

1 Supplementary materials

1.1 BoaG queries

Following are the BoaG queries that used in this work to retrieve annotations from protein sequences and clusters.

```
1 s: Sequence = input;
2 count: output collection[string] of string;

4 getTaxList := function(seq: Sequence) : string {
5   taxids := "";
6   foreach (i: int; s.annotation[i].tax_name != "")
7     taxids = taxids + s.annotation[i].tax_id + " ";
8   return taxids;
9 };

11 count[s.seqid] << getTaxList(s);
```

Figure 1: Generate list of all taxids for the entire NR database. The output of BoaG script used to generate the tree of life in Figure 3 in the main manuscript. The BoaG query and its result is publicly available on the cluster here: <http://boa.cs.iastate.edu/boag/?q=boa/job/public/80>

```
1 s: Sequence = input;
2 counts: output sum[int] of int;
3 foreach(i:int; def(s.cluster[i])){
4   if (s.cluster[i].similarity==95)
5     counts [s.cluster[i].length] << 1;}
```

Figure 2: Frequencies of proteins for different size at 95% similarity. Sequence is a top level type in BoaG language. Variable counts is an output aggregator that sums number of proteins in different size, i.e. number of amino acids. The BoaG query takes only two minutes and can be accessed here: <http://boa.cs.iastate.edu/boag/?q=boa/job/public/96>

```

1 s: Sequence = input;
2 protOut : output sum [string][string] of int;
3 distinctTax := function (seq: Sequence): int{
4 taxSet : set of string;
5 foreach(i:int; def(seq.annotation[i]))
6 add(taxSet,seq.annotation[i].tax_id);
7 return(len(taxSet));};
8 if (distinctTax(s) > 10)
9   foreach(i:int; def(s.annotation[i]))
10    protOut [s.annotation[i].define][s.seqid]<<1;

```

Figure 3: BoaG script for top protein functions for proteins that have more than 10 distinct taxonomic assignments. The output of this query is used to generate the top protein functions. The BoaG query and its result is publicly available on the cluster here: <http://boa.cs.iastate.edu/boag/?q=boa/job/public/40>

1.2 Programming efficiency in BoaG

BoaG programming performance in terms of lines of code and running time is shown in Figure 4.

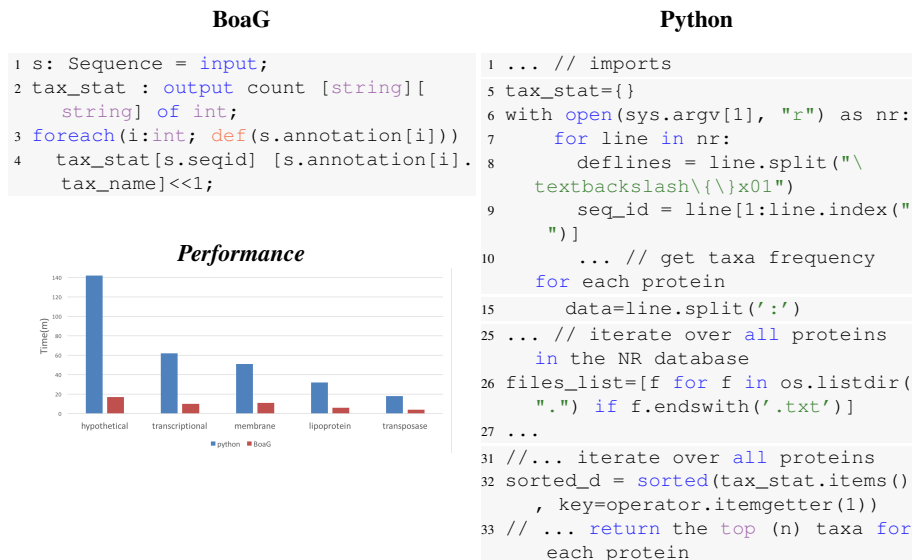


Figure 4: Comparison of Line of Code(LOC) and performance to answer query "frequency of taxonomic assignment for each protein sequence" run on the different subset of the NR dataset. On the right side, the equivalent *Boa* code needs 33 lines of code in Python whereas the BoaG script needs only four.