



# Advanced practical Programming for Scientists

**Thorsten Koch** 

Zuse Institute Berlin
TU Berlin

SS2017



#### Chapter 1 K&P



- Use descriptive names for globals, short names for locals
- Be consistent
- Use active names for functions
- Be accurate
- Indent to show structure
- Use the natural form for expressions
- Parenthesize to resolve ambiguity
- Break up complex expressions
- Be clear
- Be careful with side effects
- Use a consistent indentation and brace style
- Use idoms for consistency
- Give names to magic numbers



Integer (8, 16, 32, 64 bit)

Float (single 32 bit, double 64 bit, quad 128 bit) precision

Characters? (UTF-8)

Strings?

**Vectors** 

Compound data types



# OO style vs Imperative



```
Imperative
```

```
loc= xxx
val = yyy
If is_ok(loc) and is_valid(val) then
  add(container, val)
```

#### 00

```
If loc.is_ok() and val.is_valid() then container.add(val)
```

```
container.size() vs. size(container)

^ask object ^apply on data
```



#### Strutured vs. functional



#### Structured:

While not end of file do Process line

**Functional:** 

Process(forall lines in file)

Loops vs. Recursion



#### Polymorphism

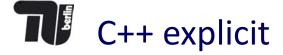


Add(5, 7)

Add(12.3, 34.7)

Add("hallo", "wie")

a **type signature** or **type annotation** defines the inputs and outputs for a <u>function</u>, <u>subroutine</u> or <u>method</u>. A type signature includes the number of arguments, the types of arguments and the order of the arguments contained by a function. A type signature is typically used during overload resolution for choosing the correct definition of a function to be called among many overloaded forms.





```
class Foo {
 public:
   Foo (int foo) : m_foo (foo) { }
     int GetFoo () { return m_foo; }
 private:
   int m_foo;
void DoBar (Foo foo)
 int i = foo.GetFoo ();
int main ()
   DoBar (42);
```





You have a

class with a constructor that constructs a string of the given size. You have a function

print(const MyString&),

and you call it with print(3).

You expect it to print "3", but it prints an empty string of length 3 instead.

http://stackoverflow.com/questions/121162/what-does-the-explicit-keyword-mean-in-c





```
if ( (country == SING) || (country == BRNI) ||
   (country == POL) || (country == ITALY) )
  If the country is Singapore, Brunei or Poland
  then the current time is the answer time
  rather than the off hook time.
  Reset answer time and set day of week.
```





```
for (theElementIndex = 0;
    theElementIndex < numberOfElements;
    theElementIndex++)
    elementArray[theElementIndex] = theElementIndex;

for (i = 0; i < nelems; i++)
    elem[i] = i;</pre>
```





```
enum { DANGER, CAUTION, CLEAR} the_signal;

If (CLEAR == the_signal)
{
  open_gates();
  start_train();
}
; = = 4
```



# **Avoid function macros**



#define isupper(c) ((c) 
$$>=$$
 'A' && (c)  $<=$  'Z')

parameter c occurs twice in the body of the macro. If isupper is called in a context like this,

while (isupper(c = getchar()))





```
int main()
  int const fixed = 20;
  int*
             var;
  int const** constptr;
   constptr = \&var;
  *constptr = &fixed;
   *var = 30;
  printf("x=%d, y=%d\n", fixed, *var);
```



#### Fun with FP Arithmetic



```
double s [2048];
 double e = 1;
  int n = 0;
4
 do \{ n = n + 1; e = e / 2; s[n] = 1 + e; \}
6
  while (e > 0);
                  /* Alternative 1 */
                  /* Alternative 2 */
  while (1 + e > 1);
                  /* Alternative 3 */
  while (s[n] > 1);
```

Does the loop terminate?

Will the program crash?

If it terminates will *n* have the save value in all alternatives? (Lines 7,8,9)



#### Fun with FP Arithmetic



```
double s [2048];
 double e = 1;
  int n = 0;
4
 do \{ n = n + 1; e = e / 2; s[n] = 1 + e; \}
6
 while (e > 0);
                  /* n = 1075
 while (1 + e > 1); /* n = 64 (Intel-P4) */
                  /* n = 53
 while (s[n] > 1);
                                             * /
```

The loop will terminate, the program will not crash and *n* is different in most cases, depending on the architecture, the compiler, and the switches.



# Introduction of C++ comments



```
a //*
//*/ b
```

In old C: a / b

In C++ : a



#### Experiment: Collecting Data at CO@Work-II



Combinatorial Optimization at Work II took place at ZIB from September 21 to October 9, 2009 with 105 participants from 23 countries.

We wanted to compute the seat allocation for the lecture hall. To do this we required ever participant to state their preferences. Everyone should send an email with a data file.

Lets see how long it took...



ASCII text with only a LF (ASCII 10) as line separator.

Fields are separated by a single space (ASCII 32)

Line 1: ParticipantNo HasLaptop EmailAddress

e.g. 67 1 koch@zib.de

0 = has no Laptop, 1 = has a Laptop

Lines 2-???: SeatNumber PreferenceValue

- Seat numbers start down at the low entrance, left to right, row by row.
- The hightest numbered seat is at the window side at the top.
- Count only seats that are physically there.
- The seat numbers in the file should be monotonically increasing.
- The preference values should be between 0 and 100.

e.g. **12 55** 

13 40

**14** 35 ...



#### Rules Regarding Preference Values



Allowed values are between 0 and 100

Only seats which are not available for the participants are allowed to get a value of 0

All numbers 1-100 have to be used at least once

The average has to be between 40-60

The difference to an adjacent seat has to be < 40

The difference to a neighboring seat has to be < 20

The data should not be randomly generated



## **Specifying Preference Offsets**



#### Lines ???-???: ParticipantNo PreferenceOffset

List indicating persons which you would like or not like to be your seat neighbor. (You have to know the ParticipantNo of the person.)

- A ParticipantNo of 0 indicates an empty seat.
- The PreferenceOffset is between -20 and 20 and will be added to your
   PreferenceValue if the person with the given ParticipantNo is your neighbor.

```
e.g. 55 17
27 -5
72 8
0 -10 ...
```

 This list can have as many entries as you like, but there should be at least 2 entries, and the occurring participant numbers have to be unique and valid.



#### **How To Submit**



#### Submission of this file is required for the course

The name of the file has to be ParticipantNo.txt

It should be <u>attached</u> to an email

Send the email to koch@zib.de

The subject of the email should be

CO@Work: SeatData for *ParticipantNo* 

Please, as soon as possible.





Mails received : 13

Different Subjects: 4 (10 1 1 1)

Wrong field spacing: 4

Seat counts : 2 (12 1)

Missing data : 1

Too much data: 1

Ok, from first view: 5 out of 13





Mails received : 23

Different Subjects: 6 (172111)

Wrong field spacing : 4

Seat counts : 4 (19 1 1)

Missing data : 2

Too much data: 0

Ok, from first view: 10

Corrected: 1

#### Add to the specification:

A seat without a desk is not allowed for the participants Seats with a 0 preference value are not relevant for the adjacency/neighboring difference rules.





Mails received : 37

Wrong subject : 11

Wrong field spacing: 8

Strange seat counts : 5

Missing data : 2

Corrected: 3





Mails received : 47

Data sets : 41 (6 corrections)

Wrong subject : 12

Wrong attachment name : 2

Wrong line separator : 29

Wrong field separator : 10

Pref. value not used : 11

Other Errors : 1

Number of seats : 153 - 181

No complains so far : 4





Mails received : 79

Data sets : 64

Wrong subject : 16

Wrong attachment name : 2

Wrong line separator : 45

Wrong field separator : 11

Pref.value not used : 22

Other Errors : 2

Number of seats : 153 - 181

No complains so far : 8





Mails received : 104

Data sets : 76

Wrong subject : 18

Wrong attachment name : 2

Pref. value not used : 19

Neighbor difference : 21

Wrong no/seq. seats: : 10

Wrong 0 seats : 20

No complains so far : 10



#### Overview of Errors in Data



	E7	E10	E11	E12	E13	E14	E16
5							Χ
6							Χ
12					Χ		Χ
13							Χ
16							X X X
18					Χ		Χ
19						Χ	Χ
20						Χ	
23					Χ		
24						Χ	
26							Χ
27						Χ	
36						Χ	
42					Χ		
45			Χ	Χ	X	Χ	Χ
47					Χ		
53					Χ		
59						X	
63			Χ		Χ	Χ	
64			Χ		X	X	Χ
71					Χ	Χ	

E7 bad seatno
E10 bad offset
E11 wrong seatno
E12 bad average
E13 prefval missing
E14 neigbour diff
E16 seat not 0

	E7	E10	E11	E12	E13	E14	E16
77							Χ
78	Χ		Χ			Χ	Χ
81				Χ	Χ		Χ
98					Χ		
99	Χ		Χ			Χ	
103					Χ	Χ	
107			Х			Χ	Χ
108			Χ			Χ	Χ
111							Χ
121							Χ
128			Х			Χ	Χ
129		Х					
134			Χ	Χ	Χ	Χ	
135					Χ		
137		Х			Χ	Χ	Χ
139					Χ		Χ
145	Χ		Χ		Χ	Χ	
160						Χ	
166					Χ	Χ	

#### **Please correct and resubmit**





Mails received : 144

Wrong subject : ~23

Wrong attachment name : 4

Data sets : 92

To be corrected : 28

Missing : 6

Pref. value not used : 14

Neighbor difference : 18

Wrong no/seq. seats : 2



#### Overview of Errors in Data



	E7	E10	E11	E12	E13	E14
12					Χ	
18					X	
23	Χ		Х			Х
24						Χ
24 27						X X X
45					Χ	Х
47					Χ	
63			Χ		X X X	Х
71					Χ	X
78	Χ		X			Х
79 103 107		Χ	Χ	Χ	Χ	
103						Х
107			X			X X X
108			Χ			Χ
110		Χ				
114						Χ
118					Χ	Х
128			X			X X X
134			Х	Х	X	X
135					Х	
136						X
136 137		Χ			Х	Χ
138					X X X	
139					Х	
160						X
166					Χ	X

E7 bad seatno

E10 bad offset

E11 wrong seatno

E12 bad average

E13 prefval missing

E14 neigbour diff

#### Please correct and resubmit





Mails received : 159

Wrong subject : ~26

Wrong attachment name : 4

Data sets : 94

To be corrected : 18

Missing : 4

Preference value not used: 9

Neighbor difference : 14

Wrong no/sequence seats: 3



#### Overview of Errors in Data



	E7	E10	E11	E12	E13	E14
18					Χ	
24						X
27						X
45					Χ	X
63					Χ	
71					Χ	X
78	Χ		X			X
79		Χ	X	X	Χ	
103						X
107			X			X
108			Χ			X
114						X
118					Χ	X
128			Χ			X
134			Χ	Χ	Χ	X
136						Х
137		Χ			Χ	X
138					Χ	

E7 bad seatno

E10 bad offset

E11 wrong seatno

E12 bad average

E13 prefval missing

E14 neigbour diff

# Please correct and resubmit





Mails received : 166

Wrong subject : ~28

Wrong attachment name : 4

Data sets : 95

To be corrected : 18

Missing : 3

Preference value not used: 7

Neighbor difference : 14

Wrong no/sequence seats: 3



#### Overview of Errors in Data



	E7	E10	E11	E12	E13	E14
24						X
27						X
45					X	X
71					X	X
78	X		X			X
79		X	X	X	X	
92					X	X
107			X			X
108			X			X
114						Χ
118					Χ	Χ
128			X			X
134			X	X	X	X
136						X
137		X			X	X

E7 bad seatno

E10 bad offset

E11 wrong seatno

E12 bad average

E13 prefval missing

E14 neigbour diff

# Please correct and resubmit



## 15 Days after the lecture – the final day



Mails received : 172

Wrong subject : ~31

Wrong attachment name : 4

Data sets : 95

To be corrected : 13

Preference value not used: 5

Neighbor difference : 13

Wrong no/sequence seats: 2



#### **Subject Variations**



#### The subject of the email should be

#### CO@Work: SeatData for *ParticipantNo*

CO@Work: SeatData for 022 CO@Work: SeatData for 222 CO@Work: SeatDatafor 222

CO@work: SeatData for 222 CO@Work: Seat Data for 222 Co@Work: SeatData for 222

CO@Work: SeatData for Participant222 CO@Work: SeatData for ParticipantNo Co@Work: SeatData for Participan222

CO@WORK: seatdata for 222 COatWork: SeatData for 222

COatWork for 222 SeatData for 222

SeatData for ParticipantNo 222

set data for participant number 222

data set participant number 222

Sitting assignment

Seats assignment



#### Overview of Errors in Data



	E7	E10	E11	E12	E13	E14
24						X
27						X
45					X	X
71					X	X
78	X		X			X
92					X	X
107			X			X
108			X			X
114						X
128			X			X
134			X	X	X	X
136						X
137		X			X	X

E7 bad seatno

E10 bad offset

E11 wrong seatno

E12 bad average

E13 prefval missing

E14 neigbour diff

Sorry, too late to correct!

Wrong line 1: 81, 129





#### You would think a ...

- ... cellular network operator knows where its base stations are located?
- ... fixed network operator can tell where the parts of its network are connected?
- ... chemical company knows how many plants they have?
- ... 5 m long pipeline cannot have a height difference from end-to-end of 100 m?
- Many companies have their data in Excel.
  There is no formal validation or referential integrality check.
- ▶ If they did formal validation, usually they found there was information they needed which they could not input and they started to "reuse" some data fields.
- If there is not at least 1 error per 100 data sets you are not looking hard enough.
- Usually the data changes all the time.
- ► They might not want to give it to you.
- ► The data might just not exist.

The first result of an optimization project is usually to improve the quality of planning data available at the company.



#### We agree mostly on the number of lines



File: ex1.dat with 100001235 lines

File: ex1.dat with 100001235 lines

File: ../../ex1.dat with 100001235 lines.

File: ex1.dat with 99999947 lines

File: ex1.dat with 100001235 lines

File ex1.dat with 239341319 lines

('File: ', '../exercise\_1/ex1.dat', 'with ', 100001235, ' lines (containing data)')

File: ../exercise\_1/ex1.dat with 100001235 lines

('File:', '../exercise\_1/ex1.dat', 'with', 100001235, 'lines')

File: ex1.dat with 100001235 lines

File: ex1.dat with 100001236 lines

File: ../exercise\_1/ex1.dat with 100001233 lines

File ex1.dat with 100001235 lines

File: ../../exercise\_1/ex1.dat with 100001236 lines

File: ex1.dat with 100001235 lines





Valid values Loc1: 49994581 with GeoMean: 36.782583

Valid values Loc1: 49994581 with GeoMean: 36.782583

Valid values Loc1: 50004466 with GeoMean: 36.781736117270

Valid values Loc1: 50004682 with GeoMean: 36.781681

Valid values Loc1: 50004706 with GeoMean: 36.781671

Valid values Loc1: 50004332 with GeoMean: 36.782305

Valid values Loc1: 50004555 with GeoMean: 36.7817

Valid Values Loc1: 50004616 with gemetric mean: ', 1.0)

Valid values Loc1: 50004851 with GeoMean: 0.000000

Valid values Loc1: 50004616 with GeoMean: 36.781761336223234

Valid values Loc1: 49504741 with GeoMean: 37.6172

Valid Values Loc1: 50004616 with GeoMean: 36.78176133614294

Valid values Loc1: 50005278 with GeoMean: 36.7817

Valid values Loc1: 50004571 with GeoMean: 50.00307664753011

Valid values Loc1: 50004620 with GeoMean: 36.78175072922403

Valid values Loc1: 50004798 with GeoMean: 36.7817

Valid values Loc1: 49994892 with GeoMean: 36.7825

Valid values Loc1: 50004777 with GeoMean: 36.7817





Valid values Loc2: 49994581 with GeoMean: 36.782583200332

Valid values Loc2: 49994783 with GeoMean: 36.782547

Valid values Loc2: 49994820 with GeoMean: 36.782505

Valid values Loc2: 49994439 with Geomean: 36.783359

Valid values Loc2: 49994670 with GeoMean: 36.7826

Valid Values Loc2: 49994716 with gemetric mean: ', 1.0)

Valid values Loc2: 49994951 with GeoMean: 0.000000

Valid values Loc2: 49994716 with GeoMean: 36.782573922686403)

Valid values Loc2: 49494352 with GeoMean: 37.6204

Valid Values Loc2: 49994716 with GeoMean: 36.782573922484126

Valid values Loc2: 49995371 with GeoMean: 36.7825

Valid values Loc2: 49994669 with GeoMean: 50.00238276443099

Valid values Loc2: 49994723 with GeoMean: 36.78255535635074

Valid values Loc2: 49994865 with GeoMean: 36.7825



# C, C++, awk, python, ada, java



Baumann	C++	75				
Wyczisk	Py3	173				
Dinsel	Ру					
Fleischer	С	30				
Gühring	Py2	203				
Hark	Py2	224				
Jeney	Ру					
Julia	C++	124				
Jung	C++					
Lewandowski	C++	68				
Luetzke	Py2	124				
Morgenroth	Java	50				
Nuernberger	С	21				
ongini	C++	116				
RayChew	C++					
Rettkowski	C++	55				
Richter	Java					
Sanny	C++	53				
Schmidt	С	41				
Seliverstov	Awk	384				
Sterling	С					
Koch:	С	17	/	15	/	9



#### **Exercise 2: Reading XML**



We provided some documentation and and an example file.

Write a program **ex2** that reads in the measured-1.0.0.2017-02-03.b0050c5c8deb1db59c7b2644414b079d.xml

And writes CSV data in format:

YYYY-MM-DD; HH; amountOfPower-Value

You will need the reading part later again.

Try to validate the XML file against the provided schema

The filename to be converted should be taken from the command line.

The output should be to stdout, any errors to stderr.

You may use whatever library to parse the XML file.