# Statistical Scientific programming: challenges in converting R to C++

Olivia Quinet
olivia.quinet @cluepoints.com

Meeting C++
BeCPP July 2019

Credit xkcd

CluePoints

Statistical
Scientific
programming

Olivia Quinet

Introduction
CluePoints
Clinical Trials
Statistical tests
SMART package

The R language
A very short introduction to R
Some examples

R to C++?

Scientific
programming
challenges

Testing
Measure
Software architecture
Data structure
Smart pointers, Pimpl,
Factories
Fail-fast/Fail-safe
Numerical errors
Accumulators
std::algorithms, boost, GSL,
BLAS, LAPACK, . . .

Conclusion

Questions,
Remarks?

# CluePoints

*CluePoints is the premier provider of Risk-Based Monitoring and Data Quality Oversight Software. Our products utilize unique statistical algorithms to determine the quality, accuracy, and integrity of clinical trial data both during and after study conduct.*

# Clinical Trials

**Why?**

- ▶ Is it working?
- ▶ Is it safe?
- ▶ Scope, dosage, . . .
- ▶ Approuval from FDA, EMA, . . .
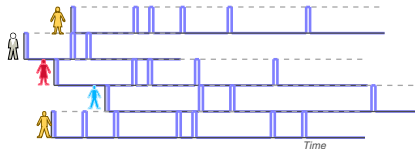
# Clinical Trials

**How?**

- ▶ Clinical protocol
- ▶ Study conducted at sites
- ▶ Patients are enrolled
- ▶ Data is collected: demographics, medical history, vital signs, adverse events, labs, patient journals, ...
- ▶ Data is verified and analyzed



Sites



Patients



Visits



Datasets

# Clinical Trials

**$$$?**

- 1.5-2.5 billion on 10-plus years
- 30% for sending investigators on sites

# Statistical tests

# SMART package

- Initially developed in the R language by the R&D team

# SMART package

- Initially developed in the R language by the R&D team
- Very good for research purposes

# SMART package

- Initially developed in the R language by the R&D team
- Very good for research purposes
- Not so much for production

# SMART package

- Initially developed in the R language by the R&D team
- Very good for research purposes
- Not so much for production
- Need for something reliable, robust and fast

# SMART package

- ▶ Initially developed in the R language by the R&D team
- ▶ Very good for research purposes
- ▶ Not so much for production
- ▶ Need for something reliable, robust and fast

# The R language

# The R language

- R is a programing language for statisticians created by statisticians
- R is weakly/dynamically typed
- R operates on named data structures: vector, matrix, array, data frame, factors, lists, objects, functions
- It is very concise
- Lot of statistical libraries

# Some examples (I)

```
1  w = !is.na(d[[field]]);
2  ctr = factor(d$center[w]);
3  npat = unclass(table(ctr));
4  v = d[w, field];
5  y = rowsum(v, ctr);
```

# Some examples (I)

```
1  w = !is.na(d[[field]]);
2  ctr = factor(d$center[w]);
3  npat = unclass(table(ctr));
4  v = d[w, field];
5  y = rowsum(v, ctr);
```

1. Select the rows where the values of the column "field" are not missing
2. Get the values as factor of the column "center" for the selected rows, i.e. the list of centers
3. Count the number of rows associated to the different centers, i.e. the number of patients per center
4. Get the values for the column "field" for the selected rows
5. Sum by center the values from (4)

# Some examples (I)

```
1  w = !is.na(d[[field]]);
2  ctr = factor(d$center[w]);
3  npat = unclass(table(ctr));
4  v = d[w, field];
5  y = rowsum(v, ctr);
```

| center | xyz |
|--------|-----|
| ctr01  |     |
| ctr02  | 1   |
| ctr01  | 2   |
| ctr03  | 3   |
| ctr05  | 4   |
| ctr02  |     |
| ctr02  | 5   |
| ctr02  |     |
| ...    |     |

| ctr   | npat | y  |
|-------|------|----|
| ctr01 | 1    | 2  |
| ctr02 | 2    | 7  |
| ctr03 | 5    | 10 |
| ...   |      |    |

# Some examples (II)

```
1  xc = x - offset;
2  v = tapply(xc, ctr, mean, na.rm=T);
3  Sn = unclass(table(ctr));
4  Sn2 = tapply(sid, ctr, function(i) sum(table(i)^2));
5  sigma = sqrt(Sn*vc[3]^2 + Sn2*vc[2]^2 + Sn^2*vc[1]^2)/Sn;
6  p = pnorm(v, sd=sigma)
```

# Some examples (II)

```
1  xc = x - offset;
2  v = tapply(xc, ctr, mean, na.rm=T);
3  Sn = unclass(table(ctr));
4  Sn2 = tapply(sid, ctr, function(i) sum(table(i)^2));
5  sigma = sqrt(Sn*vc[3]^2 + Sn2*vc[2]^2 + Sn^2*vc[1]^2)/Sn;
6  p = pnorm(v, sd=sigma)
```

1. Apply an offet to x
2. Compute per center the mean of xc
3. Number of records per center
4. Compute per center the sum of the squares of the number of values per patient
5. Compute sigma
6. Compute the p-values for each center based on a normal distribution

# Some examples (II)

```
1  xc = x - offset;
2  v = tapply(xc, ctr, mean, na.rm=T);
3  Sn = unclass(table(ctr));
4  Sn2 = tapply(sid, ctr, function(i) sum(table(i)^2));
5  sigma = sqrt(Sn*vc[3]^2 + Sn2*vc[2]^2 + Sn^2*vc[1]^2)/Sn;
6  p = pnorm(v, sd=sigma)
```

| center | subjid | x |
|--------|--------|---|
| ctr01 | s01001 | |
| ctr02 | s02001 | 1 |
| ctr01 | s01001 | 2 |
| ctr03 | s03001 | 3 |
| ctr05 | s05001 | 4 |
| ctr02 | s02002 | |
| ctr02 | s02001 | 5 |
| ctr02 | s02001 | |
| ... | | |

| ctr | Sn | Sn2 | sigma | p |
|-----|-----|-----|-------|------|
| ctr01 | 1 | 1 | 0.37 | 0.21 |
| ctr02 | 3 | 5 | 0.25 | 0.55 |
| ctr03 | 5 | 5 | 0.19 | 0.06 |
| ... | | | | |

# Some examples (III)

```
1  dd = duplicated(d$subjid);
2  v = d[[field]]
3  w = dd & c(FALSE, v[1:(length(v)-1)]==1);
4  x10 = rowsum(1-v[w], ctr[w]);
5  N10 = unclass(table(ctr[w]));
6  x10Max = rowsum(as.integer(!c(TRUE, w[1:(length(w)-1)])[w]), ctr[w]);
```

# Some examples (III)

```
1 dd = duplicated(d$subjid);
2 v = d[[field]]
3 w = dd & c(FALSE, v[1:(length(v)-1)]==1);
4 x10 = rowsum(1-v[w], ctr[w]);
5 N10 = unclass(table(ctr[w]));
6 x10Max = rowsum(as.integer(!c(TRUE, w[1:(length(w)-1)])[w]), ctr[w]);
```

1. Create a boolean vector indicating if a subjid is duplicated or not
2. Get the values of the column "field"
3. Do some wierd selection
4. Get the number of transitions $1\rightarrow0$ per center for each patient
5. Get the number of potential transitions $1\rightarrow0$ per center
6. Get the maximum number of valid transitions $1\rightarrow0$ per center

# Some examples (III)

```
1  dd = duplicated(d$subjid);
2  v = d[[field]]
3  w = dd & c(FALSE, v[1:(length(v)-1)]==1);
4  x10 = rowsum(1-v[w], ctr[w]);
5  N10 = unclass(table(ctr[w]));
6  x10Max = rowsum(as.integer(!c(TRUE, w[1:(length(w)-1)])[w]), ctr[w]);
```

| center | subjid | visit | x |
|--------|--------|-------|---|
| ctr01 | s01001 | v01 | 1 |
| ctr01 | s01001 | v02 | 0 |
| ctr01 | s01001 | v03 | 1 |
| ctr01 | s01001 | v04 | 0 |
| ctr01 | s01002 | v01 | 1 |
| ctr01 | s01002 | v02 | 1 |
| ctr02 | s02001 | v03 | 0 |
| ctr02 | s02002 | v01 | 1 |
| ctr02 | s02002 | v02 | 1 |
| ctr02 | s02002 | v03 | 0 |
| ctr03 | s03001 | v01 | 1 |
| ctr03 | s03001 | v02 | 0 |
| ctr03 | s03001 | v03 | 1 |
| ... | | | |

| ctr | X10 | N10 | x10max |
|-----|-----|-----|--------|
| ctr01 | 2 | 3 | 3 |
| ctr02 | 1 | 2 | 1 |
| ctr03 | 1 | 1 | 1 |
| ... | | | |

# How to translate R code to C++ code?

- Straightforward approach: Recode each R function in C++
  PRO C++ and R codes are similar
  CON Too many combinations of parameters/structures

# How to translate R code to C++ code?

- Straightforward approach: Recode each R function in C++
  PRO C++ and R codes are similar
  CON Too many combinations of parameters/structures

- Hard: understanding what the researcher wanted to do
  PRO Faster code
  CON C++ and R codes can be very different
      1 line in R $\longrightarrow \pm 30$ lines in C++

# How to translate R code to C++ code?

- ▶ Straightforward approach: Recode each R function in C++
  PRO C++ and R codes are similar
  CON Too many combinations of parameters/structures

- ▶ Hard: understanding what the researcher wanted to do
  PRO Faster code
  CON C++ and R codes can be very different
  1 line in R $\longrightarrow$ $\pm 30$ lines in C++

- ▶ Hardest: changing the data structure
  PRO Less ressource/faster code
  CON C++ and R codes are even more different

# How to translate R code to C++ code?

- ▶ Straightforward approach: Recode each R function in C++
  PRO  C++ and R codes are similar
  CON  Too many combinations of parameters/structures

- ▶ Hard: understanding what the researcher wanted to do
  PRO  Faster code
  CON  C++ and R codes can be very different
       1 line in R $\longrightarrow$ $\pm 30$ lines in C++

- ▶ Hardest: changing the data structure
  PRO  Less ressource/faster code
  CON  C++ and R codes are even more different

- ▶ Recoding model fitting algorithms is a huge (tremendous) task. It's easier to call the R function from the C++ code
  PRO  Updates of the fitting model code
  CON  Added dependencies

# How to translate R code to C++ code?

- ▶ Straightforward approach: Recode each R function in C++
  - PRO C++ and R codes are similar
  - CON Too many combinations of parameters/structures

- ▶ Hard: understanding what the researcher wanted to do
  - PRO Faster code
  - CON C++ and R codes can be very different
    - 1 line in R $\longrightarrow \pm 30$ lines in C++

- ▶ Hardest: changing the data structure
  - PRO Less ressource/faster code
  - CON C++ and R codes are even more different

- ▶ Recoding model fitting algorithms is a huge (tremendous) task. It's easier to call the R function from the C++ code
  - PRO Updates of the fitting model code
  - CON Added dependencies

- ▶ Beware of Numerical (in)accuracy

# How to translate R code to C++ code?

- ▶ Straightforward approach: Recode each R function in C++
  - PRO C++ and R codes are similar
  - CON Too many combinations of parameters/structures

- ▶ Hard: understanding what the researcher wanted to do
  - PRO Faster code
  - CON C++ and R codes can be very different
    1 line in R ⟶ ±30 lines in C++

- ▶ Hardest: changing the data structure
  - PRO Less ressource/faster code
  - CON C++ and R codes are even more different

- ▶ Recoding model fitting algorithms is a huge (tremendous) task. It's easier to call the R function from the C++ code
  - PRO Updates of the fitting model code
  - CON Added dependencies

- ▶ Beware of Numerical (in)accuracy

- ▶ Testing and testing and testing (no data, invalid data, NaN, Inf, . . . )

# Scientific programming challenges

# Scientific programming challenges

- Requirements include low response time and memory usage, minimizing numerical errors and error propagation.
- Testing
- Software architecture
- Data structure
- Fail-fast/Fail-safe idioms
- Exceptions
- RAII
- Pimpl idiom and smart pointers
- Factory pattern
- Iterator pattern and accumulators
- std::algorithms, boost, GSL, BLAS, LAPACK, . . .

# Testing



Credit xkcd

# Testing

- ▶ Framework
- ▶ Unit testing, Integration testing, . . .
- ▶ *Test Driven Development*
- ▶ *Behavior Driven Development* to replicate the documentation specification
- ▶ Continuous Integration

# Testing

- ▶ Framework
- ▶ Unit testing, Integration testing, . . .
- ▶ *Test Driven Development*
- ▶ *Behavior Driven Development* to replicate the documentation specification
- ▶ Continuous Integration
- ▶ Each bug must be tested

# Code for Testing

- If you cannot test your code, rewrite it
- If you cannot test your code easily, rewrite it
- If you cannot test your code independently, rewrite it
- ...

Tools like clang static analyzer and gcov/lcov code coverage are a great help

# Measure!!!



Credit xkcd

CluePoints
Statistical
Scientific
programming

Olivia Quinet

Introduction
CluePoints
Clinical Trials
Statistical tests
SMART package

The R language
A very short introduction to R
Some examples

R to C++?

Scientific programming challenges
Testing
Measure
Software architecture
Data structure
Smart pointers, Pimpl, Factories
Fail-fast/Fail-safe
Numerical errors
Accumulators
std::algorithms, boost, GSL, BLAS, LAPACK, . . .

Conclusion

Questions, Remarks?

# Measure!!!

- ▶ Select between different data structures
- ▶ Select between different algorithms
- ▶ Use generated data
- ▶ Use real data
- ▶ Use data of different sizes
- ▶ . . .

# Measure!!! – an example

**Context:** Originally, an algorithm has to be applied on vectors: $f(x, y)$

# Measure!!! – an example

CluePoints
Statistical
Scientific
programming

Olivia Quinet

Introduction
CluePoints
Clinical Trials
Statistical tests
SMART package

The R language
A very short introduction to R
Some examples

R to C++?

Scientific programming challenges
Testing
Measure
Software architecture
Data structure
Smart pointers, Pimpl, Factories
Fail-fast/Fail-safe
Numerical errors
Accumulators
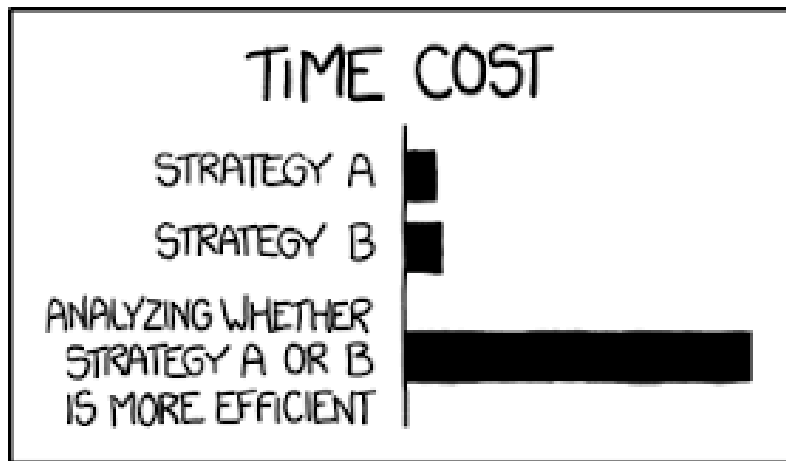std::algorithms, boost, GSL, BLAS, LAPACK, ...

Conclusion

Questions, Remarks?

**Context:** Originally, an algorithm has to be applied on vectors: $f(x, y)$
Then only on some filtered elements: $f(x, y, w)$

| X | Y | w |
|------|------|-------|
| 42.5 | 100 | true |
| ... | ... | true |
| ... | ... | false |
| ... | ... | true |
| ... | ... | false |
| ... | ... | true |
| ... | ... | true |
| ... | ... | false |
| ... | ... | ... |

# Measure!!! – an example

**Context:** Originally, an algorithm has to be applied on vectors: $f(x, y)$
Then only on some filtered elements: $f(x, y, w)$

- ▶ Modify the algorithm to take into account only the filtered vectors'elements: *filter algo*
- ▶ Create pseudo vectors with the filtered elements: *filter vector*
- ▶ Create new vectors with the filtered elements: *copy vector*

# Measure!!! – an example

**Context:** Originally, an algorithm has to be applied on vectors: $f(x, y)$
Then only on some filtered elements: $f(x, y, w)$

▶ Modify the algorithm to take into account only the filtered vectors'elements: *filter algo*

▶ Create pseudo vectors with the filtered elements: *filter vector*

▶ Create new vectors with the filtered elements: *copy vector*

| Option | Timing (s) | | | |
|---|---|---|---|---|
| | $N = 10^2$ | $N = 10^4$ | $N = 10^6$ | $N = 10^8$ |
| *filter algo* | 0.0003 | 0.008 | 0.9 | 100 |
| *filter vector* | 0.0003 | 0.006 | 0.8 | 98 |
| *copy vector* | 0.0006 | 0.015 | 4.6 | / |

# Software architecture & Data structure

# Software architecture & Data structure

**Important points to consider**

- Input/output data structure?
- *Computational units*?
- Simple but not too simple!
- Which doors are you closing?
- Expressiveness

# Software architecture & Data structure

**Important points to consider**

- ► Input/output data structure?
- ► *Computational units*?
- ► Simple but not too simple!
- ► Which doors are you closing?
- ► Expressiveness

**For this project**

- ► Data is organized in datasets, i.e. tables in which each column represents a particular variable or key variable, and each row corresponds to a given record. There may also be missing values.
- ► Statistical tests are the *computational units*.

# Levels of abstraction

- The **most** important good practice
- Divide and conquer
- *Top down* design
- *Bottom up* design
- *Separation of concerns*
- *Modularity*: *low coupling* $\longleftrightarrow$ *high cohesion*
- *Design review*

# Levels of abstraction – mathematical formula

# Levels of abstraction – mathematical formula

- ▶ Very tempting to code one mathematical formula into one function.
- ▶ Decompose the formula into meaningful steps, e.g. numerator, denominator, partial sums, . . .
- ▶ Transform the function into a class
- ▶ Transform each step into a struct

# Abstraction levels – Example

Sample variance – Standard formula

$$s_N^2 = \frac{1}{N-1} \sum_{k=1}^{N} (x_k - \bar{x})^2$$

## Abstraction levels – Example

Sample variance – Standard formula

$$s_N^2 = \underbrace{\frac{1}{N-1}}_{\text{Fraction}} \times \underbrace{\sum_{k=1}^{N} \underbrace{\left(x_k - \underbrace{\bar{x}}_{\text{Mean}}\right)^2}_{\text{Square of difference}}}_{\text{Sum}}}_{\text{Multiplication}}$$

# Abstraction levels – Example

**CluePoints**
Statistical
Scientific
programming

Olivia Quinet

Introduction
CluePoints
Clinical Trials
Statistical tests
SMART package

The R language
A very short introduction to R
Some examples

R to C++?

Scientific
programming
challenges
Testing
Measure
Software architecture
Data structure
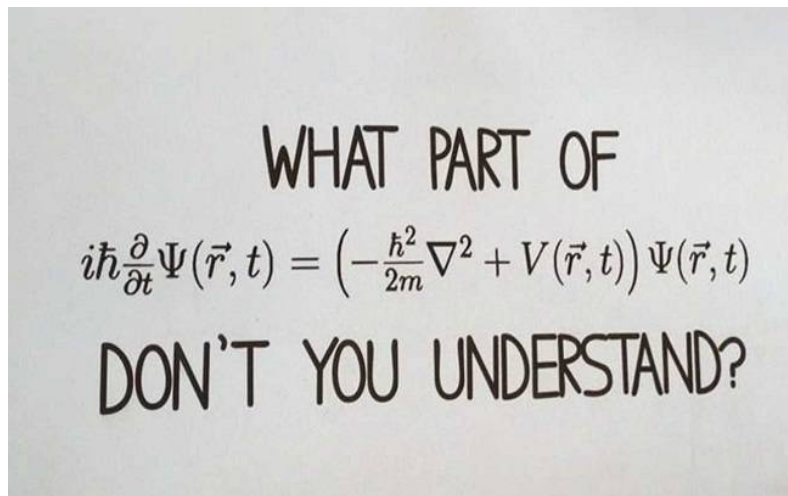Smart pointers, Pimpl, Factories
Fail-fast/Fail-safe
Numerical errors
Accumulators
std::algorithms, boost, GSL, BLAS, LAPACK, …

Conclusion

Questions,
Remarks?

```cpp
 1  namespace MATH_INTERNAL {
 2    template<typename T=double>
 3    struct sample_variance {
 4      T s2{std::numeric_limits<T>::quiet_NaN()};
 5
 6      template<typename Container>
 7      sample_variance(const Container& X) { if(X.size()>1) s2 = frac(X)*sum(X, mean(X)); }
 8
 9      operator T() const { return s2; }
10
11      template<typename Container>
12      static T frac(const Container& X) { return ONE/(X.size()-ONE); }
13
14      template<typename Container>
15      static T mean(const Container& X) { return std::accumulate(X.begin(), X.end(), ZERO)/X.size(); }
16
17      struct square_of_difference {
18        const T xbar;
19        square_of_difference(const T mean) : xbar(mean) {}
20        T operator()(const T x) const { return (x-xbar)*(x-xbar); }
21      };
22
23      template<typename Container>
24      static T sum(const Container& X, const T xbar) {
25        const square_of_difference d(xbar);
26        return std::accumulate(X.begin(), X.end(), ZERO, [d](const T s, const T x) { return s + d(x); });
27      }
28    };
29  }
30  template<typename Container>
31  inline typename Container::value_type sample_variance(const Container& X)
32  {
33    return MATH_INTERNAL::sample_variance<typename Container::value_type>(X);
34  }
```

# Data structure

- ▶ Performance requires well thought data structure
- ▶ Cache usage
- ▶ Prefetching
- ▶ Lazy evaluation
- ▶ Sparse representation
- ▶ ...

# Data structure – Example

**Duplicate patients:** comparing patient's fingerprints

# Data structure – Example

- $N$ numerical vectors of the same length $M$
  Typical cases: $N = 1000 - 40000$ and $M = 20 - 20000$

# Data structure – Example

- $N$ numerical vectors of the same length $M$
  Typical cases: $N = 1000 - 40000$ and $M = 20 - 20000$
- $N \times (N-1)/2$ scalar products

$$X \cdot Y = \sum_{i}^{M} X_i \times Y_j$$

# Data structure – Example

- $N$ numerical vectors of the same length $M$
  Typical cases: $N = 1000 - 40000$ and $M = 20 - 20000$

- $N \times (N - 1)/2$ scalar products

$$X \cdot Y = \sum_i^M X_i \times Y_j$$

- Sparse vectors!!!

## Data structure – Example

- $N$ numerical vectors of the same length $M$
  Typical cases: $N = 1000 - 40000$ and $M = 20 - 20000$
- $N \times (N-1)/2$ scalar products

$$X \cdot Y = \sum_i^M X_i \times Y_j$$

- Sparse vectors!!!
- Performance results

|  | $N = 841$ $M = 1060$ | | $N = 35613$ $M = 14304$ | |
|---|---|---|---|---|
|  | Memory | Timing | Memory | Timing |
| R |  | $\pm 1$m |  | / |
| C++ (normal vectors) | 57MB | 0.68s | 9.8GB | 29m |
| C++ (sparse vectors) | 34MB | 0.45s | 6.6GB | 38s |

# Smart pointers, Pimpl, Factory pattern

# Smart pointers, Pimpl, Factory pattern

*Pointer to implementation* or *Private implementation*

# Smart pointers, Pimpl, Factory pattern

*Pointer to implementation* or *Private implementation*

**PROS**

- ▶ Separate interface from implementation
- ▶ Decrease recompilation cycles
- ▶ Binary compatibility of shared libraries

# Smart pointers, Pimpl, Factory pattern

*Pointer to implementation* or *Private implementation*

## PROS

- ▶ Separate interface from implementation
- ▶ Decrease recompilation cycles
- ▶ Binary compatibility of shared libraries

## CONS

- ▶ Increase in memory usage
- ▶ Increase in maintenance effort
- ▶ Performance loss
- ▶ Doesn't work well with templates

# Smart pointers, Pimpl, Factory pattern

- ▶ std::unique_ptr or std::shared_ptr?
- ▶ *Mutable* or *non mutable* objects?
- ▶ Access to the objects, how often?
- ▶ Multiple inheritance, virtual inheritance (diamond problem)?
- ▶ Template member functions, template classes?
- ▶ Objects in a coherent state!

# Smart pointers, Pimpl, Factory pattern : Inheritance

# Smart pointers, Pimpl, Factory pattern: Diamond problem

# Smart pointers, Pimpl, Factory pattern: Template members

CluePoints
Statistical
Scientific
programming

Olivia Quinet

Introduction
CluePoints
Clinical Trials
Statistical tests
SMART package

The R language
A very short introduction to
R
Some examples

R to C++?

Scientific
programming
challenges
Testing
Measure
Software architecture
Data structure
Smart pointers, Pimpl,
Factories
Fail-fast/Fail-safe
Numerical errors
Accumulators
std::algorithms, boost, GSL,
BLAS, LAPACK, ...

Conclusion

Questions,
Remarks?

### api.hpp

```cpp
1 class MyPublic : public Pimpl {
2 public:
3   class impl;
4
5   MyPublic(...);
6
7   template<typename Tp> Tp as() const;
8 };
```

### api.cpp

```cpp
1 MyPublic::MyPublic(...) : Pimpl(impl::build(...)) {}
2
3 template<>
4 double MyPublic::as() const { return _valid_ptr() ? _get_ptr<impl>()->asNumber() : NaN(); }
5
6 template<>
7 std::string MyPublic::as() const { return _valid_ptr() ? _get_ptr<impl>()->asString() : std::string(); }
```

### api-impl.hpp

```cpp
1 class MyPublic::impl : public Pimpl::impl {
2 public:
3   static PTR build(...) { return make_pimpl<impl>(...); }
4
5   double asNumber() const { return ...; }
6   std::string asString() const { return ...; }
7
8 protected:
```

# Pimpl: use

```
1  typedef std::vector<Test> TESTS;
2
3  // Create the tests
4  TESTS tests;
5  tests.push_back(BetaBinomial(...));
6  tests.push_back(CountField(...));
7  tests.push_back(CountField(...));
8  tests.push_back(CNR_ByCenter(...));
9  ...
10
11 // Export the results
12 json_ostream os(...);
13 print_results(os, tests);
14
15
16 void print_results(json_ostream& os, const TESTS& tests)
17 {
18   for(const auto& test: tests) {
19     ...
20   }
21 }
```

# Pimpl: use

CluePoints

Statistical Scientific programming

Olivia Quinet

Introduction
CluePoints
Clinical Trials
Statistical tests
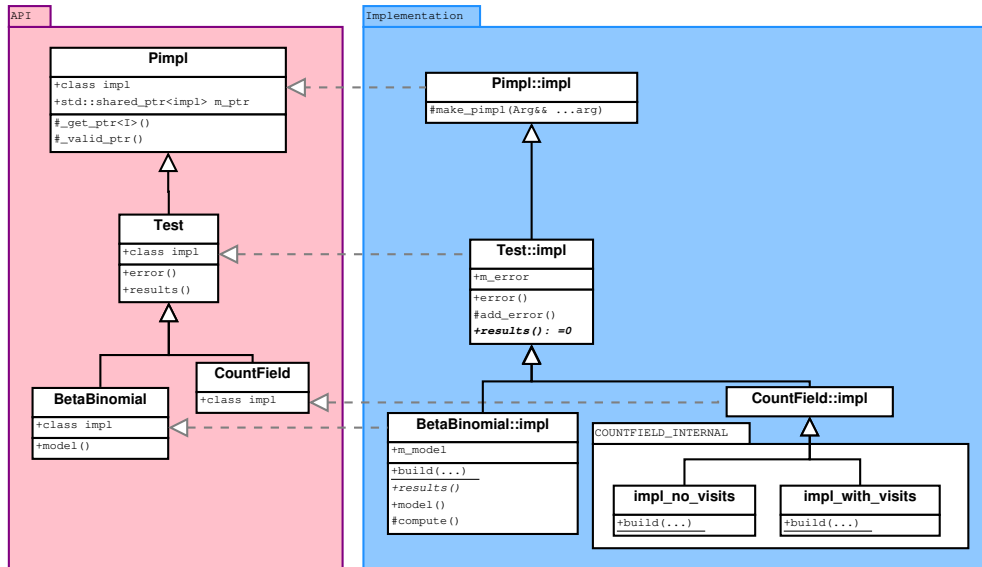SMART package

The R language
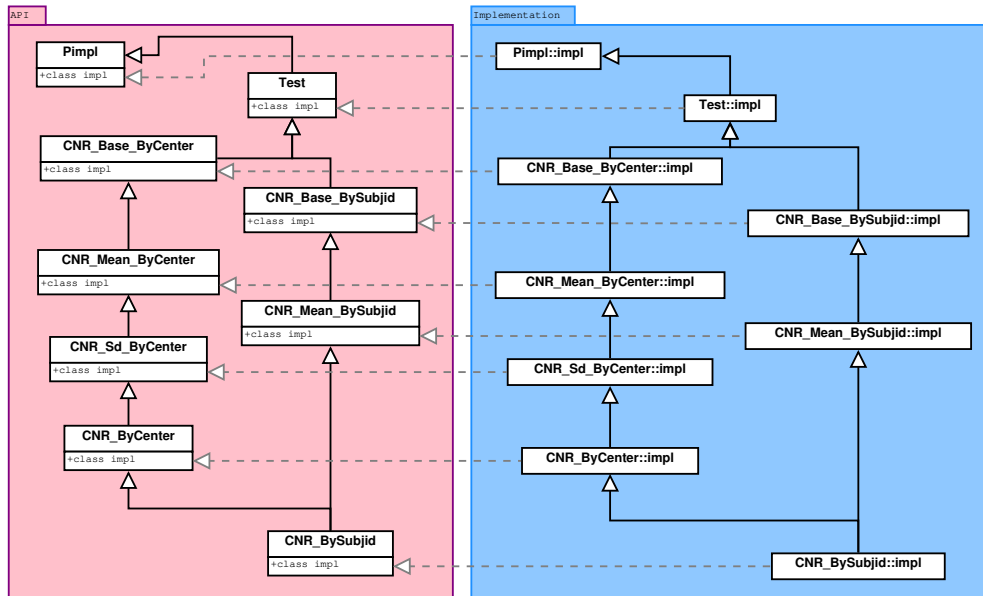A very short introduction to R
Some examples

R to C++?

Scientific programming challenges
Testing
Measure
Software architecture
Data structure
Smart pointers, Pimpl, Factories
Fail-fast/Fail-safe
Numerical errors
Accumulators
std::algorithms, boost, GSL, BLAS, LAPACK, . . .

Conclusion

Questions, Remarks?

# Fail-fast/Fail-safe idioms

# Fail-fast/Fail-safe idioms

▶ Check constraints on input/output

```
1  double foo(const std::vector<size_t>& l, const std::vector<double>& x, const std::vector<bool>& w)
2  {
3    CP_ASSERT(l.size() == x.size());
4    CP_ASSERT(l.size() == w.size());
5    // Rest of the code
6  }
```

# Fail-fast/Fail-safe idioms

- ▶ Check constraints on input/output

```
1  double foo(const std::vector<size_t>& l, const std::vector<double>& x, const std::vector<bool>& w)
2  {
3    CP_ASSERT(l.size() == x.size());
4    CP_ASSERT(l.size() == w.size());
5    // Rest of the code
6  }
```

- ▶ Fitting of statistical models might fails

```
1  try {
2    fit = vglm("cbind(a,b)~1",
3              Named("family", family),
4              Named("data", dateframe),
5              Named("control", control(Named("criterion", "coef"),
6                                       Named("stepsize", 0.5))));
7  } catch(std::exception& e) {
8    // Retry with other parameters
9  }
```

# Fail-fast/Fail-safe idioms

▶ Check constraints on input/output

```
1  double foo(const std::vector<size_t>& l, const std::vector<double>& x, const std::vector<bool>& w)
2  {
3    CP_ASSERT(l.size() == x.size());
4    CP_ASSERT(l.size() == w.size());
5    // Rest of the code
6  }
```

▶ Fitting of statistical models might fails

```
1  try {
2    fit = vglm("cbind(a,b)~1",
3              Named("family", family),
4              Named("data", dateframe),
5              Named("control", control(Named("criterion", "coef"),
6                                       Named("stepsize", 0.5))));
7  } catch(std::exception& e) {
8    // Retry with other parameters
9  }
```

▶ Propagate the error message
  ▶ Rethrow the exception
  ▶ Store the exception as an error message inside the object
  ▶ . . .

# Numerical instabilities



Credit xkcd

# Numerical instabilities

**Test on standard deviations**

- ▶ P values computed from the integration of two functions:

$$f_1(x) = pchisq(s/x^2; N, left.tail) \times dgamma(x; scale, shape)$$

$$f_2(x) = pchisq(s/x^2; N, right.tail) \times dgamma(x; scale, shape)$$



$f_1(x)$           $f_2(x)$

# Numerical instabilities

**calcPsd: test on standard deviations**

- $f_1(x)$ is unstable in case *shape* $< 1$



$f_1(x)$

# Numerical instabilities

**calcPsd: test on standard deviations**

- $f_1(x)$ is unstable in case *shape* $< 1$
- $f_1(x)$ can be rewritten by using the integration by parts theorem
  $\int_0^a u\,dv = [uv]_0^a - \int_0^a v\,du$

$$f_{1'}(x) = \frac{2s}{x^3} \times dchisq(s/x^2; N) \times pgamma(x; scale, shape, left.tail)$$

$f_1(x)$; $f_{1'}(x)$

# Minimizing numerical errors

- Sample variance – Standard formula

$$s_N^2 = \frac{1}{N-1} \sum_{k=1}^{N} (x_k - \bar{x})^2$$

# Minimizing numerical errors

- Sample variance – Standard formula

$$s_N^2 = \frac{1}{N-1} \sum_{k=1}^{N} (x_k - \bar{x})^2$$

  - Can be implemented as a 2 pass algorithm, first the mean $\bar{x}$, and the variance $s^2$ afterwards.

# Minimizing numerical errors

- Sample variance – Standard formula

$$s_N^2 = \frac{1}{N-1} \sum_{k=1}^{N} (x_k - \bar{x})^2$$

  - Can be implemented as a 2 pass algorithm, first the mean $\bar{x}$, and the variance $s^2$ afterwards.
    **BUT** the number of items $N$ can be huge

# Minimizing numerical errors

- Sample variance – Standard formula

$$s_N^2 = \frac{1}{N-1} \sum_{k=1}^{N} (x_k - \bar{x})^2$$

  - Can be implemented as a 2 pass algorithm, first the mean $\bar{x}$, and the variance $s^2$ afterwards.
    **BUT** the number of items $N$ can be huge
  - Have a one pass algorithm

# Minimizing numerical errors

- Sample variance – Standard formula

$$s_N^2 = \frac{1}{N-1} \sum_{k=1}^{N} (x_k - \bar{x})^2$$

  - Can be implemented as a 2 pass algorithm, first the mean $\bar{x}$, and the variance $s^2$ afterwards.
    **BUT** the number of items $N$ can be huge
  - Have a one pass algorithm
  - Compute the variance for increasing $N$ to observe convergence.

# Minimizing numerical errors

▶ Sample variance – One pass algorithm: *Sum of squares method*

$$s_N^2 = \frac{1}{N(N-1)} \left( N \sum_{k=1}^{N} x_k^2 - \left( \sum_{k=1}^{N} x_k \right)^2 \right)$$

One pass algorithm but the formula is unstable:

▶ *float precision:* for $\{10000f, 10001f, 10002f\}$,
the result is $-1.0666667\mathrm{e}+01$ instead of 1.

▶ *double precision:* for $\{100000000, 100000001, 100000002\}$,
the result is 0 instead of 1.

## Minimizing numerical errors

▶ Sample variance – Iterative algorithm: *Welford's recursion method*

$$
\begin{aligned}
M_k &= M_{k-1} + \frac{x_k - M_{k-1}}{k} \\
S_k &= S_{k-1} + (x_k - M_{k-1})(x_k - M_k),
\end{aligned}
$$

with $M_0 = 0$ and $S_0 = 0$, and then

$$
s_N^2 = \frac{S_N}{N-1},
$$

This stable algorithm with can be easily turned into an accumulator

# Accumulators

# Accumulators

- ▶ Separation of the operations on the elements from the iteration leads to smaller testable code.
- ▶ Statisticals tests involve operations (agregation, sum, average, variance, . . . ) on one or more variables based on one or more several key variables. E.g.: Preprocess involves taking the mean by visits or the sum by patients, the count of non missing values per center, . . .

# Accumulators as a OO pattern

Mean of the elements of a vector

▶ without accumulator

```
1    double sum{0};
2    for(const auto x: myvector) sum += x;
3    const double mean = sum / myvector.size();
```

# Accumulators as a OO pattern

Mean of the elements of a vector

- ▶ without accumulator

```
1    double sum{0};
2    for(const auto x: myvector) sum += x;
3    const double mean = sum / myvector.size();
```

- ▶ with accumulator

```
1    using namespace boost::accumulators;
2    accumulator_set<double, stats<tag::mean> > acc;
3    for(const auto x: myvector) acc(x);
4    mean(acc);
```

## Accumulator implementation

**Sample variance – Welford's recursion method**

$$
\begin{aligned}
M_k &= M_{k-1} + \frac{x_k - M_{k-1}}{k} \\
S_k &= S_{k-1} + (x_k - M_{k-1})(x_k - M_k) \\
s_k^2 &= \frac{S_k}{k-1},
\end{aligned}
$$

```cpp
1  template<typename T> class Variance {
2    size_t k; /**< Number of elements */
3    T m; /**< 0th order moment, i.e. average */
4    T s; /**< 1st order moment */
5  public:
6    Variance() : k(0), m(0), s(0) {}
7    void operator()(const T x) {
8      if(std::isnan(x)) return;
9      ++k;
10     const T pm(m);
11     m += (x-pm) * (ONE/k);
12     s += (x-pm) * (x-m);
13   }
14   T average() const noexcept { return m; }
15   T s2() const noexcept { return k > 1 ? s / (k-1) : ZERO; }
16 };
```

# std::algorithms, boost, GSL, BLAS, LAPACK, . . .

CluePoints
Statistical
Scientific
programming

Olivia Quinet

Introduction
CluePoints
Clinical Trials
Statistical tests
SMART package

The R language
A very short introduction to R
Some examples

R to C++?

Scientific
programming
challenges
Testing
Measure
Software architecture
Data structure
Smart pointers, Pimpl, Factories
Fail-fast/Fail-safe
Numerical errors
Accumulators
std::algorithms, boost, GSL, BLAS, LAPACK, . . .

Conclusion

Questions, Remarks?

Credit Big Bang Theory

# std::algorithms, boost, GSL, BLAS, LAPACK, . . .

- `<algorithm>`
- *boost*
    - Statistics, . . .
    - Logging facilities
    - System (command line arguments, . . . )
    - Thread, MPI, Serialization
    - . . .
- *GNU Scientific Library*
    - Optimization, minimization, . . .
- *BLAS* and *LAPACK*
    - Operations on matrices
- *Numerical Recipes*
    - Lots of algorithms

# Implementing Beta-Binomial distribution with boost

In probability theory and statistics, the beta-binomial distribution is a family of discrete probability distributions on a finite support of non-negative integers arising when the probability of success in each of a fixed or known number of Bernoulli trials is either unknown or random.

# Implementing Beta-Binomial distribution with boost

$$f(k; n, \alpha, \beta) = \binom{n}{k} \frac{B(k + \alpha, n - k + \beta)}{B(\alpha, \beta)}$$

- $k$ and $n$ are positive integers with $k <= n$
- $\alpha$ and $\beta$ are strictly positive numbers
- Binomial coefficient

$$\binom{n}{k} = \frac{n!}{k!(n-k)!} = \frac{\Gamma(n+1)}{\Gamma(k+1)\Gamma(n-k+1)}$$

- Beta function

$$B(x, y) = \frac{\Gamma(x) + \Gamma(y)}{\Gamma(x + y)}$$

# Implementing Beta-Binomial distribution with boost

$$f(k; n, \alpha, \beta) = \binom{n}{k} \frac{B(k + \alpha, n - k + \beta)}{B(\alpha, \beta)}$$

- Numerically fine as long as $\alpha$ and $\beta$ are small
- When $\alpha$ and $\beta$ are not small, $B(\alpha, \beta)$ tends toward zero.

# Implementing Beta-Binomial distribution with boost

$$f(k; n, \alpha, \beta) = \binom{n}{k} \frac{B(k + \alpha, n - k + \beta)}{B(\alpha, \beta)}$$

▶ Numerically fine as long as $\alpha$ and $\beta$ are small
▶ When $\alpha$ and $\beta$ are not small, $B(\alpha, \beta)$ tends toward zero.

**Trick:** Do the calculation in the log scale:

$$f(k; n, \alpha, \beta) = \exp\left( \log \binom{n}{k} + \log B(k + \alpha, n - k + \beta) - \log B(\alpha, \beta) \right)$$

$$\log \binom{n}{k} = \text{l\_binomial\_coefficient}(n, k)$$
$$= \text{lgamma}(n + 1) - \text{lgamma}(k + 1) - \text{lgamma}(-k + n + 1)$$

$$\log B(x, y) = \text{lbeta}(x, y)$$
$$= \text{lgamma}(x) + \text{lgamma}(y) - \text{lgamma}(x + y)$$

# Minimize with GSL

- ▶ GSL is a C library
- ▶ Use wrapper classes

- ▶ Pointer to the minimizer created/owned by GSL
- ▶ Pointer to the function definition struct

# Minimizers – Example

```
1    enum class M_TYPE {
2      NO_GRADIENT,
3      ...
4    };
5
6    template<M_TYPE> struct M_API; /**< Template for the C API */
7    template<M_TYPE> class M_FCT; /**< Template for defining the function to minimize */
8    template<M_TYPE> class IMinimizer; /**< Template for the minimizer */
```

# Minimizers – Example

**CluePoints**

Statistical
Scientific
programming

Olivia Quinet

Introduction
CluePoints
Clinical Trials
Statistical tests
SMART package

The R language
A very short introduction to
R
Some examples

R to C++?

Scientific
programming
challenges
Testing
Measure
Software architecture
Data structure
Smart pointers, Pimpl,
Factories
Fail-fast/Fail-safe
Numerical errors
Accumulators
std::algorithms, boost, GSL,
BLAS, LAPACK, …

Conclusion

Questions,
Remarks?

```cpp
1  /// Specialized template defining the function to minimize (no gradient)
2  template<>
3  class M_FCT<M_TYPE::NO_GRADIENT> {
4  public:
5    friend class IMinimizer<M_TYPE::NO_GRADIENT>;
6
7    /// Virtual base class for the implementation of the function to minimize
8    class Base : public boost::noncopyable {
9    public:
10   virtual ~Base() {}
11   virtual double evaluate(const double _x) = 0;
12   };
13
14   typedef std::unique_ptr<Base> PTR; /**< Type of the pointer to the instance of the function to minimize */
15   typedef gsl_function DEF; /**< Type for the definition */
16
17   M_FCT(PTR _fct, const NUMBER _minimum, const NUMBER _lower, const NUMBER _upper) : ...
18
19   double evaluate(const double _x) { return m_fct->evaluate(_x); }
20   double get_lowest_bound() const { return m_f_lower < m_f_upper ? m_lower : m_upper; }
21   double get_lowest_f_bound() const { return std::min(m_f_lower, m_f_upper); }
22
23  private:
24    ...
25  };
```

# Minimizers – Example

CluePoints
Statistical
Scientific
programming

Olivia Quinet

Introduction
CluePoints
Clinical Trials
Statistical tests
SMART package
The R language
A very short introduction to
R
Some examples
R to C++?
Scientific
programming
challenges
Testing
Measure
Software architecture
Data structure
Smart pointers, Pimpl,
Factories
Fail-fast/Fail-safe
Numerical errors
Accumulators
std::algorithms, boost, GSL,
BLAS, LAPACK, . . .
Conclusion
Questions,
Remarks?

```cpp
1  template<M_TYPE _Type> class IMinimizer : public boost::noncopyable {
2  public:
3    typedef M_FCT<_Type> FCT;
4
5    explicit IMinimizer(const std::string& _type, FCT& _fct, const double _epsabs, const double _epsrel) { ... }
6    ~IMinimizer() { ... }
7
8    bool iterate(const size_t _maxiter) {
9      if(m_can_minimize) {
10       for(size_t iter = 0; iter<_maxiter && next(); ++iter) {
11         if(converged()) return true;
12       }
13     }
14     return false;
15   }
16
17   std::string name() const
18   { return m_ptr ? M_API<_Type>::name(m_ptr) : std::string(); }
19   bool next()
20   { return m_ptr && m_can_minimize ? M_API<_Type>::iterate(m_ptr) == 0 : false; }
21   bool converged() const
22   { return m_ptr && m_can_minimize ? M_API<_Type>::converged(m_ptr, m_epsabs, m_epsrel) : false; }
23   double x() const
24   { return m_ptr ? (m_can_minimize ? M_API<_Type>::x_minimum(m_ptr) : m_fct.get_lowest_bound()) : 0; }
25   double y() const
26   { return m_ptr ? (m_can_minimize ? M_API<_Type>::f_minimum(m_ptr) : m_fct.get_lowest_f_bound()) : 0; }
27 private:
28   ...
29 };
30 typedef IMinimizer<M_TYPE::NO_GRADIENT> MinimizerNoGradient;
```

## Minimizers – Example

CluePoints
Statistical
Scientific
programming

Olivia Quinet

Introduction
CluePoints
Clinical Trials
Statistical tests
SMART package

The R language
A very short introduction to
R
Some examples

R to C++?

Scientific
programming
challenges
Testing
Measure
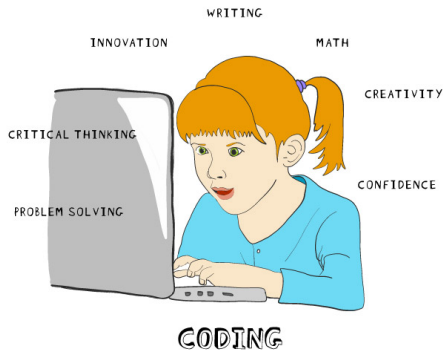Software architecture
Data structure
Smart pointers, Pimpl,
Factories
Fail-fast/Fail-safe
Numerical errors
Accumulators
std::algorithms, boost, GSL,
BLAS, LAPACK, …

Conclusion

Questions,
Remarks?

```cpp
1  template<M_TYPE> struct M_API; /// Generic template holding the API
2
3  template<>
4  struct M_API<M_TYPE::NO_GRADIENT> {
5    typedef gsl_min_fminimizer* PTR; /**< Type of the pointer to the minimizer */
6    typedef gsl_function* DEF; /**< Type of the pointer to the definition */
7
8    static PTR alloc(const STRING& _type) { return gsl_min_fminimizer_alloc(type(_type)); }
9    static void free(PTR _p) { gsl_min_fminimizer_free(_p); }
10
11   static const gsl_min_fminimizer_type* type(const std::string& _type);
12   static std::string name(PTR _p) { return gsl_min_fminimizer_name(_p); }
13
14   static bool set(PTR _p, DEF _fct, const double _minimum, const double _lower, const double _upper)
15   { return gsl_min_fminimizer_set(_p, _fct, _minimum, _lower, _upper) != GSL_EINVAL; }
16   static bool set(PTR _p, DEF _fct, const double _minimum, const double _fminimum, const double _lower, const
           double _flower, const double _upper, const double _fupper)
17   { return gsl_min_fminimizer_set_with_values(_p, _fct, _minimum, _fminimum, _lower, _flower, _upper, _fupper)
           != GSL_EINVAL; }
18
19   static int iterate(PTR _p) { return gsl_min_fminimizer_iterate(_p); }
20   static bool converged(PTR _p, const double _epsabs, const double _epsrel)
21   { return gsl_min_test_interval(x_lower(_p), x_upper(_p), _epsabs, _epsrel) == GSL_SUCCESS; }
22
23   static double x_minimum(PTR _p) { return gsl_min_fminimizer_x_minimum(_p); }
24   static double x_upper(PTR _p) { return gsl_min_fminimizer_x_upper(_p); }
25   static double x_lower(PTR _p) { return gsl_min_fminimizer_x_lower(_p); }
26   static double f_minimum(PTR _p) { return gsl_min_fminimizer_f_minimum(_p); }
27   static double f_upper(PTR _p) { return gsl_min_fminimizer_f_upper(_p); }
28   static double f_lower(PTR _p) { return gsl_min_fminimizer_f_lower(_p); }
29  };
```

# Minimizers – Example

```
1   class MyFctNoGradient : public gsl::MinimizerNoGradient::FCT::Base {
2   public:
3       MyFctNoGradient(...) { ... }
4
5       double evaluate(const double _x) override { ... }
6   };
7
8   double minimize_my_fct(...)
9   {
10      gsl::MinimizerNoGradient::FCT f(new MyFctNoGradient(...), .5*(low+hi), low, hi);
11      gsl::MinimizerNoGradient minimizer("Brent", f, 0.1, 0.1);
12      minimizer.iterate(10);
13      return minimizer.x();
14  }
```
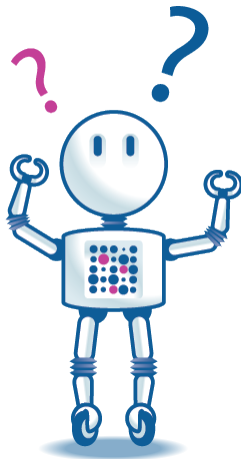
# Conclusion

- ▶ Fast production code
- ▶ No task is impossible
- ▶ Seek expertise
- ▶ Testing: Don't trust your code
- ▶ Have fun and keep learning

# Questions, Remarks?



# Thank you for your attention!