

# 第七组练习6报告

## 作业基本要求

1. 以课程大作业分组为单位
2. 下载GENBANK格式的基因组序列文件，从中提取CDS序列，保存为FASTA格式的文件
3. 整个项目至少包含两个动态共享库libgenbank.so以及libfasta.so，包含文件解读和格式的数据结构
4. 项目必须有完整的Makefile链
5. 项目必须包含测试数据和测试程序，用make test运行
6. 项目开发过程中建立dev分支和每个开发者对应的分支，必须能查询到每个开发者的开发合并路径
7. 项目最后合并到master分支并最终发布到github或者gitee的仓库中，每个开发者的账户中最后都包含一个拷贝
8. 报告必须包含核心代码的解析和测试数据的运行测试
9. 最好写gdb调试过程中怎么debug的

## libgenbank库

### 定义CDS结构体

确定需要提取的信息：基因名、蛋白id、蛋白序列等

```
1  typedef struct CDS{
2      char locus_tag[MAXWORD]; //基因名
3      char protein_id[MAXWORD];
4
5      char *translation; //蛋白序列
6
7      char left[MAXWORD];
8      char right[MAXWORD];
9
10     int end; //5表示5'部分，3表示3'部分，0表示既不是5'也不是3'
11
12     int flag_complement; //0表示正常，1表示在互补链上。
13
14 } CDS;
```

### 定义函数

1. 统计genbank文件中CDS出现次数

```
1  int countCDS (FILE *fp) //统计一共多少个CDS要提取
2  {
3      int len;
4      char line[MAXLINE];
5      int count_CDS = 0;
6      while (fgets(line, MAXLINE, fp) != NULL){
7          len = strlen(line);
8          int i=0;
9          while(isspace(line[i])) //跳过开头空格
10             i++;
```

```

11         if (line[i]=='C' && line[i+1]=='D' && line[i+2]=='S' &&
isspace(line[i+3])) //判断是否为CDS
12             {
13                 count_CDS++;
14                 continue;
15             }
16     }
17
18     return count_CDS;
19 }

```

## 2. 读取CDS中所需信息（用于写在fasta文件中第一行的注释信息中）

```

1
2 void read_CDS(CDS cdsdata[], char *filename, int *CDS_lines){
3     FILE *fp;
4     fp = fopen(filename,"r");
5     int len,CDS_line=0,nlines = 0;
6     int count = 0 ;
7     cdsdata[count].locus_tag[0] = '\0';//初始化结构体内的值
8     cdsdata[count].protein_id[0] = '\0';
9     char line[MAXLINE];
10    while (fgets(line, MAXLINE, fp) != NULL){
11        len = strlen(line);
12        nlines++;
13        int i=0;
14        while(isspace(line[i]))//跳过开头空格
15            i++;
16        if (line[i]=='C' && line[i+1]=='D' && line[i+2]=='S' &&
isspace(line[i+3]))//判断是否为CDS
17            {
18                CDS_line = nlines;
19                i+=3;
20
21                if (strcmp(line, "complement"))
22                    cdsdata[count].flag_complement = 1;
23                else cdsdata[count].flag_complement = 0;
24
25                while(!('0'<=line[i]&&line[i]<='9')){//跳过非数字
26                    if(line[i]=='<')
27                        cdsdata[count].end = 5;
28                    else if(line[i]=='>')
29                        cdsdata[count].end = 3;
30                    else cdsdata[count].end = 0;
31
32                    i++;
33                }
34                int j = 0;
35                while('0'<=line[i]&&line[i]<='9')
36                    cdsdata[count].left[j++] = line[i++];//基因起始位
置
37                cdsdata[count].left[j]='\0';
38                i+=2;
39
40                int k = 0;
41                while('0'<=line[i]&&line[i]<='9')
42                    cdsdata[count].right[k++] = line[i++];//基因终止位置

```

```

43         cdsdata[count].right[k]='\0';
44
45         count++;
46         continue;
47
48     }
49
50     if(strcmpr(line,"/locus_tag=") && (nlines == CDS_line+1) )
51     {
52         int k=0;
53         while(line[k]!='')
54             k++;
55         k++;
56         int j = 0;
57         while(line[k]!='')
58             cdsdata[count-1].locus_tag[j++]=line[k++];
59         cdsdata[count-1].locus_tag[j]='\0';
60     }
61
62     if(strcmpr(line,"/protein_id="))
63     {
64         int k=0;
65         while(line[k]!='')
66             k++;
67         k++;
68         int j = 0;
69         while(line[k]!='')
70             cdsdata[count-1].protein_id[j++]=line[k++];
71         cdsdata[count-1].protein_id[j]='\0';
72     }
73
74
75     if(strcmpr(line,"/translation="))
76         CDS_lines[count -1]=nlines;//把translation行号存在数组里
77
78
79     }
80 }

```

### 3. 提取CDS中translation序列（用于fasta文件的第二行序列）

```

1 void read_translation (FILE * fp,int count,CDS *cdsdata,char *line)
2 {
3     char Line[MAXLINE];
4     int right = atoi(cdsdata[count-1].right);
5     int left = atoi(cdsdata[count-1].left);
6     char *p,temp[right-left];
7     p = malloc(right-left);
8     int j=0,k = 0;
9     while(line[k]!='') //跳到序列处
10         k++;
11     k++;
12     while(line[k] != '\0') //存储第一行
13         temp[j++]=line[k++];
14
15     while(fgets(Line, MAXLINE, fp) != NULL){
16         int l = 0;

```

```

17     while (isspace(Line[l]))//跳过空格
18         l++;
19     while (Line[l]!='\0'){
20         temp[j++]=Line[l++];
21         if (Line[l] == ''){
22             temp[j] = '\0';
23             break;}
24     }
25     if (Line[l] == ''){
26         temp[j] = '\0';
27         break;
28     }
29 }
30 strcpy(p,temp);
31 cdsdata[count-1].translation = p;
32 printf("translation:%s\n\n",cdsdata[count-1].translation);
33 free(p);
34 }

```

## libfasta库

### writerfasta函数

将提取的信息输入到fasta格式文件中

```

1 void writerfasta(char *filename,CDS cdsdata[],int count){
2     int i;
3     FILE *fp;
4     strtok(filename,".");
5     strcat(filename,".fasta");
6     if((fp = fopen(filename,"w")) == NULL)
7     {
8         printf("Error:The file is not exist.");
9         exit(1);
10    }
11
12    for(i=0;i<=count-1;i++)
13    {
14        fputs(">",fp);
15        if(cdsdata[i].locus_tag != '\0'){
16            fputs("locus_tag:",fp);
17            fputs(cdsdata[i].locus_tag,fp);
18        }
19        if(cdsdata[i].protein_id != '\0'){
20            fputs("\tprotein_id:",fp);
21            fputs(cdsdata[i].protein_id,fp);
22        }
23        if(cdsdata[i].end != 0){
24            if(cdsdata[i].end == 3)
25                fputs("\tend:3'",fp);
26            if(cdsdata[i].end == 5)
27                fputs("\tend:5'",fp);
28        }
29        if(cdsdata[i].flag_complement == 0)
30            fputs("\tnot complementary chain",fp);
31        if(cdsdata[i].flag_complement == 1)

```

```

32         fputs("\tcomplementary chain",fp);
33
34         fputs("\n",fp);
35         fputs(cdsdata[i].translation,fp);
36         fputs("\n",fp);
37     }
38
39     fclose(fp);
40 }

```

## main函数

实现程序功能：genbank格式向fasta格式转化

```

1  int main()
2  {
3      char filename[100];
4      printf("please input your filename(path):\n");
5      scanf("%s",filename);
6      FILE *fp;
7      fp = fopen(filename,"r");
8      printf("open\n");
9
10     int count;
11     count = countCDS(fp); //统计共多少CDS
12     fclose(fp);
13     printf("count:%d\n",count);
14
15     int CDS_translation_lines[count];
16     CDS cdsdata[count];
17     read_CDS(cdsdata,filename,CDS_translation_lines); //把除translation以外的
    信息提取并存储
18     printf("read\n");
19
20     int i=0; //结果检查
21     for(i=0;i<count;i++)
22     {
23
24         printf("%d.locus_tag:%s\tprotein_id:%s\tleft:%s\tright:%s\tend:%d\tflag:%d\n",i+1,cdsdata[i].locus_tag,cdsdata[i].protein_id,cdsdata[i].left,cdsdata[i].right,cdsdata[i].end,cdsdata[i].flag_complement);
25     }
26
27     for (i=0;i<count;i++) //提取并存储translation
28     {
29         FILE *fp2;
30         char line2[MAXLINE];
31         fp2 = fopen(filename,"r");
32         int lines = 0;
33         while(fgets(line2, MAXLINE, fp2) != NULL)
34         {
35             //printf("%s\n",line2);
36             lines++;
37             if (lines == CDS_translation_lines[i])
38             {
39                 read_translation (fp2,i+1,cdsdata,line2);

```

```

39         break;
40     }
41
42     }
43     fclose(fp2);
44 }
45
46 writefasta(filename, cdsdata, count);
47
48 return 0;
49 }
50

```

## make 文件

```

1 main : main.c libgenbank.so libfasta.so fasta.h main.c
2     gcc main.c -o main
3 test :
4     ./main
5 clean :
6     rm main

```

## 代码测试

运行程序：

```

[root@localhost xiaozuoye]# make
gcc main.c -o main
[root@localhost xiaozuoye]# make test
./main
please input your filename(path):
GeneBank
open
count:1
read
1.locus_tag:   protein_id:SAL99496.1   left:9   right:1783   end:5   flag:1
translation:MEKQPNNNNMEUDUENSYQSDUARA IHSLDULKKRAVDAQGKAD
EALAADAPEEEYEALMDUYNKKWELYQRTRSNFAARFPEEGVFR TGKASGGPNQGSNK
NPUAAPALKUSDLPFLASAREASGNRALUTQNUREFATAFEALMELHQLNINDUYQR
YLPICLGKYYKTFIYSKRTLGETNIETWPMJKGWLUFTINTPRQKUKNTTAWMELTP
GSDETGEDFFHRUREFKEAHELANISADALLFFAUF TNCRFGWRNKISEAIRDTHQPF
UENFFEEMCAFASDLELNAGKPD AEDRHDNHQRSTTSSQRTNRKRSAADNNHYGNRTW
TNNNASTPRRMKGPTGYPYGGGGRYCDNGCGEKFMPPHKGUCPAI INRETNDRRSNEDR
RDSKRHQSEPD RQHQRQRDRPUSRAASEYIDQURADDUERLQSTFRGTTLDDDEDDKL
PCKLKGQKSEGNTPI LLEHEINUPLYIENKRTLALUDSGANFSSINKNFCTEHNUPIL
PHKKESNILLANAGISIKSYGYTPPITIKYNGSYTCQLEUMDLALGRMTSUWFEPSI
DAUYLQTPLKRFSIFSNAAVKLSCQ

[root@localhost xiaozuoye]#

```

查看fasta文件：

```

>locus_tag:   protein_id:SAL99496.1   end:5'   complementary chain
MEKQPNNNNMEUDUENSYQSDUARA IHSLDULKKRAVDAQGKAD
EALAADAPEEEYEALMDUYNKKWELYQRTRSNFAARFPEEGVFR TGKASGGPNQGSNK
NPUAAPALKUSDLPFLASAREASGNRALUTQNUREFATAFEALMELHQLNINDUYQR
YLPICLGKYYKTFIYSKRTLGETNIETWPMJKGWLUFTINTPRQKUKNTTAWMELTP
GSDETGEDFFHRUREFKEAHELANISADALLFFAUF TNCRFGWRNKISEAIRDTHQPF
UENFFEEMCAFASDLELNAGKPD AEDRHDNHQRSTTSSQRTNRKRSAADNNHYGNRTW
TNNNASTPRRMKGPTGYPYGGGGRYCDNGCGEKFMPPHKGUCPAI INRETNDRRSNEDR
RDSKRHQSEPD RQHQRQRDRPUSRAASEYIDQURADDUERLQSTFRGTTLDDDEDDKL
PCKLKGQKSEGNTPI LLEHEINUPLYIENKRTLALUDSGANFSSINKNFCTEHNUPIL
PHKKESNILLANAGISIKSYGYTPPITIKYNGSYTCQLEUMDLALGRMTSUWFEPSI
DAUYLQTPLKRFSIFSNAAVKLSCQ

```

第二个测试文件：CDS\_test.txt 运行后生成CDS\_test.fasta文件

命令行：

```
[root@localhost xiaozuoyel]# make test
./main
please input your filename(path):
CDS_test.txt_
```

```
1.locus_tag: protein_id:SAL99496.1 left:9 right:1783 end:5 flag:1
2.locus_tag:MAB_0174 protein_id:YP_001700928.1 left:1 right:818 end:5 flag:8
3.locus_tag:MAB_0175 protein_id:YP_001700929.1 left:1001 right:1975 end:8 flag
:8
translation:MEKQFNNNNMEUQENSYSQSDUARRAHSLSLULKKRAUDAQGXAD
EALAADAPEEYEAALMDVANKKWEYQRTKSNFAARFPEEGVFRGKASGGPNQGSNK
NPVAAAPALKUSDLPFLASARERASGNBALUTQNVREFAFAALMELHQLNINDUVR
YLPICLGKYYKTFIYSKRTLTGETNIETWPMUKGMLUEFTNTPRQKVAQNTTAAEMETP
GSDDETGEDFFHURUREFKEAHELANISADALLFFAFTNCRFGWARKKISEAIRDTHQFF
VENFFEEPCAFASDLELNAGKFDADRDHNDHQSTTSSQRTNHRKSAADNNHYGNATW
TNNNASTRPRMKGPTGFPYGGGGRYCDNGCGEKFMPPHKGUCPAIINRETNDRRSMEDR
RDSKRNHSEPDQHQRRQDRPUSRAASEYIDQURADDUERLQSTFRGTTLDDDDDKL
PCKLKGQKSEGNTPILEHEINUPLYIENKRTLALUDSGANFSSINKNFCTENHUPIL
PHKKESNILLANAGISIKSYGYTPPITIKYNGSYTCQLEDUMDLALGRTHSWFEPSI
DAUYLQTPLEKRFISNNAAUKLSQ

translation:MGPSNGLTRFTUUSULUTUTULFGWGAQRRIADDGLIULRTU
RMLLAGNGFUPNKGREUANTSTLWYLYLGGWAGGGMRLEYAALTSLULSLUGUA
LAILGTARLYAPLLAGRAAMUPAGMLUYIAIPPARDFATSGLENGULAYLGGMLM
MUKMAQAVTPTULDRGGPHQVDPRIUHAQRDQDRULSRFTUGLAFLAGLSULIR
PELALIGGGFLUMMLUARGFMSRIWIUAGGALPULYQIFRNGYYGLLUPSTAIAD
ASGSKMGQGFUYLQNMNSPYLIWIPAUILLALGUAAYQARRGQWARRQUAPGYGILA
RLUQNPTUAAUFLUSGFUQUYWIHQGGDFHARULLTPUFCMLLPISUWPLAAPDS
AAFTPKARLLTAATIGLFAGIAGWSUAMNSPGMAGDGTWYSGIUDERRFYAQT
GWAHPLTAADYLYNTPRRAULUAIDNTPDGALLLPSGNYDQMDUAPAIPPPPIPHGY
BGFHVLFTNLGLMGLMLGLDURIUDQIGLANPLAHTARITDGRIGHDKMLFPDAMI
ADGPMLEKRYPIPRYIDQWAAEAUEALKCPQTDAMLSAURKPLSPRLFUSNMLHSYE
FTTYRIDRUPREFELARCGLPMPKLDTPSYTGLPATGP

translation:MSURUKARRULSALLAUFUMPUSMAAMTINPATNAHFSREGLP
UEYLDUYSNSMGRNIRUEFQGGGPKAYLLDGLRAQDDFNGWDINTAAFEWFYQSGIS
UUMFUGGQSSFYTDWYSPSALNNKQPYTYKMETFTLTQELPAYLATNKQISATGNGUAGL
SMGGGALILAAFHPAQFRFAGSLSGFLNPSTIFMTNAIRUAMLDAGSYSUDNMAGPP
WDPAMRRNDPTUQAQALUAAGTRLYIYCAPGGSTPIDDNTDAGUALSASSLESIAWAG
NKAFQQAQYTHAGGRMANFUFPASGNHSMFYWGQQLQALKGDLIATLNG
```

CDS\_test.fasta文件内容：

```
>locus_tag: protein_id:SAL99496.1 end:5' complementary chain
MSURUKARRULSALLAUFUMPUSMAAMTINPATNAHFSREGLP
UEYLDUYSNSMGRNIRUEFQGGGPKAYLLDGLRAQDDFNGWDINTAAFEWFYQSGIS
UUMFUGGQSSFYTDWYSPSALNNKQPYTYKMETFTLTQELPAYLATNKQISATGNGUAGL
SMGGGALILAAFHPAQFRFAGSLSGFLNPSTIFMTNAIRUAMLDAGSYSUDNMAGPP
WDPAMRRNDPTUQAQALUAAGTRLYIYCAPGGSTPIDDNTDAGUALSASSLESIAWAG
NKAFQQAQYTHAGGRMANFUFPASGNHSMFYWGQQLQALKGDLIATLNG
>locus_tag:MAB_0174 protein_id:YP_001700928.1 end:5' not complementary chain
MSURUKARRULSALLAUFUMPUSMAAMTINPATNAHFSREGLP
UEYLDUYSNSMGRNIRUEFQGGGPKAYLLDGLRAQDDFNGWDINTAAFEWFYQSGIS
UUMFUGGQSSFYTDWYSPSALNNKQPYTYKMETFTLTQELPAYLATNKQISATGNGUAGL
SMGGGALILAAFHPAQFRFAGSLSGFLNPSTIFMTNAIRUAMLDAGSYSUDNMAGPP
WDPAMRRNDPTUQAQALUAAGTRLYIYCAPGGSTPIDDNTDAGUALSASSLESIAWAG
NKAFQQAQYTHAGGRMANFUFPASGNHSMFYWGQQLQALKGDLIATLNG
>locus_tag:MAB_0175 protein_id:YP_001700929.1 not complementary chain
MSURUKARRULSALLAUFUMPUSMAAMTINPATNAHFSREGLP
UEYLDUYSNSMGRNIRUEFQGGGPKAYLLDGLRAQDDFNGWDINTAAFEWFYQSGIS
UUMFUGGQSSFYTDWYSPSALNNKQPYTYKMETFTLTQELPAYLATNKQISATGNGUAGL
SMGGGALILAAFHPAQFRFAGSLSGFLNPSTIFMTNAIRUAMLDAGSYSUDNMAGPP
WDPAMRRNDPTUQAQALUAAGTRLYIYCAPGGSTPIDDNTDAGUALSASSLESIAWAG
NKAFQQAQYTHAGGRMANFUFPASGNHSMFYWGQQLQALKGDLIATLNG
```

执行make clean，自动清除main文件：

```
[root@localhost xiaozuoyel]# make clean
rm main
```

第七组小组分工

libgenbank.so --赖敏智

libfasta.so --韩德坤

main.c --赖敏智、韩德坤

makefile文件及代码测试 --王龙威

