

# Multiple Sequence Alignment

Andrew Michael





# Gene I chose to study

envZ gene from *Escherichia coli*

A histidine kinase that responds to osmolarity changes in the medium of the cell by regulating phosphorylation.

Gene sequence:

```
ATGAGGCGATTGCGCTTCTCGCCACGAAGTTCATTTGCCCGTACGTTATTGCTCATCGTCACCTTGCTGT
TCGCCAGCCTGGTGACGACTTATCTGGTGGTGCTGAACTTCGCGATTTTGCCGAGCCTCCAGCAGTTTAA
TAAAGTCCTCGCGTACGAAGTGCGTATGTTGATGACCGACAACTGCAACTGGAGGACGGCACGCAGTTG
GTTGTGCCTCCCGCTTTCCGTGCGGAGATCTACCGTGAGCTGGGGATCTCTCTACTCCAACGAGGCTG
CCGAAGAGGCAGGTCTGCGTTGGGCGCAACACTATGAATTCCTTAAGCCATCAGATGGCGCAGCAACTGGG
CGGCCCGACGGAAGTGCGCGTTGAGGTCAACAAAAGTTCGCTGTCTGCTGGCTGAAAACCTGGCTGTCTG
CCCAATATCTGGGTACGCGTGCCGCTGACCGAAATTCATCAGGGCGATTTCTCTCCGCTGTTCCGCTATA
CGCTGGCGATTATGCTATTGGCGATAGGCGGGGCGTGGCTGTTTATTGATCCAGAACCGACCGTTGGT
CGATCTCGAACACGCAGCCTTGACAGTTGGTAAAGGGATTATTCGCCGCCGCTGCGTGAGTATGGCGCT
TCGGAGGTGCGTTCCGTTACCCGTGCCTTTAACCATATGGCGGCTGGTGTTAAGCAACTGGCGGATGACC
GCACGCTGCTGATGGCGGGGGTAAGTCACGACTTGCGCACGCCGCTGACGCGTATTCGCTGGCGACTGA
GATGATGAGCGAGCAGGATGGCTATCTGGCAGAATCGATCAATAAAGATATCGAAGAGTGCAACGCCATC
ATTGAGCAGTTTATCGACTACCTGCGCACCGGGCAGGAGATGCCGATGGAATGGCGGATCTTAATGCAG
TACTCGGTGAGGTGATTGCTGCCGAAAGTGGCTATGAGCGGGAATTAAGAACCGCGCTTTACCCCGGCAG
CATTGAAGTGAAAATGCACCCGCTGTTCGATCAAACGCGCGGTGGCGAATATGGTGGTCAACGCCGCCGTG
TATGGCAATGGCTGGATCAAAGTCAGCAGCGGAACGGAGCCGAATCGCGCCTGGTTCCAGGTGGAAGATG
ACGGTCCGGGAATTGCGCCGGAACAACGTAAGCACCTGTTCAGCCGTTTGTCCGCGGCGACAGTGCGCG
CACCATTAGCGGCACGGGATTAGGGCTGGCAATTGTGCAGCGTATCGTGGATAACCATAACGGGATGCTG
GAGCTTGGCACCGAGCGAGCGGGGCGGGCTTTCCATTGCGCCTGGCTGCCAGTGCCGGTAACGCGGGCGC
AGGGCACGACAAAAGAAGGGTAA
```



# Using BLAST

BLASTx search translated nucleotide -> protein

*Citrobacter farmeri* (96.22% identity)

*Salmonella bongori* (95.56% identity)

*Yokenella regensburgei* (93.71% identity)



# Using BLAST

tBLASTn search translated protein -> translated nucleotide

*Citrobacter farmeri*

*Salmonella bongori*

*Yokenella regensburgei*



# T-Coffee

Genes:

## MSA

The multiple sequence alignment result as produced by T-coffee.

T-COFFEE, Version\_11.00 (Version\_11.00)

Cedric Notredame

SCORE=733

\*

\* **BAD** **AVG** **GOOD**

\*

Escherichia : 85

Citrobacter : 67

Salmonella : 85

Yokenella : 59

cons : 73

```
Escherichia -----ATGAGGC-----GATTGCGCTTCTCGCCACGAAG-----
Citrobacter -----ATGAGGC-----GAATGCGCTTCTCGCCACGAAG-----
Salmonella -----ATGAGGC-----GAATGCGCTTCTCGCCGCGAAG-----
Yokenella  ACTCTCTTTCTGTCCGCGATGCGGATACCGGAATCGGATAGCCAGGCGGAATCAACAAGCCACCCCG
```

```
cons          ** **          ** * * * * * ****
```

```
Escherichia TT-----CATTTGCCCGTACGTTAT--TGCTCATCGTCACCTTGCTGTTGCCAGCCTGGTGACG
Citrobacter TT-----CATTTGCCCGCACGCTGC--TGCTCATCGTCACTCTGCTGTTTGTGAGCCTGGTGACG
Salmonella  TT-----CATTTGCTCGCACGCTGT--TGCTCATCGTCACCTTGCTGTTGTTAGCCTGGTGACG
Yokenella   TTCTGCTGGTGCCAATTTCCAGCAGACCGTTATGTTATCGATAATACGCTGCACAATGCCAGACCGAG
```

```
cons          **          ** * * * * * ** * **** * **** * *
```

```
Escherichia ACTTAT-CTGGTGGTGCTGAACCTTCGCGATTTTGCCGAGCCTCCAGCAGTTTAATAAAGTCCTGCGTA
Citrobacter ACTTAT-CTGGTGGTGCTGAACCTTCGCGATTCGCCAGTCTCCAGCAGTTTAATAAGGTCTGCGCTA
Salmonella  ACCTAC-CTGGTGGTGCTGAACCTTCGCGATCTTACCGAGCCTCCAGCAGTTTAATAAGGTTCTGGCTTA
Yokenella   CCCCGTACCGCTGGTGCTGCGC---GCGCTGTGCGCCGCGAA-CAAACGGCTGGAACAGGT---GCTTA
```

```
cons          *  * * ***** * *** * **** * * * * * * * * * *
```

```
Escherichia CGAAGTGCGTATGTTGAT----GACCGACAAACTGCAACTGGAGGACGGCACGCAAGTTGGTTGTGCCT
Citrobacter CGAAGTGCGTATGCTGAT----GACCGATAAACTGCAACTGGAGGACGGCACGCAACTGGTGGTGCCT
Salmonella  TGAAGTCCGTATGCTGAT----GACCGATAAGCTGCAACTGGAGGACGGTACGCAATTAGTTGTGCCT
Yokenella   CGCTGCTCGGGT-TTGATCCCCGGGCCGTCACTTCCACCTGGAA--CCAGGCGCGATTCTGGTTCACT
```



# T-Coffee

Proteins:

*E. coli*

*C. farmeri*

*S. bongori*

*Y. regensburgae*

## MSA

The multiple sequence alignment result as produced by T-coffee.

T-COFFEE, Version\_11.00 (Version\_11.00)

Cedric Notredame

SCORE=996

\*

\* **BAD AVG GOOD**

\*

Protein : 99

Protein\_1 : 99

Protein\_2 : 99

Protein\_3 : 99

cons : 99

```
Protein      MRRLRFSRSSFARTLLLIVTLLFASLVTTYLVVLNFAILPSLQQFNKVLAYEVRMLMTDKLQLEDGTQ
Protein_1    MRRMRFSRSSFARTLLLIVTLLFASLVTTYLVVLNFAILPSLQQFNKVLAYEVRMLMTDKLQLEDGTQ
Protein_2    ---MRFSRSSFARTLLLIVTLLFVSLVTTYLVVLNFAILPSLQQFNKVLAYEVRMLMTDKLQLEDGTQ
Protein_3    MRRMRFSRSSFARTLLLIVTLLFVSLVTTYLVVLNFAILPSLQQFNKVLAYEVRMLMTDKLQLEDGTQ
```

```
cons          :*****.*****
```

```
Protein      LVVPPAFRREIYRELGISLYSNEAAEEAGLRWAQHYEFLSHQMAQQLGGPTEVRVEVNKSSPVVWLKTW
Protein_1    LVVPPAFRREIYRELGISLYSNEAAEEAGLRWAQHYEFLSHQMAQQLGGPTEVRVEVNKSSPVVWLKTW
Protein_2    LVVPPAFRREIYRELGISLYTNEAAEEAGLRWAQHYEFLSHQMAQQLGGPTEVRVEVNKSSPVVWLKTW
Protein_3    LVVPPAFRREIYRELGISLYSDEAAEDAGLRWAQHYEFLSQMAQQLGGPTEVRVEVNKSSPVVWLKTW
```

```
cons          *****:****:*****:*****
```

```
Protein      LSPNIWVRVPLTEIHQGDFSPLFRYTLAIMLLAIGGAWLFIRIQNRPLVDLEHAALQVGKGIIPPLRE
Protein_1    LSPNIWVRVPLTEIHQGDFSPLFRYTLAIMLLAIGGAWLFIRIQNRPLVDLEHAALQVGKGIIPPLRE
Protein_2    LSPNIWVRVPLTEIHQGDFSPLFRYTLAIMLLAIGGAWLFIRIQNRPLVDLEHAALQVGKGIIPPLRE
Protein_3    LSPNIWVRVPLTEIHQGDFSPLFRYTLAIMLLAIGGAWLFIRIQNRPLVDLEHAALQVGKGIIPPLRE
```

```
cons          *****
```

```
Protein      YGASEVRSVTRAFNHMAAGVKQLADDRTLLMAGVSHDLRTPLTRIRLATEMMSEODGYLAESINKDIEE
Protein_1    YGASEVRSVTRAFNHMAAGVKQLADDRTLLMAGVSHDLRTPLTRIRLATEMMGEEDGYLAESINKDIEE
Protein_2    YGASEVRSVTRAFNHMAAGVKQLADDRTLLMAGVSHDLRTPLTRIRLATEMMGEEDGYLAESINKDIEE
Protein_3    YGASEVRSVTRAFNHMAAGVKQLADDRTLLMAGVSHDLRTPLTRIRLATEMMSVEDGYLAESINKDIEE
```





# Boxshade

Boxshade Preferences

General RTF/PS/PNG Text (ASCII) output Similarities Groups

Different from consensus

Foreground Background

☒ Uppercase **ACGT**

☐ Lowercase

Identical to consensus

Foreground Background

☒ Uppercase **ACGT**

☐ Lowercase

Other prefs

Font size: 12

☐ Portrait

☒ Landscape

Similar to consensus

Foreground Background

☒ Uppercase **ACGT**

☐ Lowercase

All the same residue

Foreground Background

☒ Uppercase **ACGT**

☐ Lowercase

Protein MKRMRFSRPRSSSFARTLLLVLTLLFASLVTTYLVVLNFAILPSLQQFNKVLAYEVRMLMTD  
Protein\_1 MKRMRFSRPRSSSFARTLLLVLTLLFASLVTTYLVVLNFAILPSLQQFNKVLAYEVRMLMTD  
Protein\_2 ---MRFSRPRSSSFARTLLLVLTLLFVSLVTTYLVVLNFAILPSLQQFNKVLAYEVRMLMTD  
Protein\_3 MKRMRFSRPRSSSFARTLLLVLTLLFVSLVTTYLVVLNFAILPSLQQFNKVLAYEVRMLMTD

Protein KLQLEDGTQLVVPAPFRREIYRELGISLYSNEAAEEAGLRWAQHYEFLSHQMAQQQLGGPT  
Protein\_1 KLQLEDGTQLVVPAPFRREIYRELGISLYSNEAAEEAGLRWAQHYEFLSHQMAQQQLGGPT  
Protein\_2 KLQLEDGTQLVVPAPFRREIYRELGISLYSNEAAEEAGLRWAQHYEFLSHQMAQQQLGGPT  
Protein\_3 KLQLEDGTQLVVPAPFRREIYRELGISLYSNEAAEEAGLRWAQHYEFLSHQMAQQQLGGPT

Protein EVRVEVNKSSPVVWLKTNLSPNIWVRVPLTEIHQGDSPFLFRYT LAIMLLAIGGAWLFIR  
Protein\_1 EVRVEVNKSSPVVWLKTNLSPNIWVRVPLTEIHQGDSPFLFRYT LAIMLLAIGGAWLFIR  
Protein\_2 EVRVEVNKSSPVVWLKTNLSPNIWVRVPLTEIHQGDSPFLFRYT LAIMLLAIGGAWLFIR  
Protein\_3 EVRVEVNKSSPVVWLKTNLSPNIWVRVPLTEIHQGDSPFLFRYT LAIMLLAIGGAWLFIR

Protein IQNRPLVDLEHAALQVGKGIIPPPLEYGASEVRSVTRAFNHMAAGVKQLADDRTLLMAG  
Protein\_1 IQNRPLVDLEHAALQVGKGIIPPPLEYGASEVRSVTRAFNHMAAGVKQLADDRTLLMAG  
Protein\_2 IQNRPLVDLEHAALQVGKGIIPPPLEYGASEVRSVTRAFNHMAAGVKQLADDRTLLMAG  
Protein\_3 IQNRPLVDLEHAALQVGKGIIPPPLEYGASEVRSVTRAFNHMAAGVKQLADDRTLLMAG

Protein VSHDLRTPLTRIRLATEMMSVEEDGYLAESINKDIEECNAIEEQFIDYLR TGQEMPMEAD  
Protein\_1 VSHDLRTPLTRIRLATEMMSVEEDGYLAESINKDIEECNAIEEQFIDYLR TGQEMPMEAD  
Protein\_2 VSHDLRTPLTRIRLATEMMSVEEDGYLAESINKDIEECNAIEEQFIDYLR TGQEMPMEAD  
Protein\_3 VSHDLRTPLTRIRLATEMMSVEEDGYLAESINKDIEECNAIEEQFIDYLR TGQEMPMEAD

Protein LNAVVLGEVIAAESGYEREIEITALYPGSIQVKMHPLSIKRAVANMVVNAARYGNGWIKVSS  
Protein\_1 LNAVVLGEVIAAESGYEREIEITALYPGSIQVKMHPLSIKRAVANMVVNAARYGNGWIKVSS  
Protein\_2 LNAVVLGEVIAAESGYEREIEITALYPGSIQVKMHPLSIKRAVANMVVNAARYGNGWIKVSS  
Protein\_3 LNAVVLGEVIAAESGYEREIEITALYPGSIQVKMHPLSIKRAVANMVVNAARYGNGWIKVSS

Protein GTEFMRRAWFQVEDDGP G IAP EQRKHLFQPFVRGDSARSTISGTGLGLAIVQRIIDNHNGL  
Protein\_1 GTEFMRRAWFQVEDDGP G IAP EQRKHLFQPFVRGDSARSTISGTGLGLAIVQRIIDNHNGL  
Protein\_2 GTEFMRRAWFQVEDDGP G IAP EQRKHLFQPFVRGDSARSTISGTGLGLAIVQRIIDNHNGL  
Protein\_3 GTEFMRRAWFQVEDDGP G IAP EQRKHLFQPFVRGDSARSTISGTGLGLAIVQRIIDNHNGL

Protein EIGTSERGGLSIRAWLPVPVTRAQGTTKEG  
Protein\_1 EIGTSERGGLSIRAWLPVPVTRAQGTTKDA  
Protein\_2 EIGTSERGGLSIRAWLPVPVTRVQGTATKEA  
Protein\_3 EIGTSERGGLSIRAWLPVPVTRVQGTATKEA



# Pairwise Distance

	1	2	3	4
1. Protein sequence envZ from E. coli				
2. Protein sequence envZ from C. farmeri	0.0341			
3. Protein sequence envZ from S. bongori	36.2194	35.6468		
4. Protein sequence envZ from Y. regensburgei	0.0701	0.0628	36.2741	





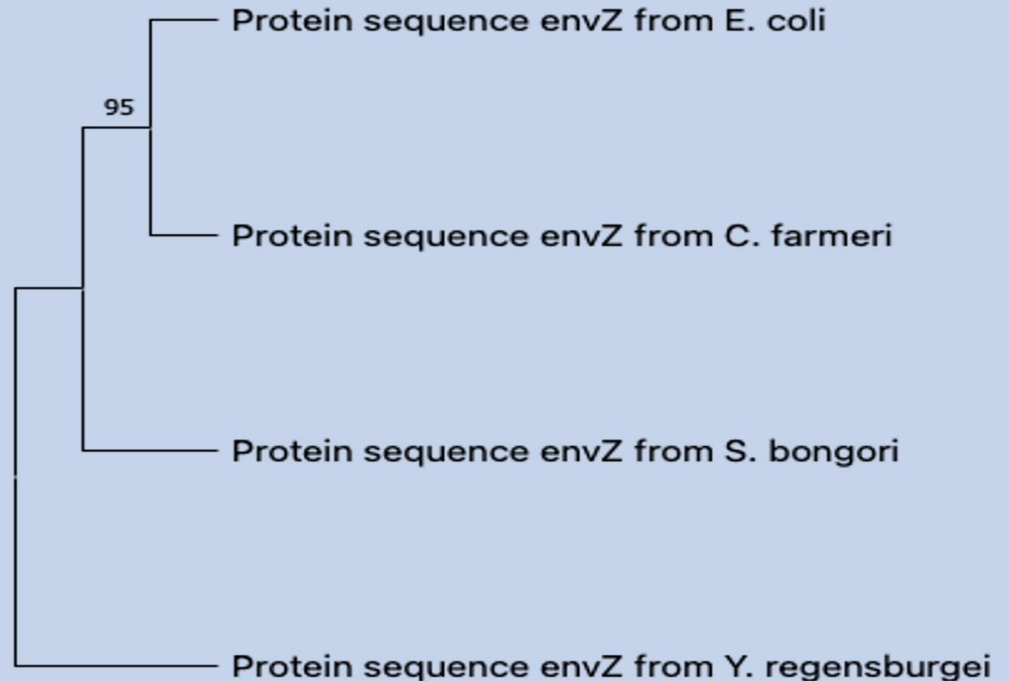
# Phylogenetic Trees

**M11: Analysis Preferences**

Phylogeny Reconstruction

Option	Setting
Statistical Method	Maximum Likelihood
Test of Phylogeny	Bootstrap method
No. of Bootstrap Replications	1000
Substitutions Type	Amino acid
Model/Method	Jones-Taylor-Thornton (JTT) model
Rates among Sites	Uniform Rates
No of Discrete Gamma Categories	Not Applicable
Gaps/Missing Data Treatment	Use all sites
Site Coverage Cutoff (%)	Not Applicable
ML Heuristic Method	Nearest-Neighbor-Interchange (NNI)
Initial Tree for ML	Make initial tree automatically (Default - NJ/Bio)
Initial Tree File	Not Applicable
Branch Swap Filter	None
Number of Threads	3

? Help X Cancel ✓ OK





# Phylogenetic Trees

**M11: Analysis Preferences**

Phylogeny Reconstruction

Option	Setting
Scope	→ All Selected Taxa
Statistical Method	→ Neighbor-joining
Test of Phylogeny	→ Bootstrap method
No. of Bootstrap Replications	→ 10000
Substitutions Type	→ Amino acid
Model/Method	→ Jones-Taylor-Thornton (JTT) model
Rates among Sites	→ Uniform Rates
Gamma Parameter	→ Not Applicable
Pattern among Lineages	→ Same (Homogeneous)
Gaps/Missing Data Treatment	→ Pairwise deletion
Site Coverage Cutoff (%)	→ Not Applicable
Number of Threads	→ 3

? Help    X Cancel    ✓ OK

