# fastx-utils之截断序列：replace

## 一、fastx-utils replace介绍

功能描述：

`fastx-utils replace` 根据 `key/value` 键值对，修改序列的名字。

命令行接口：

```
1  $ fastx-utils  replace
2
3  Usage: fastx-utils replace [options] <db> <fasta>
4
5  Options:
6    -r  remove unmapped sequence.
7    -c  remove comment.
```

可选参数：

```
1  -r        移除无映射序列；
2  -c        移除注释信息；
```

## 二、使用场景实例及其用法

示例演示：

示例文件： `protein.faa  convert.txt`

```
1  $ cat  protein.faa | head -n  8
```

```
1  >AGE80385.1
2  MKKIIIISATTIVIGITSFAYFGSKTPLHNEAKAVESQKHNNHKKEEIPAFPKADHNAKKIDNDFSVVTNPKSNLVL
   INKHRKLPDGYIPEDLTRPNVPFISPKDKEKTLLRKDAAEALENMFKAAKKEGLDLTAVSGYRSYKRQKSLHDTYVR
   RQGKAEANSVSAIPGTSEHQTGLAMDISSKSAKFQLEPIFGETAEGKWVAEHAHEFGFVIRYLEDKTDTTEYAYEPW
   HLRYVGNPYATYLYKHHLTLEEAMEDKK
3  >AGE79558.1
4  MINHELVERINFLAKKAKAEGLTEEEQRERQSLREQYLKGFRQNMLNELKGIKVVNEEGTDVTPTKLKALKKQDNAK
   LN
5  >AGE81073.1
6  MDKVISNEILQQFKDRMRLGDDEDANLRRILFASNKDLTRVCGNYDLNIDEVFKELVFERSRYVYNDALEYFDKNFL
   SQINSLSIGKALEAIKLDGD
7  >AGE81522.1
8  MKKTILTTTAALTMMGTGMGINVDHIKPAEVKADTISFYDVPNNHWATKAITNLANRNIVVGYGNGQFGFGDNVTRG
   QVARMIYNYLKPADAGNFKNPFSDIKGHMFEKEILALAKVGIIKGYGEGKFGPDDILTREQMAQVLTNAFKFEGTKK
   TSFVDVDKNSWSYKAIGALEEKGVTIGTGGNMYSPTSVVTREQYSQFLFNSINVIEKETKPEEKPNTGGEVKPEEKP
   NTGGEVKPEEKPNTGEETKPVNIPEWLETSLATNDFTFTQAWYDGSEAINKAASTNAQQIVKNINSKYGTNLKYSEV
   GAIVQLVDGAREQLWLAGMNVNDFRVTFRVSNNAMIELTKELVTLVNSDLNLDQEIQEIPSAPMKIKNVEKGDYKIR
   ISPAMADQMITIIIEKK
```

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

```
1  AGE75591.1      HD73_0001
2  AGE75592.1      HD73_0002
3  AGE75593.1      HD73_0003
4  AGE75594.1      HD73_0004
5  AGE75595.1      HD73_0005
6  AGE75596.1      HD73_0006
```

运行命令： 使用 `-r` 参数去除不能匹配的序列

```
1  $ fastx-utils  replace -r  convert.txt  protein.faa  | head -n 8
```

```
1  >HD73_4807
2  MKKIIIISATTIVIGITSFAYFGSKTPLHNEAKAVESQKHNNHKKEEIPAFPKADHNAKKIDNDFSVVTNPKSNLVL
   INKHRKLPDGYIPEDLTRPNVPFISPKDKEKTLLRKDAAEALENMFKAAKKEGLDLTAVSGYRSYKRQKSLHDTYVR
   RQGKAEANSVSAIPGTSEHQTGLAMDISSKSAKFQLEPIFGETAEGKWVAEHAHEFGFVIRYLEDKTDTTEYAYEPW
   HLRYVGNPYATYLYKHHLTLEEAMEDKK
3  >HD73_3980
4  MINHELVERINFLAKKAKAEGLTEEEQRERQSLREQYLKGFRQNMLNELKGIKVVNEEGTDVTPTKLKALKKQDNAK
   LN
5  >HD73_5496
6  MDKVISNEILQQFKDRMRLGDDEDANLRRILFASNKDLTRVCGNYDLNIDEVFKELVFERSRYVYNDALEYFDKNFL
   SQINSLSIGKALEAIKLDGD
7  >HD73_6041
8  MKKTILTTTAALTMMGTGMGINVDHIKPAEVKADTISFYDVPNNHWATKAITNLANRNIVVGYGNGQFGFGDNVTRG
   QVARMIYNYLKPADAGNFKNPFSDIKGHMFEKEILALAKVGIIKGYGEGKFGPDDILTREQMAQVLTNAFKFEGTKK
   TSFVDVDKNSWSYKAIGALEEKGVTIGTGGNMYSPTSVVTREQYSQFLFNSINVIEKETKPEEKPNTGGEVKPEEKP
   NTGGEVKPEEKPNTGEETKPVNIPEWLETSLATNDFTFTQAWYDGSEAINKAASTNAQQIVKNINSKYGTNLKYSEV
   GAIVQLVDGAREQLWLAGMNVNDFRVTFRVSNNAMIELTKELVTLVNSDLNLDQEIQEIPSAPMKIKNVEKGDYKIR
   ISPAMADQMITIIIEKK
```

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