extents extents() const noexcept;
  constexpr index_type extent(size_t r) const noexcept;
  constexpr index_type size() const noexcept;
  constexpr index_type unique_size() const noexcept;

  // [mdspan.basic.codomain], basic_mdspan observers of the codomain
  constexpr span<element_type> span() const noexcept;
  constexpr pointer data() const noexcept;

  static constexpr bool is_always_unique() noexcept;
  static constexpr bool is_always_contiguous() noexcept;
  static constexpr bool is_always_strided() noexcept;

  constexpr mapping_type mapping() const noexcept;
  constexpr bool is_unique() const noexcept;
  constexpr bool is_contiguous() const noexcept;
  constexpr bool is_strided() const noexcept;
  constexpr index_type stride(size_t r) const;
};
```

## The `mdspan` proposal: summary

### `mdspan` is the product of

- an element type `ElementType`
- an `Extents` to describe its shape
- a `LayoutPolicy` to specify the mapping between the indices and storage
- an `AccessorPolicy` to specify how individual and contiguous elements can be accessed in memory

### Simplifications. . .

- it is non-owning
- it has been designed for contiguous memory storage

# The `mdspan` proposal: specifying the shape

**Example `mdspan` declaration**

```
1  // As of P0009R10
2  using myspan = std::mdspan<double, extents<3, std::dynamic_extent, 7>>;
3
4  // With P0332: Relaxed Incomplete Multidimensional Array Type Declaration
5  using myspan = std::mdspan<double[][3][7]>; // already allowed
6  using myspan = std::mdspan<double[3][][7]>; // currently not allowed
```

**Standardization**

- The Library Evolution Working Group (LEWG) was largely in favor of `std::mdspan<double[3][][7]>`
- The Evolution Working Group (EWG) was against a language change because of unintended side effects in other contexts

**Zero-sized arrays**

- *"Its value N specifies the array bound, i.e., the number of elements in the array; N shall be greater than zero."*
- $\Rightarrow$ `double[3][7][0]` not allowed
- However `std::array<double, 0>` is a compiling special case

## The shapes of multi-dimensional arrays

### What are we looking for?

A way to specify the shapes of multi-dimensional arrays

### Goals

- Genericity: cover as much as possible of the parameter space
- Performance: both in terms of computing time and memory
- Expressivity: easy to and understand, read, and write in a concise way

## Design

1. Introduction

2. Standardization

3. **Design**

4. EDSL

5. Extents

6. Going beyond

7. Conclusion

## The shapes of multi-dimensional arrays

### What are we looking for?

A way to specify the shapes of multi-dimensional arrays

### Design goals: the GPE principle

- **G**enericity: cover as much as possible of the parameter space
- **P**erformance: both in terms of computing time and memory
- **E**xpressivity: easy to and understand, read, and write in a concise way

### However. . .

One cannot have everything at the same time

## Software design

### Software architecture

The art of balancing and compromising between genericity, performance, and expressivity.

### Design goals

- **G**enericity: cover as much as possible of the parameter space
- **P**erformance: both in terms of computing time and memory
- **E**xpressivity: easy to and understand, read, and write in a concise way

### Design complexity

- Focusing on only one axis drastically simplifies the problem (for example focusing on expressivity in a high level language)
- Focusing on two axes can still lead to reasonable level of complexity (for example performance and expressivity)
- Addressing the three axes at the same time is what can make the design a real challenge

### C++

The standard C++ library has to handle the three axes.
$$\Rightarrow \text{Still no } \texttt{std::matrix} \text{ in 2020....}$$

## What do we want (ideally) for a n-dimensional array?



### Expressivity (pseudo-code)

```
1  std::ndarray<double> a;
2  std::ndarray<double[0]> b;
3  std::ndarray<double[3][4][5]> c;
4  std::ndarray<double[3][][5]> d;
5  std::ndarray<double[3][][5], properties...> d;
```

### Genericity

- Arbitrary memory layouts
- Arbitrary indexing schemes
- Arbitrary finite ranks
- Combinations of static and dynamic dimensions
- Zero-sized case

### Performance

No computing-time/memory overhead of generic solutions compared to handcrafted solutions designed and optimized for specific scenarios

What do we need? (hint: C++20)

1. Introduction

2. Standardization

3. Design

4. **EDSL**

5. Extents

6. Going beyond

7. Conclusion

## What if. . .

> . . . the following was perfectly possible in C++20. . .
>
> - `std::ndarray<double, shape[4]()[3][5]> myarray;`
>
> - or even: `std::ndarray<contents<double>[4]()[3][5]> myarray;`

Introduction
00000

Standardization
0000000000000

Design
00000

EDSL
0000000000000000

Extents
0000000000000

Going beyond
00000

Conclusion
000
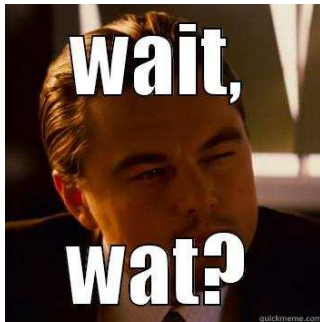
## Reverse-engineering from the expressivity goal

### With shape:

```
std::ndarray<double, shape[4]()[3][5]> myarray;
```

- shape cannot be a type, it has to be a variable
- std::ndarray should be a class template of the form `template <class, auto>`
- shape should be a variable of a type with `operator()` and `operator[]`
- `operator()`/`operator[]` should return a value of the same type so that `operator()`/`operator[]` can be called again recursively until all dimensions are specified

### With contents:

```
std::ndarray<contents<double>[4]()[3][5]> myarray;
```
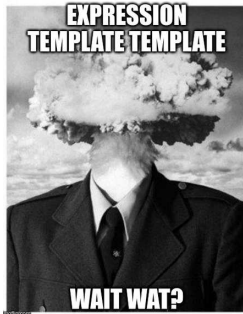
- contents cannot be a type, it has to be a variable template
- std::ndarray should be a class template of the form `template <auto>`
- contents<T> should be a variable of a type with `operator()` and `operator[]`
- `operator()`/`operator[]` should return a value of the same type so that `operator()`/`operator[]` can be called again recursively until all dimensions are specified

## Introducing expression template template



### Principle

- `shape` builds an expression template
- The result of the expression template is then injected as a template parameter

### Warning

- Takes EDSL (Embedded Domain Specific Languages) to a whole new level

Building the EDSL: tags and shaper constructors

### Tagging operators to know what operator has been called

```
1  // A tag marking the function call operator
2  template <class... Args> struct function_call_operator {};
3
4  // A tag marking the subscripting operator
5  template <class... Args> struct subscripting_operator {};
```

### Shaper class template

```
1   // The type of a shape remembering the sequence of operations it was made from
2   template <class... Ops> class shaper {
3       // Types
4       public:
5       using type = shaper<Ops...>;
6       using index_sequence = std::index_sequence_for<Ops...>;
7
8       // Lifecycle
9       public:
10      constexpr shaper() noexcept = default;
11      template <class... Args> constexpr shaper(shaper<Args...> other, std::size_t i) noexcept
12      : _data(other.to_array(i)) {}
13
14      // Implementation details
15      private:
16      std::array<std::size_t, sizeof...(Ops)> _data;
17
18      // ...
19  };
```

## Building the EDSL: shaper operations

### Shaper class template

```cpp
// The type of a shape remembering the sequence of operations it was made from
template <class... Ops> class shaper {
    // ...
    // Operators
    public:
    constexpr auto operator()() const noexcept {
        return shaper<Ops..., function_call_operator<>>(*this, -1);
    }
    constexpr auto operator()(std::size_t i) const noexcept {
        return shaper<Ops..., function_call_operator<std::size_t>>(*this, i);
    }
    constexpr auto operator[](std::size_t i) const noexcept {
        return shaper<Ops..., subscripting_operator<std::size_t>>(*this, i);
    }
    // Management
    public:
    static constexpr std::size_t rank() noexcept {
        return sizeof...(Ops);
    }
    constexpr std::size_t at(std::size_t i) const noexcept {
        return _data[i];
    }
    template <class... X> constexpr auto to_array(X... x) const noexcept {
        return to_array(index_sequence(), x...);
    }
    template <std::size_t... I, class... X> constexpr auto to_array(std::index_sequence<I...>, X... x) const noexcept {
        return std::array<std::size_t, sizeof...(Ops) + sizeof...(X)>{{std::get<I>(_data)..., x...}};
    }
};

// Default shaper instance
inline constexpr shaper shape;
```

Building the EDSL: indexed dynamic extent

## Indexing helper

```
1  // Index constant alias template
2  template <std::size_t I> using index_constant = std::integral_constant<std::size_t, I>;
3
4  // Index variable template
5  template <std::size_t I> inline constexpr index_constant<I> index;
```

## A dynamic extent

```
1  // A dynamic extent
2  template <auto...> struct indexed_dynamic_extent;
3
4  // A dynamic extent with no default value
5  template <std::size_t Index> struct indexed_dynamic_extent<Index> {
6      using is_constant = std::false_type;
7      constexpr auto operator[](index_constant<Index>) const noexcept {
8          return *this;
9      }
10 };
11
12 // A dynamic extent with a default value for initialization
13 template <std::size_t Index, std::size_t Value> struct indexed_dynamic_extent<Index, Value>
14 : std::integral_constant<std::size_t, Value> {
15     using is_constant = std::false_type;
16     constexpr auto operator[](index_constant<Index>) const noexcept {
17         return *this;
18     }
19 };
```

## Building the EDSL: indexed static extent

<div align="center">A static extent</div>

```
1  // A static extent with its value
2  template <std::size_t Index, std::size_t Value>
3  struct indexed_static_extent
4  : std::integral_constant<std::size_t, Value> {
5      using is_constant = std::true_type;
6      constexpr auto operator[](index_constant<Index>) const noexcept {
7          return *this;
8      }
9  };
```

## Building the EDSL: indexed extent maker

```
                          Indexed extent maker
 1  // A helper to build an indexed extent from a tagged operator
 2  template <class, std::size_t, auto...>
 3  struct indexed_extent_maker;
 4
 5  // Converts the function call operator into a dynamic extent
 6  template <std::size_t I, auto... Args>
 7  struct indexed_extent_maker<function_call_operator<>, I, Args...> {
 8      using type = indexed_dynamic_extent<I>;
 9  };
10
11  // Converts the function call operator into a dynamic extent with default value
12  template <std::size_t I, std::size_t V, auto... Args>
13  struct indexed_extent_maker<function_call_operator<std::size_t>, I, V, Args...> {
14      using type = indexed_dynamic_extent<I, V>;
15  };
16
17  // Converts the subscripting operator into a static extent
18  template <std::size_t I, std::size_t V, auto... Args>
19  struct indexed_extent_maker<subscripting_operator<std::size_t>, I, V, Args...> {
20      using type = indexed_static_extent<I, V>;
21  };
22
23  // Alias template to build an indexed extent from a tagged operator
24  template <class Op, std::size_t I, auto... Args>
25  using make_indexed_extent = typename indexed_extent_maker<Op, I, Args...>::type;
```

Introduction
00000

Standardization
00000000000000

Design
00000

EDSL
000000000000000

Extents
00000000000000

Going beyond
00000

Conclusion
000

## Building the EDSL: extent policy

### Extent policy class template

```
 1  // The extents, static and dynamic, of a N-dimensional array
 2  template <class... Extents>
 3  struct extent_policy
 4  : Extents... {
 5      using index_sequence = std::index_sequence_for<Extents...>;
 6      using Extents::operator[]...;
 7      static constexpr std::size_t rank() noexcept {
 8          return sizeof...(Extents);
 9      }
10  };
11
12  // Extent policy maker declaration
13  template <auto, class, class>
14  struct extent_policy_maker;
15
16  // Specialization of extent policy maker for a shape, a shaper, and an index sequence
17  template <auto S, class... Ops, std::size_t... I>
18  struct extent_policy_maker<S, shaper<Ops...>, std::index_sequence<I...>> {
19      using type = extent_policy<make_indexed_extent<Ops, I, S.at(I)>...>;
20  };
21
22  // Make an extent policy type from a shape variable
23  template <auto S, class T = decltype(S), class I = typename T::index_sequence>
24  using make_extent_policy = typename extent_policy_maker<S, typename T::type, I>::type;
```

## Building the EDSL: extend

### Extend class template

```cpp
// A class to specify the dynamic size along an axis
template <std::size_t Index>
class indexed_extender
{
    // Constants
    public:
    using is_extender = std::true_type;

    // Lifecycle
    public:
    constexpr indexed_extender() noexcept = default;
    explicit constexpr indexed_extender(std::size_t size) noexcept
    : _value(size) {
    }

    // Assignment and access
    public:
    constexpr indexed_extender operator=(std::size_t size) const noexcept {
        return indexed_extender(size);
    }
    constexpr operator std::size_t() const {
        return _value;
    }
    constexpr auto operator[](index_constant<Index>) const noexcept {
        return *this;
    }

    // Implementation details
    private:
    std::size_t _value;
};
```

Introduction
○○○○○

Standardization
○○○○○○○○○○○○○○○

Design
○○○○○

EDSL
○○○○○○○○○○○○○●○○

Extents
○○○○○○○○○○○○○○

Going beyond
○○○○○

Conclusion
○○○

## Building the EDSL: extenders

<div align="center">Extenders class template</div>

```cpp
 1  // A class grouping extenders
 2  template <class...>
 3  struct extenders;
 4
 5  // A class grouping extenders specialized for extenders
 6  template <std::size_t... I>
 7  struct extenders<indexed_extender<I>...>
 8  : indexed_extender<I>... {
 9      using indexed_extender<I>::operator[]...;
10      explicit constexpr extenders(indexed_extender<I>... i) noexcept
11      : indexed_extender<I>(i)... {
12      }
13  };
14
15  // Extenders deduction guide
16  template <std::size_t... I>
17  extenders(indexed_extender<I>... i) -> extenders<indexed_extender<I>...>;
18
19  // Extend variable template
20  template <std::size_t I>
21  inline constexpr indexed_extender<I> extend(-1);
```

Introduction
○○○○○

Standardization
○○○○○○○○○○○○○○

Design
○○○○○

EDSL
○○○○○○○○○○○○○●○

Extents
○○○○○○○○○○○○○○○

Going beyond
○○○○○

Conclusion
○○○

## Building the EDSL: ndarray

### Ndarray class template

```
 1  // Generic ndarray declaration
 2  template <class...>
 3  class basic_ndarray;
 4
 5  // Simplified ndarray definition
 6  template <class Type, class ExtentPolicy>
 7  class basic_ndarray<Type, ExtentPolicy>
 8  {
 9      // Types and constants
10      public:
11      using value_type = Type;
12      using extent_policy = ExtentPolicy;
13
14      // Lifecycle
15      public:
16      template <class... Extenders>
17      explicit constexpr basic_ndarray(Extenders... e) noexcept;
18
19      // Access
20      public:
21      constexpr extents<extent_policy::rank()> shape();
22
23      // ...
24  };
25
26  // Ndarray alias template
27  template <class Type, auto Shape, class... Args>
28  using ndarray = basic_ndarray<Type, make_extent_policy<Shape>, Args...>;
```

Introduction
○○○○○

Standardization
○○○○○○○○○○○○○○

Design
○○○○○

EDSL
○○○○○○○○○○○○○●

Extents
○○○○○○○○○○○○○○

Going beyond
○○○○○

Conclusion
○○○

## Illustration

### Usage in practice

```
 1  int main(int argc, char* argv[])
 2  {
 3      // A ndarray of:
 4      // Dimension-0: a static size of 3
 5      // Dimension-1: a dynamic size whose initial value is given in constructor
 6      // Dimension-2: a static size of 4
 7      // Dimension-3: a dynamic size with a default initial value of 2
 8      // Dimension-4: a static size of 3
 9      // And then the constructor initializes the dimension 1 to size 5
10      ndarray<int, shape[3]()[4](2)[3]> myarray(extend<1> = 5);
11      return 0;
12  }
```

Introduction    Standardization    Design    EDSL    **Extents**    Going beyond    Conclusion
○○○○○       ○○○○○○○○○○○○○    ○○○○○    ○○○○○○○○○○○○○    ●○○○○○○○○○○○○    ○○○○○       ○○○

Extents

1. Introduction

2. Standardization

3. Design

4. EDSL

5. **Extents**

6. Going beyond

7. Conclusion

Introduction
○○○○○

Standardization
○○○○○○○○○○○○○○

Design
○○○○○

EDSL
○○○○○○○○○○○○○○

**Extents**
○●○○○○○○○○○○○○

Going beyond
○○○○○

Conclusion
○○○

## The problem with extents

### Static only

When all dimensions are static, extents can just be an integer sequence ⇒ `sizeof`(extents) should be the size of an empty struct.

### Dynamic only

When all dimensions are dynamic, extents can just be a runtime array of size rank ⇒ `sizeof`(extents) = `sizeof`(std::size_t) × rank.

### Mixed static/dynamic

Well, it's a little bit more complicated...

### Optimization goal

- Memory footprint of only dynamic dimensions
- Optimal runtime access of a particular dimension
- Bonus: optimized compilation time

## Challenges

### Access

- `extents[3]`: dynamic access to the extent along axis 3
- `extents[n]`: dynamic access to the extent along axis n (runtime variable)
- `extents[index<3>]`: static access to the extent along axis 3
- `extents[index<N>]`: static access to the extent along axis N (constant expression)

### Challenges

- `extents[n]`: avoiding $\mathcal{O}(n)$ scaling for static dimensions
- `extents[index<N>]`: avoiding $\mathcal{O}(N)$ scaling in compilation-time

## Log-tuple trick: introduction

### Straightforward approach

get<N>(tuple) has to iterate over the first N types.

### Advanced approach

There is a way to exploit overload resolution to have $\mathcal{O}\left(\log\left(N\right)\right)$ compile-time access.

```
                                    Indexing
1  // Index constant type
2  template <std::size_t I>
3  struct index_constant: std::integral_constant<std::size_t, I> {};
4
5  // Index constant variable template
6  template <std::size_t I>
7  inline constexpr index_constant<I> index = {};
```

Introduction
00000

Standardization
0000000000000

Design
00000

EDSL
00000000000000

**Extents**
0000●00000000

Going beyond
00000

Conclusion
000

## Log-tuple trick: elements

<div align="center">Element wrappers</div>

```cpp
 1  // A basic element wrapper
 2  template <class T>
 3  struct tuple_element_wrapper {
 4      constexpr tuple_element_wrapper(const T& x): value(x) {}
 5      // Other constructors to be defined
 6      T value;
 7  };
 8
 9  // An indexed tuple element
10  template <std::size_t I, class T>
11  struct tuple_element: tuple_element_wrapper<T> {
12      constexpr tuple_element(const T& x): tuple_element_wrapper<T>(x) {}
13      constexpr T& operator[](index_constant<I>) {
14          return static_cast<wrapper<T>&>(*this).value;
15      }
16      constexpr const T& operator[](index_constant<I>) const {
17          return static_cast<const wrapper<T>&>(*this).value;
18      }
19  };
```

## Log-tuple trick: tuple

### Tuple

```
 1  // Base class declaration
 2  template <class Sequence, class... T>
 3  struct tuple_base;
 4
 5  // Base class specialization for index sequence
 6  template <std::size_t... I, class... T>
 7  struct tuple_base<std::index_sequence<I...>, T...>
 8  : tuple_element<I, T>... {
 9      using index_sequence = std::index_sequence<I...>;
10      using tuple_element<I, T>::operator[]...;
11      constexpr tuple_base(const T&... x): tuple_element<I, T>(x)... {}
12      // Other constructors to be defined
13  };
14
15  // Actual tuple implementation
16  template <class... T>
17  struct tuple: tuple_base<std::index_sequence_for<T...>, T...> {
18      using base = tuple_base<std::index_sequence_for<T...>, T...>;
19      using base::base;
20      using base::operator[];
21  };
22  template <class... T>
23  tuple(const T&...) -> tuple<T...>;
```

### Result

`mytuple[index<3>]` leverages overload resolution to access the element at compile-time.

Runtime element access of a static sequence: introduction

#### The problem

Considers a tuple $\mathcal{T}$ of types Types... (where, for instance each type can be either a static or dynamic extent). How to design a runtime subscript operator for this tuple, so that when provided with an index $i$ and a generic lambda $\Lambda$, it would return the result of the lambda applied to the i-th element of the tuple: $\Lambda(\mathcal{T}_i)$?

#### The trivial approach

```
1  i == 0 ? lambda(std::get<0>(tuple)) :
2    i == 1 ? lambda(std::get<1>(tuple)) :
3      i == 2 ? lambda(std::get<2>(tuple)) :
4        i == 3 ? lambda(std::get<3>(tuple)) :
5          // ...
6            lambda(std::get<N − 1>(tuple));
```

#### The trivial non-branching approach

```
1  (i == 0) * lambda(std::get<0>(tuple)) +
2    (i == 1) * lambda(std::get<1>(tuple)) +
3      (i == 2) * lambda(std::get<2>(tuple)) +
4        (i == 3) * lambda(std::get<3>(tuple)) +
5          // ...
6            (i == N − 1) * lambda(std::get<N − 1>(tuple));
```

## Runtime element access of a static sequence: complexity

---

**Access complexity**

The preceding strategies are $\mathcal{O}(N)$-access strategies. Can we do better?

---

### Dichotomic approach (for N = 6)

```
1  i < 3
2    ? i < 2
3      ? i == 0
4        ? lambda(std::get<0>(tuple))
5        : lambda(std::get<1>(tuple))
6      : lambda(std::get<2>(tuple))
7    : i < 5
8      ? i == 3
9        ? lambda(std::get<3>(tuple))
10       : lambda(std::get<4>(tuple))
11     : lambda(std::get<5>(tuple));
```

---

**Access complexity**

The preceding strategy scales as $\mathcal{O}(\log(N))$.

Runtime element access of a static sequence: variadic implementation

---

**Expressivity goal**

`overload_sequence(F...)[i](Args...)` should return the result of the i-th function F on the arguments `Args...`.

---

### Starting with the log-tuple trick

```
 1  // Index constant helper
 2  template <std::size_t I>
 3  struct index_constant: std::integral_constant<std::size_t, I> {};
 4
 5  // Index constant variable template
 6  template <std::size_t I>
 7  inline constexpr index_constant<I> index = {};
 8
 9  // Element wrapper
10  template <std::size_t I, class F>
11  struct overload_sequence_element: F {
12      static constexpr index_type<I> index = {};
13      constexpr overload_sequence_element(const F& f): F(f) {}
14      constexpr F& operator[](index_constant<I>) {return static_cast<F&>(*this);}
15      constexpr const F& operator[](index_constant<I>) const {return static_cast<const F&>(*this);}
16  };
```

Introduction
○○○○○

Standardization
○○○○○○○○○○○○○○

Design
○○○○○

EDSL
○○○○○○○○○○○○○○

**Extents**
○○○○○○○○○○●○○○○

Going beyond
○○○○○

Conclusion
○○○

Runtime element access of a static sequence: overload sequence

### Overload sequence implementation

```
 1  // Overload sequence base declaration
 2  template <class...>
 3  struct overload_sequence_base;
 4
 5  // Overload sequence base definition
 6  template <std::size_t... I, class... F>
 7  struct overload_sequence_base<std::index_sequence<I...>, F...>
 8  : overload_sequence_element<I, F>... {
 9      using index_sequence = std::index_sequence<I...>;
10      using overload_sequence_element<I, F>::operator[]...;
11      constexpr overload_sequence_base(const F&... args)
12      : overload_sequence_element<I, F>(args)... {}
13      static constexpr std::size_t size() {return sizeof...(F);}
14  };
15
16  // Overload sequence
17  template <class... F>
18  struct overload_sequence:
19  overload_sequence_base<std::index_sequence_for<F...>, F...> {
20      using base = overload_sequence_base<std::index_sequence_for<F...>, F...>;
21      using base::base;
22      using base::operator[];
23      constexpr overload_sequence_selector<overload_sequence&> operator[](std::size_t i) {
24          return overload_sequence_selector<overload_sequence&>(*this, i);
25      }
26      constexpr overload_sequence_selector<overload_sequence&> operator[](std::size_t i) const {
27          return overload_sequence_selector<const overload_sequence&>(*this, i);
28      }
29  };
30  template <class... F>
31  overload_sequence(const F&...) -> overload_sequence<F...>;
```

Runtime element access of a static sequence: dichotomic access

**Overload sequence selector implementation**

```
1  // Dichotomic access implementation
2  template <class Sequence>
3  struct overload_sequence_selector {
4      constexpr overload_sequence_selector(Sequence seq, std::size_t i): sequence(seq), index(i) {}
5      template <class... Args>
6      constexpr void operator()(Args&&... args) {
7          this->operator()<0, std::decay_t<Sequence>::size() - 1>(std::forward<Args>(args)...);
8      }
9      template <int Min, int Max, int Mid = (Min + Max) / 2, class... Args>
10     constexpr void operator()(Args&&... args) {
11         if constexpr (Min < Max) {
12             if (index < Mid) {this->operator()<Min, Mid - 1>(std::forward<Args>(args)...);}
13             else if (index > Mid) {this->operator()<Mid + 1, Max>(std::forward<Args>(args)...);}
14             else {this->operator()<Mid, Mid>(std::forward<Args>(args)...);}
15         } else {
16             sequence[index_type<Mid>{}](std::forward<Args>(args)...);
17         }
18     }
19     Sequence sequence;
20     std::size_t index;
21 };
```

Runtime element access of a static sequence: result

### Result

For a family of lambdas indexed by $i$ and each associated with an element $\mathcal{T}_i$ of a tuple $\mathcal{T}$, overload_sequence provides a runtime generic access in $\mathcal{O}\left(\log\left(N\right)\right)$.

### Application

For an heterogeneous sequence of dynamic_extent and static_extent it provides a $\mathcal{O}\left(\log\left(N\right)\right)$ runtime access with branching.

## Beyond logarithmic complexity

---

**Overload sequence selector implementation**

```
 1  // Dichotomic access implementation
 2  template <class... Extents>
 3  struct extents {
 4      static constexpr std::array<bool, sizeof...(Extents)> is_static = {Extents::is_static...};
 5      static constexpr std::array<std::size_t, sizeof...(Extents)> table = // see after;
 6      std::array<std::size_t, dynamic_rank> storage = // see after;
 7      constexpr std::size_t operator[](std::size_t i) {
 8          return is_static ? table[i] : storage[table[i]];
 9      }
10  };
```

---

**Where:**

- `storage` contains the dynamic extents known at runtime
- `table` contains:
    - `static_extent::value` if the i-th extent is static
    - the index `j` where the element is stored in `storage` if the i-th extent is dynamic

---

**Access complexity**

For an heterogeneous sequence of `dynamic_extent` and `static_extent` it provides a $\mathcal{O}(1)$ runtime access with a branching and a level of indirection.

## Optimizing extents

### Summary

- Optimizing dynamic-only extents is easy
- Optimizing static-only extents is easy
- Optimizing hybrid extents is tricky

### Axes of optimization

- In space: storing only dynamic extents at runtime
- In time: using `overload_sequence` or indirection to optimize runtime access
- In compilation time: using and reusing the log-tuple technique

### `mdspan`

The current implementation of `extents.extent(i)` seems to scale as $\mathcal{O}(N)$ (to be confirmed).

Introduction    Standardization    Design    EDSL    Extents    **Going beyond**    Conclusion
00000    0000000000000    00000    0000000000000    0000000000000    ●0000    000

Going beyond

1. Introduction

2. Standardization

3. Design

4. EDSL

5. Extents

6. **Going beyond**

7. Conclusion

## Enriching the EDSL: sky is the limit

### Layout information

- `std::ndarray<double, +shape[4]()[3][5]>`
- `std::ndarray<double, −shape[4]()[3][5]>`
- `std::ndarray<double, !shape[4]()[3][5]>`

### Symmetries

- `std::ndarray<double, shape[diagonal(4)]()[3][5]>`
- `std::ndarray<double, shape[4]()[diagonal(3)][5]>`
- `std::ndarray<double, shape<triangular>[4]()[3][5]>`

### Axis-based parallelization information

- `std::ndarray<double, shape[distributed(4)]()[3][5]>`
- `std::ndarray<double, shape[distributed(4)]()[3][vectorized(5)]>`
- `std::ndarray<double, shape[4]()[3][vectorized<gpu>(5)]>`

### Operations on shapes

- `std::ndarray<double, shape[4] * shape() * shape[5]>`

## With multidimensional subscript operator

### P2128

- P2128: Multidimensional subscript operator
- `operator[](Args&&...)`

### Evolution

- `shape[4]()[3][5]` $\Rightarrow$ `shape[4][][3][5]` (requires zero-parameter operator)
- Other opportunities to enrich the mini-language by combining `operator()` and `operator[]` for different meanings

## Towards a generic extent concept

#### Possible components

- `min`: the minimum allowed size under which it cannot shrink
- `initial`: the initial size at construction
- `threshold`: the maximum number of elements stored without external allocation
- `max`: the maximum allowed size over which is cannot be extended (can be infinite)

#### Concrete examples

- `static_extent` $\Rightarrow$ `min == initial == threshold == max`
- `dynamic_extent` $\Rightarrow$ `min == 0, max == infinity, initial == 0, threshold == 0`
- `fixed_extent` $\Rightarrow$ `min == 0, max == n, initial == 0, threshold == n`

## Generic NTTPs opened a Pandora's box

### Metalanguages

- Being able to inject objects as template parameters in C++20 is a game changer for Embedded Domain Specific Languages
- Class templates can now be seen as true mini-compilers: `template` <`auto`> `class` edsl_compiler ;

### Kind genericity

- The next revolution to come is probably the genericity regarding kinds
- Currently no way to specify universal template parameters (values, types, templates...)
- P1985: Universal template parameters
- `template` <`template auto` X> where X can be a value, a type, a template...

Conclusion

1. Introduction

2. Standardization

3. Design

4. EDSL

5. Extents

6. Going beyond

7. Conclusion

Introduction
○○○○○

Standardization
○○○○○○○○○○○○○○

Design
○○○○○

EDSL
○○○○○○○○○○○○○○

Extents
○○○○○○○○○○○○○○

Going beyond
○○○○○

Conclusion
○●○

## Summary

### Summary

- Linear algebra is everywhere
- Very complex when trying to balance between Genericity, Performance, and Expressivity
- Existing effort with `mdspan`
- C++20 with arbitrary NTTPs allows expressive ways to specify multidimensional shapes
- Expression template template to create mini-languages
- Far more powerful than relaxed incomplete multidimensional array type declaration
- Hybrid dynamic and static extent particularly tricky to implement
- Log-tuple trick useful in many situations
- Extending the EDSL: layout information, symmetries, axis-based parallelization, shape operations
- Conceptifying a generalized extent
- Class NTTPs opened a Pandora's box for EDSL
- Kind genericity will take EDSLs even further

Introduction
00000
Standardization
00000000000000
Design
00000
EDSL
0000000000000
Extents
00000000000000
Going beyond
00000
Conclusion
00●

# Thank you for your attention

## Any question?

### Acknowledgments