

## SUPPLEMENTARY MATERIAL

788

*Additional Tables*

789

Method	1000M1	1000M2	1000M3	1000M4	RNASim	RNASim2	16S.M	23S.M
SPFN								
PASTA	0.230	0.171	<b>0.049</b>	<b>0.010</b>	<b>0.109</b>	<b>0.067</b>	0.203	0.249
SEPP(F)	<b>0.197</b>	<b>0.156</b>	<b>0.052</b>	<b>0.012</b>	<b>0.104</b>	<b>0.062</b>	<b>0.174</b>	<b>0.238</b>
UPP(F)	<b>0.192</b>	<b>0.153</b>	<b>0.050</b>	<b>0.011</b>	<b>0.104</b>	<b>0.062</b>	<b>0.174</b>	<b>0.238</b>
SEPP(R)	<b>0.197</b>	<b>0.156</b>	<b>0.052</b>	<b>0.012</b>	<b>0.104</b>	<b>0.062</b>	<b>0.175</b>	<b>0.239</b>
UPP(R)	<b>0.192</b>	<b>0.152</b>	<b>0.050</b>	<b>0.011</b>	<b>0.104</b>	<b>0.062</b>	<b>0.174</b>	<b>0.238</b>
SPFP								
PASTA	0.221	0.165	<b>0.049</b>	<b>0.011</b>	<b>0.109</b>	<b>0.068</b>	0.220	<b>0.322</b>
SEPP(F)	<b>0.176</b>	<b>0.141</b>	<b>0.045</b>	<b>0.011</b>	<b>0.104</b>	<b>0.063</b>	<b>0.193</b>	<b>0.312</b>
UPP(F)	<b>0.177</b>	<b>0.142</b>	<b>0.046</b>	<b>0.011</b>	<b>0.104</b>	<b>0.063</b>	<b>0.195</b>	<b>0.312</b>
SEPP(R)	<b>0.176</b>	<b>0.141</b>	<b>0.045</b>	<b>0.011</b>	<b>0.104</b>	<b>0.063</b>	<b>0.192</b>	<b>0.312</b>
UPP(R)	<b>0.177</b>	<b>0.141</b>	<b>0.046</b>	<b>0.011</b>	<b>0.104</b>	<b>0.063</b>	<b>0.195</b>	<b>0.312</b>

Table 6. **Alignment error under low fragmentation.** Each condition has 75% full-length sequences and 25% fragmentary sequences with an average 50% length. SEPP(X) and UPP(X) operate in four steps: first they compute a PASTA alignment on the full-length sequences, then they compute a backbone tree on the backbone alignment using ML heuristic “X” (RAxML or FastTree), then they build an ensemble of profile HMMs on the backbone tree, and finally they add the fragmentary sequences into the backbone alignment. The best results for each model condition (within 1%) are shown in boldface. The error rates are averaged over 20 replicates for the simulated datasets.

Method	1000M1	1000M2	1000M3	1000M4	RNASim	RNASim2	16S.M	23S.M
SPFN								
PASTA	0.640	0.461	0.090	<b>0.016</b>	0.288	0.140	0.396	0.277
SEPP(F)	0.238	0.191	<b>0.070</b>	<b>0.014</b>	<b>0.112</b>	<b>0.069</b>	<b>0.179</b>	<b>0.239</b>
UPP(F)	<b>0.226</b>	<b>0.180</b>	<b>0.062</b>	<b>0.013</b>	<b>0.112</b>	<b>0.068</b>	<b>0.179</b>	<b>0.237</b>
SEPP(R)	0.239	0.190	<b>0.070</b>	<b>0.014</b>	<b>0.112</b>	<b>0.069</b>	<b>0.179</b>	<b>0.238</b>
UPP(R)	<b>0.226</b>	<b>0.180</b>	<b>0.062</b>	<b>0.013</b>	<b>0.112</b>	<b>0.068</b>	<b>0.178</b>	<b>0.237</b>
SPFP								
PASTA	0.613	0.444	0.088	<b>0.015</b>	0.185	0.113	<b>0.141</b>	0.338
SEPP(F)	<b>0.188</b>	<b>0.153</b>	<b>0.050</b>	<b>0.011</b>	<b>0.110</b>	<b>0.067</b>	0.197	<b>0.307</b>
UPP(F)	<b>0.188</b>	<b>0.153</b>	<b>0.049</b>	<b>0.011</b>	<b>0.110</b>	<b>0.067</b>	0.199	<b>0.306</b>
SEPP(R)	<b>0.188</b>	<b>0.153</b>	<b>0.050</b>	<b>0.011</b>	<b>0.110</b>	<b>0.067</b>	0.196	<b>0.301</b>
UPP(R)	<b>0.188</b>	<b>0.153</b>	<b>0.049</b>	<b>0.011</b>	<b>0.110</b>	<b>0.067</b>	0.199	<b>0.306</b>

Table 7. **Alignment error under high fragmentation.** PASTA is run in default mode. SEPP(X) and UPP(X) operate in four steps: first they compute a PASTA alignment on the full-length sequences, then they compute a backbone tree on the backbone alignment using ML heuristic “X” (RAxML or FastTree), then they build an ensemble of profile HMMs on the backbone tree, and finally they add the fragmentary sequences into the backbone alignment. Each condition has 50% full-length sequences and 50% fragmentary sequences with an average 25% length. The error rates are averaged over 20 replicates for the simulated datasets. The best results for each model condition (within 1%) are shown in boldface.

Method	1000M1	1000M2	1000M3	1000M4	RNASim	RNASim2	16S.M	23S.M
<b>FN Rate:</b>								
PASTA-FastTree	0.252	0.185	<b>0.103</b>	<b>0.063</b>	0.199	0.180	0.143	<b>0.143</b>
PASTA-RAXML	0.246	0.181	<b>0.095</b>	<b>0.061</b>	<b>0.186</b>	<b>0.163</b>	0.135	<b>0.137</b>
UPP(F)-FastTree	0.213	0.172	0.129	<b>0.065</b>	0.247	0.218	0.150	0.185
UPP(F)-RAXML	<b>0.159</b>	<b>0.127</b>	<b>0.095</b>	<b>0.061</b>	<b>0.185</b>	<b>0.163</b>	<b>0.121</b>	0.167
UPP(R)-FastTree	0.210	0.171	0.125	<b>0.066</b>	0.248	0.217	<b>0.112</b>	0.167
UPP(R)-RAXML	<b>0.157</b>	<b>0.128</b>	<b>0.094</b>	<b>0.061</b>	<b>0.185</b>	<b>0.163</b>	<b>0.112</b>	<b>0.143</b>
<b>FP Rate:</b>								
PASTA-FastTree	0.255	0.189	<b>0.108</b>	<b>0.085</b>	0.199	0.180	0.598	<b>0.476</b>
PASTA-RAXML	0.248	0.185	<b>0.100</b>	<b>0.083</b>	<b>0.186</b>	<b>0.163</b>	0.595	<b>0.473</b>
UPP(F)-FastTree	0.215	0.176	0.134	<b>0.087</b>	0.247	0.218	0.598	0.502
UPP(F)-RAXML	<b>0.162</b>	<b>0.131</b>	<b>0.100</b>	<b>0.083</b>	<b>0.185</b>	<b>0.163</b>	<b>0.588</b>	0.491
UPP(R)-FastTree	0.213	0.175	0.130	<b>0.088</b>	0.248	0.217	<b>0.580</b>	0.491
UPP(R)-RAXML	<b>0.160</b>	<b>0.132</b>	<b>0.099</b>	<b>0.083</b>	<b>0.185</b>	<b>0.163</b>	<b>0.584</b>	<b>0.476</b>
<b>Resolution:</b>								
PASTA-FastTree	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
PASTA-RAXML	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
UPP(F)-FastTree	1.000	1.000	1.000	1.000	1.000	1.000	0.991	1.000
UPP(F)-RAXML	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
UPP(R)-FastTree	1.000	1.000	1.000	1.000	1.000	1.000	0.991	1.000
UPP(R)-RAXML	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000

Table 8. **Tree error rates and resolution under low fragmentation for MSA-ML methods.** We show FN rates (top), FP rates (middle), and resolution (bottom); each method is given by a pair U-V where U is the MSA method and V is the tree estimation method. Each condition has 75% full-length sequences and 25% fragmentary sequences (which have an average 50% length). The best results for each model condition (within 1%) are shown in boldface. The error rates are averaged over 20 replicates for the simulated datasets.

Method	1000M1	1000M2	1000M3	1000M4	RNASim	RNASim2	16S.M	23S.M
<b>FN Rate:</b>								
PASTA-FastTree	0.760	0.615	0.377	0.205	0.541	0.495	0.485	0.381
PASTA-RAXML	0.765	0.616	0.355	<b>0.164</b>	0.436	0.362	0.409	<b>0.321</b>
UPP(F)-FastTree	0.664	0.622	0.580	0.389	0.715	0.696	0.580	0.470
UPP(F)-RAXML	<b>0.373</b>	<b>0.307</b>	<b>0.236</b>	<b>0.168</b>	<b>0.377</b>	<b>0.336</b>	0.373	0.363
UPP(R)-FastTree	0.667	0.628	0.573	0.395	0.714	0.696	0.591	0.512
UPP(R)-RAXML	<b>0.370</b>	<b>0.304</b>	<b>0.237</b>	<b>0.167</b>	<b>0.377</b>	<b>0.338</b>	<b>0.340</b>	0.363
<b>FP Rate:</b>								
PASTA-FastTree	0.761	0.617	0.381	0.224	0.541	0.495	0.758	0.622
PASTA-RAXML	0.766	0.618	0.359	<b>0.184</b>	0.436	0.362	0.723	<b>0.585</b>
UPP(F)-FastTree	0.666	0.624	0.582	0.404	0.715	0.696	0.802	0.676
UPP(F)-RAXML	<b>0.375</b>	<b>0.310</b>	<b>0.240</b>	<b>0.187</b>	<b>0.377</b>	<b>0.336</b>	0.706	0.611
UPP(R)-FastTree	0.668	0.629	0.576	0.409	0.714	0.696	0.808	0.702
UPP(R)-RAXML	<b>0.372</b>	<b>0.307</b>	<b>0.241</b>	<b>0.187</b>	<b>0.377</b>	<b>0.338</b>	<b>0.690</b>	0.611
<b>Resolution:</b>								
PASTA-FastTree	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
PASTA-RAXML	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
UPP(F)-FastTree	1.000	1.000	1.000	1.000	1.000	1.000	0.998	1.000
UPP(F)-RAXML	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000
UPP(R)-FastTree	1.000	1.000	1.000	1.000	1.000	1.000	0.998	1.000
UPP(R)-RAXML	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000

Table 9. **Tree error rates and resolution under high fragmentation for MSA-ML methods.** We show FN rates (top), FP rates (middle), and resolution (bottom); each method is given by a pair U-V where U is the MSA method and V is the tree estimation method. Each condition has 50% full-length sequences and 50% fragmentary sequences with an average 25% length. The best results for each model condition (within 1%) are shown in boldface. The error rates are averaged over 20 replicates for the simulated datasets.

Method	1000M1	1000M2	1000M3	1000M4	RNASim	RNASim2	16S.M	23S.M
<b>FN Rate:</b>								
UPP(F)-pplacer	<b>0.217</b>	<b>0.181</b>	<b>0.153</b>	<b>0.114</b>	<b>0.244</b>	<b>0.220</b>	<b>0.173</b>	<b>0.208</b>
SEPP(F)-pplacer	<b>0.222</b>	<b>0.186</b>	<b>0.156</b>	<b>0.113</b>	<b>0.243</b>	<b>0.221</b>	0.202	0.220
SEPP(F)-pplacer(c)	0.227	<b>0.189</b>	0.161	<b>0.118</b>	0.280	0.235	0.247	0.274
UPP(F)-APPLES	0.279	0.230	0.198	0.150	0.379	0.334	0.240	0.298
UPP(R)-pplacer	<b>0.215</b>	<b>0.183</b>	<b>0.150</b>	<b>0.112</b>	<b>0.243</b>	<b>0.215</b>	0.192	0.244
SEPP(R)-pplacer	<b>0.218</b>	<b>0.186</b>	<b>0.152</b>	<b>0.112</b>	<b>0.243</b>	<b>0.216</b>	0.230	0.274
SEPP(R)-pplacer(c)	<b>0.225</b>	<b>0.188</b>	<b>0.160</b>	<b>0.117</b>	0.280	0.231	0.249	0.321
UPP(R)-APPLES	0.270	0.229	0.195	0.148	0.380	0.329	0.245	0.333
<b>FP Rate:</b>								
UPP(F)-pplacer	<b>0.181</b>	<b>0.144</b>	<b>0.114</b>	<b>0.092</b>	<b>0.207</b>	<b>0.184</b>	<b>0.594</b>	<b>0.492</b>
SEPP(F)-pplacer	<b>0.187</b>	<b>0.149</b>	<b>0.116</b>	<b>0.092</b>	<b>0.206</b>	<b>0.184</b>	0.608	<b>0.500</b>
SEPP(F)-pplacer(c)	0.193	<b>0.153</b>	0.122	<b>0.098</b>	0.243	0.200	0.632	0.534
UPP(F)-APPLES	0.247	0.198	0.164	0.132	0.345	0.301	0.620	0.546
UPP(R)-pplacer	<b>0.179</b>	<b>0.146</b>	<b>0.111</b>	<b>0.091</b>	<b>0.208</b>	<b>0.178</b>	0.604	0.517
SEPP(R)-pplacer	<b>0.184</b>	<b>0.149</b>	<b>0.113</b>	<b>0.091</b>	<b>0.207</b>	<b>0.179</b>	0.624	0.534
SEPP(R)-pplacer(c)	0.190	<b>0.153</b>	0.121	<b>0.096</b>	0.244	0.195	0.632	0.567
UPP(R)-APPLES	0.238	0.197	0.161	0.129	0.345	0.295	0.622	0.576
<b>Resolution:</b>								
UPP(F)-pplacer	0.952	0.952	0.950	0.953	0.953	0.956	0.954	0.953
SEPP(F)-pplacer	0.953	0.952	0.950	0.953	0.953	0.956	0.954	0.953
SEPP(F)-pplacer(c)	0.954	0.953	0.950	0.954	0.951	0.956	0.959	0.953
UPP(F)-APPLES	0.954	0.955	0.954	0.956	0.947	0.953	0.937	0.945
UPP(R)-pplacer	0.953	0.952	0.950	0.953	0.955	0.955	0.957	0.956
SEPP(R)-pplacer	0.954	0.952	0.950	0.953	0.954	0.955	0.959	0.953
SEPP(R)-pplacer(c)	0.954	0.954	0.951	0.954	0.951	0.956	0.957	0.956
UPP(R)-APPLES	0.955	0.956	0.954	0.956	0.946	0.952	0.937	0.960

Table 10. **Tree error rates and resolution for placement-based methods under low fragmentation.** We show FN rates (top), FP rates (middle), and degree of resolution (bottom). Each method is given by a pair U(B)-V where U is the extended alignment method, B is the backbone tree method, and V is the placement method. All methods run PASTA in default mode on the full-length sequences, and differ only in how they perform the remaining steps. Each condition has 75% full-length sequences and 25% fragmentary sequences (which have an average 50% length). The best results for each model condition (within 1%) are shown in boldface. The error rates are averaged over 20 replicates for the simulated datasets.

Method	1000M1	1000M2	1000M3	1000M4	RNASim	RNASim2	16S.M	23S.M
<b>FN Rate:</b>								
UPP(F)-pplacer	<b>0.487</b>	<b>0.437</b>	<b>0.379</b>	<b>0.320</b>	<b>0.507</b>	<b>0.476</b>	<b>0.499</b>	<b>0.429</b>
SEPP(F)-pplacer	0.514	0.459	0.395	<b>0.320</b>	<b>0.509</b>	<b>0.475</b>	<b>0.496</b>	<b>0.423</b>
SEPP(F)-pplacer(c)	0.534	0.483	0.420	0.334	0.560	0.510	0.549	0.506
UPP(F)-APPLES	0.559	0.517	0.472	0.394	0.716	0.674	0.591	0.542
UPP(R)-pplacer	<b>0.488</b>	<b>0.437</b>	<b>0.380</b>	<b>0.320</b>	<b>0.507</b>	<b>0.477</b>	<b>0.496</b>	0.458
SEPP(R)-pplacer	0.520	0.452	0.392	<b>0.319</b>	<b>0.508</b>	<b>0.475</b>	0.508	0.542
SEPP(R)-pplacer(c)	0.539	0.477	0.417	0.334	0.561	0.510	0.615	0.613
UPP(R)-APPLES	0.560	0.515	0.471	0.394	0.718	0.675	0.620	0.530
<b>FP Rate:</b>								
UPP(F)-pplacer	<b>0.369</b>	<b>0.307</b>	<b>0.233</b>	<b>0.172</b>	<b>0.397</b>	<b>0.356</b>	<b>0.716</b>	<b>0.586</b>
SEPP(F)-pplacer	0.403	0.335	0.255	<b>0.172</b>	<b>0.399</b>	<b>0.356</b>	<b>0.717</b>	<b>0.578</b>
SEPP(F)-pplacer(c)	0.431	0.367	0.288	0.189	0.459	0.401	0.743	0.629
UPP(F)-APPLES	0.455	0.408	0.356	0.270	0.641	0.595	0.759	0.659
UPP(R)-pplacer	<b>0.370</b>	<b>0.307</b>	<b>0.234</b>	<b>0.170</b>	<b>0.397</b>	<b>0.362</b>	<b>0.712</b>	0.613
SEPP(R)-pplacer	0.410	0.327	0.251	<b>0.170</b>	<b>0.400</b>	<b>0.359</b>	<b>0.720</b>	0.668
SEPP(R)-pplacer(c)	0.435	0.360	0.284	0.189	0.460	0.402	0.780	0.714
UPP(R)-APPLES	0.454	0.406	0.355	0.270	0.643	0.595	0.771	0.646
<b>Resolution:</b>								
UPP(F)-pplacer	0.810	0.808	0.805	0.801	0.817	0.815	0.827	0.844
SEPP(F)-pplacer	0.811	0.809	0.807	0.801	0.817	0.815	0.833	0.836
SEPP(F)-pplacer(c)	0.816	0.813	0.811	0.801	0.812	0.818	0.822	0.815
UPP(F)-APPLES	0.806	0.812	0.816	0.811	0.790	0.804	0.796	0.822
UPP(R)-pplacer	0.810	0.809	0.805	0.800	0.818	0.819	0.821	0.855
SEPP(R)-pplacer	0.812	0.810	0.808	0.800	0.819	0.819	0.824	0.844
SEPP(R)-pplacer(c)	0.814	0.813	0.811	0.801	0.814	0.820	0.821	0.825
UPP(R)-APPLES	0.803	0.813	0.815	0.810	0.791	0.804	0.777	0.811

Table 11. **Tree error rates for placement-based methods under high fragmentation.** We show FN rates(top), FP rates (middle), and resolution (bottom). Each method is given by a pair U(B)-V where U is the extended alignment method, B is the backbone tree method, and V is the placement method. All methods run PASTA in default mode on the full-length sequences, and differ only in how they perform the remaining steps. Each condition has 50% full-length sequences and 50% fragmentary sequences (which have an average 25% length). The best results for each model condition (within 1%) are shown in boldface. The error rates are averaged over 20 replicates for the simulated datasets.

790 The full-length datasets are available at:

791 **ROSE:** <https://sites.google.com/eng.ucsd.edu/datasets/alignment/sate-i>

792 **RNASim:** <https://sites.google.com/eng.ucsd.edu/datasets/alignment/pastaupp>

793 **16S/23S:** <https://sites.google.com/eng.ucsd.edu/datasets/alignment/16s23s>

## COMMANDS

Here, we list the tools that we used and how they were run.

*PASTA 1.8.5*

```
python3 run_pasta.py -i backbone_sequences.txt -o output_dir  
--temporaries temp_dir
```

*RAxML-NG 0.9.0*

```
raxml-ng --msa sequences.txt --prefix name --threads 8  
--seed 242234 --model GTR+G --tree pars{5}
```

*FastTree 2.1*

```
FastTree -nt -gtr sequences.txt > output.tre
```

*UPP 4.3.10*

```
python3 run_upp.py -s fragment_sequences.txt  
-a backbone_alignment.txt -t backbone_tree.tre
```

*SEPP 4.3.10*

```
python3 run_sepp.py -t backbone_tree.tre -a backbone_alignment.txt  
-f fragment_sequences.txt -r raxml_info.txt
```

*pplacer 1.1*

```
pplacer -t backbone_tree.tre -r backbone_alignment.fasta  
-s raxml_info.txt -o output.json fragment_alignment.fasta
```

*APPLES 1.2.0*

813

814 APPLES requires the reference tree branch lengths to be estimated, which was done  
815 with FastTree as follows:

```
816 FastTree -nosupport -nt -nome -noml -intree backbone_tree.tre  
817 < backbone_alignment.txt > backbone_tree_me.tre
```

818 Then, APPLES was run as follows:

```
819 python3 run_apples.py -t backbone_tree_me.tre -o output.json  
820 -s backbone_alignment.txt -q fragment_alignment.txt
```

*Alignment Error*

821

822 We used FastSP 1.6.0(Mirarab and Warnow, 2011) to estimate alignment error:

```
823 java -jar FastSP_1.6.0.jar -r true_alignment.txt  
824 -e estimated_alignment.txt
```

*Tree Error*

825

826 We used the compare\_trees.py script (courtesy Erin Molloy), found in tools.zip at  
827 <https://databank.illinois.edu/datasets/IDB-1424746>, using the following command:

```
828 compare_trees.py <true_tree> <estimated_tree>
```

829 The script returns the number of false negative and positive edges, which were divided by  
830 the number of internal edges in the true tree and estimated tree, respectively.