

pan\_roary.py 的目标是根据数量超过特定阈值的每个物种，分别创建文件夹。在每个文件夹内，生成一个 txt 文件，记录该物种的所有 MAGs，一个 pan.summary.csv 文件，记录了每个物种的 Core gene 及非 Core gene 的情况。并在物种文件夹下创建指向 Prokka 注释的 gff 文件的软链接。输入文件包括 species.txt 和 tax.bac120.summary.tsv 两个文件。

其中，species.txt 记录了每个组装的 MAG 的数量：

```
1 d_Bacteria;p_Actinomycetota;c_Coribacteriia;o_Coribacteriales;f_Coribacteriaceae;g_Collinsella;Collinsella sp003436275
1 d_Bacteria;p_Bacillota_A;c_Clostridia;o_Clostridiales;f_Clostridiaceae;g_Clostridium_3s_Clostridium_3 culturomicum
6 d_Bacteria;p_Bacillota_A;c_Clostridia;o_Clostridiales;f_Clostridiaceae;g_Clostridium butyricum
1 d_Bacteria;p_Bacillota_A;c_Clostridia;o_Eubacteriales;f_Eubacteriaceae;g_Eubacterium;Eubacterium limosum
1 d_Bacteria;p_Bacillota_A;c_Clostridia;o_Lachnospirales;f_Lachnospiraceae;g_Enterocloster;Enterocloster sp021201905
2 d_Bacteria;p_Bacillota_A;c_Clostridia;o_Peptostreptococcales;f_Peptostreptococcaceae;g_Paraclostridium;P_Paraclostridium dentum
1 d_Bacteria;p_Bacillota;c_Bacilli;o_Bacillales;f_Bacillaceae;g_Bacillus;Bacillus altitudinis
1 d_Bacteria;p_Bacillota;c_Bacilli;o_Bacillales;f_Bacillaceae;g_Bacillus A;Bacillus A bombysepticus
9 d_Bacteria;p_Bacillota;c_Bacilli;o_Bacillales;f_Bacillaceae;g_Bacillus A;s_Bacillus A cereus
6 d_Bacteria;p_Bacillota;c_Bacilli;o_Bacillales;f_Bacillaceae;g_Bacillus A;s_Bacillus A paranthracis
1 d_Bacteria;p_Bacillota;c_Bacilli;o_Bacillales;f_Bacillaceae;g_Bacillus A;s_Priestia;Priestia flexa
3 d_Bacteria;p_Bacillota;c_Bacilli;o_Erysipelotrichales;f_Erysipelotrichaceae;g_Faecalicoccus;Faecalicoccus pleomorphus
35 d_Bacteria;p_Bacillota;c_Bacilli;o_Lactobacillales;f_Enterococcaceae;g_Enterococcus A;s_Enterococcus A avium
7 d_Bacteria;p_Bacillota;c_Bacilli;o_Lactobacillales;f_Enterococcaceae;g_Enterococcus B;s_Enterococcus B faecium
7 d_Bacteria;p_Bacillota;c_Bacilli;o_Lactobacillales;f_Enterococcaceae;g_Enterococcus B;s_Enterococcus B hirae
7 d_Bacteria;p_Bacillota;c_Bacilli;o_Lactobacillales;f_Enterococcaceae;g_Enterococcus C;s_Enterococcus C asiaticus
1 d_Bacteria;p_Bacillota;c_Bacilli;o_Lactobacillales;f_Enterococcaceae;g_Enterococcus C;s_Enterococcus C sp009933135
2 d_Bacteria;p_Bacillota;c_Bacilli;o_Lactobacillales;f_Enterococcaceae;g_Enterococcus D;s_Enterococcus D casseliflavus
11 d_Bacteria;p_Bacillota;c_Bacilli;o_Lactobacillales;f_Enterococcaceae;g_Enterococcus D;s_Enterococcus D gallinarum
6 d_Bacteria;p_Bacillota;c_Bacilli;o_Lactobacillales;f_Enterococcaceae;g_Enterococcus D;s_Enterococcus D innesii
50 d_Bacteria;p_Bacillota;c_Bacilli;o_Lactobacillales;f_Enterococcaceae;g_Enterococcus;Enterococcus faecalis
1 d_Bacteria;p_Bacillota;c_Bacilli;o_Lactobacillales;f_Lactobacillaceae;g_Lactocaseibacillus;Lactocaseibacillus
10 d_Bacteria;p_Bacillota;c_Bacilli;o_Lactobacillales;f_Lactobacillaceae;g_Lactocaseibacillus;s_Lactocaseibacillus paracasei
13 d_Bacteria;p_Bacillota;c_Bacilli;o_Lactobacillales;f_Lactobacillaceae;g_Ligilactobacillus;s_Ligilactobacillus salivarius
16 d_Bacteria;p_Bacillota;c_Bacilli;o_Lactobacillales;f_Lactobacillaceae;g_Limosilactobacillus;s_Limosilactobacillus fermentum
3 d_Bacteria;p_Bacillota;c_Bacilli;o_Lactobacillales;f_Lactobacillaceae;g_Pediococcus;s_Pediococcus pentosaceus
13 d_Bacteria;p_Bacillota;c_Bacilli;o_Lactobacillales;f_Streptococcaceae;g_Lactococcus;s_Lactococcus lactis
1 d_Bacteria;p_Bacillota;c_Bacilli;o_Lactobacillales;f_Streptococcaceae;g_Streptococcus;s_Streptococcus anginosus
2 d_Bacteria;p_Bacillota;c_Bacilli;o_Staphylococcales;f_Staphylococcaceae;g_Mammaliococcus;s_Mammaliococcus sciuri
1 d_Bacteria;p_Bacillota;c_Bacilli;o_Staphylococcales;f_Staphylococcaceae;g_Staphylococcus;s_Staphylococcus simulans
3 d_Bacteria;p_Bacillota;c_Negativicutes;o_Acidaminococcales;f_Acidaminococcaceae;g_Acidaminococcus;s_Acidaminococcus intestinalis
3 d_Bacteria;p_Bacteroidia;c_Bacteroidia;o_Bacteroidales;f_Bacteroidaceae;g_Bacteroides;s_Bacteroides caccae
1 d_Bacteria;p_Bacteroidia;c_Bacteroidia;o_Bacteroidales;f_Bacteroidaceae;g_Bacteroides;s_Bacteroides fragilis
3 d_Bacteria;p_Bacteroidia;c_Bacteroidia;o_Bacteroidales;f_Bacteroidaceae;g_Bacteroides;s_Bacteroides uniformis
1 d_Bacteria;p_Bacteroidia;c_Bacteroidia;o_Bacteroidales;f_Bacteroidaceae;g_Bacteroides;s_Bacteroides xylanisolvens
1 d_Bacteria;p_Bacteroidia;c_Bacteroidia;o_Bacteroidales;f_Marinifilaceae;g_Butyricimonas;s_Butyricimonas sp900184685
2 d_Bacteria;p_Bacteroidia;c_Bacteroidia;o_Bacteroidales;f_Marinifilaceae;g_Odoribacter;s_Odoribacter splanchnicus
1 d_Bacteria;p_Bacteroidia;c_Bacteroidia;o_Bacteroidales;f_Rikenellaceae;g_Alistipes A;s_Alistipes A humii
1 d_Bacteria;p_Bacteroidia;c_Bacteroidia;o_Bacteroidales;f_Rikenellaceae;g_Alistipes;s_Alistipes senegalensis
2 d_Bacteria;p_Bacteroidia;c_Bacteroidia;o_Bacteroidales;f_Tannerellaceae;g_Parabacteroides;s_Parabacteroides buchesdurhoniensis
41 d_Bacteria;p_Bacteroidia;c_Bacteroidia;o_Bacteroidales;f_Tannerellaceae;g_Parabacteroides;s_Parabacteroides distasonis
```

tax.bac120.summary.tsv，其记载了每个 MAG 所属的 species：

user_genome	classification	fastani_reference	fastani_reference_radius	fastani_taxonomy	fastani_ani	fastani_af	closest_placement	reference	c
2 X001	d_Bacteria;p_Pseudomonadota;c_Gammaproteobacteria;o_Enterobacterales;f_Enterobacteriaceae;g_Escherichia;Escherichia coli	GC_003697165.2	95.0	d_Bacter					
3 X003b1	d_Bacteria;p_Bacillota;c_Bacilli;o_Lactobacillales;f_Enterococcaceae;g_Enterococcus;Enterococcus faecalis	GC_000392875.1	95.0	d_Bacteria;p_Bacillota;					
4 X003b2	d_Bacteria;p_Bacillota;c_Bacilli;o_Lactobacillales;f_Enterobacterales;f_Enterobacteriaceae;g_Escherichia;Escherichia coli	GC_003697165.2	95.0	d_Bacter					
5 X004b1	d_Bacteria;p_Pseudomonadota;c_Gammaproteobacteria;o_Enterobacterales;f_Enterobacteriaceae;g_Escherichia;Escherichia coli	GC_003697165.2	95.0	d_Bacter					
6 X004b2	d_Bacteria;p_Bacillota;c_Bacilli;o_Lactobacillales;f_Enterococcaceae;g_Enterococcus B;s_Enterococcus B hirae	GC_000271405.2	95.0	d_Bacteria;p_Bacillota;					
7 X005b1	d_Bacteria;p_Pseudomonadota;c_Gammaproteobacteria;o_Enterobacterales;f_Enterobacteriaceae;g_Escherichia;Escherichia coli	GC_003697165.2	95.0	d_Bacter					
8 X005b2	d_Bacteria;p_Bacillota;c_Bacilli;o_Lactobacillales;f_Enterococcaceae;g_Enterococcus B;s_Enterococcus B hirae	GC_000271405.2	95.0	d_Bacteria;p_Bacillota;					
9 X006	d_Bacteria;p_Pseudomonadota;c_Gammaproteobacteria;o_Enterobacterales;f_Enterobacteriaceae;g_Escherichia;Escherichia coli	GC_003697165.2	95.0	d_Bacter					
10 X007	d_Bacteria;p_Bacillota;c_Bacilli;o_Lactobacillales;f_Enterococcaceae;g_Enterococcus B;s_Enterococcus B faecium	GC_001544255.1	95.0	d_Bacteria;p_Bacillota;					
11 X008b1	d_Bacteria;p_Bacillota;c_Bacilli;o_Lactobacillales;f_Enterococcaceae;g_Enterococcus B;s_Enterococcus B faecium	GC_001544255.1	95.0	d_Bacteria;p_Bacillota;					
12 X008b2	d_Bacteria;p_Bacillota;c_Bacilli;o_Lactobacillales;f_Enterococcaceae;g_Enterococcus;Enterococcus faecalis	GC_000392875.1	95.0	d_Bacteria;p_Bacillota;					
13 X009	d_Bacteria;p_Pseudomonadota;c_Gammaproteobacteria;o_Enterobacterales;f_Enterobacteriaceae;g_Enterobacter;Enterobacter hormaechei A	GC_001729745.1	95.0	d					
14 X010	d_Bacteria;p_Pseudomonadota;c_Gammaproteobacteria;o_Enterobacterales;f_Enterobacteriaceae;g_Enterobacter;Enterobacter hormaechei A	GC_001729745.1	95.0	d					
15 X011	d_Bacteria;p_Pseudomonadota;c_Gammaproteobacteria;o_Enterobacterales;f_Enterobacteriaceae;g_Enterobacter;Enterobacter hormaechei A	GC_001729745.1	95.0	d					
16 X013b1	d_Bacteria;p_Bacillota;c_Bacilli;o_Lactobacillales;f_Enterococcaceae;g_Enterococcus C;sp009933135	GC_009933135.1	95.0	d_Bacteria;p_Ba					
17 X013b2	d_Bacteria;p_Pseudomonadota;c_Gammaproteobacteria;o_Enterobacterales;f_Enterobacteriaceae;g_Escherichia;Escherichia coli	GC_003697165.2	95.0	d_Bacter					
18 X014b1	d_Bacteria;p_Bacillota;c_Bacilli;o_Lactobacillales;f_Enterococcaceae;g_Enterococcus;Enterococcus faecalis	GC_000392875.1	95.0	d_Bacteria;p_Bacillota;					
19 X014b2	d_Bacteria;p_Pseudomonadota;c_Gammaproteobacteria;o_Enterobacterales;f_Enterobacteriaceae;g_Enterobacter;Enterobacter hormaechei A	GC_001729745.1	95.0	d					
20 X014b3	d_Bacteria;p_Bacillota;c_Bacilli;o_Lactobacillales;f_Enterococcaceae;g_Enterococcus B;s_Enterococcus B faecium	GC_001544255.1	95.0	d_Bacteria;p_Bacillota;					
21 X014b4	d_Bacteria;p_Pseudomonadota;c_Gammaproteobacteria;o_Enterobacterales;f_Enterobacteriaceae;g_Escherichia;Escherichia coli	GC_003697165.2	95.0	d_Bacter					
22 X015	d_Bacteria;p_Pseudomonadota;c_Gammaproteobacteria;o_Enterobacterales;f_Enterobacteriaceae;g_Proteus;Proteus mirabilis	GC_000160755.1	95.0	d_Bacteria;p_Ps					
23 X016b1	d_Bacteria;p_Pseudomonadota;c_Gammaproteobacteria;o_Enterobacterales;f_Enterobacteriaceae;g_Proteus;Proteus mirabilis	GC_000160755.1	95.0	d_Bacteria;p_Ps					
24 X016b2	d_Bacteria;p_Pseudomonadota;c_Gammaproteobacteria;o_Enterobacterales;f_Enterobacteriaceae;g_Enterobacter;Enterobacter hormaechei A	GC_001729745.1	95.0	d					
25 X016b3	d_Bacteria;p_Pseudomonadota;c_Gammaproteobacteria;o_Enterobacterales;f_Enterobacteriaceae;g_Escherichia;Escherichia coli	GC_003697165.2	95.0	d_Bacter					
26 X016b4	d_Bacteria;p_Bacillota;c_Bacilli;o_Lactobacillales;f_Enterococcaceae;g_Enterococcus B;s_Enterococcus B faecium	GC_001544255.1	95.0	d_Bacteria;p_Bacillota;					
27 X018b1	d_Bacteria;p_Bacillota;c_Bacilli;o_Lactobacillales;f_Enterococcaceae;g_Enterococcus;Enterococcus faecalis	GC_000392875.1	95.0	d_Bacteria;p_Bacillota;					
28 X018b2	d_Bacteria;p_Pseudomonadota;c_Gammaproteobacteria;o_Enterobacterales;f_Enterobacteriaceae;g_Proteus;Proteus mirabilis	GC_000160755.1	95.0	d_Bacteria;p_Ps					
29 X019	d_Bacteria;p_Bacillota;c_Bacilli;o_Lactobacillales;f_Enterococcaceae;g_Enterococcus;Enterococcus faecalis	GC_000392875.1	95.0	d_Bacteria;p_Bacillota;					
30 X020b1	d_Bacteria;p_Pseudomonadota;c_Gammaproteobacteria;o_Enterobacterales;f_Enterobacteriaceae;g_Proteus;Proteus mirabilis	GC_000160755.1	95.0	d_Bacteria;p_Ps					
31 X021b1	d_Bacteria;p_Pseudomonadota;c_Gammaproteobacteria;o_Enterobacterales;f_Enterobacteriaceae;g_Proteus;Proteus mirabilis	GC_000160755.1	95.0	d_Bacteria;p_Ps					
32 X021b2	d_Bacteria;p_Pseudomonadota;c_Gammaproteobacteria;o_Enterobacterales;f_Enterobacteriaceae;g_Escherichia;Escherichia coli	GC_003697165.2	95.0	d_Bacter					
33 X021b3	d_Bacteria;p_Bacillota;c_Bacilli;o_Lactobacillales;f_Enterococcaceae;g_Enterococcus A;s_Enterococcus A avium	GC_000406065.1	95.0	d_Bacteria;p_Bacillota;					
34 X021b4	d_Bacteria;p_Bacillota;c_Bacilli;o_Lactobacillales;f_Enterococcaceae;g_Enterococcus D;s_Enterococcus D gallinarum	GC_001544275.1	95.0	d_Bacteria;p_Ba					
35 X022b1	d_Bacteria;p_Pseudomonadota;c_Gammaproteobacteria;o_Enterobacterales;f_Enterobacteriaceae;g_Enterobacter;Enterobacter hormaechei A	GC_001729745.1	95.0	d					
36 X022b2	d_Bacteria;p_Bacillota;c_Bacilli;o_Lactobacillales;f_Enterococcaceae;g_Enterococcus;Enterococcus faecalis	GC_000392875.1	95.0	d_Bacteria;p_Bacillota;					
37 X022b3	d_Bacteria;p_Pseudomonadota;c_Gammaproteobacteria;o_Enterobacterales;f_Enterobacteriaceae;g_Proteus;Proteus mirabilis	GC_000160755.1	95.0	d_Bacteria;p_Ps					
38 X024b1	d_Bacteria;p_Bacillota;c_Bacilli;o_Lactobacillales;f_Enterococcaceae;g_Enterococcus;Enterococcus faecalis	GC_000392875.1	95.0	d_Bacteria;p_Bacillota;					
39 X024b2	d_Bacteria;p_Pseudomonadota;c_Gammaproteobacteria;o_Enterobacterales;f_Enterobacteriaceae;g_Proteus;Proteus mirabilis	GC_000160755.1	95.0	d_Bacteria;p_Ps					
40 X025b1	d_Bacteria;p_Pseudomonadota;c_Gammaproteobacteria;o_Enterobacterales;f_Enterobacteriaceae;g_Proteus;Proteus mirabilis	GC_000160755.1	95.0	d_Bacteria;p_Ps					
41 X025b2	d_Bacteria;p_Bacillota;c_Bacilli;o_Lactobacillales;f_Enterococcaceae;g_Enterococcus;Enterococcus faecalis	GC_000392875.1	95.0	d_Bacteria;p_Bacillota;					
42 X026b1	d_Bacteria;p_Pseudomonadota;c_Gammaproteobacteria;o_Enterobacterales;f_Enterobacteriaceae;g_Proteus;Proteus mirabilis	GC_000160755.1	95.0	d_Bacteria;p_Ps					
43 X026b2	d_Bacteria;p_Bacillota;c_Bacilli;o_Lactobacillales;f_Enterococcaceae;g_Enterococcus D;s_Enterococcus D gallinarum	GC_001544275.1	95.0	d_Bacteria;p_Ba					
44 X028b1	d_Bacteria;p_Pseudomonadota;c_Gammaproteobacteria;o_Enterobacterales;f_Enterobacteriaceae;g_Proteus;Proteus mirabilis	GC_000160755.1	95.0	d_Bacteria;p_Ps					
45 X028b2	d_Bacteria;p_Pseudomonadota;c_Gammaproteobacteria;o_Enterobacterales;f_Enterobacteriaceae;g_Escherichia;Escherichia coli	GC_003697165.2	95.0	d_Bacter					
46 X028b3	d_Bacteria;p_Bacillota;c_Bacilli;o_Lactobacillales;f_Enterococcaceae;g_Enterococcus;Enterococcus faecalis	GC_000392875.1	95.0	d_Bacteria;p_Bacillota;					

使用方法：

python pan\_roary.py -help

-i: 输入文件 1， species.txt 文件

--tax: 输入文件 2， tax.bac120.summary.tsv 文件

-t: genome 数量的阈值，只有超过该数值的 genome 才会被创建文件夹，默认：10

--thread: roary 程序调用的线程数

- donotalign: 是否使用 MAFFT 对基因比对 (roary 内置参数)
- p: Prokka 文件夹。注意必须为绝对路径, 不能使用相对路径
- o: 输出文件, 输出每个种及其对应的 species 数目用于后期作图。

```
$python roary_split.py --help
usage: roary_split.py [-h] -i <file> --tax <file> -t <file> --thread <file>
                        --donotalign <file> -p <file> -o <file> [--version]

Roary results

optional arguments:
  -h, --help            show this help message and exit
  -i <file>, --input <file>
                        Species file: species.txt.
  --tax <file>          Tax file: tax.bac120.summary.tsv.
  -t <file>, --threshold <file>
                        Number of genome, Default: 10.
  --thread <file>       Number of threads, Default: 10.
  --donotalign <file>   Do-not-align genes, Default: T.
  -p <file>, --prokka <file>
                        Prokka directory [ABS]. this directory is usde to link
                        gff file into sub-directory and process roary.
  -o <file>, --out <file>
                        Output file.
  --version             Display version
```