



Universidade do Minho
Departamento de Informática

Github - Gestão de dados
LEI - Laboratórios de Informática 3
Grupo 9

13 de dezembro, 2021

Filipa Gomes
(a96556)

Ricardo Oliveira
(a96794)

João Loureiro
(a97257)

Introdução

Neste relatório descrevemos o nosso processo de desenvolvimento do guião 2, um sistema de gestão de usuários, commits e repositórios da plataforma Github, guardados em ficheiros .csv.

Exercicio

0.1 Estruturas de Dados

Para ser possível armazenar a informação de uma forma ordeira e organizada, foram criadas três estruturas principais de catálogos (tabelas de hash) de **Users**, **Repos** e **Commits**.

```
1 struct commit{
2     int u_id;
3     int a_id;
4     int r_id;
5     char *message;
6     int date[3];
7 };
8
9 struct commit_entry{
10     int id;
11     Commit c;
12     struct commit_entry *next;
13 };
14
15 struct catalog_c{
16     CEntry commits[MAX_TABLE];
17     int total;
18 };
```

*Exemplo do catálogo de commits.

0.2 Encapsulamento

Como pedido no guião, aplicou-se encapsulamento de dados. Para tal, foram criadas funções de *get_x* e *set_x*, de forma a que fosse possível aceder aos dados ao mesmo tempo que se cumpre o encapsulamento.

```
1
2 CEntry get_com_next (CEntry com) {
3     return com->next;
4 }
5
6 int *get_com_date (Commit commit ) {
7     int *dt=malloc(sizeof(int)*3);
8     dt[0] = commit->date[0];
9     dt[1] = commit->date[1];
10    dt[2] = commit->date[2];
11    return dt;
12 }
13
14 void set_com_next (CEntry com, CEntry cn) {
15     com->next = cn;
16 }
17
18 void set_com_date (Commit commit, int *date) {
19     commit->date[0] = date[0];
20     commit->date[1] = date[1];
21     commit->date[2] = date[2];
22 }
```

*Exemplo de `get_x` e `set_x` utilizadas para aceder aos dados relacionados com os catalogos de commits.

0.3 Importação de dados a partir de ficheiros

Foram criadas várias funções para importar os dados no projeto algumas delas são as funções `parse_hash_x`, onde `x` é o tipo a ser lido do ficheiro `.csv` e a ser colocado no catálogo respectivo.

```
1 int parse_hash_comm(CatalogC table, CatalogR r_table, CatalogU u_table, char *entrada,
   GStats gs){
2     FILE *u = fopen(entrada, "r");
3     char *buff = malloc(sizeof(char) * 500000);
4     static int id = 300000;
5
6     if (u == NULL){
7         printf("Error when opening %s! \n", entrada);
8         return 0;
9     }else{fgets(buff, 500000, u);}
10
11     while (fgets(buff, 500000, u)){
12         char *strc;
13         strc = strdup(buff);
14         int r_id = verify_int(strtok(strc, ";"));
15         int a_id = verify_int(strtok(NULL, ";"));
16         int u_id = verify_int(strtok(NULL, ";"));
17         char *cdt = strdup(strtok(NULL, ";"));
18         char *message=strdup(strtok(NULL, ";"));
19         int j = 0, data[3];
20         char *cd;
21         cd=strtok(cdt, "- :");
22         while (cd != NULL && j < 3)
23         {
24             data[j] = verify_int(cd);
25             if (data[j] == -1)
26                 j = 3;
27             j++;
28             cd = strtok(NULL, "- :");
29         }
30     }
```

*Codigo com omissões.

0.4 Estrutura de Estatísticas

Foi criada uma estrutura `GStats` que contém todos os campos necessários a guardar valores para a resolução das queries estatísticas. Inicializada toda a zeros, o preenchimento dos campos da estrutura `GStats`, é feito durante a leitura e parsing dos ficheiros.

```
1 struct global_stats {
2     int t_users;
3     int t_comms;
4     int t_repos;
5     int t_contribs;
6     int u_type[3];
7     int bot_repos;
8 };
9
10 GStats init_stats(){
11     GStats s=malloc(sizeof(struct global_stats));
12     s->t_users = 0;
13     s->t_comms = 0;
14     s->t_repos = 0;
15     s->t_contribs = 0;
16     s->bot_repos = 0;
17     s->u_type[0] = 0;
18     s->u_type[1] = 0;
19     s->u_type[2] = 0;
```

```

20     return s;
21 }

```

0.5 Queries Estatísticas

A resolução das queries estatísticas (1, 2, 3 e 4) é facilmente concretizada através do acesso dos campos da estrutura *GStats*. No caso da query dois, os colaboradores consideram-se como qualquer utilizador que seja author ou commiter id num commit de um dado repositório, podendo um dado utilizador contar várias vezes como colaborador, uma vez que, a média de colaboradores por repositório simboliza o número de users que participam no mesmo, caso se excluíssem utilizadores que já fossem colaboradores de outro repositório essa média iria diminuir de forma irrealista, pois se esse repositório só contesse esse user como colaborador seria representado como tendo apenas 0 colaboradores.

```

1 char * query_1 (GStats gs) {
2     char *o = malloc(sizeof(char)*64);
3     sprintf(o,"Bot : %d; Organization : %d; User : %d", get_type_n(gs,0), get_type_n(gs,1), get_type_n(gs,2));
4     return o;
5 }
6
7 char * query_2 (GStats gs) {
8     char *o=malloc(sizeof(char)*8);
9     double t_colab = get_t_contribs(gs);
10    double t_repos = get_t_repos(gs);
11    double x = (t_colab)/(t_repos);
12    sprintf(o,"%0.2f", x);
13    return o;
14 }

```

*Exemplo da query 1 e query 2.

0.6 Queries Parameterizadas

Houve algumas dificuldades em organizar o raciocínio de algumas queries(query 5 e 6), mas no final foi possível resolver. A dificuldade encontrava-se no tempo de execução, chegando a ser de 10 minutos ou mais. No caso da query 10, apesar de devolver commits com as maiores mensagens por repositório podem existir vários commits do mesmo user, uma vez que maior parte dos repositórios acaba por ter apenas um colaborador, sendo então mais útil retornar o top N, caso exista, das maiores mensagens e os seus users, independente de esses users poderem já ter aparecido. Foi criada uma query extra denominada de query 11, na qual é calculado o top N commits com maiores mensagens.

```

1 char *query_11(CatalogU cu, CatalogC cc, int n){
2     char *o = malloc(sizeof(char)*(32*n));
3     int k=0, j, id, sz, t1, t2;
4     CEntry re;
5     Commit r;
6     int bc[2][n];
7
8     for (int i = 0; i < MAX_TABLE; i++){
9         re = get_com_entry(cc,i);
10        while (re != NULL){
11            j = 0;
12            r = get_entry_commit(re);
13            id=get_com_u_id(r);
14            char * msg=get_com_message(r);
15            sz=utf8_strlen(msg);
16            while (j < k)
17            {
18                if (sz>bc[1][j]){
19                    t1=bc[0][j];
20                    t2=bc[1][j];

```

```

21         bc[0][j]=id;
22         bc[1][j]=sz;
23         id=t1;
24         sz=t2;
25     }
26     j++;
27 }
28 if (k < n){
29     bc[0][k]=id;
30     bc[1][k]=sz;
31     k++;
32 }
33 re= get_com_next(re);
34 }
35 }
36
37 for (int g = 0; g < k; g++)
38 {
39
40     if (g == 0)
41     {
42         char *temp = malloc(sizeof(char) * (32));
43         sprintf(temp, "%d;%s;%d\n", bc[0][g], check_user_login_by_id(cu, bc[0][g]),
44 bc[1][g] - 1);
45         strcpy(o, temp); // get_user_coms(u) );
46         free(temp);
47     }
48     else
49     {
50         char *temp = malloc(sizeof(char) * (32));
51         sprintf(temp, "%d;%s;%d\n", bc[0][g], check_user_login_by_id(cu, bc[0][g]),
52 bc[1][g] - 1);
53         strcat(o, temp); // get_user_coms(u) );
54         free(temp);
55     }
56 }
57 return o;
58 }

```

0.7 Leitor de comandos

Para a possível leitura dos ficheiros de input com os comandos, foi criada uma função *do_query_from_line* que separa a linha por um espaço (indicador da query) e seleciona a query a implementar. No caso das queries parameterizadas é feita uma separação dos parametros.

```

1 char *do_query_from_line(char *line, CatalogU u, CatalogR r, CatalogC c, GStats s){
2     char *token = strtok(line, " ");
3     int query = verify_int(token), n;
4     char *out;
5     switch (query){
6         case 1:
7             out = query_1(s);
8             break;
9
10        ...
11
12        case 5:
13            token=strtok(NULL, " ");
14            if(token!=NULL){
15                n=verify_int(token);
16                token = strtok(NULL, " ");
17                if(token!=NULL){
18                    char *d1=strdup(token);
19                    token=strtok(NULL, " ");
20                    if(token!=NULL){
21                        char *d2 =strdup( token);
22                        out = query_5(c, u, n, d1, d2);
23                    }else{out=strdup("Error Parsing Queue 5");}

```

```

24         }else{out=strdup("Error Parsing Queue 5");}
25     }else{out=strdup("Error Parsing Queue 5");}
26     break;

```

*Exemplo da leitura da query 1 e 5.

0.8 Ouput

Após a execução das queries lidas no ficheiro de input, o seu output é redirecionado e concatenado para um ficheiro.

```

1 void interprete_query_file(char * f_name, CatalogU u, CatalogR r, CatalogC c, GStats s){
2     FILE *f = fopen(f_name, "r");
3     char *tok = strtok(f_name, ".");
4     char *buff = malloc(sizeof(char) * 500000);
5     int n=1;
6     if (f == NULL){
7         printf("Error when opening %s! \n", f_name);
8     }
9     while (fgets(buff, 500000, f)){
10         char *out = malloc(sizeof(char)*64);
11         strcpy(out, tok);
12         char *op=malloc(sizeof(char)*24);
13         sprintf(op, "_command%d_output.txt", n);
14         strcat(out, op);
15         FILE *o = fopen(out, "w");
16         char *line;
17         char *temp = strdup(buff);
18         line = do_query_from_line(temp, u, r, c, s);
19         free(temp);
20         if (line!=NULL) fprintf(o, "%s", line);
21         fclose(o);
22         n++;
23     }
24     free(buff);
25     fclose(f);
26 }

```

0.9 Ficheiro com funções auxiliares

Criou-se um ficheiro com todas as funções auxiliares para a resolução das queries.

```

1 void build_dt(char *dt, int *data) {
2     int j = 0;
3     while (dt != NULL && j < 3)
4     {
5         data[j] = verify_int(dt);
6         if (data[j] == -1)
7             j = 3;
8         j++;
9         dt = strtok(NULL, "- :");
10    }
11 }

```

*Exemplo da build_dt que por sua vez transforma uma string com a data num array de ints.

0.10 Testes

Para testar o output e o tempo de execução, criou-se uma pasta *testes* e adicionou-se ficheiros de teste para cada query e um ficheiro de teste global, no qual o tempo de execução era monitorizado através do comando `time` de modo a que fosse possível melhorar o mesmo. Ficheiro de teste global:

```

1 1
2 2
3 3
4 4
5 5 100 2015-10-10 2017-10-10
6 6 100 Java
7 7 2018-01-01
8 8 100 2017-01-01
9 9 5
10 10 50
11 11 100

```

Tempo de execução

```

1 real 0m7,220s
2 user 0m6,970s
3 sysn 0m0,240s

```

0.11 Desafios na resolução do Guião 2

Ao longo do desenvolvimento deste guião depara-mo-nos com vários desafios.

Primeiramente, o encapsulamento das estruturas de dados foi uma das tarefas mais trabalhosas, sendo que era a primeira vez que nos deparávamos com tal tarefa. Isto demorou algum tempo e criou algumas dúvidas entre nós sobre qual a melhor abordagem, porém, com a ajuda disponível sentimos que fomos capazes de ultrapassar essas dificuldades.

Posteriormente, ao desenvolver as queries, surgiram mais dificuldades, sendo, desta vez, o problema da eficiência do programa. Ao início, apesar de certas queries executarem, a nosso ver, extraordinariamente rápido, outras nem tanto, causando dificuldades no desenvolvimento de melhores estratégias. A maior mais valia, no entanto, vinha da fase de Parsing dos ficheiros, uma vez que, como foi dedicada uma grande parte do tempo ao desenvolvimento das estruturas, todos os dados necessários já estavam propícios às funcionalidades das queries. Nem tudo é bom com este método, uma análise mais detalhada na leitura e armazenamento de dados no catálogo levou a um aumento do tempo de execução de queries singulares, isto é, ao executar apenas uma query o seu tempo tende a demorar cerca de 5s (para as queries estatísticas), no entanto, ao executar vários pedidos o tempo não aumenta linearmente, isto é, o custo está na leitura de dados, que permanecem na memória até ao final da execução. Dado isto, uma das estratégias, para o aumento da eficiência foi o "trade-off" entre o tempo de leitura de dados pela redução do tempo de execução das queries.

Por fim, foi necessário vários dias de trabalho árduo para que se pudesse melhorar e aperfeiçoar a resolução das queries, uma vez que, apenas o melhor tratamento de dados não é suficiente para diminuir os tempos de execução. Esta foi talvez a parte mais gratificante, pois, a cada melhoria do programa conseguia-mos ver o nosso trabalho refletido num melhor tempo. Porém, isso levou-nos a conhecer novas barreiras do nosso projeto, sendo elas a "duplicação de strings". A utilização da função `'strdup()'`, de modo a manter o encapsulamento, influencia as funções que retornam strings, uma vez que a operação de alocação de espaço pode ser demorada para grandes Data Sets.

0.12 Conclusão

O grupo sentiu que a elaboração deste guião os levou a melhorar o seu conhecimento e capacidades relativamente a programação e desenvolvimento de programas seguindo os conceitos de modularidade e modulação, aplicar conhecimentos de outras cadeiras, como por exemplo hash tables, criando uma ligação entre cadeiras diferentes como por exemplo Algoritmos e Complexidade e Programação Imperativa também como direcionada a Objetos.