

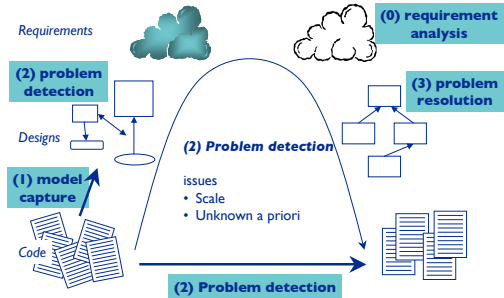
6. Code Duplication

a.k.a. *Software Cloning, Copy&Paste Programming*



- **Code Duplication**
 - + What is it?
 - + Why is it harmful?
- Detecting Code Duplication
- Approaches
- A Lightweight Approach
- Visualization (dotplots)
- Duploc

The Reengineering Life-Cycle



Code is Copied

Small Example from the Mozilla Distribution (Milestone 9)
Extract from /dom/src/base/nsLocation.cpp

```

14321 NS_IMETHODIMP                                14971 NS_IMETHODIMP                                14972 NS_IMETHODIMP
14322 LocationImpl::GetPathname(nsString&)          14981 LocationImpl::GetPathname(const nsString& aPort) 14982 LocationImpl::GetPathname(const nsString& aPort)
14323 {                                             14990 {                                             14991 {
14324   nsAutoString href;                          14999 nsAutoString href;                          14999 nsAutoString href;
14325   nsURI uri;                                   15000 nsURI uri;                                   15000 nsURI uri;
14326   mresult = NS_OK;                            14999 mresult = NS_OK;                            15000 mresult = NS_OK;
14327 }                                             15001 }                                             15001 }
14328                                           15002                                           15002
14329 result = GetHref(href);                       15004 result = GetHref(href);                       15004 result = GetHref(href);
14330 if (NS_OK == result) {                       15005 if (NS_OK == result) {                       15005 if (NS_OK == result) {
14331   #ifdef NS Gecko                               15006   #ifdef NS Gecko                               15006   #ifdef NS Gecko
14332     result = NS_NewURI(aURI, href);            14974 result = NS_NewURI(aURI, href);            15007 result = NS_NewURI(aURI, href);
14333   } else {                                       14975 } else {                                       15007 } else {
14334     result = NS_NewURI(aURI, href);            14976 result = NS_NewURI(aURI, href);            15008 result = NS_NewURI(aURI, href);
14335   }                                             14977 }                                             15008 }
14336 }                                             14978 }                                             15009 }
14337 #endif // NS Gecko                              14979 #endif // NS Gecko                              15010 #endif // NS Gecko
14338 #endif // NS Gecko                              14980 #endif // NS Gecko                              15011 #endif // NS Gecko
14339 char* file;                                    14881 char* file;                                    15012 char* file = aPathname.ToNewCString(1512);
14340 result = uri->GetPath(&file);                  14882 result = uri->GetPath(&file);                  15013 #ifdef NS Gecko
14341 const char* file;                             14883 const char* file;                             15014 #else
14342 result = uri->GetPath(&file);                  14884 result = uri->GetPath(&file);                  15015 #endif
14343 result = uri->GetPath(&file);                  14885 result = uri->GetPath(&file);                  15016 #ifdef NS Gecko
14344 result = uri->GetPath(&file);                  14886 result = uri->GetPath(&file);                  15017 #else
14345 result = uri->GetPath(&file);                  14887 result = uri->GetPath(&file);                  15018 #endif
14346 result = uri->GetPath(&file);                  14888 result = uri->GetPath(&file);                  15019 #ifdef NS Gecko
14347 result = uri->GetPath(&file);                  14889 result = uri->GetPath(&file);                  15020 #else
14348 result = uri->GetPath(&file);                  14890 result = uri->GetPath(&file);                  15021 #endif
14349 result = uri->GetPath(&file);                  14891 result = uri->GetPath(&file);                  15022 #ifdef NS Gecko
14350 result = uri->GetPath(&file);                  14892 result = uri->GetPath(&file);                  15023 #else
14351 result = uri->GetPath(&file);                  14893 result = uri->GetPath(&file);                  15024 #endif
14352 result = uri->GetPath(&file);                  14894 result = uri->GetPath(&file);                  15025 #ifdef NS Gecko
14353 result = uri->GetPath(&file);                  14895 result = uri->GetPath(&file);                  15026 #else
14354 result = uri->GetPath(&file);                  14896 result = uri->GetPath(&file);                  15027 #endif
14355 result = uri->GetPath(&file);                  14897 result = uri->GetPath(&file);                  15028 #ifdef NS Gecko
14356 result = uri->GetPath(&file);                  14898 result = uri->GetPath(&file);                  15029 #else
14357 result = uri->GetPath(&file);                  14899 result = uri->GetPath(&file);                  15030 #endif
14358 result = uri->GetPath(&file);                  14900 result = uri->GetPath(&file);                  15031 #ifdef NS Gecko
14359 result = uri->GetPath(&file);                  14901 result = uri->GetPath(&file);                  15032 #else
14360 result = uri->GetPath(&file);                  14902 result = uri->GetPath(&file);                  15033 #endif
14361 result = uri->GetPath(&file);                  14903 result = uri->GetPath(&file);                  15034 #ifdef NS Gecko
14362 result = uri->GetPath(&file);                  14904 result = uri->GetPath(&file);                  15035 #else
14363 result = uri->GetPath(&file);                  14905 result = uri->GetPath(&file);                  15036 #endif
14364 result = uri->GetPath(&file);                  14906 result = uri->GetPath(&file);                  15037 #ifdef NS Gecko
14365 result = uri->GetPath(&file);                  14907 result = uri->GetPath(&file);                  15038 #else
14366 result = uri->GetPath(&file);                  14908 result = uri->GetPath(&file);                  15039 #endif
    
```

How Much Code is Duplicated?

Usual estimates: 8 to 12% in normal industrial code
15 to 25 % is already a lot!

Case Study	LOC	Duplication without comments	with comments
gcc	460'000	8.7%	5.6%
Database Server	245'000	36.4%	23.3%
Payroll	40'000	59.3%	25.4%
Message Board	6'500	29.4%	17.4%

What Is Considered To Be Copied Code?

Duplicated Code = *Source code segments that are found in different places of a system.*

- in different files
- in the same file but in different functions
- in the same function

The segments must contain some *logic* or *structure* that can be abstracted, i.e.,

```

... computeIt(a, b, c, d); ... computeIt(w, x, y, z); is not considered duplicated code.
...                                     ...
...
... getIt(hash(tail(z))); ... getIt(hash(tail(a))); could be abstracted to a new function
...                                     ...
    
```

Copied artefacts range from expressions, to functions, to data structures, and to entire subsystems.

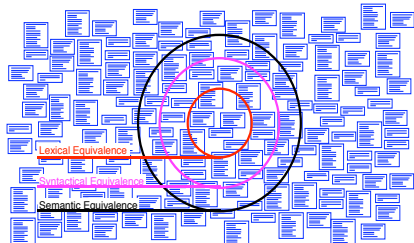
Copied Code Problems

- General negative effect:
 - + Code bloat
- Negative effects on *Software Maintenance*
 - + Copied Defects
 - + Changes take double, triple, quadruple, ... Work
 - + Dead code
 - + Add to the cognitive load of future maintainers
- Copying as additional source of defects
 - + Errors in the systematic renaming produce unintended aliasing
- Metaphorically speaking:
 - + Software Aging, "hardening of the arteries",
 - + "Software Entropy" increases even small design changes become very difficult to effect

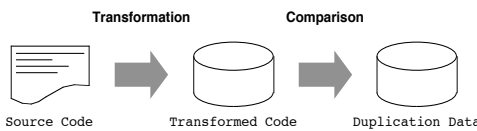
Code Duplication Detection

Nontrivial problem:

- No a priori knowledge about which code has been copied
- How to find all clone pairs among all possible pairs of segments?



General Schema of Detection Process



Author	Level	Transformed Code	Comparison Technique
[John94a]	Lexical	Substrings	String-Matching
[Duca99a]	Lexical	Normalized Strings	String-Matching
[Bake95a]	Syntactical	Parameterized Strings	String-Matching
[May96a]	Syntactical	Metric Tuples	Discrete comparison
[Kon97a]	Syntactical	Metric Tuples	Euclidean distance
[Bax98a]	Syntactical	AST	Tree-Matching

A Lightweight Approach

- Code Duplication
 - + What is it?
 - + Why is it harmful?
- Detecting Code Duplication
- Approaches
- **A Lightweight Approach**
- Visualization (dotplots)
- Duploc



Simple Detection Approach (i)

- **Assumption:**
 - Code segments are just copied and changed at a few places
- **Code Transformation Step**
 - remove white space, comments
 - remove lines that contain uninteresting code elements (e.g., just 'else' or '}")

```
...
//assign same fastid as container
fastid = NULL;
const char* fidptr = get_fastid();
if(fidptr != NULL) {
  int l = strlen(fidptr);
  fastid = newchar[ l + 1 ];
}

...
fastid=NULL;
constchar* fidptr=get_fastid();
if(fidptr!=NULL)
  intl=strlen(fidptr)
  fastid = newchar[ l + 1 ];
```

© S. Demeyer, S. Ducasse, O. Nierstrasz

Duplication.10

Simple Detection Approach (ii)

- **Code Comparison Step**
 - + Line based comparison (Assumption: Layout did not change during copying)
 - + Compare each line with each other line.
 - + Reduce search space by hashing:
 1. Preprocessing: Compute the hash value for each line
 2. Actual Comparison: Compare all lines in the same hash bucket
- **Evaluation of the Approach**
 - + Advantages: Simple, language independent
 - + Disadvantages: Difficult interpretation

© S. Demeyer, S. Ducasse, O. Nierstrasz

Duplication.11

A Perl script for C++ (i)

```
SequivalenceClassMinimalSize = 1;
$slidingWindowSize = 5;
$removeKeywords = 0;

@keywords = qw(
  then
  else
);

$keywordsRegExp = join '|', @keywords;

@unwantedLines = qw(
  return;
);

push @unwantedLines, @keywords;

while (<)> {
  chomp;
  $totalLines++;

  # remove comments of type /* */
  my $codeOnly = "";
  while($inComment && m!/* */) {
    ($inComment && m!^/) {
      $codeOnly .= $';
      $inComment = $inComment;
    }
  }
  $codeOnly = $' unless $inComment;
  $' = $';

  s!//.*!; # remove comments of type //
  s/s+//g; # remove white space
  s/$keywordsRegExp//og if
  $removeKeywords; # remove keywords
}
```

© S. Demeyer, S. Ducasse, O. Nierstrasz

Duplication.12

A Perl script for C++ (ii)

```
$codeLines++;
push @currentLines, $_;
push @currentLineNos, $_;
if($slidingWindowSize < @currentLines) {
  shift @currentLines;
  shift @currentLineNos;
}
#print STDERR "Line $totalLines >$._\n";
my $lineToBeCompared = join " ", @currentLines;
my $lineNumbersCompared = "<$ARGV>"; # append
the name of the file
$lineNumbersCompared .= join " ", @currentLineNos;
#print STDERR "$lineNumbersCompared\n";
if($bucketRef = SeqLines($lineToBeCompared)) {
  push @$bucketRef, $lineNumbersCompared;
} else {
  $seqLines($lineToBeCompared) = [
    $lineNumbersCompared ];
}
if(eof) { close ARGV } # Reset linenumber-count for next
file
```

- Handles multiple files
- Removes comments and white spaces
- Controls noise (if, {, }
- Granularity (number of lines)
- Possible to remove keywords

© S. Demeyer, S. Ducasse, O. Nierstrasz

Duplication.13

Output Sample

```
Lines:
create_property(pd,pmImplObjects,atReference,false,*ImplObjects);
create_property(pd,pmElitType,atReference,true,*ElitType);
create_property(pd,pmMinelt,stInteger,true,*lMinelt);
create_property(pd,pmMaxelt,stInteger,true,*lMaxelt);
create_property(pd,pmOwnership,stBool,true,*lOwnership);
Locations:
~/face/typesystem/SCTypesystem.C-6178/6179/6180/6181/6182
~/face/typesystem/SCTypesystem.C-6188/6189/6200/6201/6208
Lines:
create_property(pd,pmSupertype,atReference,true,*lSupertype);
create_property(pd,pmImplObjects,atReference,false,*ImplObjects);
create_property(pd,pmElitType,atReference,true,*lElitType);
create_property(pd,pmMinelt,stInteger,true,*lMinelt);
create_property(pd,pmMaxelt,stInteger,true,*lMaxelt);
Locations:
~/face/typesystem/SCTypesystem.C-6177/6178
~/face/typesystem/SCTypesystem.C-6229/6230
Lines = duplicated lines
Locations = file names and line number
```

© S. Demeyer, S. Ducasse, O. Nierstrasz

Duplication.14

Enhanced Simple Detection Approach

- **Code Comparison Step**
 - + Same as before +
 - Collect consecutive matching lines into match sequences
 - Allow holes in the match sequence
- **Evaluation of the Approach**
 - + Advantages
 - Identifies more real duplication, language independent
 - + Disadvantages
 - Less simple
 - Misses copies with (small) changes on every line

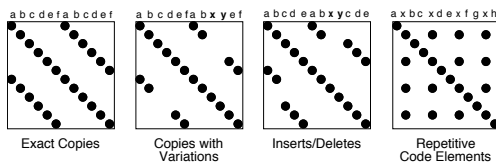
© S. Demeyer, S. Ducasse, O. Nierstrasz

Duplication.15

Visualization of Duplicated Code

- Visualization provides insights into the duplication situation
- A simple version can be implemented in three days
- Scalability issue

- Dotplots — Technique from DNA Analysis
 - Code is put on vertical as well as horizontal axis
 - A match between two elements is a dot in the matrix



© S. Demeyer, S. Ducasse, O. Nierstrasz

Duplication.16

Visualization of Copied Code Sequences

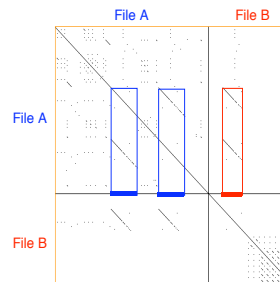
Detected Problem

File A contains two copies of a piece of code

File B contains another copy of this code

Possible Solution

Extract Method



All examples are made using Duploc from an industrial case study (1 Mio LOC C++ System)

© S. Demeyer, S. Ducasse, O. Nierstrasz

Duplication.17

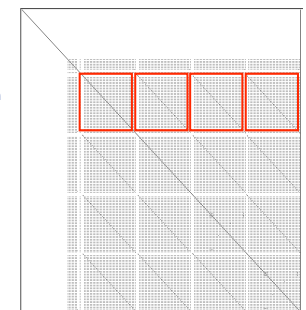
Visualization of Repetitive Structures

Detected Problem

4 Object factory clones: a switch statement over a type variable is used to call individual construction code

Possible Solution

Strategy Method



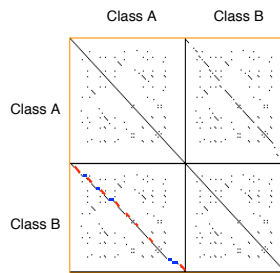
© S. Demeyer, S. Ducasse, O. Nierstrasz

Duplication.18

Visualization of Cloned Classes

Detected Problem:
Class A is an edited copy
of class B. Editing & Insertion

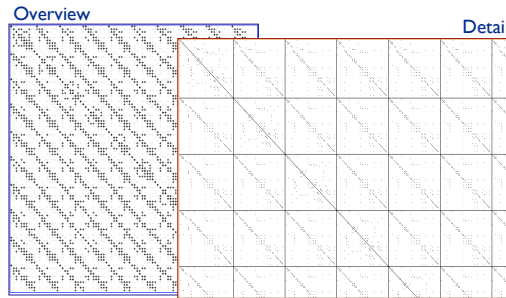
Possible Solution
Subclassing ...



© S. Demeyer, S. Ducasse, O. Nierstrasz

Duplication.19

Visualization of Clone Families



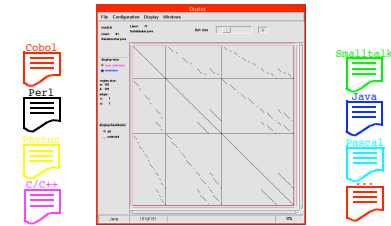
20 Classes implementing lists for different data types

© S. Demeyer, S. Ducasse, O. Nierstrasz

Duplication.20

Lightweight is sometimes not enough

Duploc is scalable, integrates detection and visualization



It runs really everywhere (Smalltalk inside)
⇒ Contact us if you need it!

© S. Demeyer, S. Ducasse, O. Nierstrasz

Duplication.21

Conclusion

- Duplicated code is a real problem
 - + makes a system progressively harder to change
- Detecting duplicated code is a hard problem
 - + some simple technique can help
 - + tool support is needed
- Visualization of code duplication is useful
 - + some basic support are easy to build
 - + one student build a simple visualization tool in three days
- Curing duplicated code is an active research area

© S. Demeyer, S. Ducasse, O. Nierstrasz

Duplication.22