

Demultiplexing Summary

Mitchell Rezzonico

October 16, 2018

Demultiplexing Script

The runtime for my script was about 14 hours. This is likely attributed to Python's sluggish rate of gzipping files. It is also a probably a consequence of checking the index lines for each filtering step.

```
#!/usr/bin/env python3
# module load python3/3.6.1
import argparse
import gzip
def get_arguments():
    parser = argparse.ArgumentParser(description="Allows users to know how to navigate this program")
    parser.add_argument("-r1", "--read1", help="Forward Read File", required=True, type=str)
    parser.add_argument("-i1", "--index1", help="Forward Index File", required=True, type=str)
    parser.add_argument("-i2", "--index2", help="Reverse Index File", required=True, type=str)
    parser.add_argument("-r2", "--read2", help="Reverse Read File", required=True, type=str)
    parser.add_argument("-qs", "--quality_score", help="Cutoff for quality score", required=True, type=int)
    parser.add_argument("-I", "--Index_library", help="Input index library. Indexes must be in the 5th column", required=True, type=str)
    return parser.parse_args()
#####
##### Argparse Variables #####
#####
args = get_arguments()
r1 = args.read1
i1 = args.index1
i2 = args.index2
r2 = args.read2
qscore = args.quality_score
indexes = args.Index_library
#####
#### Higher order functions ####
#####
def reverse_complement(seq):
    """Returns the reverse complement of a DNA sequence"""
    complement = {'A': 'T', 'C': 'G', 'G': 'C', 'T': 'A', 'N': 'N'}
    return ''.join([complement[base] for base in seq[::-1]])
def convert_phred(letter):
    """Converts a single character into a phred score"""
    return(ord(letter)-33)
#####
### Isolate Indexes from file ###
#####
index_list = []
with open(indexes) as i:
    for col in i:
        col = col.strip("\n").split("\t")
        index_list.append(col[4])
index_list = index_list[1:] # remove the header line
#####
# Open and Create Output Files #
#####
ctr_dict = {}
F_index_dict = {}
R_index_dict = {}
for index in index_list:
    Forward = gzip.open(index + "_forward.gz", "wt")
    Reverse = gzip.open(index + "_reverse.gz", "wt")
    F_index_dict[index] = Forward
    R_index_dict[index] = Reverse
    ctr_dict[index] = 0
bad_ctr = 0
LQ_Fz = gzip.open("bad_reads_forward.gz", "wt")
```

```

LQ_Rz = gzip.open("bad_reads_reverse.gz", "wt")
#####
### Bad Read Filter Counters ###
#####
undetermined = 0
low_qscore = 0
false_read = 0
index_hopping = 0
#####
# Open Input Files, Sort Reads #
#####
total_count = 0
with gzip.open(r1, "rt") as r1, gzip.open(i1, "rt") as i1, gzip.open(i2, "rt") as i2, gzip.open(r2, "rt") as r2:
    while True: # keep looping until the break

        r_header1 = r1.readline()
        r_seq1 = r1.readline()
        r_plus1 = r1.readline().strip()
        r_qscore1 = r1.readline()

        # set condition for when the loop finally ends so that there isn't an error
        if r_header1 == '':
            break

        total_count += 1

        i_header1 = i1.readline()
        i_seq1 = i1.readline()
        i_plus1 = i1.readline()
        i_qscore1 = i1.readline()

        i_header2 = i2.readline()
        i_seq2 = i2.readline()
        i_plus2 = i2.readline()
        i_qscore2 = i2.readline()

        r_header2 = r2.readline()
        r_seq2 = r2.readline()
        r_plus2 = r2.readline().strip()
        r_qscore2 = r2.readline()

        isBad = False

        # filter out indexes with undetermined nucleotides
        if "N" in i_seq1 or "N" in i_seq2:
            isBad = True
            undetermined += 1
            r_plus1 += "Undetermined | "
            r_plus2 += "Undetermined | "

        # calculate the mean qscore of the index sequence
        converted_NT = []
        for NT in i_qscore1:
            converted_NT.append(convert_phred(NT))
        converted_NT_mean = sum(converted_NT) / len(converted_NT)

        # filter out indexes with a low mean qscore
        if converted_NT_mean < qscore:
            isBad = True
            low_qscore += 1
            r_plus1 += "Low Quality | "
            r_plus2 += "Low Quality | "

        # filter out indexes that do not appear in the index library
        if i_seq1.strip() not in index_list and reverse_complement(i_seq2.strip()) not in index_list:
            isBad = True
            false_read += 1

```

```

        r_plus1 += "Index not recognized | "
        r_plus2 += "Index not recognized | "

# filter and track occurrences of index hopping
    else:
        if reverse_complement(i_seq1.strip()) != i_seq2.strip():
            isBad = True
            index_hopping += 1
            r_plus1 += "Index hopping | "
            r_plus2 += "Index hopping | "

# build final bad output files
    if isBad:
        LQ_Fz.write(r_header1.strip()+':'+i_seq1.strip()+'\n')
        LQ_Fz.write(r_seq1)
        LQ_Fz.write(r_plus1 + '\n')
        LQ_Fz.write(r_qscore1)

        LQ_Rz.write(r_header2.strip()+':'+i_seq2.strip()+'\n')
        LQ_Rz.write(r_seq2)
        LQ_Rz.write(r_plus2 + '\n')
        LQ_Rz.write(r_qscore2)

        bad_ctr += 1

# once all the bad reads have been filtered out, write good reads to their respective files
    if isBad == False:

        F_index_dict[i_seq1.strip()].write(r_header1.strip()+':'+i_seq1.strip()+'\n')
        F_index_dict[i_seq1.strip()].write(r_seq1)
        F_index_dict[i_seq1.strip()].write(r_plus1 + '\n')
        F_index_dict[i_seq1.strip()].write(r_qscore1)

        R_index_dict[i_seq1.strip()].write(r_header2.strip()+':'+i_seq2.strip()+'\n')
        R_index_dict[i_seq1.strip()].write(r_seq2)
        R_index_dict[i_seq1.strip()].write(r_plus2 + '\n')
        R_index_dict[i_seq1.strip()].write(r_qscore2)

        ctr_dict[i_seq1.strip()] += 1
#####
##### Close Output Files #####
#####

for index in index_list:
    F_index_dict[index].close()
    R_index_dict[index].close()
LQ_Fz.close()
LQ_Rz.close()
#####
## Relevant Output Statements ##
#####
print("INDEX", "\t", "RAW COUNT", "\t", "PROPORTION (%)")
for k in ctr_dict:
    print(k, "\t", ctr_dict[k], "\t", (ctr_dict[k]/total_count) * 100)

bad_proportion = (bad_ctr / total_count) * 100
print("Bad reads compose {}% of the entire file.".format(bad_proportion))
print("There are {} instances of indexes with undetermined NTs in their sequences.".format(undetermined))
print("There are {} instances of indexes with a mean phred quality score below {}".format(low_qscore, qscore))
print("There are {} instances of indexes that are not included in the list of indexes.".format(false_read))
print("There are {} instances of index hopping.".format(index_hopping))

```

Sbatch script submitted to Talapas

```
#!/usr/bin/env bash
```

Test files

4

[illegible]

Output from test files

```

# INDEX      RAW COUNT      PROPORTION (%)
# AAAAAAAAA  1    20.0
# CGATCGAT   0    0.0
# GATCAAGG   0    0.0
# AACAGCGA   0    0.0
# TAGCCATG   0    0.0
# CGGTAATC   0    0.0
# CTCTGGAT   0    0.0
# TACCGGAT   0    0.0
# CTAGCTCA   0    0.0
# CACTTCAC   0    0.0
# GCTACTCT   0    0.0
# ACGATCAG   0    0.0
# TATGGCAC   0    0.0
# TGTTCCGT   0    0.0
# GTCCTAAG   0    0.0
# TCGACAAG   0    0.0
# TCTTCGAC   0    0.0
# ATCATGCG   0    0.0
# ATCGTGGT   0    0.0
# TCGAGAGT   0    0.0
# TCGGATTG   0    0.0
# GATCTTGC   0    0.0
# AGAGTCCA   0    0.0
# AGGATAGC   0    0.0
# Bad reads compose 80.0% of the entire file.
# There are 1 instances of indexes with undetermined NTs in their sequences.
# There are 4 instances of indexes with a mean phred quality score below 30.
# There are 0 instances of indexes that are not included in the list of indexes.
# There are 1 instances of index hopping.

```

Output from assigned samples

#	INDEX	RAW COUNT	PROPORTION (%)
#	GTAGCGTA	6423487	1.768353678388878
#	CGATCGAT	4370218	1.203099045060928
#	GATCAAGG	4877454	1.3427385658401032
#	AACAGCGA	7247905	1.9953118092031854
#	TAGCCATG	8048298	2.2156559782980567
#	CGGTAATC	3693609	1.0168319888683928
#	CTCTGGAT	28060203	7.7248328192130895
#	TACCGGAT	61542468	16.942332048765696
#	CTAGCTCA	13196603	3.632958462792515
#	CACTTCAC	3197050	0.880131792512877
#	GCTACTCT	5162621	1.4212436073238208
#	ACGATCAG	6186925	1.7032293490538877
#	TATGGCAC	8622836	2.3738234013307786
#	TGTTCCGT	10905086	3.0021153528055797
#	GTCTAAG	6461033	1.7786899034343697
#	TCGACAAG	2947217	0.811354023595009
#	TCTTCGAC	34754032	9.56761029111521
#	ATCATGCG	7264829	1.999970901321384
#	ATCGTGGT	4814723	1.3254690369068285
#	TCGAGAGT	5608742	1.544058475845626

```
# TCGGATTC      2840779      0.7820521772893567
# GATCTTGC      2773588      0.7635548327777811
# AGAGTCCA      7709563      2.122403935716036
# AGGATAGC      6771157      1.8640654815520914
# Bad reads compose 30.218113040988513% of the entire file.
# There are 4205183 instances of indexes with undetermined NTs in their sequences.
# There are 99116240 instances of indexes with a mean phred quality score below 30.
# There are 7867981 instances of indexes that are not included in the list of indexes.
# There are 23623721 instances of index hopping.
```

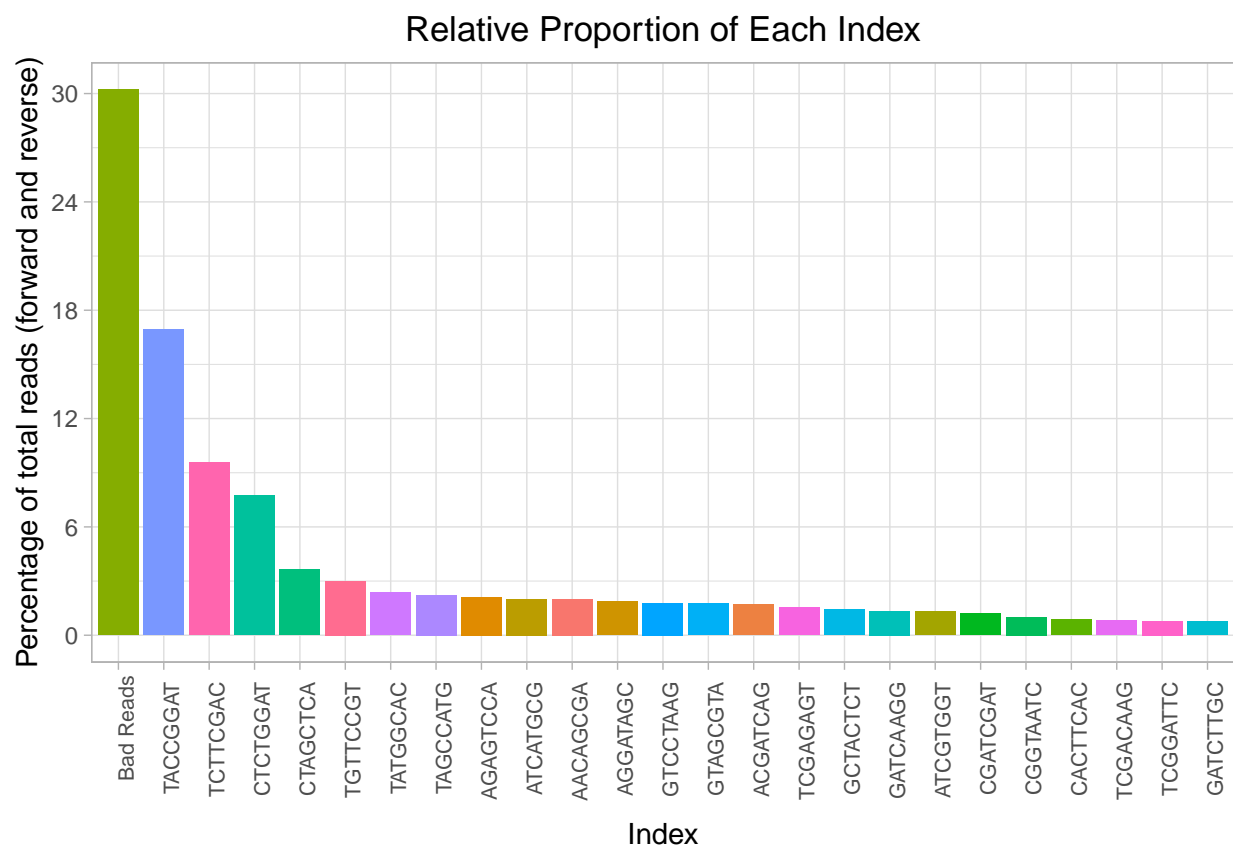
Files created after script completed

Head -8 output for AACAGCGA and bad reads

[illegible]

+Undetermined / Low Quality / Index not recognized /
#A-AFFJJFJJJJJJJJJJJJJ<JAFJJJJF<JFJJAJJJJJJJJJJJJJJJFJJJAJJFJJFJJF<JJA-JJJ-<AF AF--FF<JAFJF

Output plots



Instances* of Reads Written to Bad Output Files

*Instances are not mutually exclusive

