

# Phylotyper: *In silico* predictor of molecular subtypes from gene sequences

## Supplementary Information

Matthew D. Whiteside, Chad R. Laing and Victor P.J. Gannon

### Contents

<b>1</b>	<b>Supplementary Figures</b>	<b>1</b>
1.1	Supplementary Figure S1 . . . . .	1
<b>2</b>	<b>Supplementary Tables</b>	<b>2</b>
2.1	Supplementary Table S1 . . . . .	2
2.2	Supplementary Table S2 . . . . .	2

## 1 Supplementary Figures

### 1.1 Supplementary Figure S1

Figure S1: Stx2 Phylogenetic tree showing the Phylotyper marginal likelihoods as pie charts. This output is provided to the user.

## 2 Supplementary Tables

### 2.1 Supplementary Table S1

Table S1: Subtype Schemes in Phylotyper

Name	Description	Species	Loci
stx1	Shiga-toxin 1 subtype	<i>E. coli</i>	2
stx2	Shiga-toxin 1 subtype	<i>E. coli</i>	2
eae	Intimin subtype	<i>E. coli</i>	1
flic	H-serotype based on flagellin gene; fliC	<i>E. coli</i>	1
wz	O-serotype based on the wzy and wzx genes	<i>E. coli</i>	2

### 2.2 Supplementary Table S2

Table S2 contains performance metrics from a leave-one-out cross-validation analysis comparing Phylotyper and a top-BLAST hit approach. The analysis examines the four *E. coli* schemes available in Phylotyper. In this multi-class analysis, precision, recall and  $F_1$  score are calculated for each individual class provided that at least one instance of the class is in the training set. The individual class positive and negatives are summed to calculate an overall precision, recall and  $F_1$  score for the scheme.

Table S2: Leave-One-Out Cross Validation Results

Scheme	Phylotyper			Sequence-similarity		
	Precision	Recall	F <sub>1</sub> Score	Precision	Recall	F <sub>1</sub> Score
<i>E. coli</i> Stx1	1.00	0.94	0.97	0.94	0.94	0.94
<i>E. coli</i> Stx2	1.00	0.99	0.99	0.93	0.93	0.93
<i>E. coli</i> Intimin	1.00	0.98	0.99	0.99	0.98	0.99
<i>E. coli</i> H-serotype	0.99	0.98	0.98	0.96	0.85	0.90
<i>E. coli</i> O-serotype	row4	row4	row4	row4	row4	row4

Formula:

1.  $Precision = TP / (TP + FP)$

2.  $Recall = TP / (TP + FN)$

3.  $F_1 \text{ score} = 2 * Precision * Recall / (Precision + Recall)$

$TP = \text{True Positive}$ ,  $FP = \text{False Positive}$ ,  $FN = \text{False Negative}$