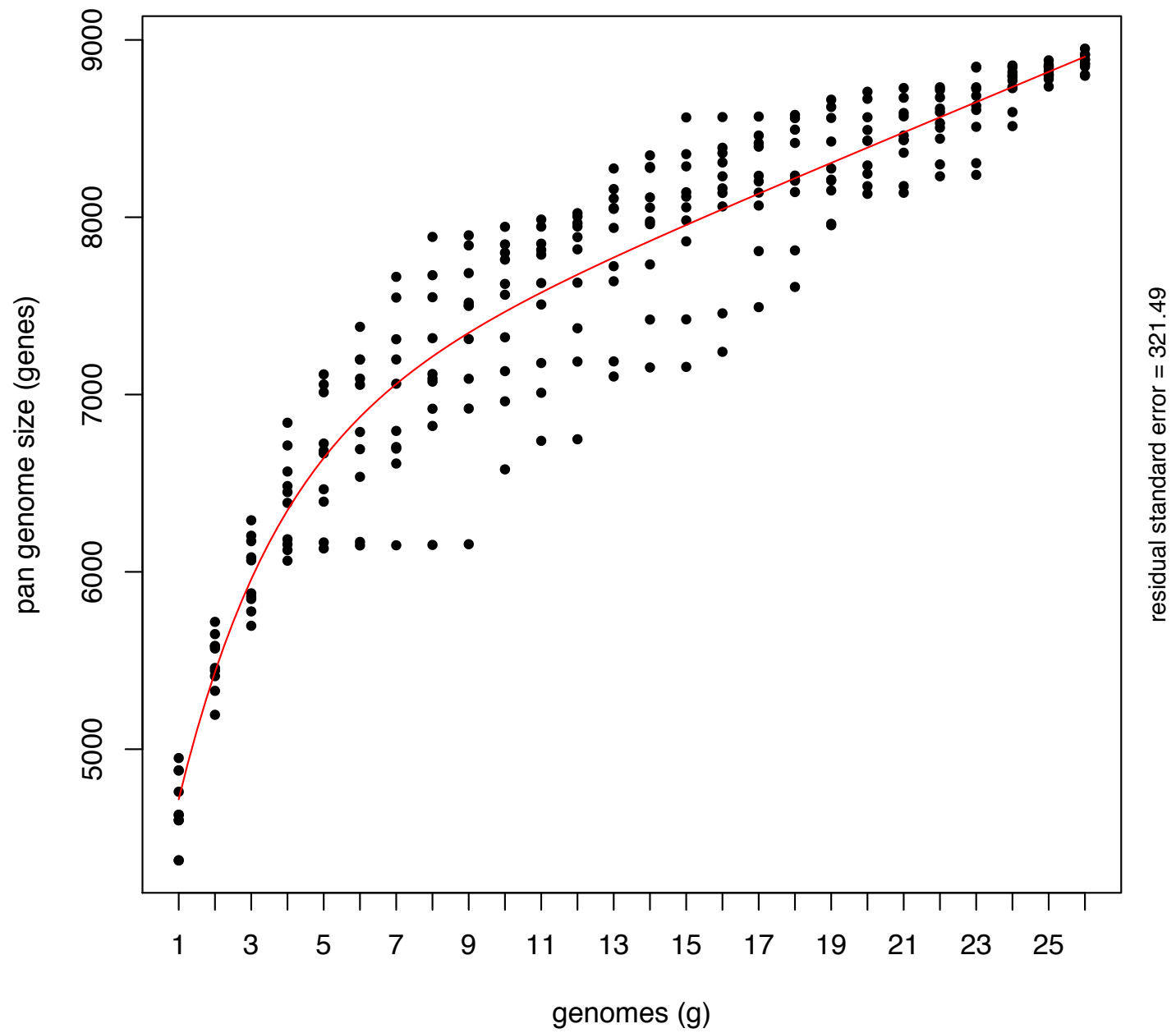**Supp. Fig. 1**

$$\text{pangenes}(g) = 4717 + 85.2(g-1) + 1317 \exp\left(\frac{-2}{2.72}\right) \frac{1 - \exp\left(\frac{-(g-1)}{2.72}\right)}{1 - \exp\left(\frac{-1}{2.72}\right)}$$

pan genome size (genes)

genomes (g)

residual standard error = 321.49

**Supp. Fig. 2**

**Supplementary File 1:** Quality evaluation for 24 UPEC de novo draft genome assemblies.
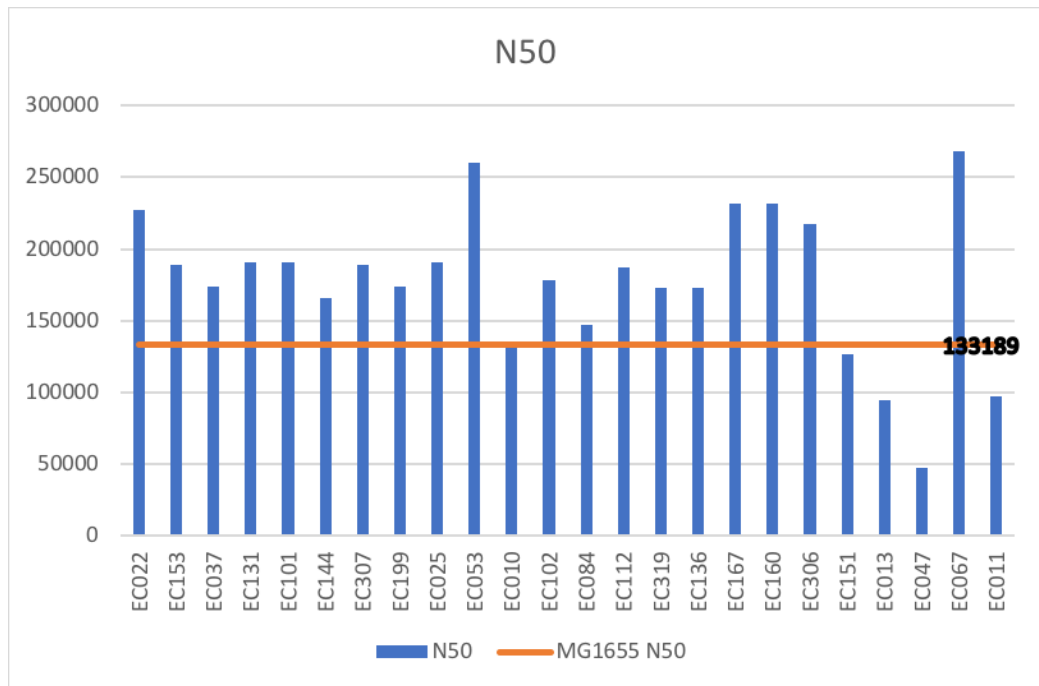
We calculated the total assembly length and the N50 value (the length of the smallest contig which when added to the other larger contigs of the assembly yields at least 50% of the genome assembly) of our 24 *de novo* draft genome assemblies. **Table S1-A** shows the obtained values for these statistics. **Figure S1-1** shows the total Assembly length of our assemblies in comparison with genome length of *E. coli* K.12 MG1655 and UPEC CFT073. **Figure S1-2** shows the N50 values of our assemblies in comparison with the N50 value obtained by the SPAdes developers in a *de novo* genome assembly with 26 million reads of an Illumina 2x100 paired end library of *E. coli* K-12 MG1655 (http://cab.spbu.ru/software/spades/).

To further examine the completeness of our assemblies we aligned the forward reads of each of our 24 UPEC libraries (included in the assembly input) against the corresponding genome assemblies with bowtie. and measured the percentage of aligned reads. We expect that if an assembly is nearly complete (represents most of the genome of a strain) most of the library reads will align to it, as the reads are expected to represent the complete genome in fragments. **Table S1-B** shows the percentage of each library reads that align to the corresponding *de novo* draft genome assembly.
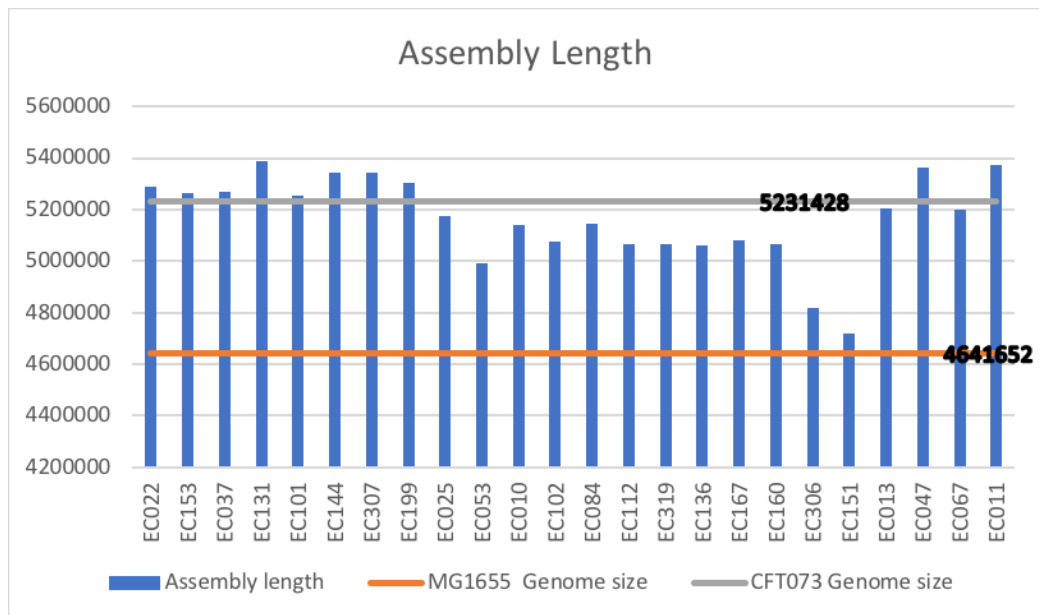
**Table S1-A. N50 and Assembly Length values for our de-novo draft genome assemblies.**

| Serotype-ST-Phylogroup | Sample | N50 | Assembly length |
|---|---|---|---|
| 025:H4 (ST 131) B2 | EC022 | 227017 | 5290048 |
| 025:H4 (ST 131) B2 | EC153 | 188599 | 5266057 |
| 025:H4 (ST 131) B2 | EC037 | 173941 | 5268188 |
| 025:H4 (ST 131) B2 | EC131 | 190541 | 5388762 |
| 025:H4 (ST 131) B2 | EC101 | 191226 | 5256127 |
| 025:H4 (ST 131) B2 | EC144 | 165794 | 5343505 |
| 025:H4 (ST 131 )B2 | EC307 | 189333 | 5344076 |
| 025:H4 (ST 131) B2 | EC199 | 173944 | 5303066 |
| 025:H4 (ST 131) B2 | EC025 | 191227 | 5174605 |
| 016:H5 (ST-131) B2 | EC053 | 260028 | 4990640 |
| 075:H5 (ST-1193) B2 | EC010 | 134538 | 5141539 |
| 075:H5 (ST-1193) B2 | EC102 | 178414 | 5075017 |
| 075:H5 (ST-14) B2 | EC084 | 147108 | 5147666 |
| 08:H9 (ST-423) B1 | EC112 | 187031 | 5065990 |
| 08:H9 (ST-423) B1 | EC319 | 172754 | 5066920 |
| 08:H9 (ST-423) B1 | EC136 | 172754 | 5060390 |
| 08:H9 (ST-423) B1 | EC167 | 231780 | 5080520 |
| 08:H9 (ST-423) B1 | EC160 | 231544 | 5066067 |
| 045:H11 (ST-297) B1 | EC306 | 217310 | 4817787 |
| 0116:H48 (ST-351) B1 | EC151 | 126473 | 4717866 |
| 016:H4 (ST10) A | EC013 | 94645 | 5203374 |
| NA (ST-69) D | EC047 | 47191 | 5362429 |
| NA (ST-69) D | EC067 | 267993 | 5201365 |
| 01:H6 (ST-648) - | EC011 | 97802 | 5374532 |

**Supplementary File 1:** Quality evaluation for 24 UPEC de novo draft genome assemblies.



**Figure S1-1**. N50 values for our 24 de novo genome assemblies. The orange line indicates de N50 value reported by the SPAdes developers for a Illumina paired-end library of *E. coli* K-12 MG1655.
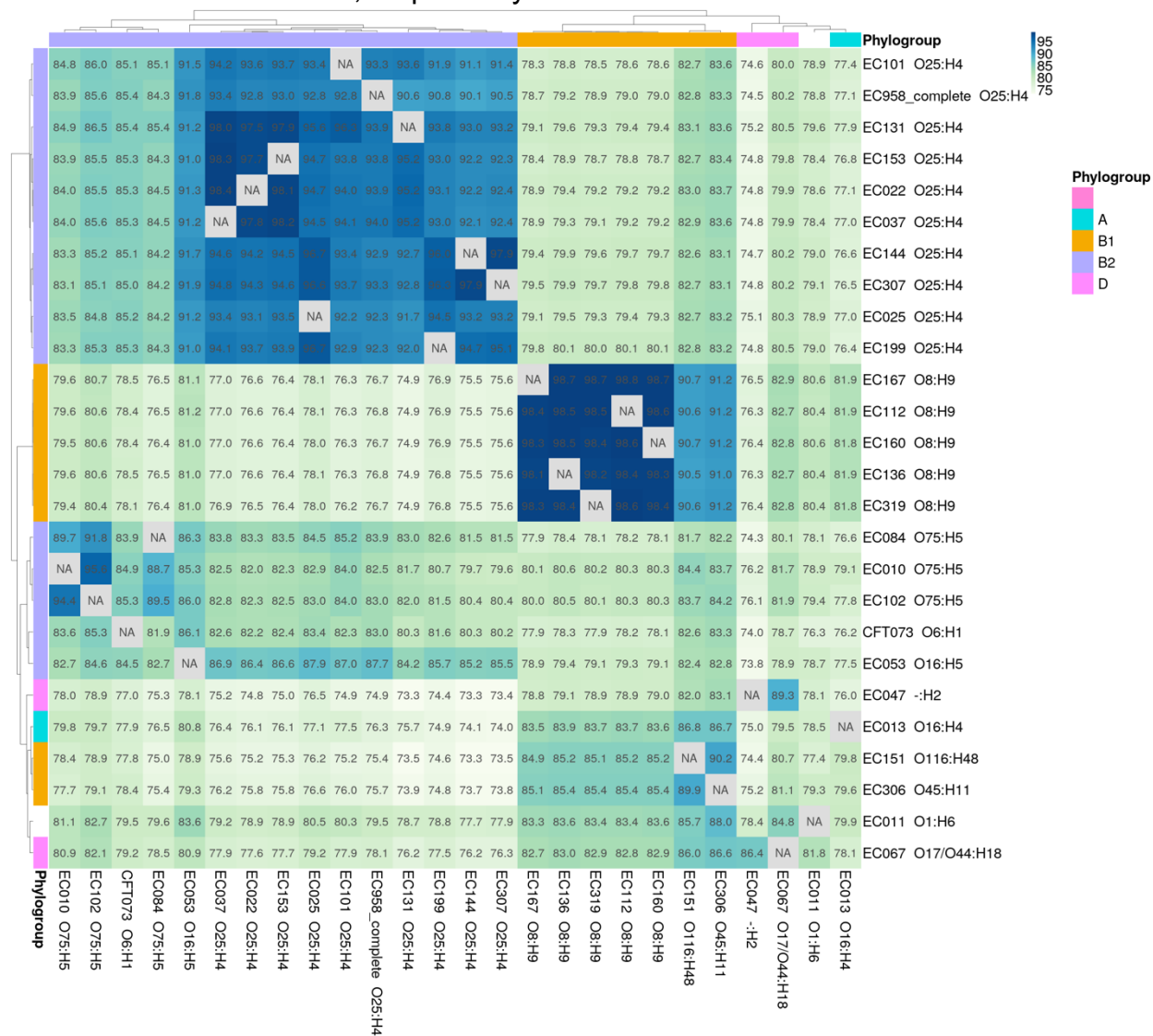


**Figure S1-2**. Assembly length values for our 24 de novo genome assemblies. The orange and gray lines indicate de genome size of *E. coli* K-12 MG1655 and UPEC CFT073, respectively.

**Supplementary File 2:** Additional synteny assessment statistics for the 24 UPEC strains and the model strains CFT073 and EC958 ST-131.

We identified all of the regions comprising at least 5 contiguous genes with a maximum gap of 1 gene present in only one of the strains (see Materials and Methods). **Figures S2-1, S2-2 and S2-3** show different statistics obtained from the i-ADHoRe results.

**Figure S2-1.** Synteny percentage for each strain comparison. The percentage of the total genes of each strain a (column labels) that appear in syntenic regions in each strain b (row labels) is shown. The strains were hierarchically clustered according to these values. The phylogroups and serotypes of the strains are shown at the heatmap margins and in the row and column labels, respectively.

**Figure S2-2.** Synteny percentage for the genes shared between each pair of strains. The percentage of the shared genes that appear in regions with the same synteny in both strains is shown. The strains were hierarchically clustered according to these values. The phylogroups and serotypes of the strains are shown at the heatmap margins and in the row and column labels, respectively.