



Diferença estrutural entre versões de um programa

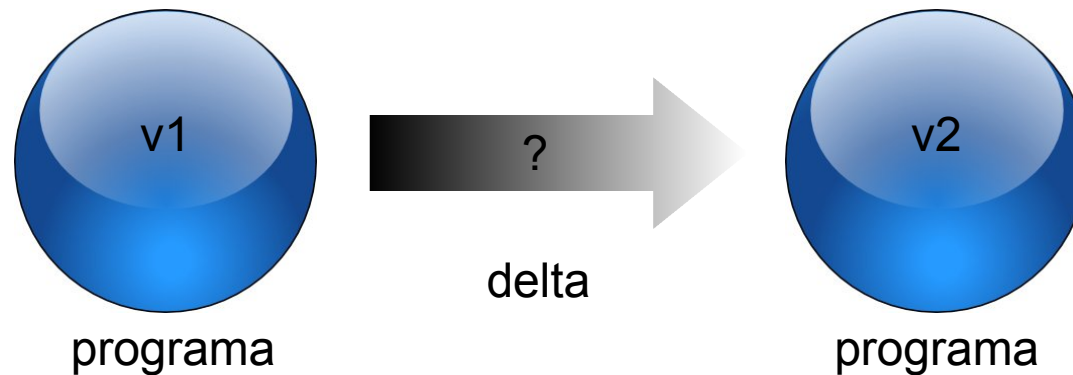
Evolução de Software @ UFCG

6 de agosto de 2008

Rodrigo Rocha <rodrigorgs@gmail.com>

Questão de Pesquisa

- A partir duas versões de um programa, como descobrir quais mudanças foram feitas?



- Importante ferramenta para entender a evolução de programas

Questão de Pesquisa



- GNU diff: arquivo, linha, caractere

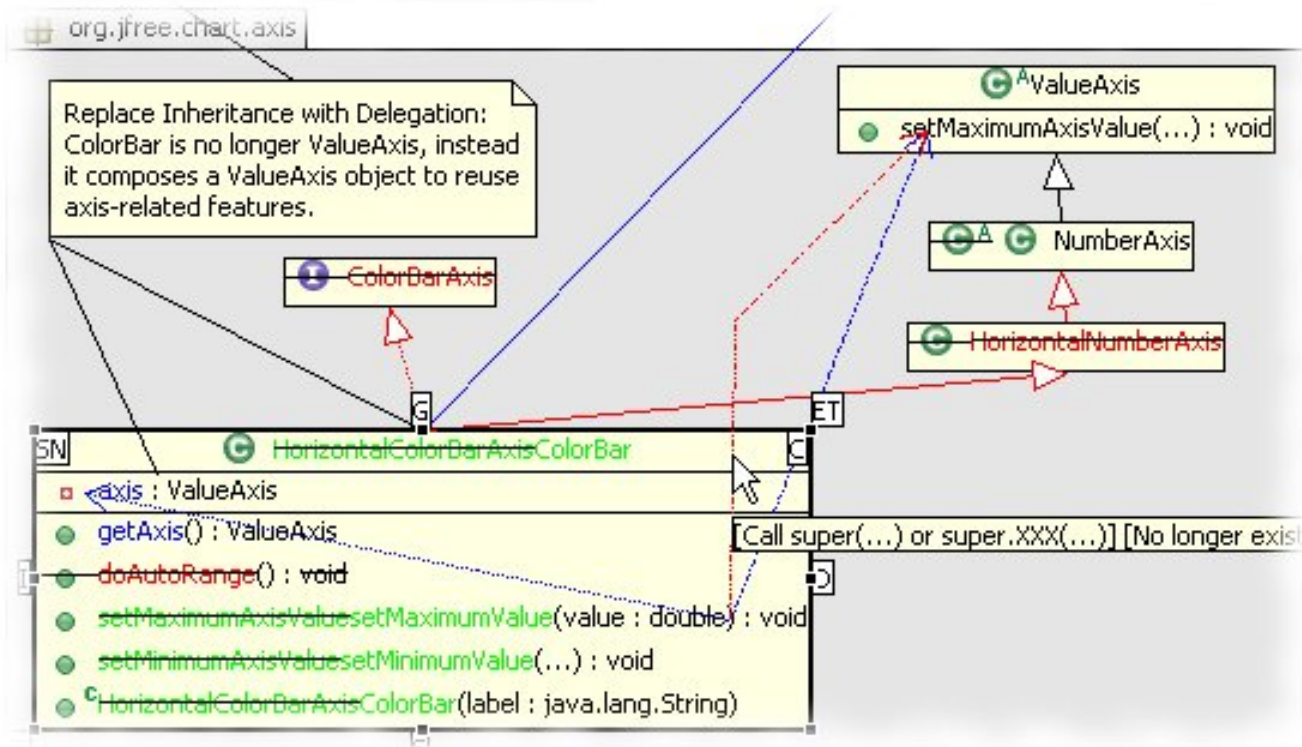
```
13 my ($mail_dir, $id, $message_file, $pa
14 my ($output, $success, $folder, @folder
15
16 !
17 !
18 # define some settings we will need lat
19 $mail_dir = $ENV{HOME} . '/mail';
20 $message_file = 'Message_file'; # save
21 $done_file = 'Sent_already.txt'; # Stor
22
23 $pause = '10'; # pause this many second
24
25 open(FILE, "<$done_file") || die "Can't
26 foreach $line (split ("\n", <FILE>)) {
27     next unless ($line =~ /^(.*)\s+(\d+)
28     $id = $1; $success = $2;
29     $Done_hash{$id} = $success;
30 }
31
32 # Go through the folders, forward the
33 @folders = <$mail_dir/*>;
34 foreach $folder (@folders){
35
36     $mailbox = Mail::MboxParser->new($fol
37     foreach $message ($mailbox->get_messa
38
39         if ($message =~ /Message-ID:\s*<
40
41
```

http://en.wikipedia.org/wiki/Image:Tkdifff_screenshot.png

Questão de Pesquisa



- Estrutural: pacote, classe, método, atributo...



<http://www.cs.ualberta.ca/~xing/jdevanviewer.html>

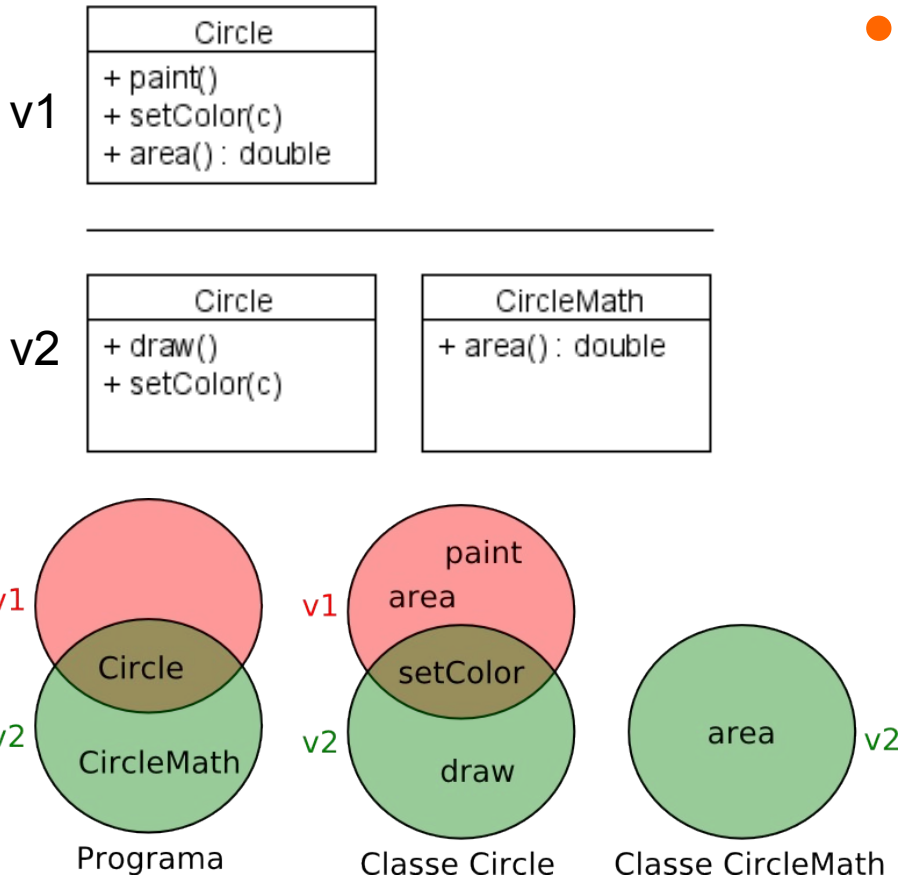
Sumário

- Definições
- Exemplo
- Programação estruturada: BEAGLE
 - Análise de Bertillonage
 - Análise de dependência
- Programação OO: UMLDiff
 - Similaridade de nome
 - Similaridade estrutural
- Considerações finais

Definições

- **Entidade** (de um programa):
 - (normalmente possui um identificador)
 - classe, método, atributo... (OO)
 - função, variável, arquivo... (prog. estruturada)
- **Mudança** (de uma entidade):
 - adicionar, remover, manter
 - renomear, mover
- **Identidade**
 - O que torna uma entidade reconhecível e diferenciável das outras

Exemplo



- Ações:
 - Adicionou CircleMath
 - Manteve Circle
 - renomeou paint para draw
 - manteve setColor
 - moveu area para CircleMath

Problema: identificar a **origem** de entidades **aparentemente novas**



BEAGLE

[Tu e Godfrey, 2002]

BEAGLE



- Ferramenta para estudo de evolução
- Programação estruturada
- Entidades: funções
- Métricas de similaridade entre funções
 - Análise de Bertillonage
 - Análise de dependência



Análise de Bertillonage

- Inicialmente usada para detectar clones
- Seleciona rapidamente prováveis entidades originais (triagem inicial)
- Caracteriza um trecho de código através de 5 métricas (dimensões):
 - {S,D,Ciclomatic}-Complexity
 - Albrecht
 - Kafura
- Similaridade = distância euclidiana

Análise de Bertillonage

v1 → `build_binary_op_noddefault in cp-typeck.c`

v2 →

1. `combine in fold-const.c:`
d=1005745.47
2. `recog_4 in insn-recog.c:`
d=2496769.23
3. `insn-recog.c in recog_5.c:`
d=7294066.05
4. `fprop in hard-params.c:`
d=8444858.78
5. `build_binary_op_noddefault in c-typeck.c:`
d=8928753.44

- Selecciona as 5 entidades mais prováveis
- Escolhe a que tiver o mesmo nome, se existir

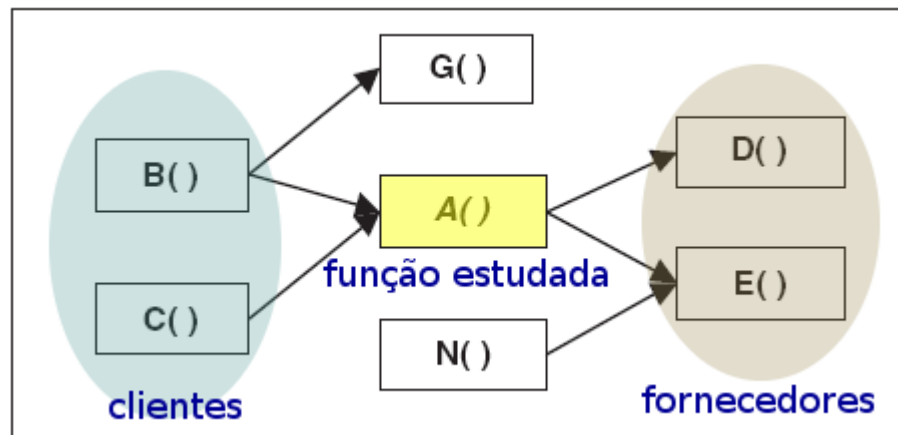
Análise de Bertillonage



- Ponto forte
 - Bom tratamento de renomear e mover
- Pontos fracos
 - Se confunde quando há funções parecidas
 - Mudanças na implementação das funções prejudicam a análise

Análise de dependência

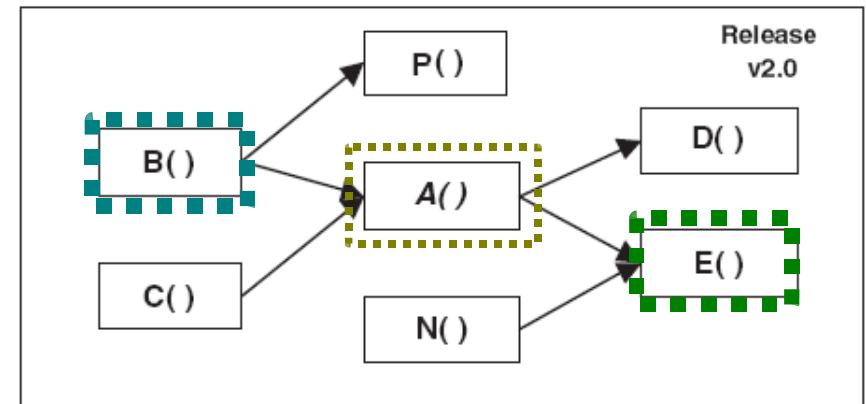
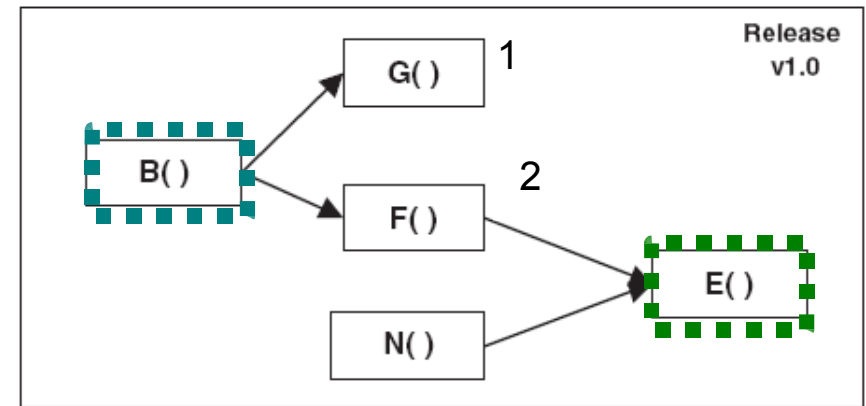
- Uma função...
 - ... chama outras funções (**forneecedores**) e
 - ... é chamada por outras funções (**clientes**)
- Uma função pode ser caracterizada pelo seu conjunto de forneecedores e clientes



Análise de dependência



- Funções “adicionadas” (v2):
{**A**, C, D, P}
- Funções “removidas” (v1):
{G, F}
- Clientes de A: {**B**, C} (v2)
 - Chama {**G**, **F**} (v1)
- Fornecedores de A: (v2)
{D, **E**}
 - Chamada por {**F**, N} (v1)
- Qual a origem de A?
 - **F** (2 pontos) ou G (1)



Análise de dependência

- Ponto forte
 - Pouco afetada por mudanças na implementação das funções
- Ponto fraco
 - Assume que clientes e fornecedores são estáveis
 - Pior caso: todas as funções são renomeadas ou movidas

Algoritmo



-
- Como combinar Bertillonage e dependência?
 - (O artigo não explica o algoritmo)

Avaliação



- Estudo de caso: GCC 2.7.2.3 => EGCS 1.0
- Resultados
 - Não avalia o algoritmo em termos de precisão e revocação!
 - Relatos específicos de sucesso
 - Avalia quanto do código do EGCS 1.0 é novo

File	Func	New	Old	Type
<code>gcc/cplus-dem.c</code>	36	36	0	Mostly New
<code>gcc/crtstuff.c</code>	5	5	0	Mostly New
<code>gcc/insn-output.c</code>	107	95	12	Mostly New
<code>gcc/final.c</code>	33	20	13	Half-Half
<code>gcc/regclass.c</code>	20	12	8	Half-Half

UMLDiff

[Xing e Stroulia, 2005]



UMLDiff

- Parte da ferramenta JDEvAn (Java Design Evolution and Analysis)
- Programação OO
- Entidades: classes, métodos, atributos...
- Métricas de similaridade
 - Similaridade de nome
 - Similaridade estrutural

Similaridade de nome

- Aplicada a identificadores
- Métrica baseada em LCS (maior subsequência comum)

VerticalDrawAction $\left\{ \begin{array}{l} \text{AddVerticalAction} \\ ? \\ \text{DrawVerticalAction} \end{array} \right.$

$$\frac{2 |\text{LCS}(s_1, s_2)|}{|s_1| + |s_2|}$$

$$\frac{\text{VerticalDrawAction} \quad \text{AddVerticalAction}}{\text{VerticalAction} \quad \text{DrawAction}} \frac{2 \times 14}{35} = 0.8$$

Sensível a permutações de nomes!

$$\frac{\text{VerticalDrawAction} \quad \text{DrawVerticalAction}}{\text{VerticalAction} \quad \text{DrawAction}} \frac{2 \times 14}{36} = 0.77$$

Similaridade de nome

- Métrica baseada em pares de caracteres adjacentes (*case-insensitive*)

- $P(\text{"VerticalDrawAction"}) = \{(v,e), (e,r), (r,t), (t,i), (i,c), \dots, (o,n)\}$

$$\frac{2|P(s_1) \cap P(s_2)|}{|P(s_1)| + |P(s_2)|}$$

- DrawVerticalAction $\Rightarrow 0.88$
- AddVerticalAction $\Rightarrow 0.73$

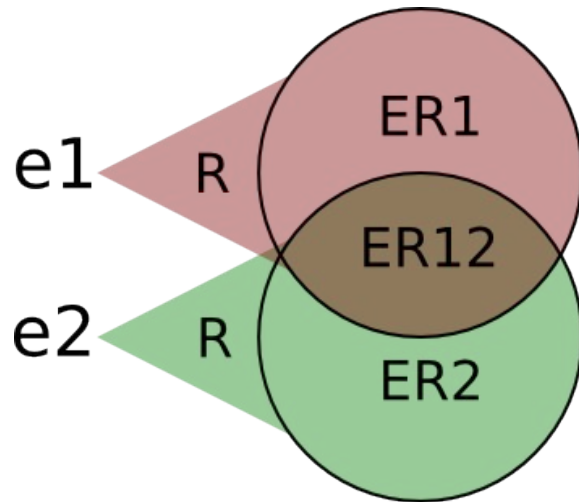
Similaridade estrutural

- relativa a um tipo de relacionamento R)

Table 2. Facts for computing structure similarity

Entity Type	Entity and relationship facts
Package	The top-level classes and interfaces it <i>contains</i>
Named class and interface	The fields, methods, constructors, inner classes and interfaces it <i>contains</i> ; classes/interfaces it <i>uses</i> and is <i>used by</i>
Field	The methods/constructors that <i>read</i> and <i>write</i> it
Method and constructor	Type of parameters it <i>declares</i> ; fields it <i>reads</i> and <i>writes</i> ; methods/constructors it <i>calls</i> and is <i>called by</i>

Similaridade estrutural



$$\approx \frac{|ER_{12}|}{|ER_1 \cup ER_2|}$$

- ERn são multi-conjuntos
- Se ER12 é vazio, usa (similaridade de nome)^k
- Para calcular ER12, usa um critério arbitrário de comparação de entidades

Algoritmo UMLDiff

- Similaridade geral: média de todas as métricas calculadas (mais ou menos)
- Identifica entidades mantidas => renomeadas => movidas
- Parâmetros: limiar de *renomear*, limiar de *mover*
- Limitação: não identifica *renomear* + *mover*

- Estudo de caso: 31 releases do JFreeChart em 4 anos, analisadas par-a-par (até 800 classes)
- Performance
 - Intel Centrino 1.6GHz, 768M RAM, Windows XP, PostgreSQL 7.4.5
 - Tempo de um diff: entre 1 e 60 min (média: 10 a 12 min)

- Metodologia
 - Inspeção manual das correspondências encontradas para calcular a precisão
 - Para obter uma idéia da revocação, inspeção manual de entidades renomeadas e movidas encontrados pelo UMLDiff com limiar 1%
- Resultados (limiar de 30%)
 - Revocação de renomear: 96,4%
 - Revocação de mover: 97,1%

- Resultados (precisão, limiar de 30%)

Table 5. *UMLDiff* results at threshold 30%

Type of change	#Correct	#Reported	Precision
Renamed package	29	29	100%
Renamed class/interface	121	128	94.5%
Moved class/interface	306	306	100%
Renamed field/method	1927	2024	95.2%
Moved field/method	630	721	87.3%
Data type and return type	677	710	95.4%
Visibility modifier	845	855	98.8%
Non-visibility modifier	299	303	98.7%
Class inheritance	180	185	97.3%
Interface inheritance	970	1025	94.6%
Total	5894	6286	95.2%

Considerações finais

- Resultados melhores
 - Incrementos pequenos
 - Uso disciplinado do SCM: separar reengenharia de mudanças funcionais
- Entidades com poucos relacionamentos são de difícil análise (ex.: getters, métodos abstratos)
- As técnicas não consideram mudanças do tipo união ou separação



Referências

- Qiang Tu, Michael W. Godfrey, "**An Integrated Approach for Studying Architectural Evolution,**" in Proc. 2002 Intl. Workshop on Program Comprehension (IWPC 2002), Paris, June 2002.
- Zhenchang Xing, Eleni Stroulia. **UMLDiff: An Algorithm for Object-Oriented Design Differencing.** Proceedings of the 20th IEEE/ACM international Conference on Automated software engineering, 2005.