

fastx-utils之获取序列集合的子集或补集：subseq

一、fastx-utils subseq介绍

功能描述：

`fastx-utils subseq` 根据提供的序列 `ID`，获取序列集合的交集序列或者补集序列。

命令行接口：

```
1 $ fastx-utils subseq
2
3 Usage: fastx-utils subseq [options] <in.fa/fq> <name.list>
4
5 Options:
6     -s    supplementary set option.
```

可选参数：

```
1 -s    补充设置选项；
```

二、使用场景实例及其用法

示例演示：

示例演示：

示例文件： `protein.faa`, `list.txt`

```
1 $ cat protein.faa | head -n 6
```

```
1 >AGE80385.1
2 MKKIIIIISATTIVIGITSFAYFGSKTPLHNEAKAVESQKHNNHKKEEIPAFPKADHNAKKIDNDFS VVTNPKSNLVL
  INKHRKLPGYIPEDLTRPNVPFISPKDKEKTLRKDAEAELENMFKA AKKEGLDLTAVSGYRSYKRQKSLHDTYVR
  RQGKAEANSVSAIPGTSEHQTGLAMDISSKSAKFQLEPIFGETAEGKWAEHAHEFGFVIRYLEDKTDTTEYAYEPW
  HLRVYGNPYATYLYKHHLTLEEAMEDKK
3 >AGE79558.1
4 MINHELVERINFLAKKAKAEGLTEEEQRRERQSLREQYLKGFRQNMLNELKGIKVVNEEGTDVTPTKLKALKKQDNAK
  LN
5 >AGE81073.1
6 MDKVISNEILQQFKDRMRLGDDDEDANLRRILFASNKDLTRVCGNYDLNIDEVFKELVFERSRYVYNDALEYFDKNFL
  SQINSLSIGKALEAIKLDGD
```

```
1 $ cat list.txt | head -n 6
```

1	AGE75591.1	chromosomal replication initiation protein [Bacillus thuringiensis serovar kurstaki str. HD73]
2	AGE75592.1	DNA polymerase III, beta subunit [Bacillus thuringiensis serovar kurstaki str. HD73]
3	AGE75593.1	hypothetical protein HD73_0003 [Bacillus thuringiensis serovar kurstaki str. HD73]
4	AGE75594.1	recombination protein F [Bacillus thuringiensis serovar kurstaki str. HD73]
5	AGE75595.1	DNA gyrase subunit B [Bacillus thuringiensis serovar kurstaki str. HD73]
6	AGE75596.1	DNA gyrase, A subunit [Bacillus thuringiensis serovar kurstaki str. HD73]

运行命令：

参数选项**1**： 默认参数, 求列表序列的交集。

```
1 $ fastx-utils subseq protein.faa list.txt | head -n 6
```

```
1 >AGE79558.1
2 MINHELVERINFLAKKAKAEGLTEEEQRERQSLREQYLKGFQNMNLKGIKVVNEEGTDVTPTKLKALKKQDNAK
  LN
3 >AGE76491.1
4 MKEMTAKELEEKLLRKEAVNIVDVREVEEVAEGKIPEACNIPLGLLEFRMHeldKKKEYIIVCRSGGRSARAVQFLE
  SYGFQAINMVGMLAWEGKV
5 >AGE77971.1
6 MEFQLLVTCILQEGNAYFLVTKVDDVITLKVPIAGVAGLFLALGVPRCS
```

参数选项**1**： 设置 **-s**, 求列表序列的补集。

```
1 $ fastx-utils subseq protein.faa -s list.txt | head -n 6
```

```
1 >AGE80385.1
2 MKKIIISATTIVIGITSFAYFGSKTPLHNEAKAVESQKHNNHKKEEIPAFPKADHNAKKIDNDFS SVVTNPKSNLVL
  INKHKRLPDGYIPEDLTRPNVPFISPKDKEKTLRKDAEALENMFKA AKKEGLDLTAVSGYRSYKRQKSLHDTYVR
  RQGKAEANSVSAIPGTSEHQTGLAMDISSKSAKFQLEPIFGETAEGKWAEHAHEFGFVIRYLEDKTDTTEYAYEPW
  HLRYVGNPYATYLYKHHLTLEEAMEDKK
3 >AGE81073.1
4 MDKVISNEILQQFKDRMRLGDDEDANLRRILFASNKDLTRVCGNYDLNIDEVFKELVFERSRYVYNDALEYFDKNFL
  SQINLSIGKALEAIKLDGD
5 >AGE81522.1
6 MKKTILTTAALTMMGTGMGINVDHIKPAEVKADTISFYDVPNNHWATKAITNLANRNIVVGYGNGQFGFGDNVTRG
  QVARMYNYLKPADAGNFKNPFSDIKGHMFEKEILALAKVGIKGYGEGKFGPDDILTREQMAQVL TNAFKFEGTKK
  TSFVDVDKNSWSYKAIGALEEKGVITGTGGNMYSPSVVTRQYSQFLFNSINVIEKETKPEEKPN TGGEVKPEEK
  NTGGEVKPEEKPN TGGEVKPN IPEWLETSLATNDFTFTQAWYDGSEAINKAASTNAQQIVKNINSKYGTNLKYSEV
  GAIVQLVDGAREQLWLAGMNVNDFRVTFRVSNAMIETKELVTLVNSDLNLDQEIQEIPSAPMKIKNVEKG DYKIR
  ISPAMADQMITIIIEKK
```

本文材料为 **BASE (Biostack Applied bioinformatic SEies)** 课程 **Linux Command Line Tools for Life Scientists** 材料，版权归上海逻捷信息科技有限公司 所有。

Last Update: 2020-08-10 11:56 AM

