# Alignment

The main goal is to compare some multimple alignment algorithms.

## 1. Materials and methods

I work with nucleotide sequences of gene of translation termination factor (SUP35) from different species of yeasts
In this project I use the following programms:
- CLUSTAL 2.1
- MUSCLE v3.8.1551
- MAFFT v7.475
- Kalign (3.3)
- T-COFFEE Version_13.41.0.28bdc39
- PRANK v.170427
- transeq (Version: EMBOSS:6.6.0.0)
- getorf (Version: EMBOSS:6.6.0.0)

## 2. Using 6 various alignment algorithms on 10 DNA sequences

### clustalw

By default we generate the alignment in clustalw format

```
clustalw -INFILE=./data/raw_data/SUP35_10seqs.fa
```

time:

```
real    0m3,747s
user    0m3,732s
sys     0m0,015s
```

Determining the output file:

```
clustalw -INFILE=./data/raw_data/SUP35_10seqs.fa -OUTFILE=./data/processed_data/SUP35_10seqs.clustalw
```

Generating the output in fasta format:

```
clustalw -INFILE=./data/raw_data/SUP35_10seqs.fa -OUTPUT=FASTA -OUTFILE=./data/processed_data/SUP35_10seqs.clustalw.fa
```

### muscle

```
muscle   -in ./data/raw_data/SUP35_10seqs.fa -out ./data/processed_data/SUP35_10seqs_muscle.fa
```

time:

```
real    0m4,742s
user    0m4,712s
sys     0m0,022s
```

### mafft

```
mafft --auto ./data/raw_data/SUP35_10seqs.fa >./data/raw_data/SUP35_10seqs_mafft.fa
```

time:

```
real    0m4,486s
user    0m4,351s
sys     0m0,156s
```

### kalign

```
kalign <./data/raw_data/SUP35_10seqs.fa >./data/processed_data/SUP35_10seqs_kalign.fa
```

time:

```
  real     0m0,242s
  user     0m0,357s
  sys      0m0,008s
```

## t_coffee

```
  t_coffee -infile=./data/raw_data/SUP35_10seqs.fa -outfile=./data/processed_data/SUP35_10seqs_tcoffee.fa
```

time:

```
  real     1m31,287s
  user     1m29,863s
  sys      0m1,477s
```

## prank

```
  prank -d=./data/raw_data/SUP35_10seqs.fa -o=./data/processed_data/SUP35_10seqs_prank.fa
```

time:

```
  real     0m12,179s
  user     0m12,006s
  sys      0m0,163s
```

# 3. Comparation of the results on 10 DNA sequences

The table containing real times and grapical representation of the results can be found in Supplementary (Table 1, Figures 1-5)
I find mafft and t-coffee the best ones

# 4. Reverse complement problem

While opening the file ./data/raw_data/SUP35_10seqs_strange_aln.fa we can see that alignment for 1 sequence is not satisfactory (Supplementary, fig. 6). But if we BLAST it, we see that it is the same gene of the same organism. The answer is to find the reverse complement sequence and realign it

# 5. Using 6 various alignment algorithms on 250 DNA sequences

## clustalw

```
  clustalw -INFILE=./data/raw_data/SUP35_250seqs.fa OUTPUT=FASTA -OUTFILE=./data/processed_data/SUP35_250seqs.clustalw.fa
```

time:

```
  real     29m3,728s
  user     29m2,200s
  sys      0m1,483s
```

## muscle

```
  muscle -in ./data/raw_data/SUP35_250seqs.fa -out ./data/processed_data/SUP35_250seqs_muscle.fa
```

time

```
  real     1m20,447s
  user     1m20,169s
  sys      0m0,276s
```

## mafft

```
mafft --auto ./data/raw_data/SUP35_250seqs.fa >./data/processed_data/SUP35_250seqs_mafft.fa
```

time:

```
  real    0m39,807s
  user    0m39,104s
  sys     0m0,744s
```

## kalign

```
kalign <./data/raw_data/SUP35_250seqs.fa >./data/processed_data/SUP35_250seqs_kalign.fa
```

time:

```
  real    0m3,912s
  user    0m7,910s
  sys     0m0,105s
```

## t_coffee

```
t_coffee -infile=./data/raw_data/SUP35_250seqs.fa -outfile=./data/processed_data/SUP35_250seqs_tcoffee.fa
```

I decided to interrupt the process after:

```
  real    39m1,756s
  user    0m0,761s
  sys     0m0,275s
```

## prank

```
prank -d=./data/raw_data/SUP35_250seqs.fa -o=./data/processed_data/SUP35_250seqs_prank.fa
```

time:

```
  real    2m39,690s
  user    2m35,851s
  sys     0m1,994s
```

# 6. Comparation of the results on 250 DNA sequences

The table containing real times and grapical representation of the results can be found in Supplementary (Table 1, Figures 7-10)
I find CLUSTAL to be the best when you can wait or muscle when you want to make it faster.

# 7. Translation

Let's translate the 10 DNA sequences:

```
transeq -sequence ./data/raw_data/SUP35_10seqs.fa -outseq ./data/processed_data/SUP35_10seqs.t.faa
```

Another way is to use getorf. getorf gives you the ORF and its coordinates in nucleotides. We shold give the minsize near the protein size (in nucleotides) to get its sequence witout other small peptides.

```
getorf -sequence ./data/raw_data/SUP35_10seqs.fa -outseq ./data/processed_data/SUP35_10seqs.g.faa -noreverse -minsize 500
```

# 8. Using 6 various alignment algorithms on 10 protein sequences

We use protein sequences of the same gene from the same organisms

## clustalw
```

```
clustalw -INFILE=./data/processed_data/SUP35_10seqs.g.faa -OUTFILE=./data/raw_data/SUP35_10seqs.clustalw.faa -OUTPUT=FASTA -TYPE=pro
```

time:

```
  real    0m0,715s
  user    0m0,715s
  sys     0m0,000s
```

```
clustalo --infile=./data/processed_data/SUP35_10seqs.g.faa --outfile=./data/processed_data/SUP35_10seqs.clustalo.faa --verbose
```

time:

```
  real    0m0,625s
  user    0m1,072s
  sys     0m0,109s
```

## muscle

```
muscle -in ./data/processed_data/SUP35_10seqs.g.faa -out ./data/processed_data/SUP35_10seqs_muscle.faa
```

time:

```
  real    0m0,358s
  user    0m0,342s
  sys     0m0,016s
```

## mafft

```
mafft --auto ./data/processed_data/SUP35_10seqs.g.faa >./data/processed_data/SUP35_10seqs_mafft.faa
```

time:

```
  real    0m0,618s
  user    0m0,576s
  sys     0m0,094s
```

## kalign

```
kalign <./data/processed_data/SUP35_10seqs.g.faa >./data/processed_data/SUP35_10seqs_kalign.faa
```

time:

```
  real    0m0,076s
  user    0m0,176s
  sys     0m0,016s
```

## t_coffee

```
t_coffee -infile=./data/processed_data/SUP35_10seqs.g.faa -outfile=./data/processed_data/SUP35_10seqs_tcoffee.faa
```

time:

```
  real    0m15,731s
  user    0m15,329s
  sys     0m0,402s
```

## prank
```

```
prank -d=./data/processed_data/SUP35_10seqs.g.faa -o=./data/processed_data/SUP35_10seqs_prank.faa
```

time:

```
real    0m12,729s
user    0m12,461s
sys     0m0,317s
```

## 9. Comparation of the results on 10 protein sequences

The table containing real times and grapical representation of the results can be found in Supplementary (Table 1, Figures 11-15)
I think muscle is the best for this goal.

## 10. Add more alignments using muscle/mafft:

Here I align 2 more DNA sequences and add them to the files with 250 aligned DNA sequences using muscle of mafft.

```
muscle -in ./data/raw_data/SUP35_2addseqs.fa -out ./data/processed_data/SUP35_2addseqs_muscle.fa
muscle -profile -in1 ./data/processed_data/SUP35_250seqs_muscle.fa -in2 ./data/processed_data/SUP35_2addseqs_muscle.fa -out ./data/p
```

```
mafft --auto ./data/raw_data/SUP35_2addseqs.fa > ./data/processed_data/SUP35_2addseqs_mafft.fa
mafft --add ./data/processed_data/SUP35_2addseqs_mafft.fa ./data/processed_data/SUP35_250seqs_mafft.fa > ./data/processed_data/SUP35
```

## Supplementary

| | CLUSTAL | MUSCLE | MAFFT | KALIGN | T-COFFEE | PRANK |
|---|---|---|---|---|---|---|
| **10 DNA SEQUENCES** | 3,747s | 4,742s | 4,486s | 0,242s | 1m 31,287s | 12,179s |
| **250 DNA SEQUENCES** | 29m 3,728s | 1m 20,447s | 39,807s | 3,912s | > 39m | 2m 39,690s |
| **10 PROTEIN SEQUENCES** | 0,715s | 0,358s | 0,618s | 0,076s | 15,731s | 12,729s |

Table 1. Real times for the alignments



*Figure 1. Alignment of 10 DNA sequences using CLUSTAL*



*Figure 2. Alignment of 10 DNA sequences using MUSCLE*

*Figure 3. Alignment of 10 DNA sequences using MAFFT*



*Figure 4. Alignment of 10 DNA sequences using T-COFFEE*

*Figure 5. Alignment of 10 DNA sequences using PRANK*



*Figure 6. Strange alignment*

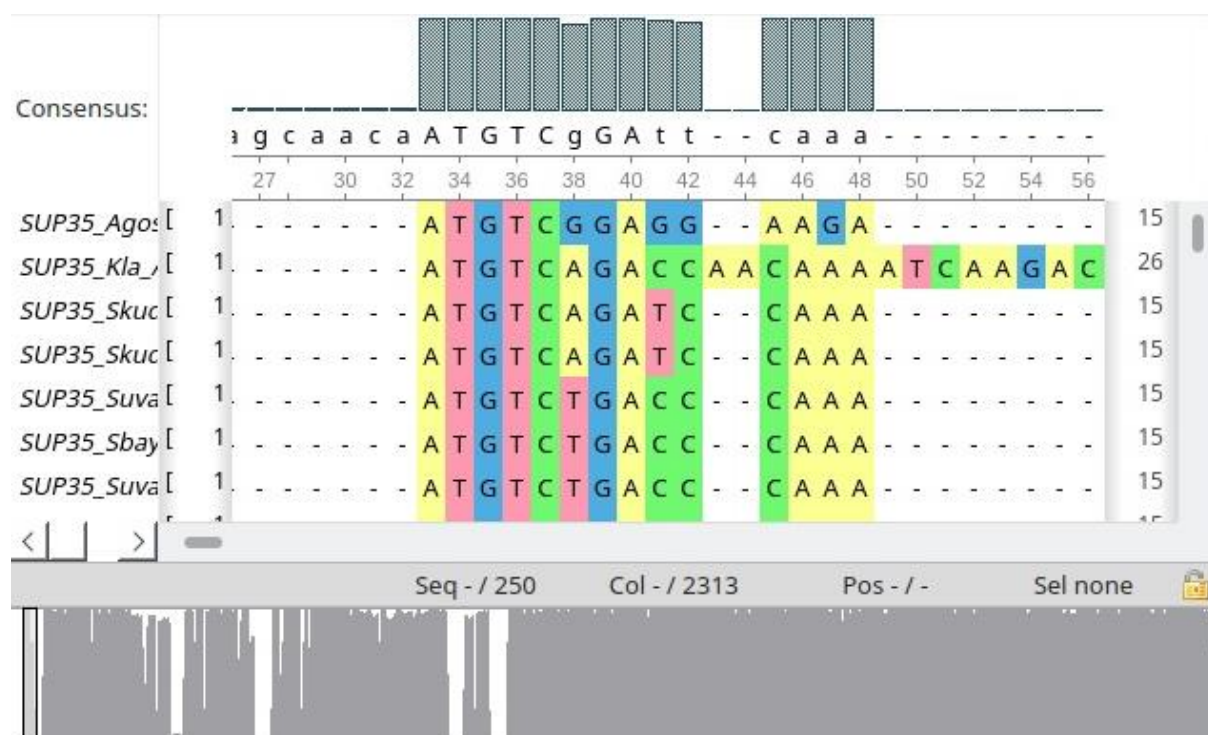*Figure 7. Alignment of 250 DNA sequences using CLUSTAL*



*Figure 8. Alignment of 250 DNA sequences using MUSCLE*

*Figure 9. Alignment of 250 DNA sequences using MAFFT*



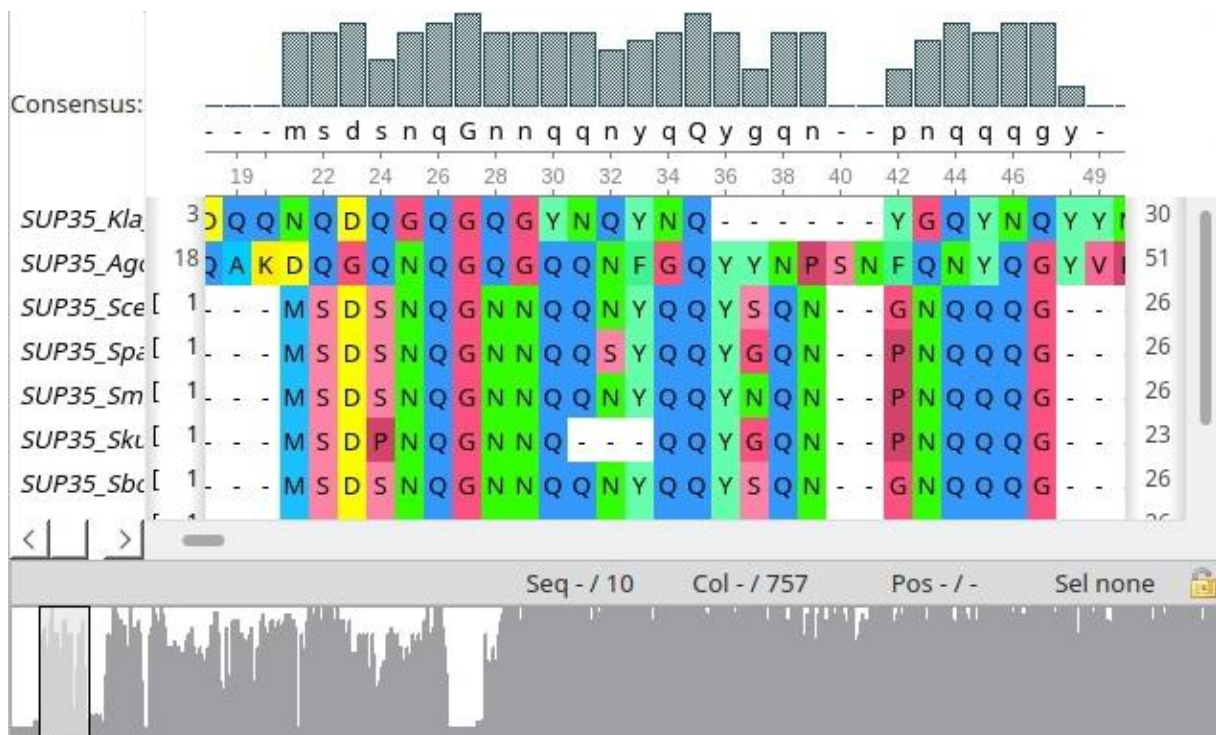*Figure 10. Alignment of 250 DNA sequences using PRANK*

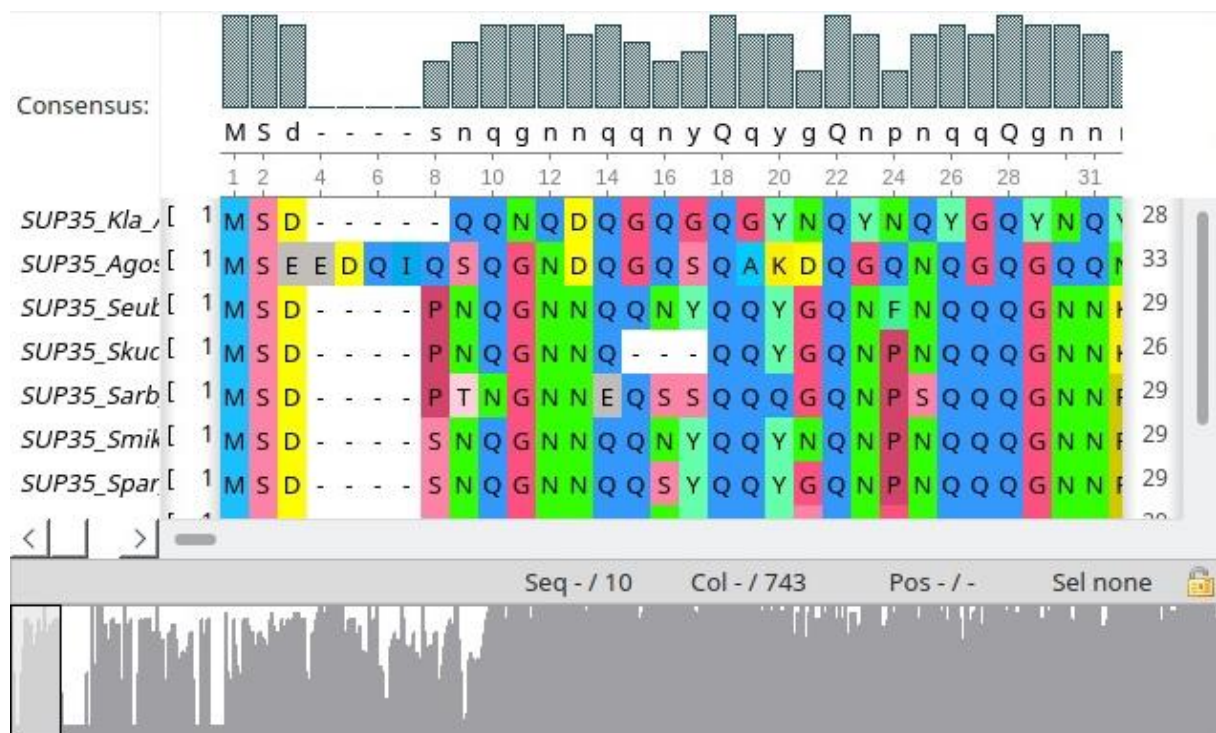*Figure 11. Alignment of 10 protein sequences using CLUSTAL*



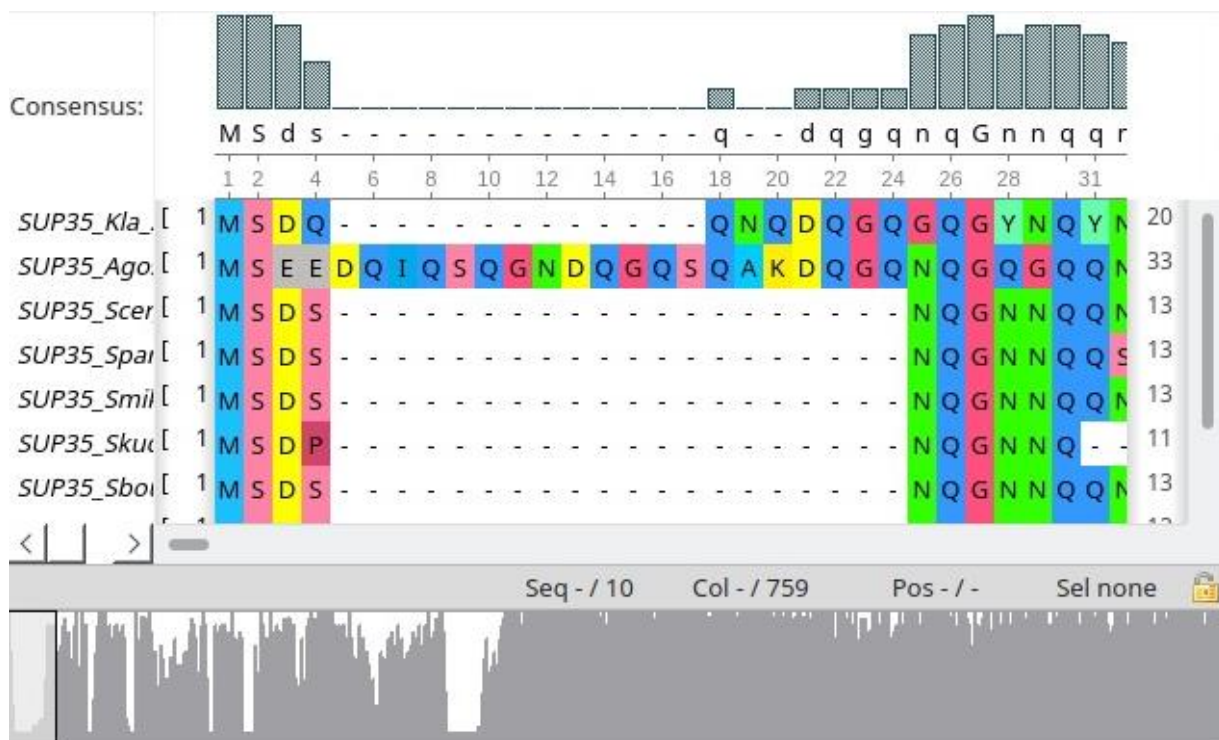*Figure 12. Alignment of 10 protein sequences using MUSCLE*

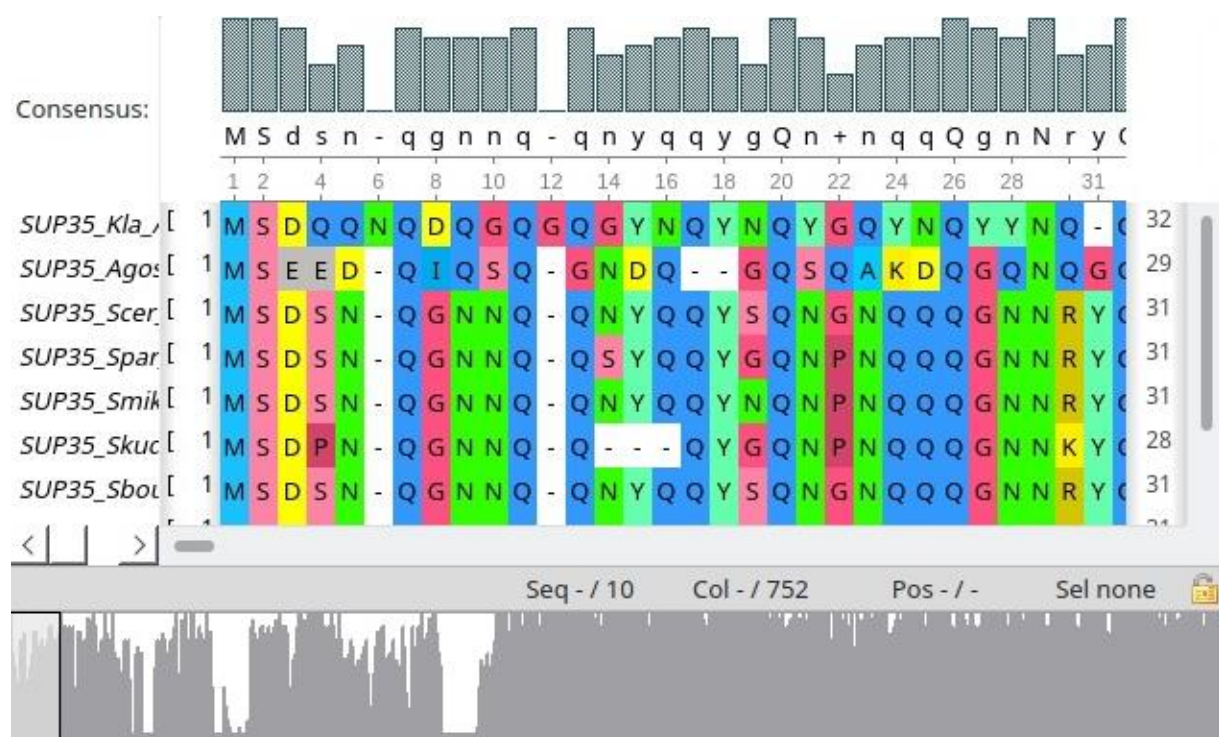*Figure 13. Alignment of 10 protein sequences using MAFFT*



*Figure 14. Alignment of 10 protein sequences using T-COFFEE*

*Figure 15. Alignment of 10 protein sequences using PRANK*