

Ejercicio de parseo de archivos FASTA

Pablo Vinuesa

2019-02-22

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	2
Edición de las cabeceras FASTA mediante herramientas de filtrado de UNIX	7
Generación automática de archivos FASTA especie-específicos (avanzado)	7

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

mkdir sesion1_parseo_fastas

cd sesion1_parseo_fastas

cp -r /home/vinuesa/cursos/intro2genomics/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
```

```
## 125
```

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

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

```
## >EU574327.1_B._liaoningense_strain_ViHaR5
## >EU574326.1_B._liaoningense_strain_ViHaR4
## >EU574325.1_B._liaoningense_strain_ViHaR3
## >EU574324.1_B._liaoningense_strain_ViHaR2
## >EU574323.1_B._liaoningense_strain_ViHaR1
```

3. Cuenta el número de géneros y especies que contiene el archivo FASTA

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

```
##      1 >AY591540.1_B._alpha_bv._genistearum_strain_BC-C1
##      1 >AY591541.1_B._canariense_bv._genistearum_strain_BC-C2
##      1 >AY591542.1_B._canariense_bv._genistearum_strain_BC-P5
##      1 >AY591543.1_B._beta_strain_BC-P6
##      1 >AY591544.1_B._japonicum_bv._genistearum_strain_BC-P14
##      1 >AY591545.1_B._canariense_bv._genistearum_strain_BC-P22
##      1 >AY591546.1_B._canariense_bv._genistearum_strain_BC-MAM1
##      1 >AY591547.1_B._canariense_bv._genistearum_strain_BC-MAM5
##      1 >AY591548.1_B._canariense_bv._genistearum_strain_BES-1
##      1 >AY591549.1_B._canariense_bv._genistearum_strain_BES-2
##      1 >AY591550.1_B._canariense_bv._genistearum_strain_BC0-1
##      1 >AY591551.1_B._beta_strain_BRE-1
##      1 >AY591552.1_B._canariense_bv._genistearum_strain_BRE-4
##      1 >AY591553.1_B._canariense_bv._genistearum_strain_BTA-1
##      1 >AY591554.1_B._beta_strain_BC-MK6
##      1 >AY591555.1_B._japonicum_bv._glycinearum_strain_DSMZ30131
##      1 >AY591556.1_B._japonicum_bv._glycinearum_strain_X3-1
##      1 >AY591557.1_B._japonicum_bv._glycinearum_strain_X6-9
##      1 >AY591558.1_B._japonicum_bv._genistearum_strain_BGA-1
##      1 >AY591559.1_B._japonicum_bv._genistearum_strain_BLup-MR1
##      1 >AY591560.1_B._japonicum_bv._genistearum_strain_FN13
##      1 >AY591561.1_B._sp._CICS70
##      1 >AY591562.1_B._japonicum_bv._glycinearum_strain_USDA122
##      1 >AY591563.1_B._japonicum_bv._glycinearum_strain_Nep1
##      1 >AY591564.1_B._liaoningense_bv._glycinearum_strain_LMG18230
##      1 >AY591565.1_B._yuanmingense_strain_TAL760
##      1 >AY591566.1_B._yuanmingense_strain_CCBAU_10071
##      1 >AY591567.1_B._alpha_strain_CIAT3101
##      1 >AY591568.1_B._elkanii_strain_USDA76
##      1 >AY591569.1_B._elkanii_strain_USDA94
##      1 >AY591570.1_B._sp._BTai1
##      1 >AY591571.1_B._sp._IRBG127
##      1 >AY591572.1_B._sp._IRBG231
##      1 >AY591573.1_B._yuanmingense_strain_LMTR28
##      1 >AY591574.1_B._liaoningense_strain_Spr3-7
```

```

##      1 >AY591575.1_B._elkanii_strain_USDA46
##      1 >AY591576.1_B._canariense_strain_ISLU16
##      1 >AY653743.1_B._canariense_strain_BC-P24
##      1 >AY653744.1_B._canariense_strain_BC-MAM2
##      1 >AY653745.1_B._canariense_strain_BC-MAM6
##      1 >AY653746.1_B._canariense_strain_BC-MAM8
##      1 >AY653747.1_B._canariense_strain_BC-MAM9
##      1 >AY653748.1_B._canariense_strain_BC-MAM11
##      1 >AY653749.1_B._canariense_strain_BC-MAM12
##      1 >AY653750.1_B._beta_strain_BC-MK1
##      1 >EU574248.1_B._yuanmingense_strain_BuCeG2
##      1 >EU574249.1_B._yuanmingense_strain_BuCeG3
##      1 >EU574250.1_B._yuanmingense_strain_BuCeG4
##      1 >EU574251.1_B._sp._BuCeR1
##      1 >EU574252.1_B._sp._BuCeR2
##      1 >EU574253.1_B._yuanmingense_strain_BuCeR3
##      1 >EU574254.1_B._yuanmingense_strain_BuCeR4
##      1 >EU574255.1_B._yuanmingense_strain_BuCeR5
##      1 >EU574256.1_B._elkanii_strain_BuMiN1
##      1 >EU574257.1_B._elkanii_strain_BuMiN2
##      1 >EU574258.1_B._elkanii_strain_BuMiN3
##      1 >EU574259.1_B._elkanii_strain_BuMiN4
##      1 >EU574260.1_B._liaoningense_strain_BuMiN6
##      1 >EU574261.1_B._elkanii_strain_BuMiT1
##      1 >EU574262.1_B._sp._BuMiT10
##      1 >EU574263.1_B._liaoningense_strain_BuMiT3
##      1 >EU574264.1_B._liaoningense_strain_BuMiT4
##      1 >EU574265.1_B._liaoningense_strain_BuMiT5
##      1 >EU574266.1_B._elkanii_strain_BuMiT6
##      1 >EU574267.1_B._elkanii_strain_BuMiT7
##      1 >EU574268.1_B._elkanii_strain_BuMiT8
##      1 >EU574269.1_B._elkanii_strain_BuMiT9
##      1 >EU574270.1_B._elkanii_strain_BuNoG1
##      1 >EU574271.1_B._elkanii_strain_BuNoG4
##      1 >EU574272.1_B._sp._BuNoG5
##      1 >EU574273.1_B._elkanii_strain_BuNoR1
##      1 >EU574274.1_B._elkanii_strain_BuNoR2
##      1 >EU574275.1_B._elkanii_strain_BuNoR3
##      1 >EU574276.1_B._elkanii_strain_BuNoR4
##      1 >EU574277.1_B._yuanmingense_strain_InBu02
##      1 >EU574278.1_B._yuanmingense_strain_InIn01
##      1 >EU574279.1_B._yuanmingense_strain_InIn02
##      1 >EU574280.1_B._yuanmingense_strain_InIn03
##      1 >EU574281.1_B._yuanmingense_strain_InIn04
##      1 >EU574282.1_B._yuanmingense_strain_InIn05
##      1 >EU574283.1_B._yuanmingense_strain_InIn08
##      1 >EU574284.1_B._yuanmingense_strain_InIn09
##      1 >EU574285.1_B._yuanmingense_strain_InIn10
##      1 >EU574286.1_B._yuanmingense_strain_InJa01
##      1 >EU574287.1_B._yuanmingense_strain_InJa02
##      1 >EU574288.1_B._yuanmingense_strain_InJa03
##      1 >EU574289.1_B._yuanmingense_strain_InJa04
##      1 >EU574290.1_B._yuanmingense_strain_InJa05
##      1 >EU574291.1_B._yuanmingense_strain_InJa06

```

```

##      1 >EU574292.1_B._yuanmingense_strain_InJa07
##      1 >EU574293.1_B._yuanmingense_strain_InJa08
##      1 >EU574294.1_B._yuanmingense_strain_InJa09
##      1 >EU574295.1_B._yuanmingense_strain_InKo01
##      1 >EU574296.1_B._yuanmingense_strain_InKo02
##      1 >EU574297.1_B._yuanmingense_strain_InRo02
##      1 >EU574298.1_B._japonicum_strain_NeMa01
##      1 >EU574299.1_B._japonicum_strain_NeMa02
##      1 >EU574300.1_B._japonicum_strain_NeMa10
##      1 >EU574301.1_B._japonicum_strain_NeMa11
##      1 >EU574302.1_B._japonicum_strain_NeMa12
##      1 >EU574303.1_B._japonicum_strain_NeMa16
##      1 >EU574304.1_B._japonicum_strain_NeRa01
##      1 >EU574305.1_B._japonicum_strain_NeRa02
##      1 >EU574306.1_B._japonicum_strain_NeRa03
##      1 >EU574307.1_B._japonicum_strain_NeRa04
##      1 >EU574308.1_B._japonicum_strain_NeRa05
##      1 >EU574309.1_B._japonicum_strain_NeRa06
##      1 >EU574310.1_B._japonicum_strain_NeRa07
##      1 >EU574311.1_B._japonicum_strain_NeRa08
##      1 >EU574312.1_B._japonicum_strain_NeRa11
##      1 >EU574313.1_B._japonicum_strain_NeRa12
##      1 >EU574314.1_B._japonicum_strain_NeRa14
##      1 >EU574315.1_B._japonicum_strain_NeRa15
##      1 >EU574316.1_B._japonicum_strain_NeRa16
##      1 >EU574317.1_B._liaoningense_strain_ViHaG3
##      1 >EU574318.1_B._yuanmingense_strain_ViHaG4
##      1 >EU574319.1_B._yuanmingense_strain_ViHaG5
##      1 >EU574320.1_B._liaoningense_strain_ViHaG6
##      1 >EU574321.1_B._liaoningense_strain_ViHaG7
##      1 >EU574322.1_B._liaoningense_strain_ViHaG8
##      1 >EU574323.1_B._liaoningense_strain_ViHaR1
##      1 >EU574324.1_B._liaoningense_strain_ViHaR2
##      1 >EU574325.1_B._liaoningense_strain_ViHaR3
##      1 >EU574326.1_B._liaoningense_strain_ViHaR4
##      1 >EU574327.1_B._liaoningense_strain_ViHaR5

```

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
```

```

##      1 >EU574327.1_B._liaoningense_strain_ViHaR5
##      1 >EU574326.1_B._liaoningense_strain_ViHaR4
##      1 >EU574325.1_B._liaoningense_strain_ViHaR3
##      1 >EU574324.1_B._liaoningense_strain_ViHaR2
##      1 >EU574323.1_B._liaoningense_strain_ViHaR1
##      1 >EU574322.1_B._liaoningense_strain_ViHaG8
##      1 >EU574321.1_B._liaoningense_strain_ViHaG7
##      1 >EU574320.1_B._liaoningense_strain_ViHaG6
##      1 >EU574319.1_B._yuanmingense_strain_ViHaG5
##      1 >EU574318.1_B._yuanmingense_strain_ViHaG4
##      1 >EU574317.1_B._liaoningense_strain_ViHaG3
##      1 >EU574316.1_B._japonicum_strain_NeRa16
##      1 >EU574315.1_B._japonicum_strain_NeRa15
##      1 >EU574314.1_B._japonicum_strain_NeRa14

```

```

##      1 >EU574313.1_B._japonicum_strain_NeRa12
##      1 >EU574312.1_B._japonicum_strain_NeRa11
##      1 >EU574311.1_B._japonicum_strain_NeRa08
##      1 >EU574310.1_B._japonicum_strain_NeRa07
##      1 >EU574309.1_B._japonicum_strain_NeRa06
##      1 >EU574308.1_B._japonicum_strain_NeRa05
##      1 >EU574307.1_B._japonicum_strain_NeRa04
##      1 >EU574306.1_B._japonicum_strain_NeRa03
##      1 >EU574305.1_B._japonicum_strain_NeRa02
##      1 >EU574304.1_B._japonicum_strain_NeRa01
##      1 >EU574303.1_B._japonicum_strain_NeMa16
##      1 >EU574302.1_B._japonicum_strain_NeMa12
##      1 >EU574301.1_B._japonicum_strain_NeMa11
##      1 >EU574300.1_B._japonicum_strain_NeMa10
##      1 >EU574299.1_B._japonicum_strain_NeMa02
##      1 >EU574298.1_B._japonicum_strain_NeMa01
##      1 >EU574297.1_B._yuanmingense_strain_InRo02
##      1 >EU574296.1_B._yuanmingense_strain_InKo02
##      1 >EU574295.1_B._yuanmingense_strain_InKo01
##      1 >EU574294.1_B._yuanmingense_strain_InJa09
##      1 >EU574293.1_B._yuanmingense_strain_InJa08
##      1 >EU574292.1_B._yuanmingense_strain_InJa07
##      1 >EU574291.1_B._yuanmingense_strain_InJa06
##      1 >EU574290.1_B._yuanmingense_strain_InJa05
##      1 >EU574289.1_B._yuanmingense_strain_InJa04
##      1 >EU574288.1_B._yuanmingense_strain_InJa03
##      1 >EU574287.1_B._yuanmingense_strain_InJa02
##      1 >EU574286.1_B._yuanmingense_strain_InJa01
##      1 >EU574285.1_B._yuanmingense_strain_InIn10
##      1 >EU574284.1_B._yuanmingense_strain_InIn09
##      1 >EU574283.1_B._yuanmingense_strain_InIn08
##      1 >EU574282.1_B._yuanmingense_strain_InIn05
##      1 >EU574281.1_B._yuanmingense_strain_InIn04
##      1 >EU574280.1_B._yuanmingense_strain_InIn03
##      1 >EU574279.1_B._yuanmingense_strain_InIn02
##      1 >EU574278.1_B._yuanmingense_strain_InIn01
##      1 >EU574277.1_B._yuanmingense_strain_InBu02
##      1 >EU574276.1_B._elkanii_strain_BuNoR4
##      1 >EU574275.1_B._elkanii_strain_BuNoR3
##      1 >EU574274.1_B._elkanii_strain_BuNoR2
##      1 >EU574273.1_B._elkanii_strain_BuNoR1
##      1 >EU574272.1_B._sp._BuNoG5
##      1 >EU574271.1_B._elkanii_strain_BuNoG4
##      1 >EU574270.1_B._elkanii_strain_BuNoG1
##      1 >EU574269.1_B._elkanii_strain_BuMiT9
##      1 >EU574268.1_B._elkanii_strain_BuMiT8
##      1 >EU574267.1_B._elkanii_strain_BuMiT7
##      1 >EU574266.1_B._elkanii_strain_BuMiT6
##      1 >EU574265.1_B._liaoningense_strain_BuMiT5
##      1 >EU574264.1_B._liaoningense_strain_BuMiT4
##      1 >EU574263.1_B._liaoningense_strain_BuMiT3
##      1 >EU574262.1_B._sp._BuMiT10
##      1 >EU574261.1_B._elkanii_strain_BuMiT1
##      1 >EU574260.1_B._liaoningense_strain_BuMiN6

```

```

##      1 >EU574259.1_B._elkanii_strain_BuMiN4
##      1 >EU574258.1_B._elkanii_strain_BuMiN3
##      1 >EU574257.1_B._elkanii_strain_BuMiN2
##      1 >EU574256.1_B._elkanii_strain_BuMiN1
##      1 >EU574255.1_B._yuanmingense_strain_BuCeR5
##      1 >EU574254.1_B._yuanmingense_strain_BuCeR4
##      1 >EU574253.1_B._yuanmingense_strain_BuCeR3
##      1 >EU574252.1_B._sp._BuCeR2
##      1 >EU574251.1_B._sp._BuCeR1
##      1 >EU574250.1_B._yuanmingense_strain_BuCeG4
##      1 >EU574249.1_B._yuanmingense_strain_BuCeG3
##      1 >EU574248.1_B._yuanmingense_strain_BuCeG2
##      1 >AY653750.1_B._beta_strain_BC-MK1
##      1 >AY653749.1_B._canariense_strain_BC-MAM12
##      1 >AY653748.1_B._canariense_strain_BC-MAM11
##      1 >AY653747.1_B._canariense_strain_BC-MAM9
##      1 >AY653746.1_B._canariense_strain_BC-MAM8
##      1 >AY653745.1_B._canariense_strain_BC-MAM6
##      1 >AY653744.1_B._canariense_strain_BC-MAM2
##      1 >AY653743.1_B._canariense_strain_BC-P24
##      1 >AY591576.1_B._canariense_strain_ISLU16
##      1 >AY591575.1_B._elkanii_strain_USDA46
##      1 >AY591574.1_B._liaoningense_strain_Spr3-7
##      1 >AY591573.1_B._yuanmingense_strain_LMTR28
##      1 >AY591572.1_B._sp._IRBG231
##      1 >AY591571.1_B._sp._IRBG127
##      1 >AY591570.1_B._sp._BTai1
##      1 >AY591569.1_B._elkanii_strain_USDA94
##      1 >AY591568.1_B._elkanii_strain_USDA76
##      1 >AY591567.1_B._alpha_strain_CIAT3101
##      1 >AY591566.1_B._yuanmingense_strain_CCAU_10071
##      1 >AY591565.1_B._yuanmingense_strain_TAL760
##      1 >AY591564.1_B._liaoningense_bv._glycinearum_strain_LMG18230
##      1 >AY591563.1_B._japonicum_bv._glycinearum_strain_Nep1
##      1 >AY591562.1_B._japonicum_bv._glycinearum_strain_USDA122
##      1 >AY591561.1_B._sp._CICS70
##      1 >AY591560.1_B._japonicum_bv._genistearum_strain_FN13
##      1 >AY591559.1_B._japonicum_bv._genistearum_strain_BLup-MR1
##      1 >AY591558.1_B._japonicum_bv._genistearum_strain_BGA-1
##      1 >AY591557.1_B._japonicum_bv._glycinearum_strain_X6-9
##      1 >AY591556.1_B._japonicum_bv._glycinearum_strain_X3-1
##      1 >AY591555.1_B._japonicum_bv._glycinearum_strain_DSMZ30131
##      1 >AY591554.1_B._beta_strain_BC-MK6
##      1 >AY591553.1_B._canariense_bv._genistearum_strain_BTA-1
##      1 >AY591552.1_B._canariense_bv._genistearum_strain_BRE-4
##      1 >AY591551.1_B._beta_strain_BRE-1
##      1 >AY591550.1_B._canariense_bv._genistearum_strain_BC0-1
##      1 >AY591549.1_B._canariense_bv._genistearum_strain_BES-2
##      1 >AY591548.1_B._canariense_bv._genistearum_strain_BES-1
##      1 >AY591547.1_B._canariense_bv._genistearum_strain_BC-MAM5
##      1 >AY591546.1_B._canariense_bv._genistearum_strain_BC-MAM1
##      1 >AY591545.1_B._canariense_bv._genistearum_strain_BC-P22
##      1 >AY591544.1_B._japonicum_bv._genistearum_strain_BC-P14
##      1 >AY591543.1_B._beta_strain_BC-P6

```

```
##      1 >AY591542.1_B._canariense_bv._genistearum_strain_BC-P5
##      1 >AY591541.1_B._canariense_bv._genistearum_strain_BC-C2
##      1 >AY591540.1_B._alpha_bv._genistearum_strain_BC-C1
```

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
```

```
## >EU574327.1_B._liaoningense_strain_ViHaR5
## >EU574326.1_B._liaoningense_strain_ViHaR4
## >EU574325.1_B._liaoningense_strain_ViHaR3
## >EU574324.1_B._liaoningense_strain_ViHaR2
## >EU574323.1_B._liaoningense_strain_ViHaR1
## >EU574322.1_B._liaoningense_strain_ViHaG8
## >EU574321.1_B._liaoningense_strain_ViHaG7
## >EU574320.1_B._liaoningense_strain_ViHaG6
## >EU574319.1_B._yuanmingense_strain_ViHaG5
## >EU574318.1_B._yuanmingense_strain_ViHaG4
```

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 número 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;
```

```
## >EU574327.1_B._liaoningense_strain_ViHaR5
## >EU574326.1_B._liaoningense_strain_ViHaR4
## >EU574325.1_B._liaoningense_strain_ViHaR3
## >EU574324.1_B._liaoningense_strain_ViHaR2
## >EU574323.1_B._liaoningense_strain_ViHaR1
```

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.*//; s/
```

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
```

```
##
## >EU574327.1_B._liaoningense_strain_ViHaR5      ATGAAGCTCGGCAAGAACGACCGGTCCATGGACATCGAGGCGGTGTCCTCCGGCT
```

```
## >EU574326.1_B._liaoningense_strain_ViHaR4 ATGAAGCTCGGCAAGAACGACCGGTCCATGGACATCGAGGCGGTGTCTCCGGCT
## >EU574325.1_B._liaoningense_strain_ViHaR3 ATGAAGCTCGGCAAGAACGACCGGTCCATGGACATCGAGGCGGTGTCTCCGGCT
## >EU574324.1_B._liaoningense_strain_ViHaR2 ATGAAGCTCGGCAAGAACGACCGGTCCATGGACATCGAGGCGGTGTCTCCGGCT
## >EU574323.1_B._liaoningense_strain_ViHaR1 ATGAAGCTCGGCAAGAACGACCGGTCCATGGACATCGAGGCGGTGTCTCCGGCT
## >EU574322.1_B._liaoningense_strain_ViHaG8 ATGAAGCTCGGCAAGAACGACCGGTCCATGGACATCGAGGCGGTGTCTCCGGCT
## >EU574321.1_B._liaoningense_strain_ViHaG7 ATGAAGCTCGGCAAGAACGACCGGTCCATGGACATCGAGGCGGTGTCTCCGGCT
## >EU574320.1_B._liaoningense_strain_ViHaG6 ATGAAGCTCGGCAAGAACGACCGGTCCATGGACATCGAGGCGGTGTCTCCGGCT
## >EU574319.1_B._yuanmingense_strain_ViHaG5 ATGAAGCTCGGCAAGAACGACCGCTCCATGGACATCGAGGCGGTGTCTCCGGCT

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

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
```

```
## >EU574319.1_B._yuanmingense_strain_ViHaG5 ATGAAGCTCGGCAAGAACGACCGCTCCATGGACATCGAGGCGGTGTCTCCGGCT
## >EU574318.1_B._yuanmingense_strain_ViHaG4 ATGAAGCTCGGCAAGAACGACCGCTCCATGGACATCGAGGCGGTGTCTCCGGCT
## >EU574297.1_B._yuanmingense_strain_InRo02 ATGAAGCTCGGCAAGAACGACCGCTCCATGGACATCGAGGCGGTGTCTCCGGCT
## >EU574296.1_B._yuanmingense_strain_InKo02 ATGAAGCTCGGCAAGAACGATCGTCCATGGACATCGAGGCGGTCTCTCCGGCT
## >EU574295.1_B._yuanmingense_strain_InKo01 ATGAAGCTCGGCAAGAACGATCGTCCATGGACATCGAGGCGGTGTCTCCGGCT

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
```

```
## >EU574319.1_B._yuanmingense_strain_ViHaG5
## ATGAAGCTCGGCAAGAACGACCGCTCCATGGACATCGAGGCGGTGTCTCCGGCTCGCTCGGGCTCGATATCGCGCTCGGCATCGGCGGCTTGCCCAAGG
## >EU574318.1_B._yuanmingense_strain_ViHaG4
## ATGAAGCTCGGCAAGAACGACCGCTCCATGGACATCGAGGCGGTGTCTCCGGCTCGCTCGGGCTCGATATCGCGCTCGGCATCGGCGGCTTGCCCAAGG
## >EU574297.1_B._yuanmingense_strain_InRo02
## ATGAAGCTCGGCAAGAACGACCGCTCCATGGACATCGAGGCGGTGTCTCCGGCTCGCTCGGGCTCGATATCGCGCTCGGCATCGGCGGCTTGCCCAAGG

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 $(grep '>' recA_Bradyrhizobium_vinuesa.fnaed | cut -d_ -f3); do
  grep "$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
```

```
## >EU574316.1_B._japonicum_strain_NeRa16 ATGAAGCTCGGCAAGAACGACCGGTTCGATGGATGTCGAGGCGGTGTCTCCGGTTCTCT
```



```
## >EU574315.1_B._japonicum_strain_NeRa15    ATGAAGCTCGGCAAGAACGACCGGTCGATGGATGTCGAGGCGGTGTCCTCCGGTTCTCT
## >EU574314.1_B._japonicum_strain_NeRa14    ATGAAGCTCGGCAAGAACGACCGGTCGATGGATGTCGAGGCGGTGTCCTCCGGTTCTCT
## >EU574313.1_B._japonicum_strain_NeRa12    ATGAAGCTCGGCAAGAACGACCGGTCGATGGATGTCGAGGCGGTGTCCTCCGGTTCTCT
## >EU574312.1_B._japonicum_strain_NeRa11    ATGAAGCTCGGCAAGAACGACCGGTCGATGGATGTCGAGGCGGTGTCCTCCGGTTCTCT
```

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
```

```
## >EU574316.1_B._japonicum_strain_NeRa16
## >EU574315.1_B._japonicum_strain_NeRa15
## >EU574314.1_B._japonicum_strain_NeRa14
## >EU574313.1_B._japonicum_strain_NeRa12
## >EU574312.1_B._japonicum_strain_NeRa11
```

16. y confirmemos que son fastas regulares

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

```
## >EU574316.1_B._japonicum_strain_NeRa16
## ATGAAGCTCGGCAAGAACGACCGGTCGATGGATGTCGAGGCGGTGTCCTCGGGTCTCTCGGGCTCGACATTGCACTGGGGATCGGCGGTCTGCCAAGG
## >EU574315.1_B._japonicum_strain_NeRa15
## ATGAAGCTCGGCAAGAACGACCGGTCGATGGATGTCGAGGCGGTGTCCTCCGGTTCTCTCGGGCTCGACATTGCACTGGGGATCGGCGGTCTGCCAAGG
## >EU574314.1_B._japonicum_strain_NeRa14
## ATGAAGCTCGGCAAGAACGACCGGTCGATGGATGTCGAGGCGGTGTCCTCCGGTTCTCTCGGGCTCGACATTGCGCTGGGGATCGGCGGTCTGCCAAGG
```