



Universidade do Minho
Departamento de Informática

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

10 de Novembro, 2021

Filipa Gomes
(a96556)

Ricardo Oliveira
(a96794)

João Loureiro
(a97257)

Introdução

Neste relatório descrevemos o nosso processo de desenvolvimento do guião 1, um sistema de gestão de usuários, commits e repositórios da plataforma Github, guardados em ficheiros .csv.

Exercicio 1

0.1 Estruturas de Dados

Para ser possível organizar a informação de uma forma ordeira e organizada, foram criadas três estruturas principais de **Users**, **Repos** e **Commits**, que depois são organizadas em arrays dinamicos para fácil e rápida implementação.

```
1 typedef struct Repos{           typedef struct Commits{
2     int cmt_id;                  int repo_id;
3     int owner_id;               int author_id;
4     char *full_name;            int committer_id;
5     char *license;              int cmt_date[6];
6     char *has_wiki;             char *message;
7     char *description;          }Commit;
8     char *language;             typedef struct Users{
9     char *default_branch;        int id;
10    int cr_date[6];              char *login;
11    int updated_at[6];           char *type;
12    int forks_count;             int c_date[6];
13    int open_issues;            int followers;
14    int stargazers_count;        int *flr_id;
15    int size;                   int following;
16 }Repo;                          int *flg_id;
17                                int pub_gist;
18                                int pub_repo;
19                                }User;
```

0.2 Importação de dados a partir de ficheiros

Foram criadas várias funções para importar os dados no projeto algumas delas são as funções *parse_x* e com o auxílio de uma função *build_x*, onde x é o tipo a ser lido do ficheiro .csv.

```
1 Commit *parse_commits(int *i, char *entrada){
2     Commit *commits= malloc(sizeof(Commit));
3     FILE *r = fopen(entrada, "r");
4     char *buff = malloc(sizeof(char) * 200000);
5
6     while (fgets(buff, 200000, r)){
7         build_commits(&commits[*i], buff, i);}
8 }
9
10 int build_commits(Commit *data, char *line, int *i){
11     char *dt;
12     data->repo_id = verify_int(strsep(&line, ";"));
13     data->message = strdup(strsep(&line, "\n"));
14 }
```

*Codigo com omissões.

0.3 Verificação de campos

A verificação dos campos do ficheiro e a sua validação é efetuada em simultâneo com a sua construção do seu *type* usando as funções *verify_int*, *build_date* e usando a função especifica a cada tipo *verify_x*, onde x é o tipo a ser construido no momento.

0.4 Exportação de dados validados

Sendo o objetivo do exercício 1 a validação dos dados guardados nos ficheiros .csv e a sua exportação para ficheiros "corretos", foi necessário criar a função `save_x`, onde `x` é o tipo a ser guardado, que recebe um array de `*x` e o escreve num ficheiro com o nome "`ok-x.csv`".

```
1 int save_commits(Commit *u, int i, char *saida){
2     FILE *us = fopen(saida, "w");
3     fprintf(us, "repo_id;author_id;committer_id;commit_at;message\n");
4     int j = 0;
5     while (j < i){
6         fprintf(us, "%d;%d;%d;%d-%02d-%02d %02d:%02d:%02d;%s\n", u[j].repo_id, u[j].
7             author_id, u[j].committer_id, u[j].cmt_date[0], u[j].cmt_date[1], u[j].cmt_date[2],
8             u[j].cmt_date[3], u[j].cmt_date[4], u[j].cmt_date[5], u[j].message);
9         j++;
10    }
11    fclose(us);
12    return 0;
13 }
```

*exemplo do `save_commits` que é guardado com o formato "`repo_id;author_id;committer_id;commit_at;message`"

Exercicio 2

0.5 Combater o problema

No exercicio 2 foi nos pedido que fosse feita a confirmação de dados entre ficheiros, por exemplo verificar se um commit tinha sido efetuado por um utilizador real e se não para eliminar o mesmo. Mas isso apresentava um outro problema por detrás, as nossas estruturas de dados não foram desenhadas para serem acedidas assim com tanta diferença de "*indices*", o que tornaria o nosso programa extremamente lento e custoso em termos de *CPU Time*. A nossa solução foi a utilização de *Hash_tables*. As *Hash_tables* utilizadas são apenas arrays de *int* com listas ligadas para evitar perda de dados nas colisões.

```
1 #define MAX_TABLE 300000
2
3 /** Struct para representar commits na hashtable*/
4 typedef struct commit_entry{
5     int id;
6     struct commit_entry *next;
7 } commit_entry;
8
9
10 /** Representa uma hashtable */
11 typedef struct hash_table {
12     user_entry *users[MAX_TABLE];
13     repo_entry *repos[MAX_TABLE];
14     commit_entry *commits[MAX_TABLE];
15 } Hash_table;
```

0.6 Gestão dos dados

Chegamos também á conclusão de que o ficheiro `users-ok.csv` não ia sofrer alterações neste exercício, pelo que o qual bastava apenas guardar os registos na *Hash_table*, fazendo do `users-final.csv` o mais rápido de exportar. Foram implementadas várias funções com o intuito de resolver o problemas de repositórios e commits inválidos.

De modo a inserir os registos de ids nas *Hash_tables*, criamos uma função `insert_x`, onde `x` representa o tipo a ser inserido que fizesse esse mesmo `insert`, mais tarde esses dados seriam depois utilizados para fazer a construção do tipo de struct, repositório ou commit. A `parse_hash_x`,

que semelhante a *parse_x*, interpreta os ficheiros de entrada (desta vez os ficheiros resultantes do exercício 1), e a *build_hash_x*, onde x representa o tipo a ser construído, que guarda os registos repartidos nos seus respetivos campos da struct, verificando antes com o auxílio da *Hash_table*, se o registo continua a ser válido de acordo com as novas restrições do exercício 2. Os dados são depois guardados da mesma forma que no exercício 1, porém com a função *save_hash_x*, onde x é o tipo a guardar.

```

1 void insert_repo(int repo, Hash_table *table){
2     size_t pos = hash(repo);
3     repo_entry *new_repo = table->repos[pos];
4     repo_entry *x = calloc(1, sizeof(struct repo_entry));
5     x->id = repo;
6     if (new_repo == NULL) {
7         table->repos[pos] = x;
8     } else {
9         while ( new_repo->next != NULL) {
10             new_repo = new_repo->next;
11         }
12         new_repo->next = x;
13     }
14 }
15
16 Repo *parse_hash_repos(Hash_table *table, int *i, char *entrada)
17 {
18     Repo *repos= malloc(sizeof(Repo));
19     int at = 1;
20     FILE *r = fopen(entrada, "r");
21     char *buff = malloc(sizeof(char) * 500000);
22
23     if (r == NULL){
24         printf("Error when opening %s!\n", entrada);
25         return NULL;
26     }else{fgets(buff, 500000, r);}
27
28     while (fgets(buff, 500000, r)){
29         build_hash_repos(&repos[*i], table, buff, i);
30         if ((*i) == at){
31             repos = realloc(repos, sizeof(Repo) * (*i)*2);
32             at = 2 * (*i);
33         }
34     }
35     free(buff);
36     fclose(r);
37
38     return repos;
39 }

```

*exemplo do *insert_repos* e do *parse_hash_repos*.

0.7 Desafios na resolução do Guião 1

Como grupo encontramos muitas dificuldades em implementar as Hash Tables da biblioteca *glib*, logo tentamos procurar outra solução, o que nos custou algum tempo de trabalho e dedicação. A conclusão do exercício 2 foi a mais complexa em termos de compreensão do que era pedido e como implementar eficientemente. De resto, o Guião 1 foi relativamente fácil de implementar.