### Comparison to BLAST

```
Basic Local Alignment Search Tool
                                                                               BLAST+ 2.15.0 is here!
BLAST finds regions of similarity between biological sequences. The
                                                                                We have included two exciting new features in the latest
program compares nucleotide or protein sequences to sequence
                                                                                BLAST+ release
databases and calculates the statistical significance.
                                                          Learn more
                                                                                Tue, 28 Nov 2023
                                                                                                                             More BLAST news...
```

## Install Blast tools

Out[3]:

• https://blast.ncbi.nlm.nih.gov/doc/blast-help/downloadblastdata.html https://www.ncbi.nlm.nih.gov/books/NBK279690/

id

## In [1]: import pandas as pd

Compare our prediction with BLAST Tools

```
from Bio.Seq import Seq
        from Bio.SeqRecord import SeqRecord
        from Bio import SeqIO
        import json
        from tqdm import tqdm
        import sklearn.metrics as metrics
        from BLAST.utils import metric
In [2]: | df = pd.read_csv('./BLAST/src/testset.csv')
```

```
In [3]: df.head()
```

sequence length class

```
0
                 Thhalv10001938m -- EIL MMMFNEMGMYGKMDFFSSTSLGEIDVCPLPQAEQDHPVVEEDYTDD...
                                                                                                  582
                                                                                                         EIL
         1
                       Tp1g37880 -- EIL
                                         MGVFFSDPDSIAEDDGYTDNELDVDELEKRIWKQEMRLRRLKEHRK...
                                                                                                   211
                                                                                                         EIL
         2
                      Pbr004574.1 -- EIL
                                           MGDVGEIGPDISSDIEEDLRCDNIAEKDVSDEEIEAEDLERRMWKD...
                                                                                                  608
                                                                                                         EIL
              Sme2.5_02278.1_g00005.1 --
        3
                                       MNNEVVEENQEFDDEEINYDDLKRRMWKDRMRMQILKGKKRDMMIE...
                                                                                                         EIL
                                                                                                  494
        4
                                                                                                         EIL
                  ONIVA11G15950.1 -- EIL
                                         MDASKKSVMTKEEQQLSPAASPAAAVMTAEADAINEEQDKAAAATT...
                                                                                                  466
In [4]: records = []
```

```
for i, row in df.iterrows():
     record = SeqRecord(Seq(row['sequence']), id=str(i+1), description=" | " + row['class'])
     records.append(record)
 # Save to a FASTA file
 output_file = "./BLAST/testset.fasta"
 SeqIO.write(records, output_file, "fasta")
 print(f"FASTA file saved as {output_file}")
FASTA file saved as ./BLAST/testset.fasta
```

Build reference(train) database for blast

```
In [5]: ref_pah = "../data/mix_data/trainset"
```

```
In [6]: gene_info_path = "../data/gene_info.json"
        with open(gene_info_path, 'r') as json_file:
            gene_info = json.load(json_file)
In [7]: records = []
        for gene, info in tqdm(gene_info.items()):
            path = f"{ref_pah}/{info['file_code']}.csv"
            df_gene = pd.read_csv(path)
            df_gene = df_gene[df_gene['class'] == 1]
            # write
            for _, row in df_gene.iterrows():
                record = SeqRecord(Seq(row['sequence']), id=row['id'] + "--" + gene, description="")
                records.append(record)
        # Save to a FASTA file
        output_file = "./BLAST/trainset.fasta"
        SeqIO.write(records, output_file, "fasta")
        print(f"FASTA file saved as {output_file}")
```

**BUILD DATABASE FOR BLAST** 

FASTA file saved as ./BLAST/trainset.fasta

100%

```
Building a new DB, current time: 07/07/2024 09:42:23
New DB name: /Users/genereux/Documents/UM6P/COURS-S2/S2-PROJECT/pygenomics/test/BLAST/database/py
```

In [8]: !makeblastdb -in ./BLAST/trainset.fasta -dbtype prot -out ./BLAST/database/pygenomics\_ref\_db

```
genomics_ref_db
New DB title: ./BLAST/trainset.fasta
Sequence type: Protein
Deleted existing Protein BLAST database named /Users/genereux/Documents/UM6P/COURS-S2/S2-PROJECT/py
genomics/test/BLAST/database/pygenomics_ref_db
Keep MBits: T
Maximum file size: 300000000B
Adding sequences from FASTA; added 256279 sequences in 4.78068 seconds.
 RUN BLAST SEARCH TOOL
```

58/58 [00:09<00:00, 6.35it/s]

# View Blast Analysis results

```
In [18]: columns = [
             "query_id", "subject_id", "perc_identity", "alignment_length",
```

"mismatches", "gap\_opens", "q\_start", "q\_end", "s\_start", "s\_end",

In [16]: #!blastp -query ./BLAST/testset.fasta -db ./BLAST/database/pygenomics\_ref\_db -out ./BLAST/blast\_res

```
"evalue", "bit_score"
blast_results = pd.read_csv("./BLAST/blast_result.txt", sep="\t", names=columns)
```

3

0

1

2

Accuracy

R**hetzikt**ison Score

**01968**6

019606

In [26]: **import** importlib

**Confusion Matrix** 

EIL

EIL

```
#blast_results.to_csv("./BLAST/blast_results.csv", index=False)
          # Display the first few rows of the DataFrame
          blast_results.head()
Out[18]:
             query_id
                                   subject_id perc_identity alignment_length mismatches gap_opens q_start q_end
          0
                    1
                              676714038--EIL
                                                    89.696
                                                                                     49
                                                                       592
                                                                                                 6
                                                                                                          1
                                                                                                               581
                          XP_013629559.1--EIL
                    1
                                                    84.628
                                                                                     59
                                                                                                          1
          1
                                                                       592
                                                                                                 10
                                                                                                               581
                      GSBRNA2T00066199001-
          2
                                                    83.277
                                                                       592
                                                                                     59
                                                                                                 10
                                                                                                          1
                                                                                                               581
```

85.135

XP\_009133916.1--EIL 4 1 85.135 592 56 10 1 581 Proccessing to make final inference base on threshold In [19]: voting\_threshold = 50 unique\_query\_ids = sorted(blast\_results['query\_id'].unique()) predictions = ["Unknown" for \_ in range(len(df))]

best\_match = filtered\_df[filtered\_df['perc\_identity'] > voting\_threshold].sort\_values(by='perc\_

592

56

10

581

0.9686

**Metrics Descriptio** 

```
if not best_match.empty:
    subject_id = best_match.iloc[0]['subject_id']
    subject_id_part = subject_id.split('--')[-1]
```

1

1

1

for query\_id in tqdm(unique\_query\_ids, desc="Max Voting : "):

filtered\_df = blast\_results[blast\_results['query\_id'] == query\_id]

# Find the row with the highest perc\_identity greater than voting\_threshold

XP\_009133914.1--EIL

```
predictions[query_id-1] = (subject_id_part)
         blast_prediction = pd.DataFrame({'prediction': predictions})
                                             64037/64037 [15:59<00:00, 66.74it/s]
        Max Voting : 100%|
In [20]: with open("./BLAST/src/class_mapping.json", 'r') as json_file:
             class_mapping = json.load(json_file)
         true_label = pd.read_csv("./BLAST/src/true_labels.csv")
In [21]: blast_prediction["predicted_label"] = blast_prediction.prediction.map(class_mapping)
         blast_prediction.head()
            prediction predicted_label
Out[21]:
```

	3	EIL	1
	4	EIL	1
	Display performance		
In [22]:	<pre>metric.show_metrics(true_label, blast_prediction.predicted_label)</pre>		
Overall Score			

```
importlib.reload(metric)
Out[26]: <module 'BLAST.utils.metric' from '/Users/genereux/Documents/UM6P/COURS-S2/S2-PROJECT/pygenomics/t
          est/BLAST/utils/metric.py'>
In [27]: |metric.show_confusion(true_label, blast_prediction.predicted_label, class_mapping)
```

Webbased tes metrics for each class independently and then takes the average, weighted by the number of instances of each clas

Calculates metrics globally by counting the total true positives, false negatives, and false positive

Calculates metrics for each class independently and then takes the average (treating all classes equally

```
LBD -
СЗН
STAT
```

