《Linux 生物信息基础》小组讨论总结报告

第4组,第5次讨论

组长: 陈奕晗

执笔:朱瑾煜、邹济平、陈奕晗、高培翔

1时间

2021年4月23日, 15:00~17:00

2地点

3W 咖啡厅

3人员

陈奕晗、邹济平、朱瑾煜、高培翔

4方式

线下讨论

5 主题

5.1 上次课涉及的 BLAST 内容复习

(朱瑾煜主要执笔,陈奕晗补充)

5.2 vi 文本编辑

(邹济平执笔)

5.3 关于 MyNCBI

(高培翔执笔)

6内容

6.1 上次课涉及的 BLAST 内容

6.1.1 建立index—— makeblastdb 命令

使用 makeblastdb 命令可以读取 fasta 格式的序列文件,生成 *.phr、*.psq、*.pin 格式的索引文件。

例 1: 将从 Swiss-Prot 上获取的人类蛋白质序列文件转换成本地 BLAST 可读取的索引文件

makeblastdb -dbtype prot -in sp_human.fasta -out sp_human

其中

- -dbtype [prot/nucl] 即 database type,数据库类型。包括 prot (protein,蛋白质) 和 nucl (nucleic acid,核酸) 两种。
- -in [filename] 输入的 fasta 格式文件。
- -out [filename] 输出的索引文件的名称(不含扩展名)。默认是输入文件的名称。

执行上述命令后,可以看到生成了 sp_human.phr、sp_human.pin、sp_human.psq 三个索引文件。

```
<leb4c@bbt> 16:04:26 ~/MyBlastDB
$ makeblastdb -in sp_human.fasta -dbtype prot -out sp_human
Building a new DB, current time: 04/23/2021 16:11:16
New DB name: /rd1/home/leb4c/MyBlastDB/sp_human
New DB title: sp_human.fasta
Sequence type: Protein
Deleted existing Protein BLAST database named /rd1/home/leb4c/MyBlastDB/sp_human
Keep MBits: T
Maximum file size: 100000000B
Adding sequences from FASTA; added 20395 sequences in 0.432522 seconds.
<leb4c@bbt> 16:11:17 ~/MyBlastDB
$ 1s
sp_human.fasta sp_human.pin uniprot_sprot.phr uniprot_sprot.psq
sp_human.phr
             sp_human.psq uniprot_sprot.pin
<leb4c@bbt> 16:11:19 ~/MyBlastDB
$
```

image-20210423161149614

6.1.2 不使用index,直接从目标序列(subject) 中进行BLAST 搜索—— - subject

例 2: 以 17 个拟南芥 SBP (Squamosa promoter-binding-like protein) 转录因子搜索 19 个水稻 SBP 转录因子,输出结果到 ARATH-ORYSJ.xls 文件中。

blastp -query 17SPL_ARATH.FASTA -subject 19SPL_ORYSJ.FASTA -outfmt 7 -e
value 0.1 -out ARATH-ORYSJ.xls

其中:

- -query [filename]用于搜索的序列。
- -subject [filename] 目标序列,即被搜索的序列。如果使用 index,则用 -db [filename] 取代。
- -outfmt [int] 输出格式,详见 help 中的内容。这里使用 -outfmt 7,即带注释的表格格式。
- -evalue [float] 允许的期望值 (expectation value) 阈值。期望值越低表示相似度越高,其阈值越低也就使得搜索结果越严格。默认值为 10,这里设成 0.1。
- -out [filename] 输出文件名。

结果如下:

```
<leb4c@bbt> 16:40:28 ~/BLAST
$ 1s
17SPL_ARATH.FASTA
                      HBA_SW_2.txt
                                        SPL3D_CDS.FASTA
                                                           ZMTF_CDS.nin
                     HBA_SW.txt
19SPL_ORYSJ.FASTA
                                        SPL3D_PEP.FASTA
                                                           ZMTF_CDS.nsq
CEA21_HUMAN_C2.FASTA input
                                        spl7_arath.needle ZMTF_PEP.FASTA
                                                           ZMTF_PEP.phr
CEA21_HUMAN.FASTA
                      output
                                        SPL_BLAST.txt
CEAM5 HUMAN.FASTA
                      out.txt
                                       ZMTF_CDS.FASTA
                                                           ZMTF_PEP.pin
HBA HUMAN.FASTA
                     SPL3_ARATH.FASTA ZMTF_CDS.nhr
                                                           ZMTF_PEP.psq
<leb4c@bbt> 16:40:28 ~/BLAST
$ blastp -query 17SPL_ARATH.FASTA -subject 19SPL_ORYSJ.FASTA -outfmt 7 -evalue 0.1 -out ARAT
H-ORYSJ.xls
<leb4c@bbt> 16:40:41 ~/BLAST
$ 1s
                     HBA_SW_2.txt
                                       SPL3D_PEP.FASTA
                                                           ZMTF_PEP.FASTA
17SPL_ARATH.FASTA
19SPL_ORYSJ.FASTA
                     HBA SW.txt
                                        spl7_arath.needle ZMTF_PEP.phr
                                                           ZMTF_PEP.pin
ARATH-ORYSJ.xls
                      input
                                        SPL_BLAST.txt
CEA21_HUMAN_C2.FASTA output
                                       ZMTF_CDS.FASTA
                                                           ZMTF_PEP.psq
CEA21_HUMAN.FASTA
                                       ZMTF_CDS.nhr
                      out.txt
CEAM5_HUMAN.FASTA
                      SPL3_ARATH.FASTA ZMTF_CDS.nin
HBA_HUMAN.FASTA
                      SPL3D_CDS.FASTA ZMTF_CDS.nsq
<leb4c@bbt> 16:40:43 ~/BLAST
$
```

image-20210423164139737

在本地终端输入如下命令,输入密码后,可以将这个.xls 文件转移到本地,进而可以用 Excel 软件打开:

scp leb4c@117.78.18.116:~/BLAST/ARATH-ORYSJ.xls /Users/jinyuzhu/Desktop 用 Excel 查看结果如下: (仅显示部分)

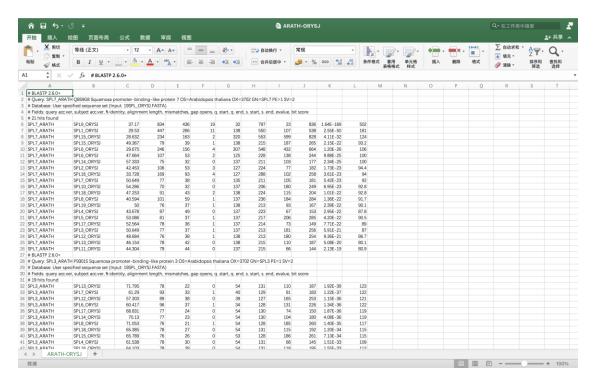


image-20210423164728585

6.1.3 指定搜索的起止位点—— -query_loc 和 -subject_loc

例 3: 以人癌胚抗原 CEA21*HUMAN 中恒定结构域 (147~231 位) 搜索 CEAM5*HUMAN 中 6 个恒定结构域,将结果保存到 output 目录下的 CEA_BLAST.txt 中。(癌胚抗原的相关数据可以从癌胚抗原家族结构域信息网站 http://www.carcinoembryonic-antigen.de/human/index.html 中获得)

blastp -query CEA21_HUMAN.FASTA -query_loc 147-231 -subject CEAM5_HUMAN.
FASTA -subject_loc 145-675 -out output/CEA_BLAST.txt

其中:

- -query_loc [ini]-[ter] 指定 query 序列参与比对的部分,从起始位点 [ini] 到终止位点 [ter]。
- -subject_loc [ini]-[ter] 指定 subject 序列参与比对的部分,从起始位点 [ini] 到终止位点 [ter]。

使用 less 命令逐屏读取输出文件,可以看到序列比对的结果。一共搜索出 6 段。 现截取一部分如下。

```
Score = 61.2 bits (147), Expect = 5e-17, Method: Compositional matrix adjust.
 Identities = 34/85 (40%), Positives = 49/85 (58%), Gaps = 1/85 (1%)
Query 147 PSIQAS-STTVTEKGSVVLTCHTNNTGTSFQWIFNNQRLQVTKRMKLSWFNHMLTIDPIR 205
             PSI ++ S V +K +V TC
                                           T++ W N Q L V+ R++LS N LT+
Sbjct 503 PSISSNNSKPVEDKDAVAFTCEPEAQNTTYLWWVNGQSLPVSPRLQLSNGNRTLTLFNVT 562
Query 206 QEDAGEYQCEVSNPVSSNRSDPLKL 230
              + DA Y C + N VS+NRSDP+ L
Sbjct 563 RNDARAYVCGIQNSVSANRSDPVTL 587
 Score = 60.8 bits (146), Expect = 7e-17, Method: Compositional matrix adjust. Identities = 33/85 (39%), Positives = 48/85 (56%), Gaps = 1/85 (1%)
Query 147 PSIQAS-STTVTEKGSVVLTCHTNNTGTSFQWIFNNQRLQVTKRMKLSWFNHMLTIDPIR 205
PSI ++ S V +K +V TC ++ W NNQ L V+ R++LS N LT+ +
Sbjct 147 PSISSNNSKPVEDKDAVAFTCEPETQDATYLWVNNQSLPVSPRLQLSNGNRTLTLFNVT 206
Query 206 QEDAGEYQCEVSNPVSSNRSDPLKL 230
Sbict 207 RNDTASYKCETQNPVSARRSDSVIL 231
 Score = 57.8 bits (138), Expect = 8e-16, Method: Compositional matrix adjust.
 Identities = 29/78 (37%), Positives = 45/78 (58%), Gaps = 0/78 (0%)
Query 153 STTVTEKGSVVLTCHTNNTGTSFQWIFNNQRLQVTKRMKLSWFNHMLTIDPIRQEDAGEY 212
S V ++ +V LTC T++ W NNQ L V+ R++LS N LT+ ++ D G Y Sbjct 332 SNPVEDEDAVALTCEPEIQNTTYLWWVNNQSLPVSPRLQLSNDNRTLTLLSVTRNDVGPY 391
Query 213 QCEVSNPVSSNRSDPLKL 230
              +C + N +S + SDP+ L
Sbjct 392 ECGIQNKLSVDHSDPVIL 409
: ||
```

image-20210423170031529

6.1.4 几种 BLAST 的区别和选择

BLAST 套件的 blastn、blastp、blastx、tblastn 和 tblastx 子工具的用途分别如下:

- 1. **blastn** 是将给定的核酸序列与核酸数据库中的序列进行比较。**(核酸搜核 酸)**
- 2. **blastp** 是使用蛋白质序列与蛋白质数据库中的序列进行比较,可以寻找较远的关系。**搜蛋白质序列时,应首选 blastp。(蛋白搜蛋白)**
- 3. **blastx** 则是将给定的核酸序列按照六种阅读框架将其翻译成蛋白质与蛋白质数据库中的序列进行比对,**在转录组分析中**对分析新序列和 EST (expression sequence tag) 很有用。(**核酸搜蛋白**)
- 4. **tblastn** 将给定的氨基酸序列与核酸数据库中的序列(双链)按不同的阅读框进行比对,**在转录组分析中**对于寻找数据库中序列没有标注的新编码区很有用。**(蛋白搜核酸)**
- 5. **tblastx** 只在特殊情况下使用,它将 DNA 被检索的序列和核酸序列数据库中的序列按不同的阅读框全部翻译成蛋白质序列,然后进行蛋白质序列比对。 (核酸翻译之后搜核酸)

在下面的例子中,我们分别采用这5种方法进行BLAST。

例 4: 以拟南芥 SPL3 转录因子搜索 764 个玉米转录因子数据集:

已有 SPL3D_PEP.FASTA、SPL3D_CDS.FASTA 和索引文件 ZMTF_PEP.*、ZMTF_CDS.*,期望值阈值设为 0.01。

1. **blastp** 以拟南芥 SPL3 转录因子 DNA 结合结构域搜索 764 个玉米转录因子 数据集蛋白质序列

blastp -query SPL3D_PEP.FASTA -db ZMTF_PEP -outfmt 6 -evalue 0.01 在玉米转录因子数据集中搜出 4 个结果。

<pre><!--eb4c@bbt--> 17:11:03 ~/BLAST</pre>									
<pre>\$ blastp -query</pre>	SPL3D_PEP.FASTA	-db ZMT	F_PEP	-outfmt	6 -evalue	0.01			
AT2G33810.1	PTZm00608.1	64.865	74	26	0	3	76	171	2442.16e-31 107
AT2G33810.1	PTZm00605.1	64.865	74	26	0	3	76	144	2173.88e-31 105
AT2G33810.1	PTZm00606.1	56.757	37	16	0	3	39	180	2168.94e-13 55.5
AT2G33810.1	PTZm00607.1	51,220	41	20	0	3	43	171	2111.26e-10 50.1

image-20210423172102619

2. **blastn** 以拟南芥 SPL3 转录因子编码区核苷酸序列搜索 764 个玉米转录因子数据集编码区核苷酸序列

blastn -query SPL3D_CDS.FASTA -db ZMTF_CDS -outfmt 6 -evalue 0.01 没有搜出结果。说明此方法灵敏度不高。

3. **tblastn** 以拟南芥 SPL3 转录因子 DNA 结合结构域搜索 764 个玉米转录因子 数据集编码区核苷酸序列

tblastn -query SPL3D_PEP.FASTA -db ZMTF_CDS -outfmt 6 -evalue 0.0 1

在玉米转录因子数据集中搜出4个结果。

```
<leb4c@bbt> 17:21:15 ~/BLAST
$ tblastn -query SPL3D_PEP.FASTA -db ZMTF_CDS -outfmt 6 -evalue 0.01
                               64.865 74
64.865 74
AT2G33810.1
               PTZm00608.1
                                               26
AT2G33810.1
               PTZm00605.1
                                                       0
                                                                                                2.86e-30
                                                                                                                105
                                                                                1279
AT2G33810.1
               PTZm00607.1
                                58.000 50
                                                21
                                                                                        1428
                                                                                               2.44e-17
                                                                                                                69.7
AT2G33810.1
               PTZm00606.1
                               58,140 43
                                               18
                                                                                630
                                                                                        758
                                                                                                5.86e-16
                                                                                                                65.1
```

image-20210423172338501

4. **blastx** 以拟南芥 SPL3 转录因子编码区核苷酸序列搜索 764 个玉米转录因子数据集蛋白质序列

blastx -query SPL3D_CDS.FASTA -db ZMTF_PEP -outfmt 6 -evalue 0.01 在玉米转录因子数据集中搜出 4 个结果。

<le>Acount 17:22:55 ~/RLAST \$ blastx -query SPL3D_CDS.FASTA -db ZMTF_PEP -outfmt 6 -evalue 0.01 AT2G33810.1 PTZm00608.1 64.865 74 228 2.40e-31 107 171 244 26 AT2G33810.1 PTZm00605.1 64.865 74 26 217 4.31e-31 AT2G33810.1 PTZm00606.1 180 9.93e-13 55.5 56.757 37 117 AT2G33810.1 PTZm00607.1 51.220 41 1.40e-10 50.1

image-20210423172406342

5. **tblastx** 以拟南芥 SPL3 转录因子编码区核苷酸序列(翻译后)搜索 764 个 玉米转录因子数据集编码区核苷酸序列

tblastx -query SPL3D_CDS.FASTA -db ZMTF_CDS -outfmt 6 -evalue 0.0 1

在玉米转录因子数据集中搜出5个结果。最后两项来自同一条序列。

```
$ tblastx -query SPL3D_CDS.FASTA -db ZMTF_CDS -outfmt 6 -evalue 0.01
AT2G33810.1
                PT7m00608.1
                                64.865 74
                                                 26
                                                         0
                                                                         228
                                                                                 549
                                                                                          328
                                                                                                  4.31e-32
                                                                                                                  128
AT2G33810.1
                PTZm00605.1
                                64.865 74
                                                 26
                                                         0
                                                                         228
                                                                                 727
                                                                                          948
                                                                                                  5.91e-32
                                                                                                                  128
AT2G33810.1
                PTZm00607.1
                                58.000 50
                                                 21
                                                         0
                                                                         156
                                                                                 1279
                                                                                         1428
                                                                                                  5.05e-18
                                                                                                                  82.2
AT2G33810.1
                PTZm00606.1
                                58.140 43
                                                 18
                                                                         135
                                                                                 630
                                                                                          758
                                                                                                  4.31e-16
                                                                                                                  75.8
                                                                 137
AT2G33810.1
                PTZm00606.1
                                56.818
                                                                                                  1.39e-06
```

image-20210423172616898

BLAST 的主要特点就是: 速度快, 共线性输出结果简单易读。对于比较小的序列 (如 cDNA 等) 对大基因组的比对, BLAST 无疑是首选。

BLAST 虽然性能优异,但是它自身也存在着一定的局限性,对于特殊的任务需要注意选择合适的软件。例如 BLAST 用于远亲缘物种间的核酸序列比对时,比对精度就不够高,建议使用专门为此用途开发的 Blastz 软件。

6.2 vi 文本编辑

vi/vim 共分为三种模式,分别是命令模式 (Command mode),输入模式 (Insert mode)和底线命令模式 (Last line mode)。

6.2.1 命令模式

此状态下敲击键盘动作会被 vim 识别为命令,而非输入字符。

以下是常用的几个命令:

- i——切换到输入模式,以输入字符。
- x ——删除当前光标所在处的字符。
- :——切换到底线命令模式,以在最底一行输入命令。

命令模式只有一些最基本的命令,因此仍要依靠底线命令模式输入更多命令。

6.2.2 输入模式

在命令模式下按下 i 就进入了输入模式。

在输入模式中,可以使用以下按键:

字符按键以及 Shift 组合——输入字符

ENTER (回车键) ——换行

BACKSPACE (退格键)——删除光标前一个字符

DEL (删除键) ——删除光标后一个字符

方向键——在文本中移动光标

HOME/END——移动光标到行首/行尾

Page Up/Page Down——上/下翻页

Insert——切换光标为输入/替换模式,光标将变成竖线/下划线

ESC——退出输入模式,切换到命令模式

6.2.3 底线命令模式

在命令模式下按下:(英文冒号)就进入了底线命令模式。

底线命令模式可以输入单个或多个字符的命令,可用的命令非常多。

在底线命令模式中,基本的命令有(已经省略了冒号):

q——退出程序

w——保存文件

下面列举一些 vi 中常用的指令。

6.2.4 光标移动的方法

操作 说明

h 或 向左 光标向左移动一个字符

箭头键(←)

j 或 向下 光标向下移动一个字符

箭头键(↓)

k 或 向上 光标向上移动一个字符

箭头键(↑)

1或 向右 光标向右移动一个字符

箭头键(→)

如果你将右手放在键盘上的话,你会发现 hjkl 是排列在一起的,因此可以使用这四个按钮来移动光标。如果想要进行多次移动的话,例如向下移动 30 行,可以使用 "30j" 或 "30↓" 的组合按键,亦即加上想要进行的次数(数字)后,按下动作。

[Ctrl] + [f] 屏幕『向下』移动一页,相当于 [Page Down]按键

[Ctrl] + [b] 屏幕『向上』移动一页,相当于 [Page Up] 按键

[Ctrl] + [d] 屏幕『向下』移动半页

[Ctrl] + [u] 屏幕『向上』移动半页

n 那个 n 表示『数字』,例如 20。按下数字后再按空格键,光标会向 右移动这一行的 n 个字符。例如 20 则光标会向后面移动 20 个字符距 离。

0 或功能 这是数字『0』: 移动到这一行的最前面字符处

键[Home]

\$或功能 移动到这一行的最后面字符处

键[End]

H 光标移动到这个屏幕的最上方那一行的第一个字符

L 光标移动到这个屏幕的最下方那一行的第一个字符

G 移动到这个档案的最后一行

nG n 为数字。移动到这个档案的第 n 行。例如 20G 则会移动到这个档案

的第 20 行(可配合:set nu)

gg 移动到这个档案的第一行,相当于 1G。

n enter 为回车。n 为数字。光标向下移动 n 行

6.2.5 搜索替换的方法

命令	说明
/word	向光标之下寻找一个名称为 word 的字符串。
?word	向光标之上寻找一个字符串名称为 word 的字符串。
	使用 /word 配合 n 及 N 是非常有帮助的! 可以让你 重复的找到一些你搜寻的关键词!

n1 与 n2 为数字。在第 n1 与 n2 行之间寻找 word1 :n1,n2s/word1/word2/g

> 这个字符串,并将该字符串取代为word2! 举例来 说, 在 100 到 200 行之间搜寻 line 并取代为 LINE

则: 『:100,200s/line/LINE/g』。

:1,\$s/word1/word2/g 从第一行到最后一行寻找 word1 字符串,并将该字

或:%s/word1/word2/g 符串取代为 word2

:1,\$s/word1/word2/gc 从第一行到最后一行寻找 word1 字符串,并将该字 或:%s/word1/word2/gc

符串取代为 word2! 且在取代前显示提示字符给用

户确认 (confirm) 是否需要取代

6.2.6 删除、复制、粘贴的方法

命令	说明
x, X	在一行字当中, x 为向后删除一个字符 (相当于 [del] 按键), X 为向前删除一个字符(相当于 [backspace] 亦即 是退格键)
nx	n 为数字,连续向后删除 n 个字符。举例来说,我要连续删除 10 个字符, $\lceil 10x \rceil$ 。
dd	删除游标所在的那一整行
ndd	n 为数字。删除光标所在的向下 n 行,例如 20dd 则是删除 20 行
уу	复制游标所在的那一行(常用)
nyy	n 为数字。复制光标所在的向下 n 行,例如 20yy 则是复制 20 行(常用)
y0	复制光标所在的那个字符到该行行首的所有数据
y\$	复制光标所在的那个字符到该行行尾的所有数据
p, P	p 为将已复制的数据在光标下一行贴上, P 则为贴在游标上一行!
J	将光标所在行与下一行的数据结合成同一行
С	重复删除多个数据,例如向下删除 10 行,[10cj]
u	复原前一个动作。(常用)
[Ctrl]+r	重做上一个动作。(常用)
•	小数点,意思是重复前一个动作的意思。 如果你想要重复删除、重复贴上等等动作,按下小数点『.』就好了

6.2.7 储存、离开等指令

命令	说明
: W	将编辑的数据写入硬盘档案中(常用)
: q	离开 vi (常用)

:q! 若曾修改过档案,又不想储存,使用!为强制离开不储存档案。

:wq 储存后离开, 若为:wq! 则为强制储存后离开 (常用)

:w [filename] 将编辑的数据储存成另一个档案(类似另存新档)

:r [filename] 在编辑的数据中,读入另一个档案的数据。亦即将 『filename』

这个档案内容加到游标所在行后面

:n1,n2 w
[filename]

将 n1 到 n2 的内容储存成 filename 这个档案。

:! command 暂时离开 vi 到指令行模式下执行 command 的显示结果! 例如

『:! ls /home』即可在 vi 当中察看 /home 底下 以 ls 输出的档案

信息

6.3 关于 MyNCBI

MyNCBI 的功能之一:保存的搜索策略 (Saved Search Strategies)

保存方法: 在使用个人账户登录到 NCBI 并进行 BLAST 之后,采取以下两种操作之一:

1) 在想要保存的 BLAST 报告页面点击 Save Search 按钮;

2) 在 Recent Results 页面,找到想要保存的检索条目,点击对应的 Save 按钮。

用途:保存一次检索使用的程序类型,输入序列和各种参数,以后可以随时查看此次检索的结果。

7问题与建议

无。