

Ejercicio de parseo de archivos FASTA

Pablo Vinuesa

2018-02-23

Contents

Presentación	1
Búsqueda y descarga de secuencias en GenBank usando el sistema ENTREZ	1
Práctica de parseo de archivos FASTA descargados de NCBI mediante ENTREZ	1
Acceso a las secuencias	1
Inspección y estadísticas básicas de las secuencias descargadas	1
Edición de las cabeceras FASTA mediante herramientas de filtrado de UNIX	2
Generación automática de archivos FASTA especie-específicos (avanzado)	3

Presentación

Este código corresponde a unas prácticas escritas por Pablo Vinuesa para el manual de Bioinformática y Sistemática Molecular de la Facultad de Ciencias - UNAM, Abril 2015.

Para correr los ejercicios, asegúrate de tener el archivo `recA_Bradyrhizobium_vinuesa.fna` en el directorio actual de trabajo.

Búsqueda y descarga de secuencias en GenBank usando el sistema ENTREZ

El archivo `recA_Bradyrhizobium_vinuesa.fna` contiene secuencias del gen *recA* de bacterias del género *Bradyrhizobium* depositadas en GenBank por P. Vinuesa. Este bloque muestra el comando usado para descargarlas. El comando debe pegarse en la ventana superior del sistema ENTREZ.

```
# pega esta sentencia en la ventana de captura para interrogar la base de datos de nucleótidos
# de NCBI mediante el sistema ENTREZ
'Bradyrhizobium[orgn] AND vinuesa[auth] AND recA[gene]'
```

Práctica de parseo de archivos FASTA descargados de NCBI mediante ENTREZ

Acceso a las secuencias

```
cd $HOME/intro2genomics
cp -r /home/vinuesa/cursos/intro2genomics/sesion1_parseo_fastas .
cd sesion1_parseo_fastas
```

Inspección y estadísticas básicas de las secuencias descargadas

1. ¿Cuántas secuencias hay en el archivo `recA_Bradyrhizobium_vinuesa.fna`?

```
grep -c '>' recA_Bradyrhizobium_vinuesa.fna
```

```
## 117
```

2. Veamos las 5 primeras líneas de cabeceras fasta usando **grep** y **head**

```
grep '>' recA_Bradyrhizobium_vinuesa.fna | head -5
```

```
## >gi|145926563|gb|EF190191.1| Bradyrhizobium japonicum strain X6-9 RNA polymerase beta subunit (rpoB)
## >gi|145926559|gb|EF190189.1| Bradyrhizobium elkanii strain USDA 94 RNA polymerase beta subunit (rpoB)
## >gi|145926555|gb|EF190187.1| Bradyrhizobium elkanii strain USDA 46 RNA polymerase beta subunit (rpoB)
## >gi|145926551|gb|EF190185.1| Bradyrhizobium yuanmingense strain TAL760 RNA polymerase beta subunit (rpoB)
## >gi|145926547|gb|EF190183.1| Bradyrhizobium japonicum strain Nep1 RNA polymerase beta subunit (rpoB)
```

3. Cuenta el numero de generos y especies que contiene el archivo FASTA

```
grep '>' recA_Bradyrhizobium_vinuesa.fna | cut -d' ' -f3 | sort | uniq -c
```

```
##      12 canariense
##      18 elkanii
##       3 genosp.
##      28 japonicum
##      15 liaoningense
##       9 sp.
##      32 yuanmingense
```

4. Imprime una lista ordenada de mayor a menor, del numero de especies que contiene el archivo FASTA

```
grep '>' recA_Bradyrhizobium_vinuesa.fna | cut -d' ' -f2,3 | sort | uniq -c | sort -nrk1
```

```
##      32 Bradyrhizobium yuanmingense
##      28 Bradyrhizobium japonicum
##      18 Bradyrhizobium elkanii
##      15 Bradyrhizobium liaoningense
##      12 Bradyrhizobium canariense
##       9 Bradyrhizobium sp.
##       3 Bradyrhizobium genosp.
```

Edición de las cabeceras FASTA mediante herramientas de filtrado de UNIX

5. Exploremos todas las cabeceras FASTA del archivo recA_Bradyrhizobium_vinuesa.fna usando **grep**

```
# grep '>' recA_Bradyrhizobium_vinuesa.fna | less # para verlas por página
grep '>' recA_Bradyrhizobium_vinuesa.fna | head # para no hacer muy extensa la salida
```

```
## >gi|145926563|gb|EF190191.1| Bradyrhizobium japonicum strain X6-9 RNA polymerase beta subunit (rpoB)
## >gi|145926559|gb|EF190189.1| Bradyrhizobium elkanii strain USDA 94 RNA polymerase beta subunit (rpoB)
## >gi|145926555|gb|EF190187.1| Bradyrhizobium elkanii strain USDA 46 RNA polymerase beta subunit (rpoB)
## >gi|145926551|gb|EF190185.1| Bradyrhizobium yuanmingense strain TAL760 RNA polymerase beta subunit (rpoB)
## >gi|145926547|gb|EF190183.1| Bradyrhizobium japonicum strain Nep1 RNA polymerase beta subunit (rpoB)
## >gi|145926543|gb|EF190181.1| Bradyrhizobium liaoningense strain LMG 18230 RNA polymerase beta subunit (rpoB)
## >gi|145926539|gb|EF190179.1| Bradyrhizobium sp. IRBG 131 RNA polymerase beta subunit (rpoB) gene, partial
## >gi|145926535|gb|EF190177.1| Bradyrhizobium japonicum strain FN13 RNA polymerase beta subunit (rpoB)
## >gi|145926531|gb|EF190175.1| Bradyrhizobium sp. CIAT 3101 RNA polymerase beta subunit (rpoB) gene, partial
## >gi|145926527|gb|EF190173.1| Bradyrhizobium sp. BTAi1 RNA polymerase beta subunit (rpoB) gene, partial
```

6. simplifiquemos las cabeceras FASTA usando el comando **sed** (stream editor)

El objetivo es eliminar redundancia y los campos gb|no.de.acceso, así como todos los caracteres ‘(, ; :)’ que impedirían el despliegue de un árbol filogenético, al tratarse de caracteres reservados del formato NEWICK. Dejar solo el numero GI, así como el género, especie y cepa indicados entre corchetes.

Es decir vamos a: - reducir Bradyrhizobium a ‘B.’ - eliminar ‘RNA poly ...’ y reemplazarlo por ‘]’ - eliminar ‘genosp.’ - sustituir espacios por guiones bajos

Noten el uso de expresiones regulares como `.*'y'[[[:space:]]]`

```
sed 's/ Bra/ [Bra/; s/|gb.*| /|/; s/Bradyrhizobium /B./; s/genosp\. //; s/ RNA.*//; s/[[[:space:]]/_/g;

## >gi|145926563| [B.japonicum_strain_X6-9]
## >gi|145926559| [B.elkanii_strain_USDA_94]
## >gi|145926555| [B.elkanii_strain_USDA_46]
## >gi|145926551| [B.yuanmingense_strain_TAL760]
## >gi|145926547| [B.japonicum_strain_Nep1]
```

8. Cuando estamos satisfechos con el resultado, guardamos la salida del comando en un archivo usando `>` para redirigir el flujo de STDOUT a un archivo de texto

```
sed 's/ recom.*cds//; s/ Bra/ [Bra/; s/|gb.*| /|/; s/Bradyrhizobium /B. /; s/genosp\. //; s/ RNA.*\/\|/
```

Generación automática de archivos FASTA especie-específicos (avanzado)

9. Convertir archivos FASTA a formato “FASTAB” usando `perl` 1-liners.

Vamos a transformar los FASTAS de tal manera que las secuencias queden en la misma línea que su cabecera, separada de ésta por un tabulador. Esto puede ser muy útil para filtrar el archivo resultante con `grep`. Veamos un ejemplo:

```
perl -pe 'unless(/^>){s/\n//g}; if(/^>){s/\n\t/g}; s/>/\n>/' recA_Bradyrhizobium_vinuesa.fnaed | head

##
## >gi|145926563| [B._japonicum_strain_X6-9] TCGGGCAAGCTGCTGTTTCGCCGCCCGCGTCATCCCGTATCGCGGCTCCTGGCTCGACAT
## >gi|145926559| [B._elkanii_strain_USDA_94] TCCGGCAAGCTGCTGTTTGCCGCGCGCGTGATTCCGTATCGCGGTTCTGGCTCG
## >gi|145926555| [B._elkanii_strain_USDA_46] TCCGGCAAGCTGCTGTTTGCCGCGCGCGTGATTCCGTATCGCGGTTCTGGCTCG
## >gi|145926551| [B._yuanmingense_strain_TAL760] TCGGGCAAGCTGCTGTTTCGCCGCCCGCGTCATCCCGTATCGCGGCTCCTGG
## >gi|145926547| [B._japonicum_strain_Nep1] TCGGGCAAGCTGCTGTTTCGCCGCCCGCGTGATTCCGTATCGCGGCTCCTGGCTCGACAT
## >gi|145926543| [B._liaoningense_strain_LMG_18230] TCGGGCAAGCTGCTGTTTCGCCGCCCGCGTCATCCCGTATCGCGGCTCCTGG
## >gi|145926539| [B._sp._IRBG_131] TCGGGCAAGCTGCTGTTTCGCCGCCCGCGTCATTCCGTATCGCGGCTCCTGGCTCGACATCGAGTTCG
## >gi|145926535| [B._japonicum_strain_FN13] TCGGGCAAGCTGCTGTTTCGCCGCCCGCGTCATCCCGTATCGCGGCTCCTGGCTCGACAT
## >gi|145926531| [B._sp._CIAT_3101] TCGGGCAAGCTGCTGTTTCGCTGCCCGCGTGATCCCGTATCGCGGCTCCTGGCTCGACATCGAGTTCG

perl -pe 'unless(/^>){s/\n//g}; if(/^>){s/\n\t/g}; s/>/\n>/' recA_Bradyrhizobium_vinuesa.fnaed > recA
```

10. Filtrar el archivo `fnaedtab` generado en 9 para obtener solo las secuencias de `B._yuanmingense` del mismo, guardarlo en un archivo y convertirlo de nuevo a formato FASTA.

```
grep yuanmingense recA_Bradyrhizobium_vinuesa.fnaedtab | head -5

## >gi|145926551| [B._yuanmingense_strain_TAL760] TCGGGCAAGCTGCTGTTTCGCCGCCCGCGTCATCCCGTATCGCGGCTCCTGG
## >gi|145926545| [B._yuanmingense_strain_LMTR28] TCGGGCAAGCTGCTGTTTCGCCGCCCGCGTCATCCCGTATCGCGGCTCCTGG
## >gi|145926529| [B._yuanmingense_strain_CCBAU10071] TCGGGCAAGCTGCTGTTTCGCCGCCCGCGTCATCCCGTATCGCGGTTCT
## >gi|190612281| [B._yuanmingense_strain_ViHaG5] TCGGGCAAGCTGCTGTTTCGCCGCCCGCGTCATCCCGTATCGCGGCTCCTGG
## >gi|190612237| [B._yuanmingense_strain_InRo02] TCGGGCAAGCTGCTGTTTCGCCGCCCGCGTCATCCCGTATCGCGGCTCCTGG

grep yuanmingense recA_Bradyrhizobium_vinuesa.fnaedtab > recA_Byuanmingense.fnaedtab
```

11. Estas dos líneas no contienen nada nuevo en cuanto a sintaxis. Simplemente llamamos a `perl` para sustituir los tabuladores por saltos de línea y así reconstituir el FASTA.

```
perl -pe 'if(/^>){s/\t/\n/}' recA_Byuanmingense.fnaedtab | head -5

## >gi|145926551| [B._yuanmingense_strain_TAL760]
## TCGGGCAAGCTGCTGTTTCGCCGCCCGCGTCATCCCGTATCGCGGCTCCTGGCTCGACATCGAGTTTCGACGCCAAGGACATCGTCTATGCGCGTATCGAC
## >gi|145926545| [B._yuanmingense_strain_LMTR28]
```

```
## TCGGGCAAGCTGCTGTTCGCCGCCCCGCGTCATCCCGTATCGCGGCTCCTGGCTCGACATCGAGTTCGACGCCAAGGACATCGTCTATGCGCGTATCGACCG
## >gi|145926529|[B._yuanmingense_strain_CCBau10071]
```

```
perl -pe 'if(</>){s/\t/\n/}' recA_Byuanmingense.fnaedtab > recA_Byuanmingense.fna
```

12. Llamar a un bucle for de shell para generar archivos fastab para todas las especies

```
for sp in $(grep '>' recA_Bradyrhizobium_vinuesa.fnaed | cut -d_ -f2); do
  grep "$sp" recA_Bradyrhizobium_vinuesa.fnaedtab > recA_B${sp}.fnaedtab
done
```

13. Veamos el resultado

```
ls *fnaedtab
```

```
## recA_Balpha.fnaedtab
## recA_Bbeta.fnaedtab
## recA_Bcanariense.fnaedtab
## recA_Belkanii.fnaedtab
## recA_Bjaponicum.fnaedtab
## recA_Bliaoningense.fnaedtab
## recA_Bradyrhizobium_vinuesa.fnaedtab
## recA_Bsp..fnaedtab
## recA_Byuanmingense.fnaedtab
```

```
head -5 recA_Bjaponicum.fnaedtab
```

```
## >gi|145926563|[B._japonicum_strain_X6-9] TCGGGCAAGCTGCTGTTCGCCGCCCCGCGTCATCCCGTATCGCGGCTCCTGGCTCGACAT
## >gi|145926547|[B._japonicum_strain_Nep1] TCGGGCAAGCTGCTGTTCGCCGCCCCGCGTGATTCCGTATCGCGGCTCCTGGCTCGACAT
## >gi|145926535|[B._japonicum_strain_FN13] TCGGGCAAGCTGCTGTTCGCCGCCCCGCGTCATCCCGTATCGCGGCTCCTGGCTCGACAT
## >gi|145926519|[B._japonicum_strain_BGA-1] TCGGGCAAGCTGCTGTTCGCCGCCCCGCGTCATCCCGTATCGCGGCTCCTGGCTCG
## >gi|145926561|[B._japonicum_strain_X3-1] TCGGGCAAGCTGCTGTTCGCCGCCCCGCGTCATCCCGTATCGCGGCTCCTGGCTCGACAT
```

14. Finalmente convertimos todos los archivos fnatabed a FASTA con el siguiente bucle for:

```
for file in *fnaedtab; do perl -pe 'if(</>){s/\t/\n/}' $file > ${file%.*}.fna; done
```

15. Visualizemos las cabeceras de dos archivos FASTA especie-específicos

```
grep '>' recA_Bjaponicum.fna | head -5
```

```
## >gi|145926563|[B._japonicum_strain_X6-9]
## >gi|145926547|[B._japonicum_strain_Nep1]
## >gi|145926535|[B._japonicum_strain_FN13]
## >gi|145926519|[B._japonicum_strain_BGA-1]
## >gi|145926561|[B._japonicum_strain_X3-1]
```

16. y confirmemos que son fastas regulares

```
head -6 recA_Bjaponicum.fna
```

```
## >gi|145926563|[B._japonicum_strain_X6-9]
## TCGGGCAAGCTGCTGTTCGCCGCCCCGCGTCATCCCGTATCGCGGCTCCTGGCTCGACATCGAGTTCGACGCCAAGGACATCGTCTATGCGCGTATCGACCG
## >gi|145926547|[B._japonicum_strain_Nep1]
## TCGGGCAAGCTGCTGTTCGCCGCCCCGCGTGATTCCGTATCGCGGCTCCTGGCTCGACATCGAGTTCGATGCCAAGGACATCGTCTATGCGCGTATCGACCG
## >gi|145926535|[B._japonicum_strain_FN13]
## TCGGGCAAGCTGCTGTTCGCCGCCCCGCGTCATCCCGTATCGCGGCTCCTGGCTCGACATCGAGTTCGACGCCAAGGACATCGTCTATGCGCGTATCGACCG
```