1      **Supplementary information**

2

3      **Extended insight into the *Mycobacterium chelonae-abscessus* complex through**

4      **whole genome sequencing of *Mycobacterium salmoniphilum* outbreak and**

5      ***Mycobacterium salmoniphilum*-like strains**

6

7      Phani Rama Krishna Behra[1#], Sarbashis Das[1#], B. M. Fredrik Pettersson[1], Lisa Shirreff[2], Tanner

8      DuCote[2], Karl-Gustav Jacobsson[3], Don G. Ennis[2] and Leif A. Kirsebom[1*]

9

10     [1]Department of Cell and Molecular Biology
11     Box 596, Biomedical Centre
12     SE-751 24 Uppsala, Sweden
13
14     [2]Department of Biology,
15     University of Louisiana,
16     Lafayette, Louisiana, USA
17
18     [3]Department of Neuroscience
19     Box 593, Biomedical Centre
20     SE-751 24 Uppsala, Sweden
21
22
23
24     [#]These authors contributed equally to this work
25
26     *corresponding author
27      Leif.kirsebom@icm.uu.se

28

29

30

31

32   **Supplementary Methods:** DNA isolation, Genome assembly, annotation, plasmid, phage,

33   identification of IS elements, horizontal gene transfer (HGT) analysis and identification of SNV

34   and mutational hotspots.

35   **Supplementary Table S1:** Compilation of mycobacterial species/strains, and genomes used in

36   the present study.

37   **Supplementary Table S2:** Summary of predicted ncRNA genes in MCAC members.

38   **Supplementary Table S3:** Summary of predicted phages in the different MCAC members.

39   **Supplementary Table S4:** Summary of IS-elemnets in the different MCAC members.

40   **Supplementary Table S5a-e:** List of core and unique genes in $Msal^{T}$, $Mche^{T}$, $Msal$-like$^{CCUG64054}$

41   and $Mabs^{ATCC19977}$.

42   **Supplementary Table S6:** List of hotspot genes, annotation and function.

43   **Supplementary Table S7:** Horizontal gene transfer analysis. Sheet 1, summary of predicted

44   HGT genes in $Msal$ and $Msal$-like strains and other mycobacteria. Subsequent sheets contain

45   detail information of predicted HGT genes in individual strains.

46   **Supplementary Table S8:** Virulence factor analysis. Sheet 1, list of genomes used in VF

47   analysis and sheet 2 list predicted virulence genes along with functional classification for $Msal$

48   and $Msal$-like strains and other mycobacteria.

49   **Supplementary Table S9a, b:** (a) List of genes encoding ribosomal proteins in MCAC-

50   members. (b) List of genes encoding translation factors in MCAC-members.

51

52 **Supplementary Methods**

53 *DNA isolation*

54 For Illumina sequencing, cells were lyzed by bead beating (2 x 1 min, 6.5 m/s, 5 min on ice

55 between runs, 0.1 mm silica/zirconium beads) in equal volumes of TE-buffer (10 mM Tris-HCl,

56 pH 7.5; 1 mM EDTA) and DNAZol reagent (Invitrogen) using a FastPrep24 device (MP

57 Biotech). This was followed by chloroform extraction and ethanol precipitation, resuspension in

58 1 x TE-buffer and removal of RNA and proteins using RNase A and Proteinase K treatment for

59 1 h each. Chromosomal DNA was retrieved by phenol/chloroform extraction and ethanol

60 precipitation.

61   For PacBio sequencing, 500 mL of exponentially growing culture was pelleted and

62 resuspended in 11 mL of Qiagen buffer B1 (containing 1 mg/mL RNase A) and transferred to a

63 tube containing 2 g ($\geq$ 60000 U) Lipase (product number 80612, Sigma-Aldrich). After

64 dissolving the Lipase, the tubes were incubated for 2 h at $37^{o}$C in a waterbath followed by:

65 addition of 600 µL lysozyme (100 mg/mL) and 3 h of incubation, addition of 500 µL of

66 Proteinase K (20 mg/mL) and incubation for 1.5 h, and addition of Qiagen buffer B2 and

67 incubation for 16 h at $50^{o}$C. The cell lysate was cleared by centrifugation. The DNA was

68 recovered using Qiagen Genomic-tip 500/G following the protocol supplied by the manufacturer

69 and further purified using the MoBio PowerClean Pro DNA Clean-Up Kit.

70   Before submitting DNA for sequencing (PacBio or Illumina), DNA size and quality was

71 estimated using spectrophotometry and agarose gel electrophoresis.

72

73 *Genome assembly, annotation, plasmid, phage, identification of IS elements, identification of*

74 *SNV and mutational hotspots and horizontal gene transfer (HGT) analysis*

75 *Genome assembly*: The PacBio-generated reads were assembled using the SMRT-analysis

76 HGAP3 assembly pipeline (Chin *et al*. 2013) and polished using Quiver (Pacific Biosciences,

77 Menlo Park, CA, USA). Assembly of the Illumina generated reads was performed using the A5-

78 Assembly pipeline (version 20140604) with a minimum contig size of 200 bases (Tritt *et al*.

79 2012). The MAUVE program (Darling et al. 2004) was used for genome reordering and whole

80 genome alignment. This alignment was plotted using genoplotR (Guy et al. 2010). The

81 RNAmmer (Lagesen *et al*. 2007) and tRNAScan-SE (Lowe and Eddy 1997) programs were used

82 to predict the rRNA and tRNA genes. All the genomes were annotated and functionally

83 classified using Prokka [version 1.11] (Seemann 2014) and RAST server (http://rast.nmpdr.org/)

84 (Aziz *et al*. 2008), respectively.

85 *Plasmid, phage and IS element predictions*: Assembled scaffolds were subjected to BLAST

86 search using the NCBI plasmid database (downloaded March 2016). We considered a scaffold

87 belonging to a plasmid if more than 90% of the scaffold sequence aligned with a plasmid

88 sequence from the plasmid database with more than 90% identity.

89  Phage sequences were predicted using the PHAST server (Zhou *et al*. 2011), while prediction

90 of IS elements was done using the ISsaga webserver (Varani *et al*. 2011).

91 *Identifications of mutational hotspots and SNVs*: Mutational hotspots were identified using

92 Shewhart Control Chart, as described by Das *et al*. (2012). Briefly, SNVs were identified

93 between $Msal^T$ and other *Msal* strains in a pairwise manner using MUMmer (Delcher *et al*.

94 1999). The reference genome $Msal^T$ was divided into non-overlapping windows of 2000 bases

95 and the average number of SNVs in each of the windows was determined. The average SNV

96 values were subsequently used in Shewhart Control Chart for the prediction of hotspots.

97 *Horizontal Gene Transfer* (*HGT*): Putative horizontal gene transfer events were predicted using

98 the tool HGTector v0.2.2 (Zhu *et al*. 2014). Briefly, this approach is a combination of BLASTp

99 and taxonomy searches. For the BLASTp search analysis, we used the DIAMOND v0.9.10

100 tool and build the database file (NCBI NR database downloaded from the HGTector

101 source v2017-6-30) (Buchfink et al. 2015). The parameters we used for the BLASTp analysis,

102 percentage identity = >60% and query coverage = >70%, e-value = <1e-100. For the taxonomy

103 search (using HGTector), we used three parameters, "self" group, "close" group, and "distal"

104 group hierarchical classification to predict putative horizontal genes where "self =

105 Mycobacteriaceae", and "close = Corynebacteriales" (as of Feb 2018, NCBI taxonomy; Sayers *et*

106 *al*. 2009). The "distal" group = all other organisms except the "self" and "close" groups (Zhu *et*

107 *al*. 2014; see Supplementary Table S7 for further information). Finally, the predicted HGTs were

108 analysed by performing a Mann–Whitney-Wilcoxon test (in R ver 3.2.2, 2015-08-14, on

109 platform x86_64-pc-linux-gnu) for GC-content of genome-encoded protein-coding gene

110 sequences (CDS, excluding horizontal transferred genes) and GC-content for candidate

111 horizontal transferred genes.

112 To identify common and unique genes in $Msal^T$, $Mche^T$, $Msal$-like$^{CCUG64054}$ and $Mabs^{ATCC19977}$

113 we combined PanOCT ortholog clustering (Fouts *et al*. 2012) with BLASTp search (Boratyn *et*

114 *al*. 2013) and cut-off e-value = <1e-05, percent identity = >45% and query coverage = >70%.

115

116 **References**

117 Aziz, R. K. *et al*. The RAST Server: rapid annotations using subsystems technology. *BMC*

118 *Genomics* **9**, 75 (2008).

119 Buchfink, B., Xie, C. & Huson, D.H. Fast and sensitive protein alignment using DIAMOND.

120 *Nat Meth* **12**, 59–60 (2015).

121 Boratyn, G. M. *et al*. BLAST: a more efficient report with usability improvements. *Nucl Acids*

122 *Res* **41**, W29–W33 (2013).

123 Chin, C-S. *et al*. Nonhybrid, finished microbial genome assemblies from long-read SMRT

124 sequencing data. *Nat Meth* **10**, 563–569 (2013).

125 Darling, A. C. E., Mau, B., Blattner, F. R. & Perna, N. T. Mauve: multiple alignment of

126 conserved genomic sequence with rearrangements. *Genome Res* 14, 1394–1403 (2004).

127    Das, S. *et al*. Identification of hot and cold spots in genome of Mycobacterium tuberculosis using

128    Shewhart Control Charts. *Sci Rep* **2**, 297 (2012).

129    Delcher, A. L., Kasif, S., Fleischmann, R. D., Peterson, J., White, O. & Salzberg, S. L.

130    Alignment of whole genomes. *Nucl Acids Res* **27**, 2369-2376 (1999).

131    Fouts, D. E., Brinkac, L., Beck, E., Inman, J. & Sutton, G. PanOCT: automated clustering of

132    orthologs using conserved gene neighborhood for pan-genomic analysis of bacterial strains and

133    closely related species. *Nucl Acids Res* **40**, e172–e172 (2012).

134    Guy, L., Kultima, J. R. & Andersson, S. G. E. genoPlotR: comparative gene and genome

135    visualization in R. *Bioinformatics* (Oxford, England) **26**, 2334–2335 (2010).

136    Lagesen, K. *et al*. RNAmmer: consistent and rapid annotation of ribosomal RNA genes. *Nucl*

137    *Acids Res* **35**, 3100-3108 (2007).

138    Lowe, T. M. & Eddy, S. R. tRNAscan-SE: a program for improved detection of transfer RNA

139    genes in genomic sequence. *Nucl Acids Res* **25**, 955-964 (1997).

140    Sayers, E. W. *et al*. Database resources of the National Center for Biotechnology Information.

141    *Nucl Acids Res* **37**, D5-15 (2009).

142    Seemann, T. Prokka: rapid prokaryotic genome annotation. *Bioinformatics* (Oxford, England).

143    **30**, 2068–2069 (2014).

144    Tritt, A., Eisen, J. A., Facciotti, M. T. & Darling, A. E. An integrated pipeline for de novo

145    assembly of microbial genomes. *PLoS ONE*. **7**, e42304 (2012).

146    Varani, A.M., Siguier, P., Gourbeyre, E., Charneau, V. & Chandler, M. ISsaga is an ensemble of

147    web-based methods for high throughput identification and semi-automatic annotation of

148    insertion sequences in prokaryotic genomes. *Genome Biol* **12**, R30 (2011).

149    Zhou, Y., Liang, Y., Lynch, K. H., Dennis, J. J. & Wishart, D. S. PHAST: a fast phage search

150    tool. *Nucl Acids Re*s **39**, W347–52 (2011).

151    Zhu, Q., Kosoy, M. & Dittmar, K. HGTector: an automated method facilitating genome-wide

152    discovery of putative horizontal gene transfers. *BMC Genomics* **15**, 717 (2014).

153

154

155

156

157

158

159

160

161

162

163     **Figure S1** Genome alignment and genome-wide distribution of tRNA genes.

164     (a) Whole-genome alignment for the complete genomes $Mche^{T}$, $Msal^{T}$, $Mabs^{ATCC19977}$ and the

165     $Msal$-like$^{CCUG64054}$ draft genome. Each horizontal block represents one genome and vertical lines

166     between the genomes correspond to homologous regions whereas blue diagonal lines correspond

167     to genomic inversions. Of note, we cannot conclusively state that the indicated inversion in

168     $Msal$-like$^{CCUG64054}$ is real due to draft genome status. White gaps correspond to the absence of

169     genes while regions in black represent phage sequences; red stars mark intact phages while black

170     stars mark incomplete/questionable phage sequences (see text for details).

171     (b) $Msal^{T}$ complete genome, blue and red marked tRNA genes refers to transcription from the

172     positive and negative strands, respectively.

173     (c) $Mche^{T}$ complete genome, tRNA genes in red and blue as in (b).
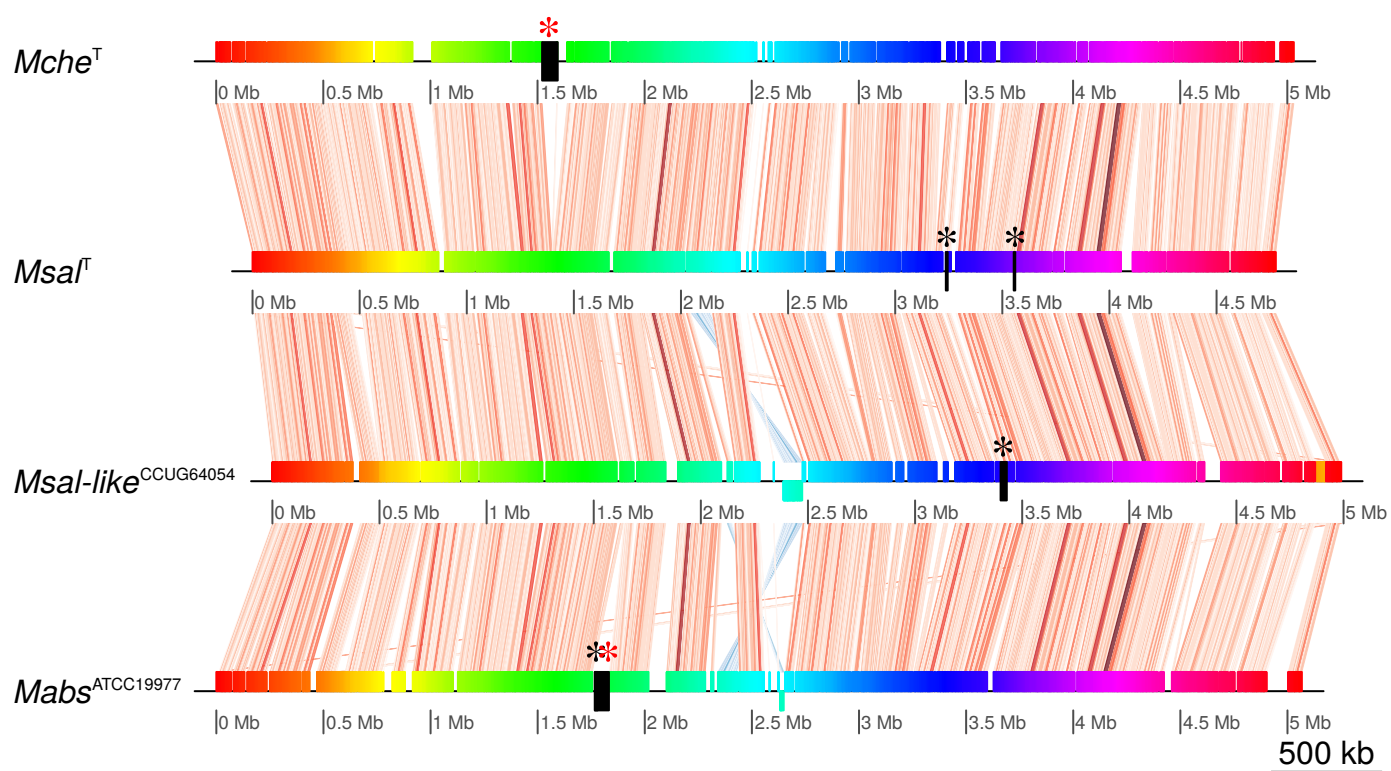
174

175

Figure-S1:

a)

b)



$Msal^T$

4776625 bp

56 tRNAs organized in to 41 operons.
(tRNA gene with in 200 nucleotides region)

c)



$Mche^T$

5030282 bp

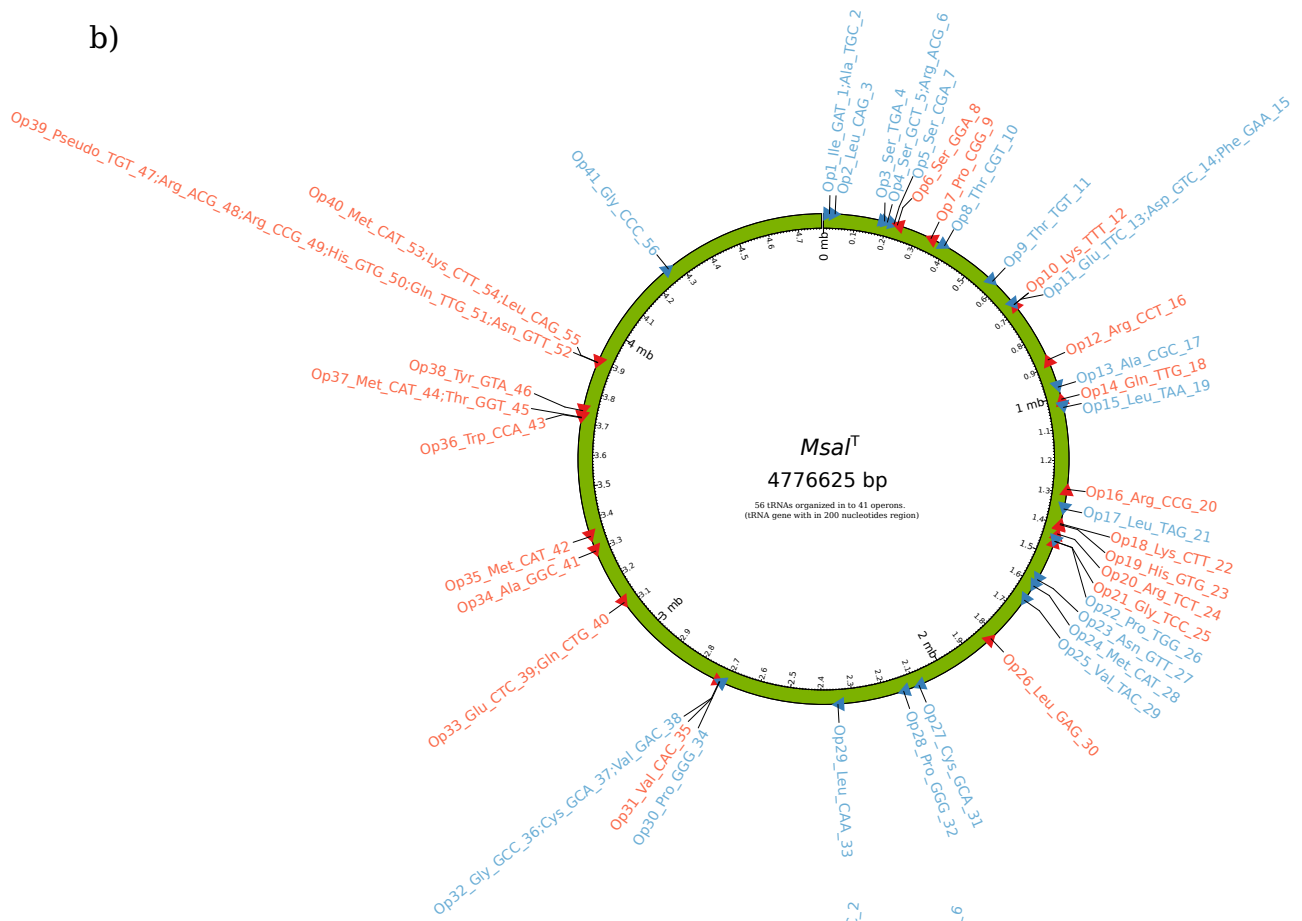47 tRNAs organized in to 39 operons.
(tRNA gene with in 200 nucleotides region)

176    **Figure S2** Overview of genome assembly of *Msal* and *Msal*-like strains, and *Mfra*[DSM45524T]

177    (Illumina derived sequences/reads).

178

179

Figure-S2:

a.

180    **Figure S3** Bar plot showing predicted cumultative phage sequence lengths and classification as

181    indicated (see also main text). X and Y axis indicates species/strain names and length of the

182    phage sequence in kilo bases (Kb), respectively.

183

184

Figure-S3:

185 **Figure S4** (a) Amino-acid percentage identity plots for the different MCAC-members shown in
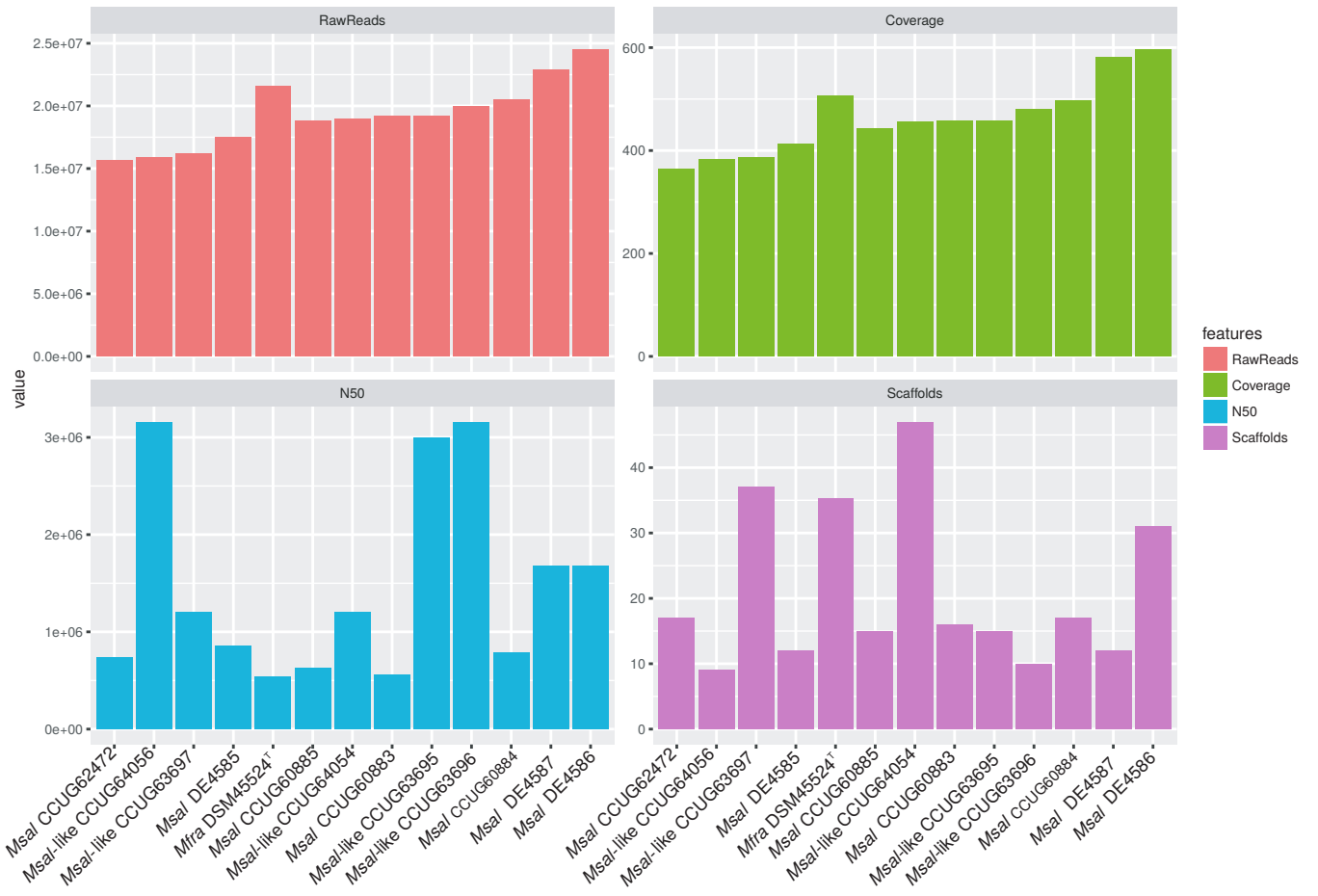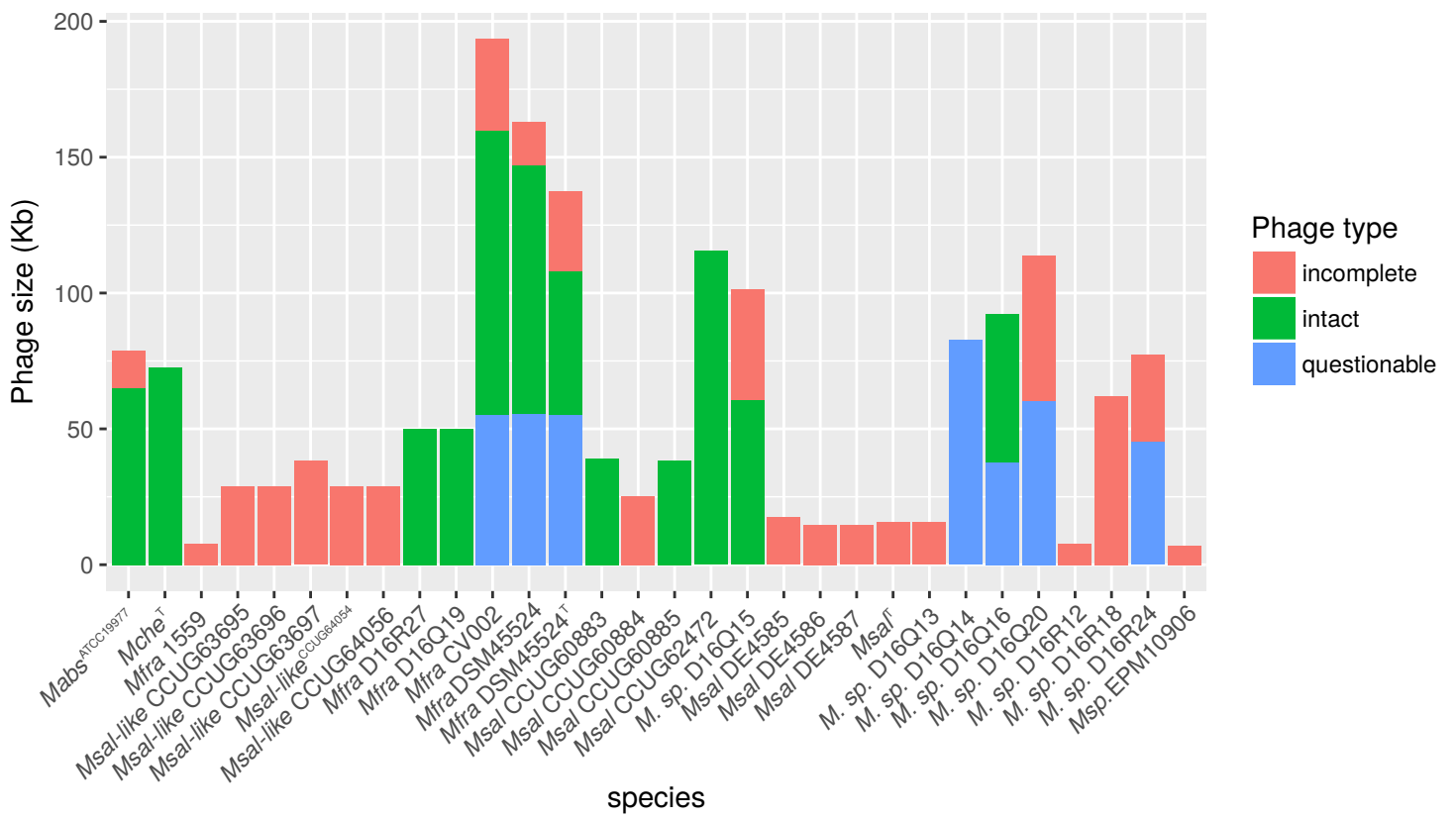
186 Fig 3a (see main text) and as indicated. Y-axis (count) refers to the number of genes while the

187 X-axis gives the percentage identity (PID).

188 (b) Alignment of the *ileS* sequences upstream of the predicted translational start sites (marked

189 ATG and GTG codons) in MCAC-members, *Mtb*H37Rv and *Msmeg*MC$^2$155 as indicated. Red

190 and blue residues mark the "T-box" signatures (see main text for references).

191 (c) A secondary structure model of the *Mabs*[ATCC19977] T-box using the *Msmeg*MC$^2$155 T-box as

192 template (see Ref 23 main text). The highligthed boxes (dashed lines) marks K-turn and putative

193 S-turn motifs while • mark conserved residues in the K-turn, S-turn and T-box. The inset

194 highlight the S-turn region in *Msmeg*MC$^2$155 and the putative S-turn structure in *Msal*[T].

195

196

Figure-S4:

a.



CDS percentage identitycoverage plot

Figure-S4:

b.

Figure-S4:

C.



$Mabs^{ATTC19977}$

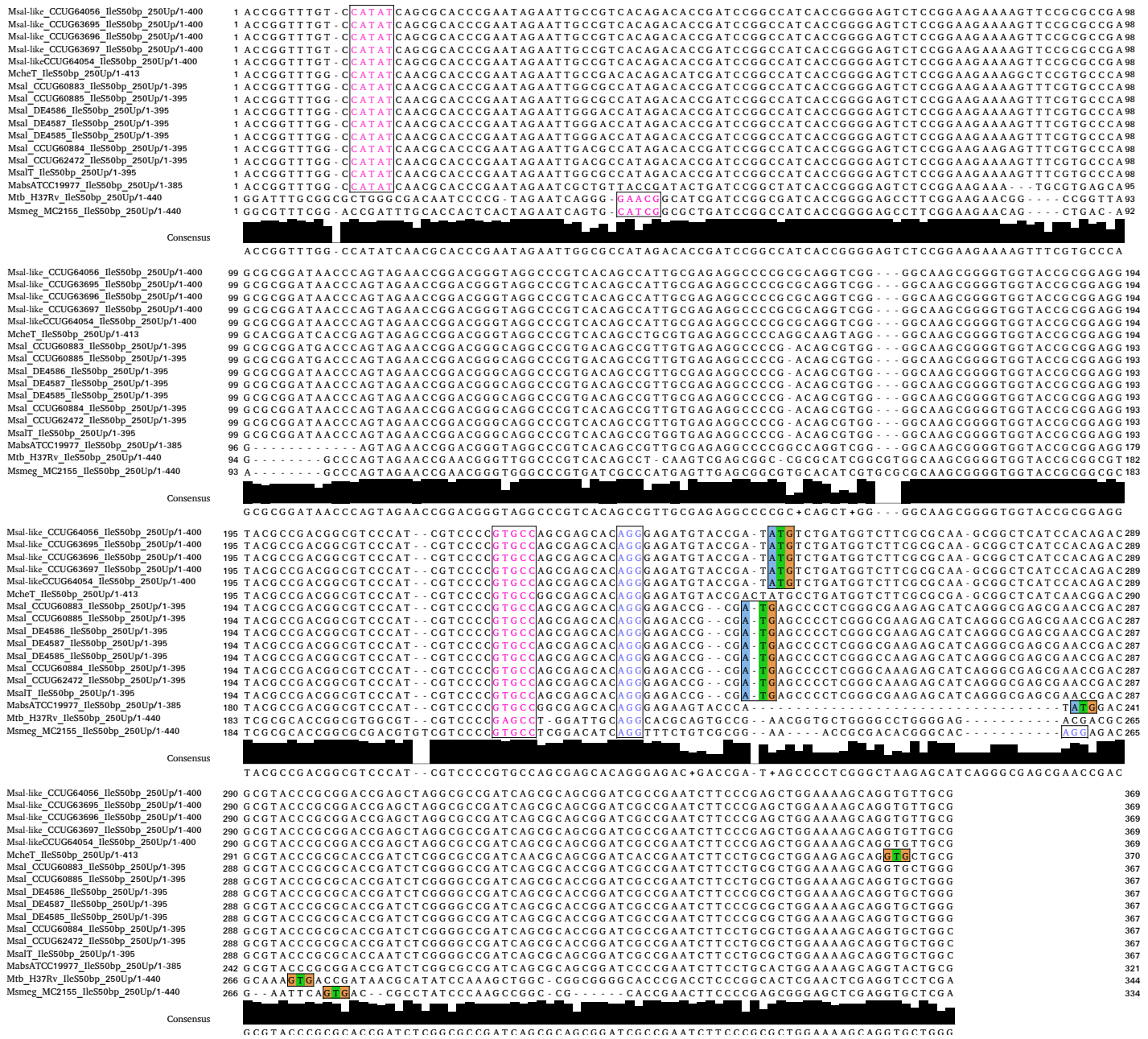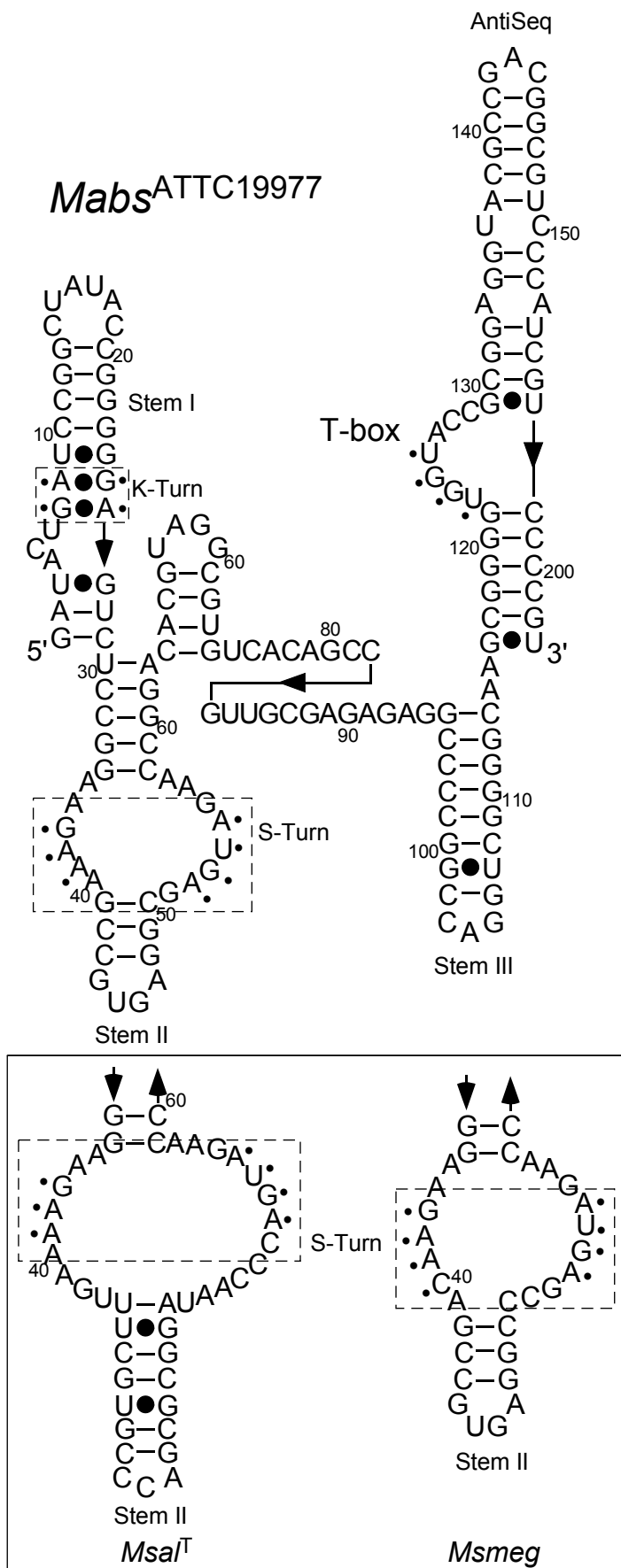197 **Figure S5** Phylogenetic relationship of MCAC-members.

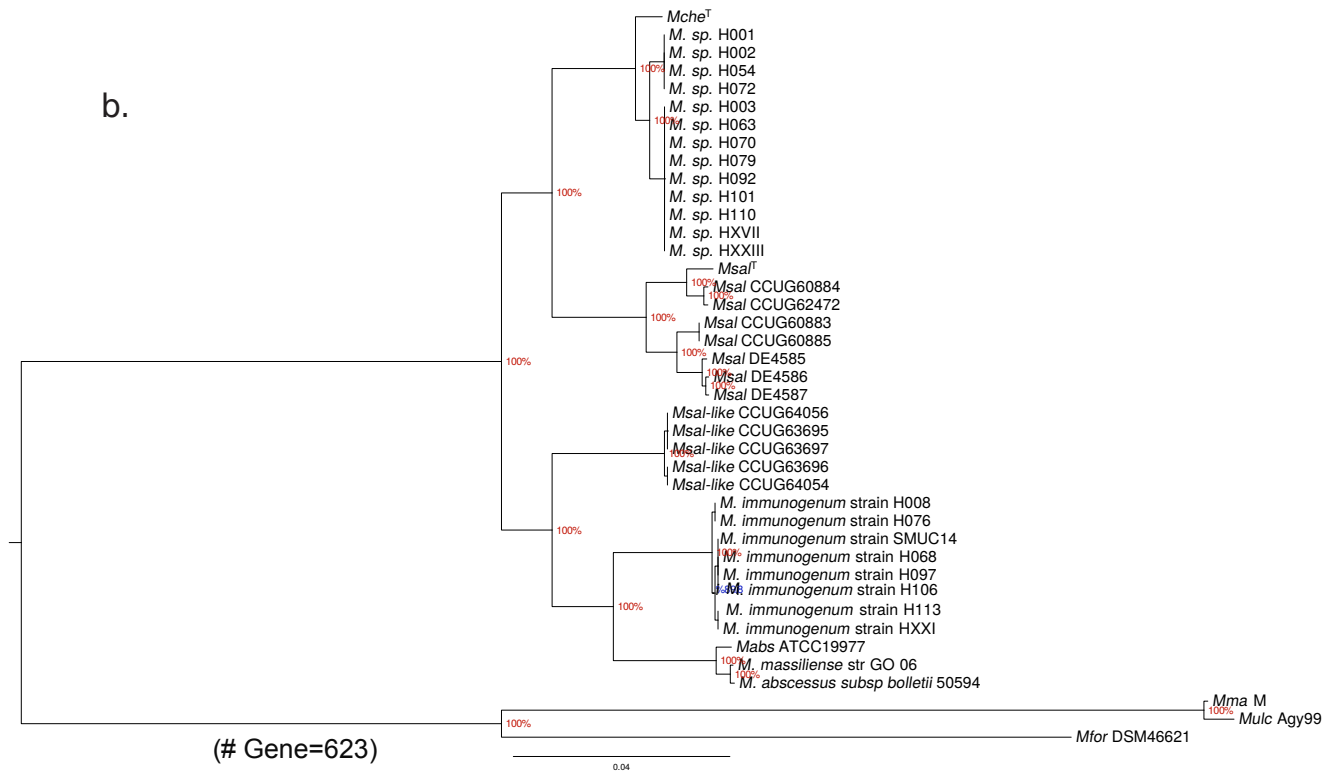198 (a) Phylogenetic tree based on 16S rDNA for MCAC-members as indicated.

199 (b) Core gene phylogenetic tree (n=623) for MCAC-members as indicated.

200 The phylogenetic trees were generated as described in Methods using *Mfor*$^{DSM46621}$, *Mulc*$^{Agy99}$

201 and *Mma* M strain as outgroups.

202

203

Figure-S5:

a.

(16S rDNA)

b.

(# Gene=623)

204    **Figure S6** Functional classification of genes in *Mabs*<sup>ATCC19977</sup>, *Mche*<sup>T</sup>*, Msal*<sup>T</sup> and *Msal-*

205    like<sup>CCUG64054</sup> into subsystem as indicated.

206    (a) Subsystem classification of genes predicted to be present in *Mabs*<sup>ATCC19977</sup> (2796 genes),

207    *Mche*<sup>T</sup> (2552 genes), *Msal*<sup>T</sup> (2544 genes) and *Msal*-like<sup>CCUG64054</sup> (2756 genes).

208    (b) Subsystem classification of unique genes in *Mabs*<sup>ATCC19977</sup> (299 genes), *Mche*<sup>T</sup> (130 genes),

209    *Msal*<sup>T</sup> (134 genes) and *Msal*-like<sup>CCUG64054</sup> (247 genes).

210    (c) Classification of unique genes present in *Mabs*<sup>ATCC19977</sup>, *Mche*<sup>T</sup>, *Msal*<sup>T</sup> and *Msal*-like<sup>CCUG64054</sup>

211    in the subcategory "Amino Acids and Derivatives".

212    (d) Classification of unique genes present in *Mabs*<sup>ATCC19977</sup>, *Mche*<sup>T</sup>, *Msal*<sup>T</sup> and *Msal-*

213    like<sup>CCUG64054</sup> in the subcategory "Carbohydrates".

214    (e) Classification of unique genes present in *Mabs*<sup>ATCC19977</sup>, *Mche*<sup>T</sup>, *Msal*<sup>T</sup> and *Msal*-like<sup>CCUG64054</sup>

215    in the subcategory "Fatty Acids, Lipids and Isoprenoids".

216    For (a) – (e), one gene can be classified in more than one subsystem.

217    (f) Shewart control chart showing the average SNVs frequencies in *Msal* strains (n = 8). Black

218    and red dots mark in-control SNV and out of control (hotspots) frequencies, respectively.

219

220

# Figure-S6:



**a.** Functional classification of *Mycobacterium spp.*

**b.** Subsystem classification of unique genes

Figure-S6:

c.

**Unique genes – subcategory of Amino Acids and Derivates**

d.

**Unique genes – subcategory of Carbohydrates**

e.

**Unique genes – subcategory of Fatty Acids, Lipids, and Isoprenoids**

Figure-S6:

f.



Number of groups = 2357
Center = 97.41699
StdDev = 22.87435
LCL = 0
UCL = 166.04
Number beyond limits = 69
Number violating runs = 258

221 **Figure S7** Horizontal gene transfer analysis in *Msal* and *Msal*-like strains, *Mche*[T], *Mabs*[ATCC19977]

222 and *Mfra*[DSM45524T].

223 (a) Bar plot showing number of predicted horizontally transferred genes. Y and X axis represent

224 mycobacterial strains/species and predicted number of genes, respectively.

225 (b) Venn diagram showing common and predicted unique horizontally transferred genes in

226 *Mabs*[ATCC19977], *Mche*[T], *Msal*[T] and *Msal*-like[CCUG64054].

227 (c) Heat map showing the probable source of the HGT genes for *Msal* and *Msal*-like strains,

228 *Mche*[T], *Mabs*[ATCC19977] and *Mfra*[DSM45524T]. The vertical tree represents the heat map clustering of

229 the column wise dendogram. Color code, see top left corner of the plot.

230

231

Figure-S7

a.



b.



c.

232 **Figure S8** Heat map showing distribution of virulence genes in $Mabs^{ATCC19977}$, $Mche^{T}$, $Msal^{T}$

233 and $Msal$-like$^{CCUG64054}$ and other mycobacteria as indicated. The vertical tree represents the heat

234 map clustering of the column wise dendogram. Green = present and gray = absent.

235

236

Figure-S8:

237    **Figure S9** Compilation of tRNA genes predicted in MCAC-members.

238    (a) Heat map showing presence (green) and absence (light grey) of tRNA genes in MCAC-

239    members indicated below. The clustered tRNA gene names with mycobacteria strain/species and

240    tRNA isoacceptor name, *e.g.*, *Msal*-DSM43276_7tRNA_Ser_CGA, are listed on the right. The

241    horizontal and vertical trees represent the heat map clustering of the column and row wise

242    dendograms.

243    (b) *Mche*$^T$ *vs Msal*$^T$. Blue and red marked tRNA genes refers to transcription from the positive

244    and negative strands, respectively. Blue lines mark that the locations of the tRNA genes have not

245    shifted.

246    (c) *Mabs*$^{ATCC19977}$ *vs Msal*-like$^{CCUG64054}$. Blue and red marked tRNA genes, and blue lines as in

247    (b; see above), while red lines mark tRNA genes that have shifted position on the chromosome.

248    Of note, the *Msal*-like$^{CCUG64054}$ is a draft genome while *Mabs*$^{ATCC19977}$ is a complete genome.

249    (d) *Mabs*$^{ATCC19977}$ *vs Mche*$^T$. Blue and red marked tRNA genes, and blue lines as in (b; see

250    above).

251    (e) Analysis of the gene synteny for a tRNA gene cluster encompassing nine genes in *Msal* and

252    *Msal*-like strains, and *Mche*$^T$ and *Mabs*$^{ATCC19977}$ as indicated. The tRNA genes are marked in red

253    and the vertical boxes marked in brown highlight homologous genes. For further details see

254    main text and Figs 6c, S1a and S9a.

255    (f) Sequence alignments of the common and "extra" tRNA genes as indicated. With respect to

256    tRNA$^{Leu}$CAG the arrows mark residues forming the amino acid acceptor-stem, D-stem,

257    anticodon-stem and T-stem. For details see the main text.

258

259

Figure-S9

a)

Figure-S9

b)

$Mche^T$

$Msal^T$

Figure-S9

c)

Figure-S9

d)

Figure-S9

e)

# Figure-S9 f)

**Thr_TGT**

```
                    **  ********  **  *  ****  **   *******  *    *  *  ****  *  *   *  **  ***
MsalT_Reg_11tRNA_Thr_TGT     -CCTCCTTAGCTCAGTGGTA-GAGCACCGCTCTTGTAAAGCGAAGGTCGTCAGTTCAATCCTGACAGGGGGCTCAA   74
MsalT_Extra_47tRNA_Pseudo_TGT GCCGAATTAGCTCAACGGGAAGAGCGCCTGCCTTGTAAA-CAGGTGGAGCGGGTTCGAGTCCTGCCGGAGGC----   71
                    1.......10........20........30........40........50........60........70......
```

**Arg_ACG**

```
                    *  ***************  ****   *  *  ****  *  *  **  *   *   ******  **     ****  *  ***
MsalT_Reg_6tRNA_Arg_ACG   CGCCCGTAGCTCAACGGATAGAGCATCTGACTACGGATCAGAAGGTTAGGGGTTCGAATCCCTTCGGGCGCACCA   75
MsalT_Extra_48tRNA_Arg_ACG  CCCCCGTAGCTCAACGGACAGAGTGACAGCCTACGGAGCTGAGGAT-GGAAGTTCGATTCTTCCCGGGGG-ACC-   72
                    1.......10........20........30........40........50........60........70......
```
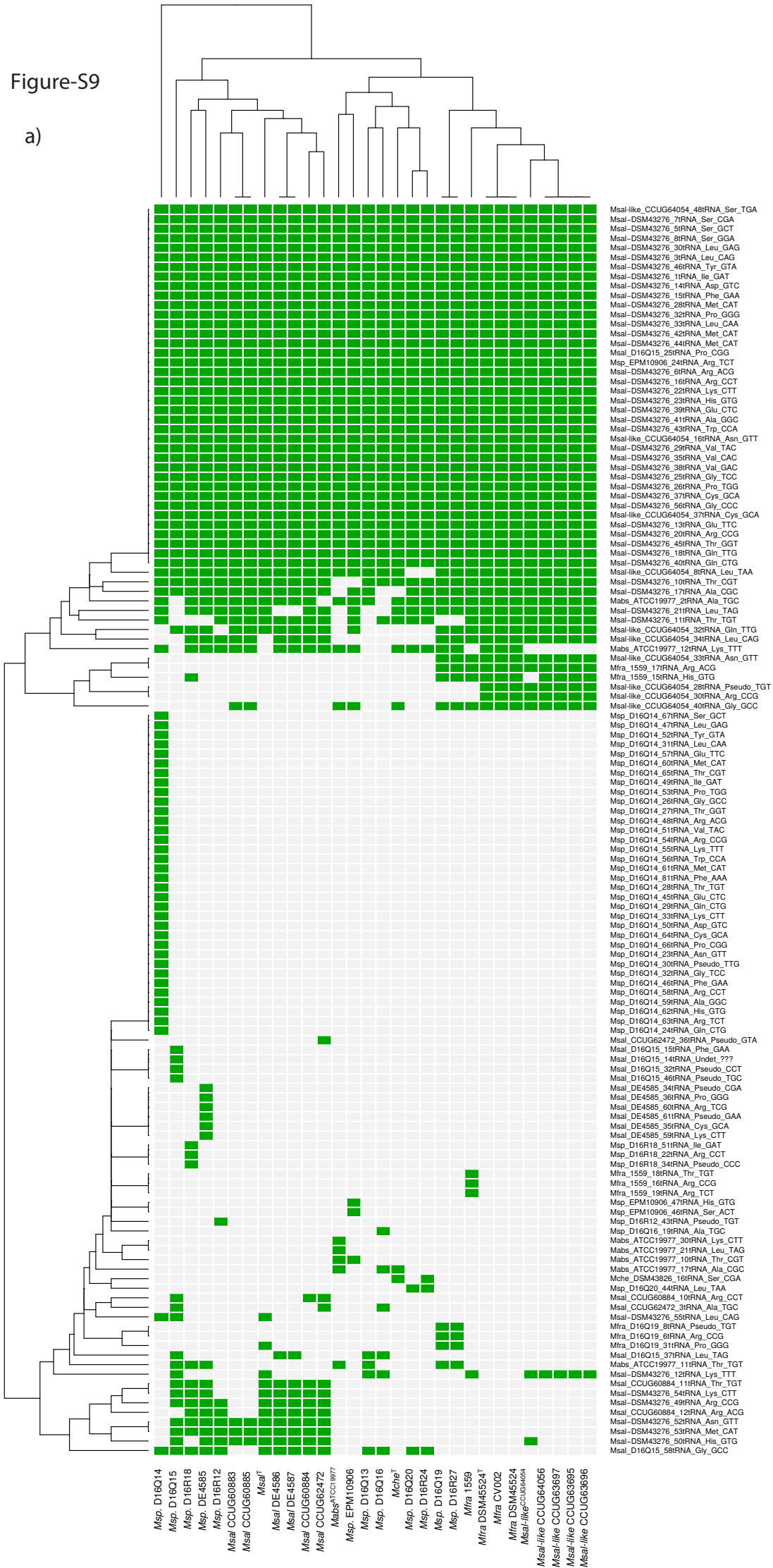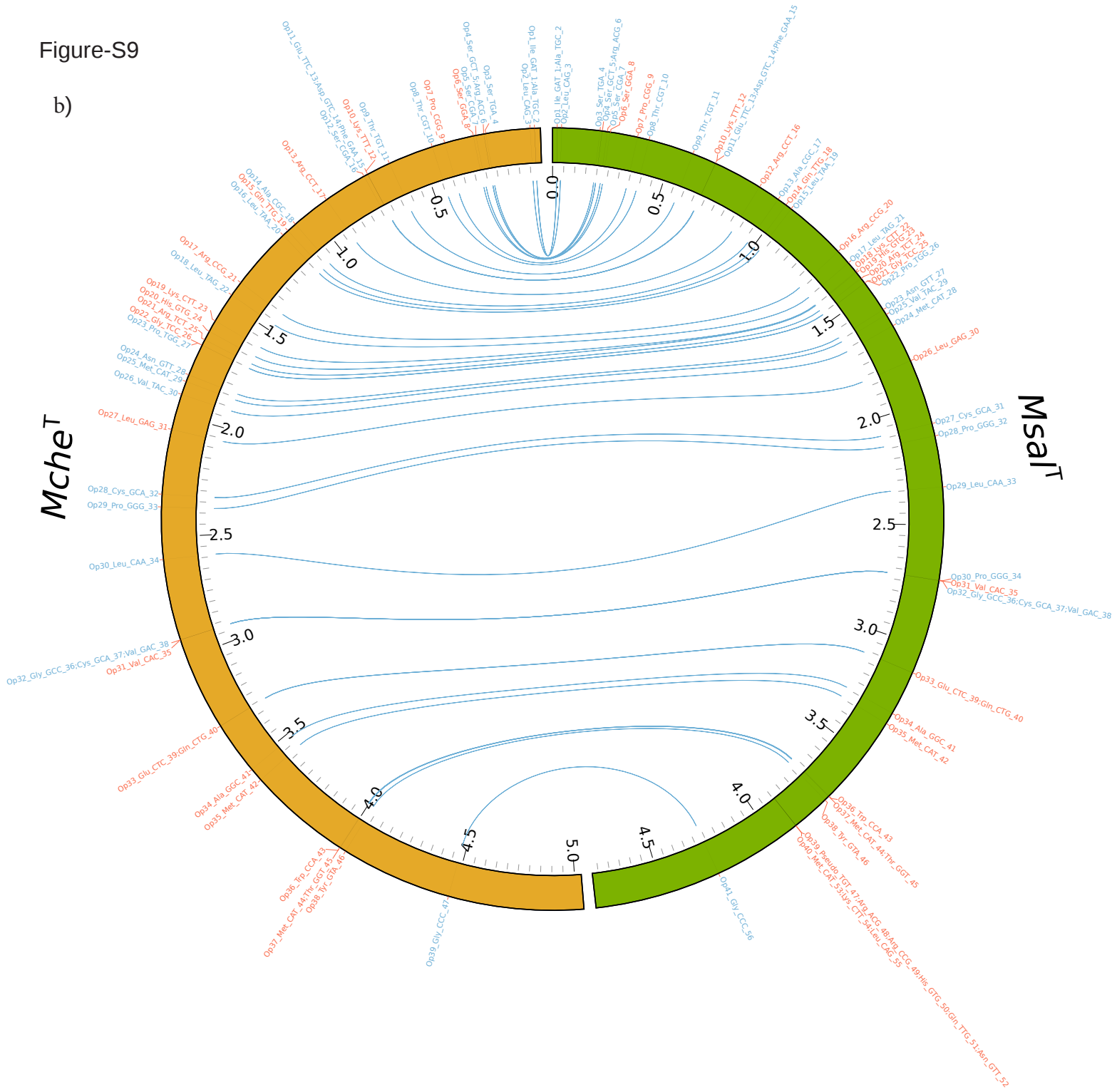
**Arg_CCG**

```
                    ***  ******  **   ******  **   *  *   ******  *  *  ***   **  *******    **  ****  ***
MsalT_Reg_20tRNA_Arg_CCG   GCCCTCGTAGCTCAGGGGATAGAGCGTCTGCCTCCGGAGCAGAAGGCCGCAGGTTCGATTCCTGCCGAGGGC---   72
MsalT_Extra_49tRNA_Arg_CCG  ACCCCGTAGCGCAATGGATAGCGCACCGGTTTCCGGAACCGGAGGTTGCGGGTTCGAACCCCGCCGGGGGTGCC   75
                    1.......10........20........30........40........50........60........70......
```

**His_GTG**

```
                    *  ***  **  **   *    **  *    **  *   ****   *  *  *  ****  ******  *  ****  **  *  *  ****
MsalT_Reg_23tRNA_His_GTG   GTGGCTGTAGTTCAGTTGGTAGAGCACCAGGTTGTGATCCTGGCTGTCGCGGGTTCGAGTCCCGTCAGCCACCCC   75
MsalT_Extra_50tRNA_His_GTG  GCGGCAGTGGTGTAAC-GGAAACACAGCGGTCTGTGACTCCG-CAGTCGAGGGTTCAACTCCCTTCGGTCGCCCC   73
                    1.......10........20........30........40........50........60........70......
```

**Gln_TTG**

```
                    *  *  *  ******  *******  ***   *    *  **  ****    **  *   *******************  **  *  *
MsalT_Reg_18tRNA_Gln_TTG   TCCGTCGTGGTGTAATCGGCAGCACCTCTGATTTTGGTTCAGATAGTTCAGGTTCGAGTCCTGGCGACGGA   71
MsalT_Extra_51tRNA_Gln_TTG  TGCCTCGTGGGGTAATCGGTAGCCCACCGGACTTTGAATCCGGGAGTTCAGGTTCGAGTCCTGGTGAGGCA   71
                    1.......10........20........30........40........50........60........70.
```

**Asn_GTT**

```
                    ****  **  *  *  **   **  *  *  *   *  **  ********  *  **   *    *******  ***    *****
MsalT_Reg_27tRNA_Asn_GTT   -CCCCTGTAGCTCAATTGGCAGAGCTTCCGACTGTTAATCGGACGGTTCTTGGTTCGAGTCCAAGCGGGGGAG   72
MsalT_Extra_52tRNA_Asn_GTT  TCCCCGGTTGGTTAACCGGTAAACCAGCGGATTGTTAATCCG-CGACTGCAGGTTCGAATCCTGCTCGGGGA-   71
                    1.......10........20........30........40........50........60........70...
```

**Met_CAT**

```
                    *  *  **  *  ***    *  ***  *    *  *****  *  *    ****    *  **  *  *    *   *  *  **
MsalT_Reg_42tRNA_Met_CAT    CGCGGGGTGGAGCAGCTCGGTAGCTCGCTGGGCTCATAACCCAGAGGTCGCAGGTTCAAATCCTGTCCCCGCTACT-   76
MsalT_Extra_53tRNA_Met_CAT  TGCGCGGTAGGACAGTCCGGTAGTCCAGCGGGCTCATGACCCGAAAGTCGCAGGTTCGAATCCTGCCCGCGCGACC-   76
MsalT_Reg_28tRNA_Met_CAT    -GGGCGGTAGCTCAGTTGGTTAGAGCCGTGGACTCATAATCCATTGGTCGCGGGTTCGAGTCCCGCCCGCCCCTACAA   76
MsalT_Reg_44tRNA_Met_CAT    GGCGGGTGTAGCTCAGTGGTTAGAGCGCACGACTCATAATCGTGAGGTCGGGGATCGAGCTCCCCACCGCTACC-   76
                    1.......10........20........30........40........50........60........70......
```

**Lys_CTT**

```
                    ************  ***  **  *  ****  *  **  *******  *  ***  *  ****  ****   **   ******  **
MsalT_Reg_22tRNA_Lys_CTT   GCGCCGTTAGCTCAGTTGGTAGAGCAGCTGACTCTTAATCAGCGGGTCCGGGGTTCGAGCCCCTGACGGCGCACC   75
MsalT_Extra_54tRNA_Lys_CTT  GCGCCGTTAGCTGAGTAGGGATAGCACCGGATTCTTAATCCG-GGGACGGGGGATCGAAACCTCCACGGCGTACG   74
                    1.......10........20........30........40........50........60........70.....
```

**Leu_CAG**

```
MsalT_Reg_3tRNA_Leu_CAG    GGCGAGTGGCGGAATGGCaGACGCGCTGGCTTCAGGTGCCAGTGTCCTTCGGGACGTGGGGGTTCAAGTCCCCCTTCGCccaCA
                    >>>>>>..>>>..........<<<.>>>>>.......<<<<<.>>>>....<<<<..>>>>>.......<<<<<<<<<<...<.

MsalT_Extra_55tRNA_Leu_CAG  GCCCCTCTGGCCCAACTGGAaGAGGCGTCCCGTTCAGGGCGGGAAGGtTCCTGGTTCGAATCCAGGGAGGGGTA
                    >>>>>>..>>>..........<<<.>>>>>.......<<<<<.....>>>>>.......<<<<<<<<<<<<.
```
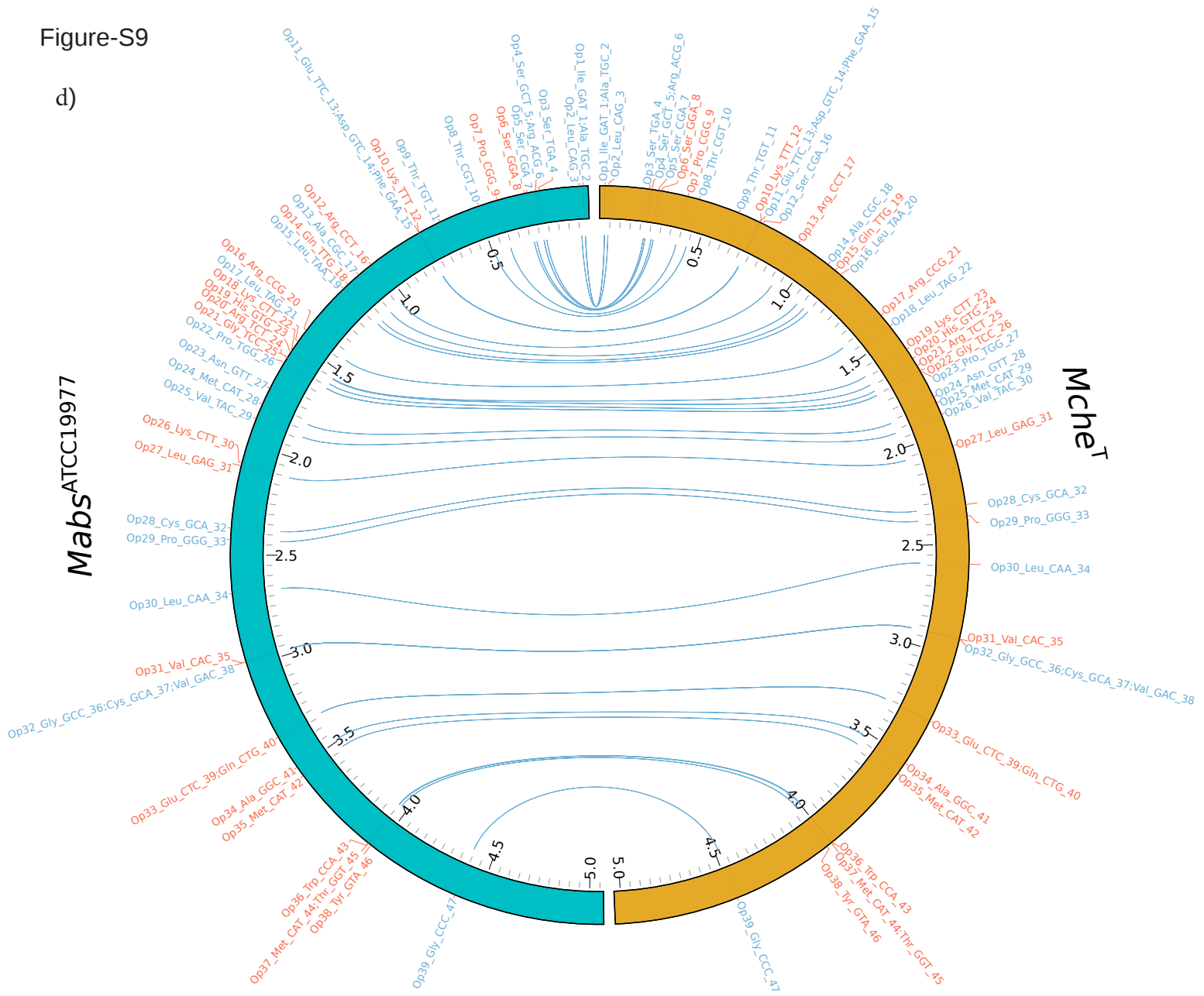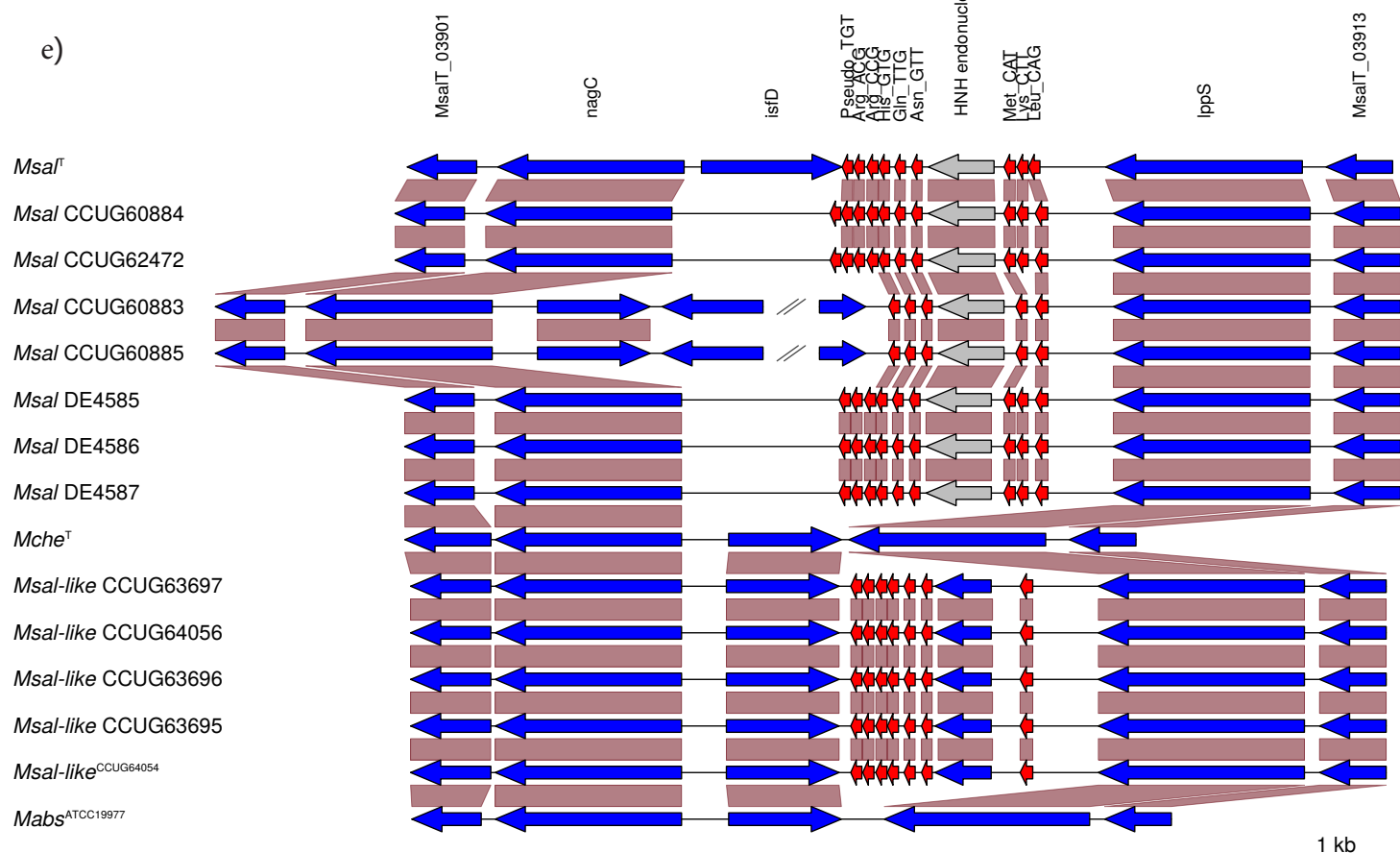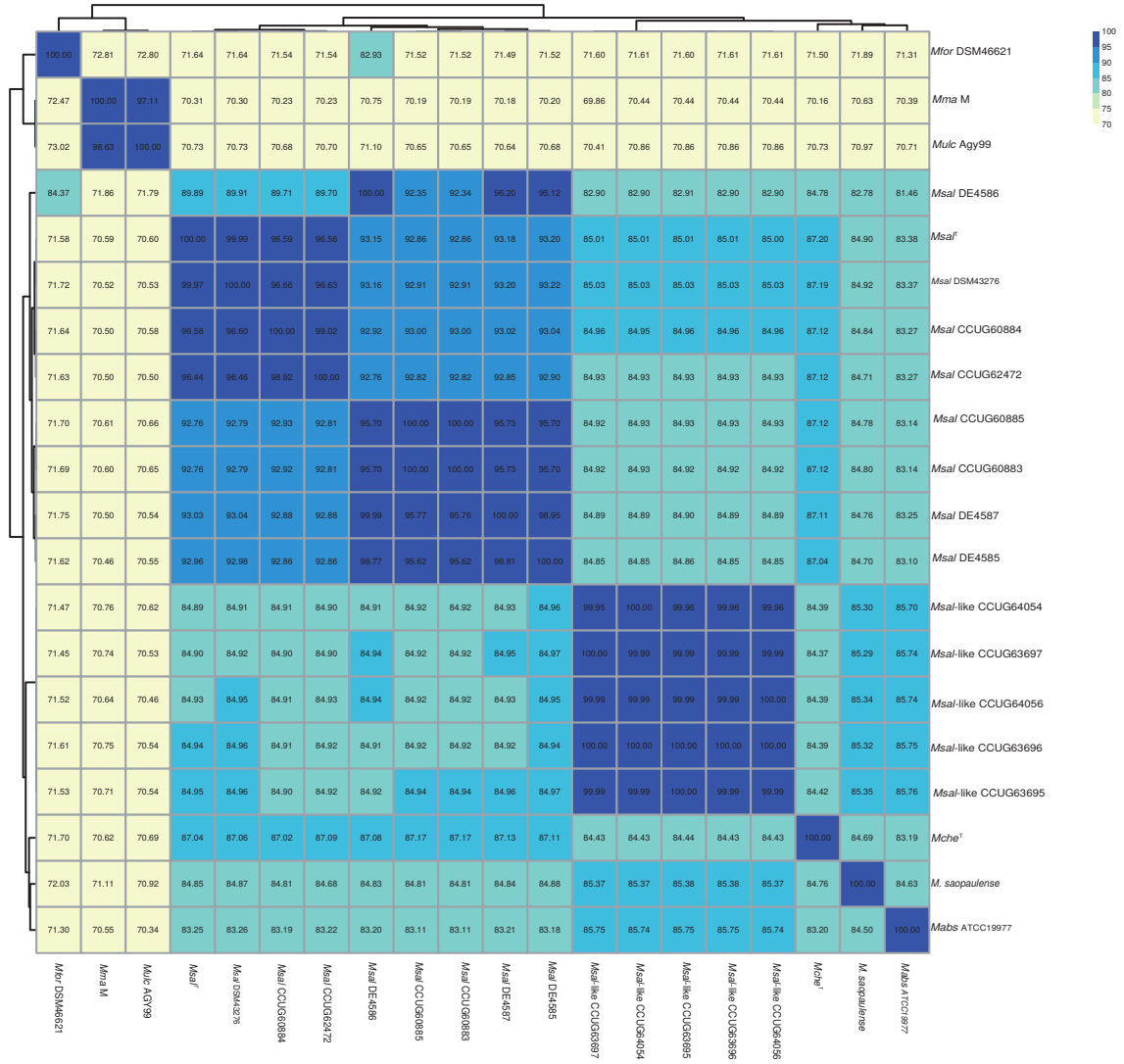
260 **Figure S10** Extended ANI analysis including *Mycobacterium saopaulense*.

261 (a) Heat map showing ANI values for "all-versus-all" *Msal* and *Msal*-like strains, *Mabs*<sup>ATCC19977</sup>,

262 *Mche*<sup>T</sup>, *M. saopaulense*, *Mma* M strain, *Mulc*<sup>Agy99</sup> and *Mfor*<sup>DSM46621</sup> as indicated. ANI values

263 were clustered based on unsupervised hierarchical clustering (see Methods, main text and Fig 2).

264 (b) Dendogram, extracted from the heat map shown in (a), displaying clustering of different

265 strains / based on ANI values.

266

Figure-S10:

a.



b.