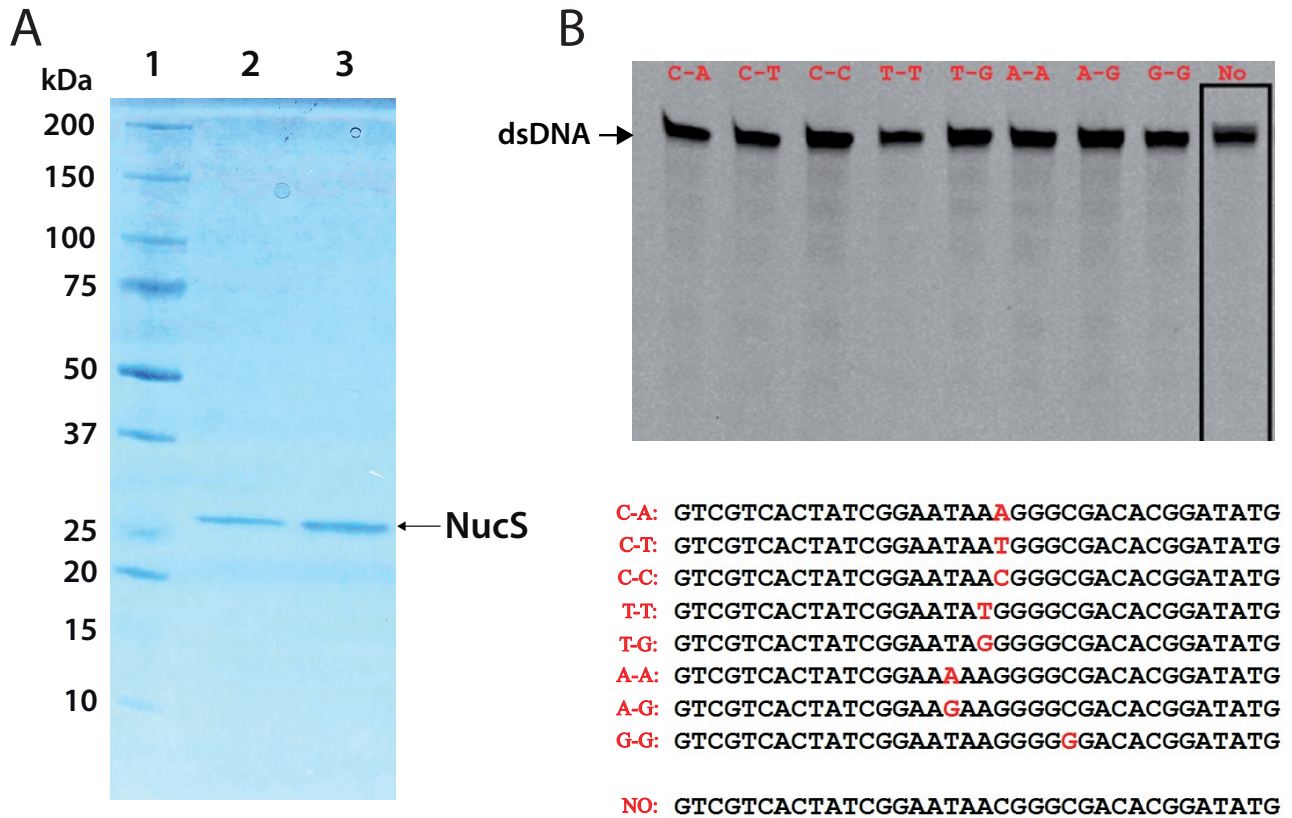


## Supplementary Information

### Supplementary Figures.



**Supplementary Figure 1. Purified mycobacterial NucS does not cleave mismatch DNA substrates in vitro.** A) Native *M. smegmatis* NucS protein purification. Recombinant NucS protein was expressed in *E. coli*, purified and concentrated, as described (see Methods). Native NucS protein was analysed in a SDS-PAGE gel and stained with Coomassie Brilliant Blue. Lane 1, molecular marker; Lane 2: purified native *M. smegmatis* NucS; Lane 3: concentrated *M. smegmatis* NucS protein (used for EMSA and nuclease assays). B) Double stranded DNA (30 nM) with or without the indicated internal mismatches, was incubated with 300 nM NucS. Sequence consensus is shown at the bottom and mismatch positions are indicated in red. No significant cleavage activity was observed with the shown 36-mers or with longer 75-mer (not displayed here) mismatch substrates.

**A**

```

100% t g a t c a a g c t g t t c g g c g a g c a c t g g t g c g g t c c g g a g a g c c t c g c g t c g g a g t c g g a g g c g t a c g c g g t c c t g 74
95% t g a t c a a g c t g t t c g g G g a g c a c t g g t g c g g t c c c g a g a g c c t c g c g t c g g a g t c G g a g g c g t a c g c g g t c c t G 74
90% t g a t c a a g c t G t t c g g G g a A c a c t g g t g c g g C c c G a g a g c c t G c g t c g g a g t c G g a g g c g t a c g c G t c c t G 74
85% t g a t c a a g c t C t t c g g G g a A c a c t g g t g c g g C c c G a A a g c c t G c g t c G g a g t c G g a g g c G t a c g c G t c c t G 74

100% g c g g a c g c c c c g g t g c c g g t g c c c c g c c t c c t c g g c g c g g c g a g c t g c g g c c c g g c a c c g g a g c c t g g c c g t g 148
95% g c g g a c g c c c c g g t g c c g g t C c c c c g c c t c c t c g g c g c g G g g c g a g c t g c g g c c c g g c a c c G g a g c c t g g c c g t g 148
90% g c g g a c g c c c c G g t g c c g g t C c c c c g c c t G c t c g g c g G g g c g a g c t g c g g c G g g c a c c G g a g c c t g g c c G t g 148
85% g c g g a c g c c c c G c t g t c c g g t C c c c c g c c t G c t c g g c g G g g c g a g c t C c g g c G g g c a c G g a g c G t g g c c G t g 148

100% g c c c t a c c t g g t g a t g a g c c g g a t g a c c g g c a c c a c c t g g c g g t c c g c g a t g g a c g g c a c g a c c g a c c g g a a c g 222
95% g c c G t a c c t g g t g a t g a g c c G C a t g a c c g g c a c c a c c t g g c g C t c c g c g a t g g a c g g c a c g a c G a c c g g a a c g 222
90% g c c G t a c c t g g t C a t g a g c c G C a t g a c c g g G a c c a c c t g g c g C t c c g c g a t g g a c g g G a c g a c G a c c g g a a c g 222
85% g c c G t a c c t g g t C a t g a g c c G C a t g a c c g g G a c c a c c G t g g c g C t c c g C a t g g a c g g G a c g a c G a c c g C a a c g 222

100% c g c t g c t c g c c c t g g c c c g c g a a c t c g g c c g g g t g t c t g g c c g g c t g c a c a g g g t g c c g c t g a c c g g g a a c a c c 296
95% c g c t g c t c g c c c t C c c c c g c g a a c t c g g c c g g g t g t G g g c c g g c t g c a c a g g g t g c c g c t C a c c g g g a a c a c c 296
90% c C t g c t c g c c c t C c c c c g c g a a c t G g g c c g g t g t C t G g g c c g g c t C a c a g g g t g c c g c t C a c c g g a a c a c c 296
85% c C t g c t G g c c c t C g c c c g c g a G c t G g g c c g G t g t G g g c c g g c t C a c a g g g t C c g c t C a c c g g C a a c a c G 296

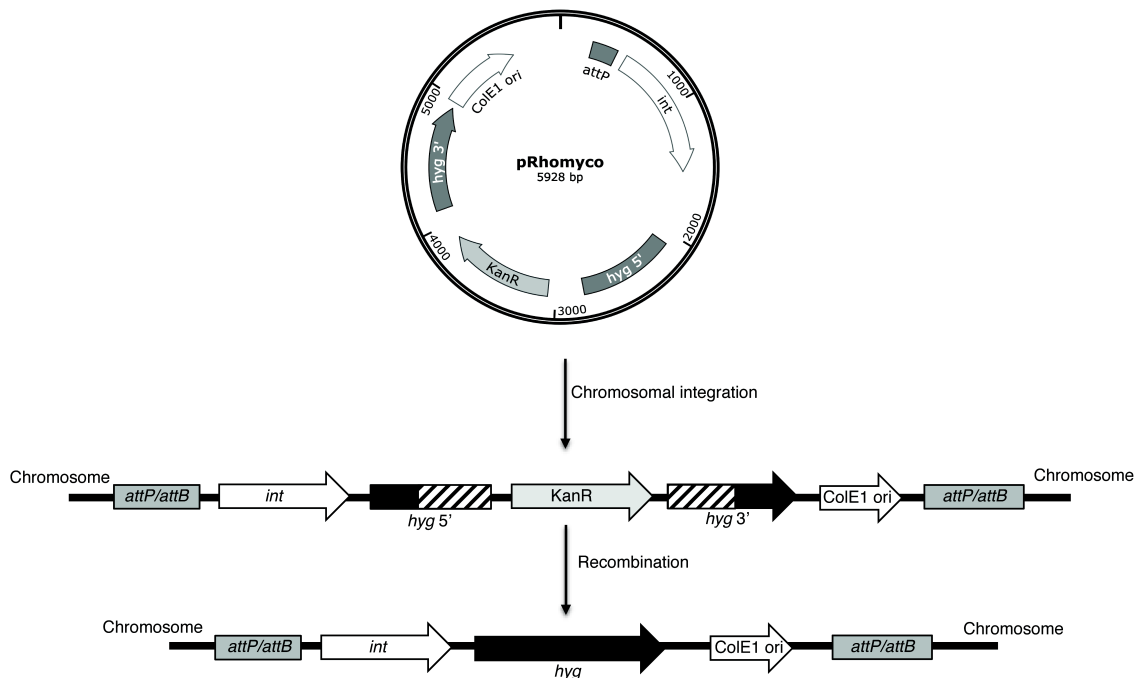
100% g t g c t c a c c c c c a t t c c g a g g t c t t c c c g g a a c t g c t g c g g g a a c g c c g c g c g g c g a c c g t c g a g g a c c a c c g 370
95% g t g c t c a c c c c G c a t t c c g a g g t c t t c c c g g a a c t C t g c g g g a a c g c c g c g c g g c g a c G g t c g a g g a c c a c c g 370
90% g t g c t c a c c c c G c a t t c c g a g g t G t t c c c g g a a c t C t g c g g g a a c g G c g c g c g g c g a c G t c g a g g a c c a c c g 370
85% g t g c t G a c c c c G c a t t c G a g g t G t t c c c G a a c t C t g c g G a a c g G c g c g C g c g a c G t c g a G a c c a c c g 370

100% c g g g t g g g g c t a c c t c t c g c c c c g g c t g c t g g a c c g c c t g g a g g a c t g g c t g c c g g a c g t g g a c a c g c t g c t g g 444
95% c g g g t g g g g t a c c t c t c g c c c c g g c t g c t G a a c c c c t g g a g g a c t g g c t g c c c g a c g t g g a c a c g c t g c t G 444
90% G g g g t g g g g t a c c t c t c g c c c c g g c t g c t G a a c c c c t g g a g g a c t g g c t g c c c g a c g t g g a c a c G c t g c t G 444
85% G g g g t g g g g t a c c t G t c g c c G c g g c t g c t G a a c c c c t G a g g a c t g g c t C c c G a c g t G a c a c C t g c t G 444

100% c c g g c g c g a a c c c c g g t t c g t c c a c g g c g a c c t g c a c g g g a c c a a c a t c t t c g t g g a c c t g g c c g c g a c c g a 517
95% c c g g c g c g a a c c c c g C t t c g t c c a c g g c g a c c t C a c g g g a c c a a c a t c t t c g t G a c c t g g c c g c g a c G g a 517
90% c c g g c g G g a a c c c c G t t c g t c c a c g g G g a c c t C a c g g g a c G a a c a t c t t c g t G a c c t g g c G g c g a c G g a 517
85% c c g g c g G g a a c c c c G t t c g t G a c g g G g a c c t C a c g g C a c G a a c a t c t t c g t G a c c t G g c G g c g a c G g a 517

```

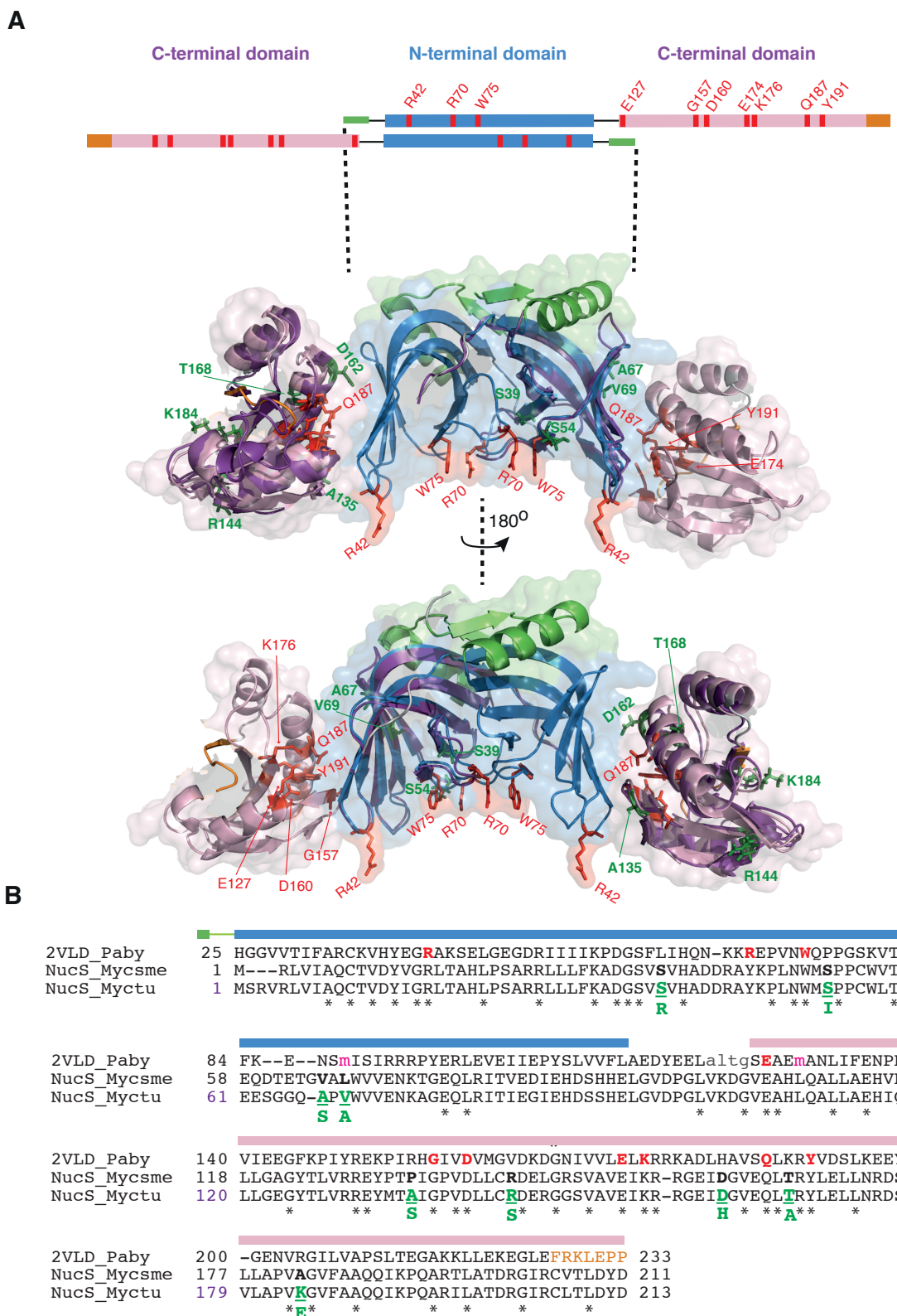
**B**



**Supplementary Figure 2. Tools for measuring recombination rates.**

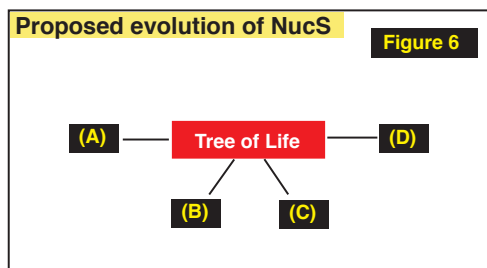
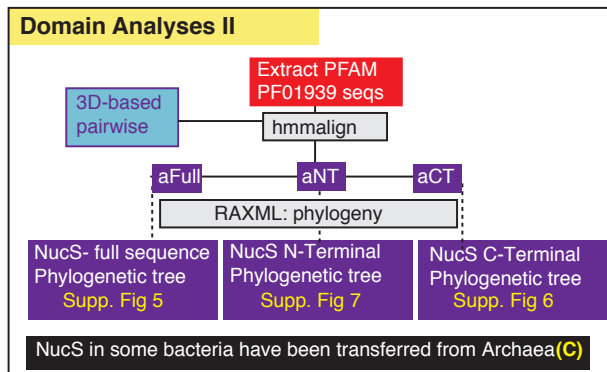
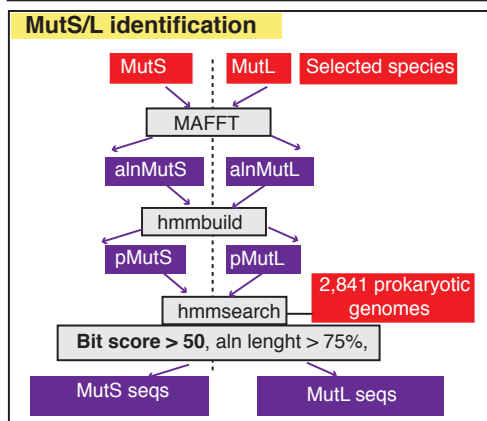
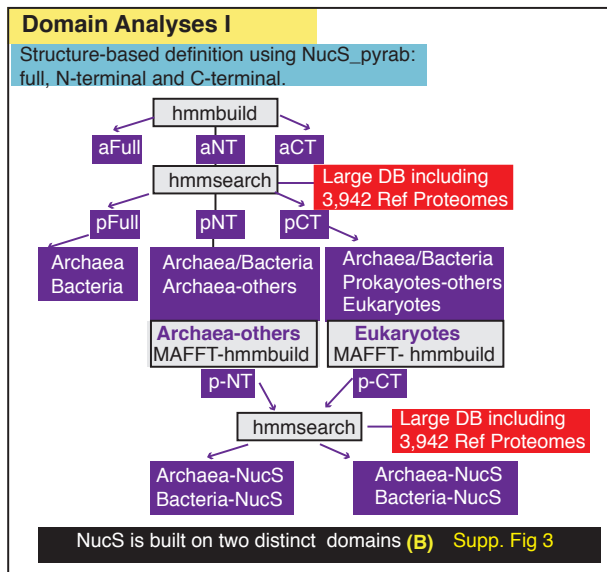
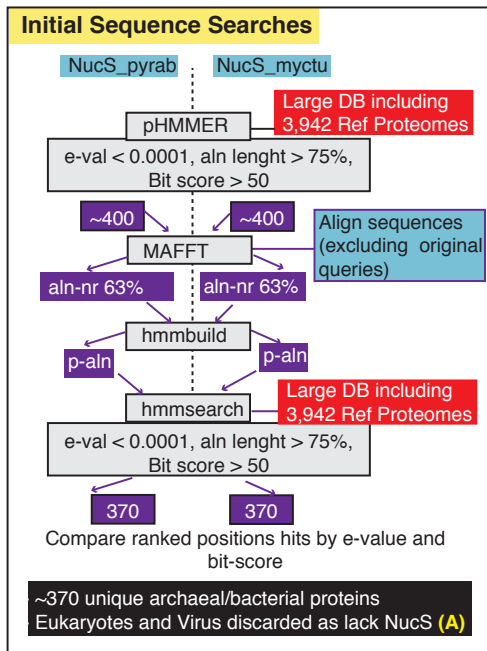
A) Alignment of the 517 bp overlapping sequences, with 100%, 95%, 90% and 85% identities, shared by *hyg 5'* and *hyg 3'* regions, used to construct the different pRhomyco versions. Changes

introduced in the original sequence are highlighted. B) Cartoon of plasmid pRhomyco before and after recombination. *hyg 5'* and *hyg 3'* are truncated alleles of the *hyg* gene carrying a 3'-terminal or 5'-terminal deletion, respectively. Both alleles overlap 517 bp (striped regions) and are separated by a 1,200 bp region containing *aph3* gene (Kan-R). Four versions of pRhomyco were generated carrying *hyg 5'* alleles 100%, 95%, 90% or 85% identical, in its overlapping region, to the *hyg 3'*. Plasmids integrate in the chromosome of *M. smegmatis* mc<sup>2</sup> 155 and its  $\Delta$ *nucS* derivative. Site-specific recombination between the *attP* from the plasmid and the unique bacterial *attB* is promoted by the plasmid-encoded integrase.



**Supplementary Figure 3. Domain characterization of *M. tuberculosis* NucS and polymorphic residues.** A) The upper part is a schematic representation of the homodimeric structure of *P. abyssi* NucS (2VLD) <sup>1</sup>. The two distinct domains of NucS, N-terminal and C-terminal, are in blue and

pink, respectively. The first residues (1-25) of the NucS *P. abyssi* N-terminal domain, missing in the *M. tuberculosis* model, are depicted as green cartoon. The putative  $\beta$ -clamp binding motif in *P. abyssi* NucS is shown in orange. The important catalytic and DNA binding sites are depicted over the *P. abyssi* structure as vertical red bars with the residues numbered above. Back and front views of the structural superimposition of the homodimeric resolved structure of NucS from *P. abyssi*<sup>1</sup> and the *M. tuberculosis* model (purple) is presented below the schema (colour code as described for the upper scheme). The important catalytic and DNA binding sites are depicted as red sticks and the residues where polymorphisms described in this work have been detected are shown as green sticks over the structural imposition. B) Multiple structure-based alignment of NucS from *P. abyssi* (2VDL), *M. smegmatis* mc<sup>2</sup> 155 and *M. tuberculosis* CDC155. Colour code is the same as in panel A. The amino acid change of each naturally occurring *M. tuberculosis* polymorphism is depicted in green. Important catalytic and DNA binding residues are in red. Grey low-case indicates regions without structural information; magenta “m” indicates Seleno-Met modifications in the structure of *P. abyssi*. Residues of the putative  $\beta$ -clamp binding motif in *P. abyssi* NucS are shown in orange. Only regions that align among the three proteins are shown. Asterisks indicate identical residues in the three aligned sequences.



- External info
- Program
- Results
- Input file
- Output file

### Phylogenetic Profiling

Supplementary data file 1  
Figure 5

**Procedure:**

```

If NO NucS is identified
and "Unassembled WGS" ..... Undef NucS
else, translated searches vs genome(*)
  if NucS found ..... NucS
  if NucS not found ..... LikelyNot

if MutS/L is found ..... MutS/L
if MutS/L not found:
and "Unassembled WGS" ..... Undef MutS/L
else, translated searches vs genome(*)
  if MutS/L found ..... MutS/L
  if MutS/L not found ..... LikelyNot

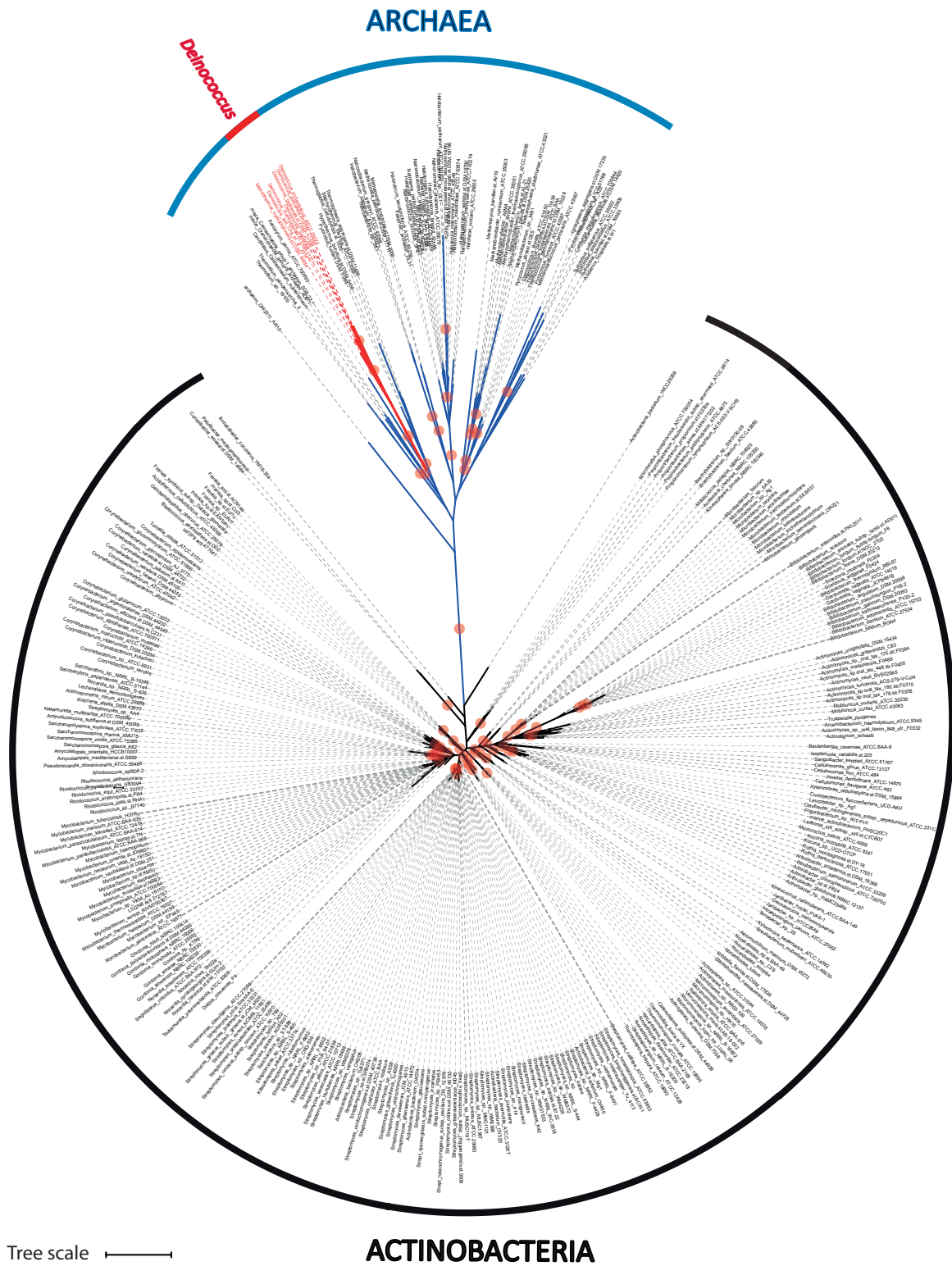
If NucS IS identified
if MutS/L is found ..... NucS-MutS/L
if MutS/L not found:
and "Unassembled WGS" ..... Undef MutS/L
else, translated searches vs genome(*)
  if MutS/L found ..... NucS-MutS/L
  if MutS/L not found ..... NucS-only**
  
```

(\*) Only complete genomes Actinobacteria phylum and Archaea  
 \*\* Confident absence

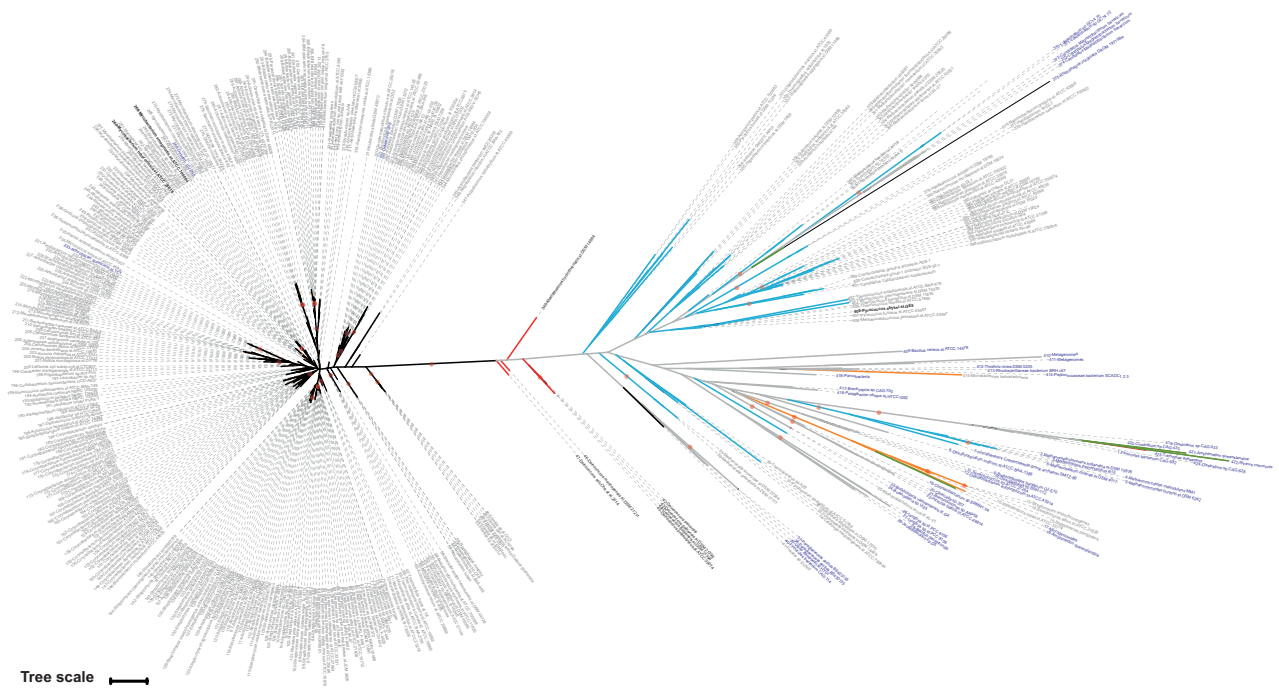
NucS shows a disperse distribution pattern (D)

Supplementary Figure 4. Computational procedures to analyse NucS distribution and evolution. Each panel depicts a different protocol. Yellow fonts indicate figures and/or tables. Yellow uppercase letters are pieces of evidence used in the model presented in Fig. 6 (main text).



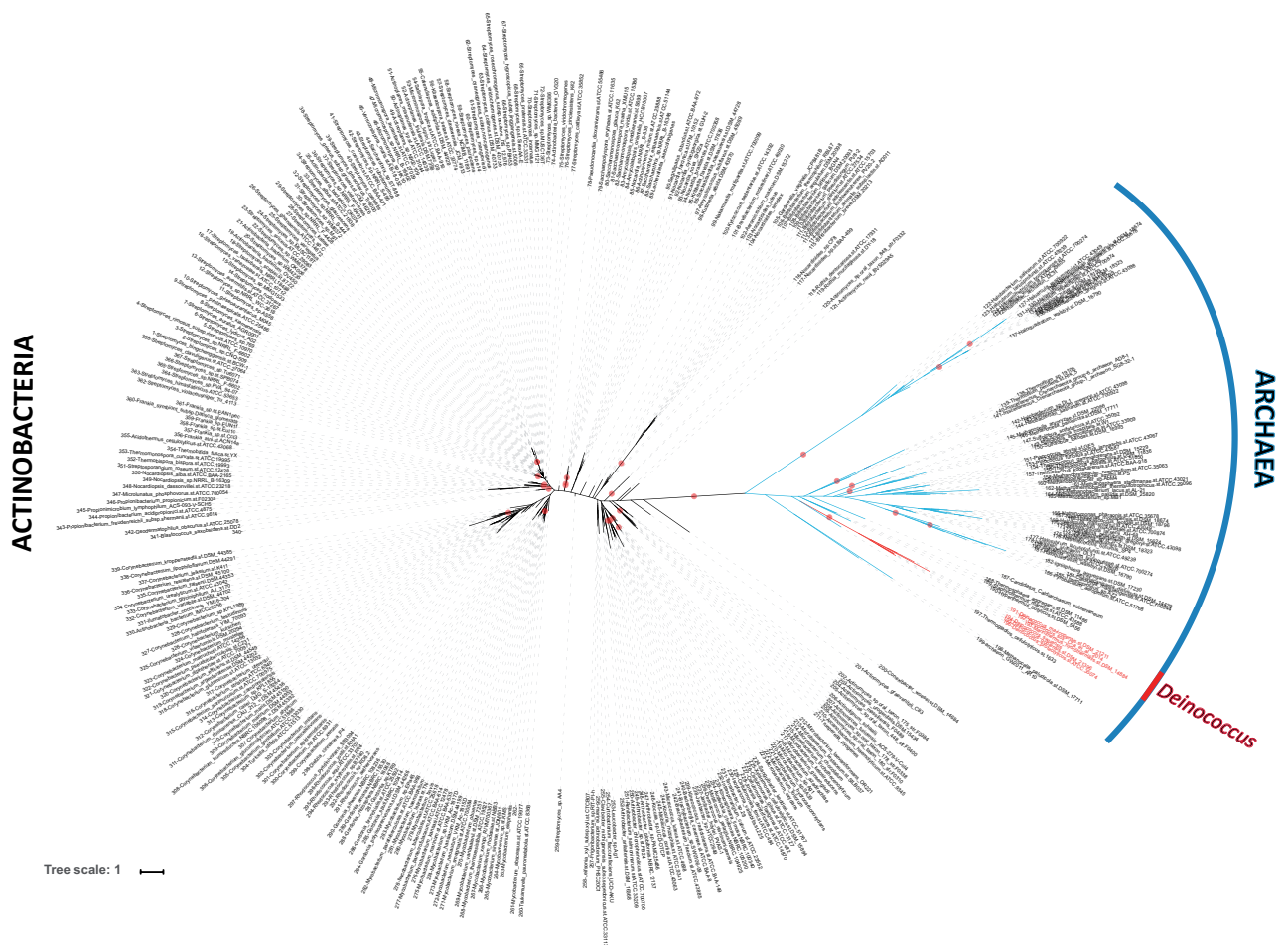


**Supplementary Figure 5. Phylogenetic analyses of NucS.** Unrooted ML tree of full NucS sequences from the PFAM PF01939. Black branches are Actinobacteria; blue branches are Archaea; red branches are species from the Deinococcus-Thermus group. Labels indicate species name. Red circles represent >80% bootstrap in 1,500 replicates.



**Supplementary Figure 6. Phylogenetic analyses of the NucS C-terminal (CT) region.** Unrooted ML tree of CT sequences from the PFAM PF01939 domain. Grey font indicates that this domain is in NucS and blue font that is found outside NucS. Black branches, Actinobacteria; red Deinococcus-Thermus; grey, other Bacteria; light blue, Archaea; green, eukaryotic sequences. Bold fonts indicate the main species used in this study. Labels indicate species name. The circles represent >80% bootstrap in 1,500 replicates.





**Supplementary Figure 7. Phylogenetic analyses of NucS N-terminal (NT) region.** Unrooted ML tree of NT sequences from the PFAM PF01939 domain. Black branches, Actinobacteria; red Deinococcus-Thermus; light blue, Archaea. Labels indicate species name. The circles represent >80% bootstrap in 1,500 replicates.

## Supplementary Tables

**Supplementary Table 1. Mutation rates of *M. smegmatis* and its  $\Delta nucS$  derivatives.** Rates of spontaneous mutations conferring rifampicin (Rif-R) and streptomycin resistance (Str-R) of *M. smegmatis* mc<sup>2</sup> 155 (WT), *M. smegmatis*  $\Delta nucS$  and *M. smegmatis*  $\Delta nucS$  complemented with the wild-type *nucS* from *M. smegmatis* mc<sup>2</sup> 155 (*nucS*<sub>Sm</sub>). Fold change indicates the increase in mutation rate with respect to the strain *M. smegmatis* mc<sup>2</sup> 155 (set to 1). Mut Rate: mutation rate (mutations per cell per generation).

Strain	Genotype	Mut Rate Rif	95% CI	Fold	Mut Rate Str	95% CI	Fold
mc <sup>2</sup> 155	WT	2.07x10 <sup>-9</sup>	1.47-2.74x10 <sup>-9</sup>	1	5.58x10 <sup>-10</sup>	2.25-11.1x10 <sup>-10</sup>	1
$\Delta nucS$	$\Delta nucS$	3.11x10 <sup>-7</sup>	2.81-3.39x10 <sup>-7</sup>	150.2	4.82x10 <sup>-8</sup>	3.53-6.17x10 <sup>-8</sup>	86.3
$\Delta nucS/nucS_{Sm}$	Complemented	3.02x10 <sup>-9</sup>	2.2-3.97x10 <sup>-9</sup>	1.46	1.43x10 <sup>-9</sup>	0.77-2.33x10 <sup>-9</sup>	2.56

**Supplementary Table 2. Mutational spectrum of *M. smegmatis* mc<sup>2</sup> 155 and its  $\Delta$ nucS derivative.**

A) Spontaneous mutations conferring rifampicin resistance found in 41 independent Rif-R mutants in the WT and  $\Delta$ nucS strains. The columns show the position of the mutations in the *rpoB* gene sequence, the codon change (modified bases are in bold), the amino acid change caused by each mutation and the number of independent Rif-R mutants isolated from mc<sup>2</sup> 155 (WT) or its  $\Delta$ nucS derivative. B) Specificity of base substitutions, summarized by class, produced in the WT and  $\Delta$ nucS strains.

A)

Position	Codon change	Amino acid change	WT	$\Delta$ nucS
A 1295 G	GAC→GGC	Asp 432 Gly	0	1
C 1324 T	CAC→TAC	His 442 Tyr	6	6
C 1324 G	CAC→GAC	His 442 Asp	2	0
A 1325 G	CAC→CGC	His 442 Arg	13	20
A 1325 C	CAC→CCC	His 442 Pro	4	0
G 1334 A	CGT →CAT	Arg 445 His	0	6
G 1334 C	CGT→CCT	Arg 445 Pro	1	0
G 1334 T	CGT→CTT	Arg 445 Leu	2	0
C 1340 T	TCG→TTG	Ser 447 Leu	9	6
C 1340 G	TCG→TGG	Ser 447 Trp	2	0
T 1346 C	CTG→CCG	Leu 449 Pro	0	2
G 1348 A	GGC→AGC	Gly 450 Ser	1	0
$\Delta$ CCAGCTGTC (1275-1283)	CCAGCTGTC	Ser 425 Arg + $\Delta$ GlnLeuSer (426-428)	1	0
TOTAL			41	41

B)

<b>Mutation</b>	<b>WT (%) n=41</b>	<b><i>ΔnucS</i> (%) n=41</b>
G:C→A:T	16 (39.0)	18 (43.9)
A:T→G:C	13 (31.7)	23 (56.1)
G:C→T:A	2 (4.9)	0
A:T→T:A	0	0
G:C→C:G	5 (12.2)	0
A:T→C:G	4 (9.7)	0
Deletions	1 (2.4)	0

**Supplementary Table 3. Mutation rates of *S. coelicolor* A3(2) M145 and its  $\Delta nucS$  derivatives.** Rates of spontaneous mutations conferring rifampicin (Rif-R) and streptomycin resistance (Str-R) of *S. coelicolor* A3(2) M145, its  $\Delta nucS$  derivative and *S. coelicolor*  $\Delta nucS$  complemented with the wild-type *nucS* from *S. coelicolor* (*nucS*<sub>SCO</sub>). 95% confidence intervals (CI) are indicated. Mut Rate: mutation rate (mutations per cell per generation).

Strain	Genotype	Mut Rate Rif-R	95% CI	Fold	Mut Rate Str-R	95% CI	Fold
<i>S. coelicolor</i> A3(2) M145	WT	$5.11 \times 10^{-9}$	2.87-8.01 $\times 10^{-9}$	1	$4.30 \times 10^{-10}$	0.24-18.9 $\times 10^{-10}$	1
$\Delta nucS$	$\Delta nucS$	$5.54 \times 10^{-7}$	4.20-6.97 $\times 10^{-7}$	108.4	$8.49 \times 10^{-8}$	5.78-11.4 $\times 10^{-8}$	197.4
$\Delta nucS/nucS_{SCO}$	Complemented	$2.01 \times 10^{-9}$	0.83-3.88 $\times 10^{-9}$	0.5	$4.30 \times 10^{-10}$	0.24-18.9 $\times 10^{-10}$	1



**Supplementary Table 4. Characteristics of the *M. tuberculosis* representative strains containing NucS polymorphisms.** Genome ID/common name, NucS polymorphism, resistance profile, lineage and origin of the strains are shown. Resistance profile indicates not detected/unknown antibiotic resistances (Susceptible) or MDR, multidrug-resistant strain (expressing at least rifampicin and isoniazid resistance).

<b>GenomeID/ name</b>	<b>Polymorphism</b>	<b>Resistance profile</b>	<b>Lineage</b>	<b>Origin</b>
CDC1551	WT	Susceptible	4	North America
TKK_02_0079	S39R	MDR	4	South Africa
MTB_N1057	S54I	Susceptible	4	South Asia
KT-0040	A67S	Susceptible	2	S. Korea (Broad Inst)
ERR036236	V69A	Susceptible	1	Unknown
BTB 04-388	A135S	MDR	3	Sweden (Broad Inst)
BTB 07-246	R144S	MDR	4	Sweden (Broad Inst)
TKK_03_0044	D162H	Susceptible	4	South Africa
HN2738	T168A	Unknown	Unknown	Unknown (Broad Inst)
MTB_X632	K184E	MDR	4	Central America

**Supplementary Table 5. Effect of *M. tuberculosis* NucS naturally occurring polymorphisms on mutation rates.** Rates of spontaneous mutations conferring rifampicin resistance (Mut rate, mutations per cell per generation) of *M. smegmatis*  $\Delta$ *nucS* complemented with the wild-type *nucS* from *M. tuberculosis* (*nucS*<sub>TB</sub>) or the 9 polymorphic alleles. 95% confidence intervals (CI) are shown. Fold change indicates the increase in mutation rate with respect to the strain *M. smegmatis*  $\Delta$ *nucS* complemented with wild-type *nucS*<sub>TB</sub> ( $\Delta$ *nucS*/*nucS*<sub>TB</sub>), set to 1.

<b>Amino acid change</b>	<b>Codon change</b>	<b>Mut rate</b>	<b>95% CI</b>	<b>Fold change</b>
<i>M. smegmatis</i> $\Delta$ <i>nucS</i> / <i>nucS</i> <sub>TB</sub>	<i>nucS</i> <sub>TB</sub> from CDC1551	4.17x10 <sup>-9</sup>	3.29-5.12x10 <sup>-9</sup>	1
S39R	AGC→AG <u>G</u>	3.49x10 <sup>-7</sup>	2.93-4.0x10 <sup>-7</sup>	83.7
S54I	AGT→A <u>T</u> T	7.94x10 <sup>-9</sup>	5.49-11.07x10 <sup>-9</sup>	1.9
A67S	GCG→ <u>T</u> CG	2.78x10 <sup>-9</sup>	1.87-3.84x10 <sup>-9</sup>	0.7
V69A	GTG→G <u>C</u> G	8.77x10 <sup>-9</sup>	6.50-11.2x10 <sup>-9</sup>	2.1
A135S	GCG→ <u>T</u> CG	3.57x10 <sup>-8</sup>	2.77-4.38x10 <sup>-8</sup>	8.6
R144S	CGC→ <u>A</u> GC	3.16x10 <sup>-8</sup>	2.40-4.0x10 <sup>-8</sup>	7.6
D162H	GAC→ <u>C</u> AC	8.98x10 <sup>-9</sup>	6.42-11.80x10 <sup>-9</sup>	2.2
T168A	ACC→ <u>G</u> CC	3.96x10 <sup>-8</sup>	3.09-4.81x10 <sup>-8</sup>	9.5
K184E	AAG→ <u>G</u> AG	3.06x10 <sup>-8</sup>	2.37-3.75x10 <sup>-8</sup>	7.3

**Supplementary Table 6.** Sequence of oligonucleotides used in this work.

Oligonucleotide	Sequence (5'-3')	Purpose
SecMycomar	CCC <span style="text-decoration: underline;">G</span> AAAAGTGCCACCTAAATTGTAAGCG	Localization of <i>Tn</i> insertion site
DelnucSm5F	CCC <span style="text-decoration: underline;">GCTGCAG</span> CTGGCCGAGTTCCG	<i>nucS<sub>Sm</sub></i> <i>M. smegmatis</i> deletion
DelnucSm5R	CGGT <span style="text-decoration: underline;">AAGCTT</span> GGCTATCACGAGGCGCACCC	<i>nucS<sub>Sm</sub></i> <i>M. smegmatis</i> deletion
DelnucSm3F	CGGA <span style="text-decoration: underline;">AAGCTT</span> AGCGACGAGTACCGGCTCTT	<i>nucS<sub>Sm</sub></i> <i>M. smegmatis</i> deletion
DelnucSm3R	AATC <span style="text-decoration: underline;">GTGCAC</span> CGAACCCATCAACTTACCGA	<i>nucS<sub>Sm</sub></i> <i>M. smegmatis</i> deletion
CompnucTBF	ACT <span style="text-decoration: underline;">GGAATTC</span> TCGAGTGGTGGCCTTCTCGGATGGCAT	<i>nucS<sub>TB</sub></i> for $\Delta$ <i>nucS<sub>Sm</sub></i> complementation
CompnucTBR	ACTG <span style="text-decoration: underline;">AAGCTT</span> TCAGAACAGCCGGTACTCGCCGCT	<i>nucS<sub>TB</sub></i> for $\Delta$ <i>nucS<sub>Sm</sub></i> complementation
CompnucSmF	ACT <span style="text-decoration: underline;">GGAATTC</span> CCGCGCCAGCGAATTGTCGGCGTTCAT	<i>nucS<sub>Sm</sub></i> for $\Delta$ <i>nucS<sub>Sm</sub></i> complementation
CompnucSmR	CGGC <span style="text-decoration: underline;">AAGCTT</span> TCAGAAGAGCCGGTACTCGTCGCT	<i>nucS<sub>Sm</sub></i> for $\Delta$ <i>nucS<sub>Sm</sub></i> complementation
RifRRDRrpoBF	GTGGCGGCGATCAAGGAGTTCTTC	RRDR- <i>rpoB<sub>Sm</sub></i> amplification
RifRRDRrpoBR	GGCGACCGACACCATCTGGCGCGG	RRDR- <i>rpoB<sub>Sm</sub></i> amplification
DelnucSco5F	GTGT <span style="text-decoration: underline;">AAGCTT</span> CGGCACCGCGGTGAGTGTGC	<i>nucS<sub>Sco</sub></i> <i>S. coelicolor</i> deletion
DelnucSco5R	CGAC <span style="text-decoration: underline;">GGATCC</span> GCGGGCAATGACGAGACGCA	<i>nucS<sub>Sco</sub></i> <i>S. coelicolor</i> deletion
DelnucSco3F	GCTG <span style="text-decoration: underline;">GGATCC</span> TTCTGAGGGCGGACGCGTC	<i>nucS<sub>Sco</sub></i> <i>S. coelicolor</i> deletion
DelnucSco3R	CAGG <span style="text-decoration: underline;">GATATC</span> TGCCCGCCCTGGTCGGCGAG	<i>nucS<sub>Sco</sub></i> <i>S. coelicolor</i> deletion
CompnucScoF	TACG <span style="text-decoration: underline;">GAATTC</span> GGGTTCTCCTCTCGCACCCCGACAGCAGGGG	<i>nucS<sub>Sco</sub></i> for $\Delta$ <i>nucS<sub>Sco</sub></i> complementation
CompnucScoR	TACG <span style="text-decoration: underline;">GGATCC</span> TCAGAACAGCCCGAGCTTGTCTCCTCGATGCC	<i>nucS<sub>Sco</sub></i> for $\Delta$ <i>nucS<sub>Sco</sub></i> complementation
Rechyg3F	CAGTTTCATTTGATGCTCGATGAG	Recombination. pRhomyco 100%
Rechyg3R	GACTA <span style="text-decoration: underline;">ACTAGT</span> CAGGCGCCGGGGCGGTGT	Recombination. pRhomyco 100%
Rechyg5F	GATG <span style="text-decoration: underline;">CAGCTG</span> GGAGTGGCTGTGACACAAGAATCC	Recombination. pRhomyco 100%
Rechyg5R	CGATA <span style="text-decoration: underline;">AAGCTT</span> CGAATTCTGCAGCTCG	Recombi. pRhomyco 100%, 95%, 90% and 85%
Rechyg5intR	GATG <span style="text-decoration: underline;">TGATCA</span> CCGGGTCGGGCTCG	Recombi. pRhomyco 95%, 90% and 85%
Rec95-BclIF	GATG <span style="text-decoration: underline;">TGATCA</span> AGCTGTTCCGGGAGCACTG	Recombination. pRhomyco 95%
Rec95-NheIR	GATG <span style="text-decoration: underline;">GCTAGC</span> GGAAGTCGACGATCCCGGTGA	Recombination. pRhomyco 95%
Rec90-85-BclIF	GATG <span style="text-decoration: underline;">TGATCA</span> AGCTCTTCGGGGAACACTG	Recombination. pRhomyco 90% and 85%
Rec90-85-PciIR	GATG <span style="text-decoration: underline;">ACATGT</span> GGAAGTCGACGATCCCGGTGA	Recombination. pRhomyco 90% and 85%
S39RF	GGCCGACGGATCGGT <span style="text-decoration: underline;">CAG</span> GGTACATGCTGACGACCG	Site-directed mutagenesis <i>nucS<sub>TB</sub></i>
S39RR	CGGTCGTGACGATG <span style="text-decoration: underline;">TACC</span> CTGACCGATCCGTGGCC	Site-directed mutagenesis <i>nucS<sub>TB</sub></i>
S54IF	CCGTTGAACTGGATG <span style="text-decoration: underline;">AT</span> TCCCGCTGCTGGTTG	Site-directed mutagenesis <i>nucS<sub>TB</sub></i>

S54IR	CAACCAGCACGGCGGA <u>A</u> TTCATCCAGTTCAACGG	Site-directed mutagenesis <i>nucS<sub>1B</sub></i>
A67SF	AGAGTCGGCGGCCAG <u>T</u> CGCCAGTGTGGTGGTCCG	Site-directed mutagenesis <i>nucS<sub>1B</sub></i>
A67SR	CGACCACCACACTGGCG <u>A</u> CTGGCCGCGGACTCTTC	Site-directed mutagenesis <i>nucS<sub>1B</sub></i>
V69AF	CCGGCGGCCAGGCGCCAG <u>C</u> GTGGGTGGTCGAGAAC	Site-directed mutagenesis <i>nucS<sub>1B</sub></i>
V69AR	GTTCTCGACCACCAC <u>G</u> CTGGCGCCTGGCCGCGG	Site-directed mutagenesis <i>nucS<sub>1B</sub></i>
A135SF	GCCGCGAGTACATGAC <u>C</u> TCGATCGACCCGTCGAC	Site-directed mutagenesis <i>nucS<sub>1B</sub></i>
A135SR	GTCGACGGTCCGATCG <u>A</u> GGTCATGTACTCGCGG	Site-directed mutagenesis <i>nucS<sub>1B</sub></i>
A162HF	GCGGCGTGGCGAGAT <u>C</u> ACGCGGTGGAGCAGCTGAC	Site-directed mutagenesis <i>nucS<sub>1B</sub></i>
A162HR	GTCAGCTGCTCCACG <u>C</u> CGTGATCTGCCACGCCGC	Site-directed mutagenesis <i>nucS<sub>1B</sub></i>
T168AF	GCGTGGAGCAGCT <u>G</u> CCCGCTACCTCGAGTTGC	Site-directed mutagenesis <i>nucS<sub>1B</sub></i>
T168AR	GCAACTCGAGGTAGCGGG <u>C</u> CAGCTGCTCCACGC	Site-directed mutagenesis <i>nucS<sub>1B</sub></i>
K184EF	GTGCTCGCGCCGGT <u>C</u> GAGGGGTGTTGCCG	Site-directed mutagenesis <i>nucS<sub>1B</sub></i>
K184ER	CGGCAACACCC <u>C</u> CTCGACCGGCGGAGCAC	Site-directed mutagenesis <i>nucS<sub>1B</sub></i>
pvv16-seq-fwd	CGGTGAGTCGTAGGTCGGGACGG	PCR amplification and sequencing
pvv16-seq-rev	TGCCTGGCAGTCGATCGTACGCTAG	PCR amplification and sequencing

## Supplementary Methods.

### Generation of a $\Delta nucS$ knockout mutant in *M. smegmatis*.

To generate a  $\Delta nucS$  knockout mutant in *M. smegmatis*, a 1.0 kb PstI-HindIII upstream and HindIII-SalI downstream fragments of *nucS* were amplified by PCR, using delnucSm 5F, 5R, 3F and 3R primers, and cloned in-frame into the p2NIL vector. Then, a pGOAL19 PacI-cassette, carrying a  $\beta$ -galactosidase *lacZ* gene, a hygromycin-resistance *hyg* gene and a *sacB* gene, that confers sucrose sensitivity, was also inserted into the p2NIL vector. The resulting plasmid p2NIL- $\Delta nucS_{Sm}$ , harbouring an in-frame deletion of the target gene (lacking  $\approx 95\%$  of the gene sequence), was electroporated into *M. smegmatis* mc<sup>2</sup> 155. Cells were plated on Middlebrook 7H10 agar-X-gal (100  $\mu\text{g ml}^{-1}$ ) plus kanamycin (25  $\mu\text{g ml}^{-1}$ ) and hygromycin (50  $\mu\text{g ml}^{-1}$ ) and incubated for 4-5 days at 37°C. Single-crossover merodiploid clones were grown in 7H9 broth without antibiotics to allow a second crossover event. Cultures were diluted and counter-selected on Middlebrook 7H10-X-gal (100  $\mu\text{g ml}^{-1}$ ) plates containing 2% sucrose and incubated 4-5 days at 37°C. Double-crossover clones were tested for kanamycin and hygromycin susceptibility to confirm the loss of the plasmid. Finally, *M. smegmatis* mc<sup>2</sup> 155  $\Delta nucS$  colonies were tested by PCR and sequencing to verify the unmarked *nucS* deletion.

### Complementation of the *M. smegmatis* $\Delta nucS$ mutant.

For complementation with *nucS* from *M. smegmatis* mc<sup>2</sup> 155, the *MSMEG\_4923* gene (*nucS<sub>Sm</sub>*), including its own 46-bp promoter region, was amplified by PCR with primers compnucSmF and compnucSmR primers, digested with EcoRI and HindIII and cloned into the integrative vector pMV361, rendering the complementation vector pMV-*nucS<sub>Sm</sub>*. Similarly, for complementation with *nucS* from *M. tuberculosis*, the wild-type full-length *MT\_1321* gene from CDC1551 control strain (*nucS<sub>TB</sub>*), with its own 61-pb upstream promoter region, was amplified by PCR, using compnucTBF and compnucTBR primers, cloned into the pMV361 vector to generate the complementation vector, pMV-*nucS<sub>TB</sub>*. Putative complemented mutants were obtained upon electroporation of the plasmids into *M. smegmatis* mc<sup>2</sup> 155  $\Delta nucS$  and incubation of the plated samples on Middlebrook 7H10 agar plus kanamycin (25  $\mu\text{g ml}^{-1}$ ) for 3-5 days at 37°C. Finally, *M. smegmatis* mc<sup>2</sup> 155  $\Delta nucS$  complemented mutants were analysed by PCR and sequencing to verify the proper insertion of the genes.



### **Generation of $\Delta nucS$ *S. coelicolor* knockout mutant and complementation.**

pIJ6650 vector was used to clone in-frame two 1.5 kb DNA fragments, one HindIII-BamHI *nucS* upstream fragment plus one BamHI-EcoRV *nucS* downstream fragment, previously amplified by PCR (primers delnucSco 5F, 5R, 3F and 3R). *E. coli* ET12567 (pUZ8002) was transformed with pIJ- $\Delta nucS_{Sco}$  (containing the in-frame deletion of the *nucS*<sub>Sco</sub> gene) and conjugated with *S. coelicolor* A3(2) M145. Following the isolation of apramycin resistant single-crossover, putative double-crossover mutants were isolated and the unmarked deletion of the *nucS*<sub>Sco</sub> gene was verified by PCR. *S. coelicolor*  $\Delta nucS$  was complemented with a wild-type copy of *nucS*<sub>Sco</sub> cloned in pSET152. pSET152 is a non-replicative plasmid in *Streptomyces* that carries the *attP* site and the integrase gene of the C31 phage and consequently can integrate into the *attB* site<sup>2</sup> in the chromosome of *Streptomyces*. *nucS*<sub>Sco</sub> gene with its own promoter was amplified by PCR with primers compnucScoF and compnucScoR, cloned into pSET152 using EcoRI and BamHI sites to give the pSET-*nucS*<sub>Sco</sub> plasmid that was introduced into *E. coli* ET12567 (pUZ8002) by transformation. Finally, pSET-*nucS*<sub>Sco</sub> from *E. coli* ET12567 (pUZ8002, pSET-*nucS*<sub>Sco</sub>) was introduced into *S. coelicolor*  $\Delta nucS$  by conjugation in order to integrate the construction into the chromosome.

### **Construction of pRhomyco plasmids.**

To create a template for measuring homologous and homeologous recombination in *M. smegmatis*, we designed a recombination assay based on the hygromycin-resistance (Hyg-R) gene *hyg*, using the integrative plasmid pMV361. For pRhomyco 100%, a fragment denominated *hyg* 3' (from nucleotide 195 of the coding sequence to the stop codon) of the *hyg* gene, was PCR amplified from pRAM vector, using rechyg3F and rechyg3R primers and digested with EcoRV/SpeI to be cloned in StuI/SpeI targets of pMV361. Additionally, a fragment denominated *hyg* 5' containing 711 bp (from the start codon of the *hyg* gene to nucleotide 711) was PCR amplified and cloned using PvuII/HindIII with rechyg5F and rechyg5R primers. Both fragments share two overlapping *hyg* fragments of 517 bp (nucleotide 195 to 711). The homeologous recombination vectors, named pRhomyco 95%, 90% and 85%, were constructed by replacing the overlapping 517-bp fragment of *hyg* 5' for synthetic fragments that have 95%, 90% or 85% sequence similarity to the original *hyg* fragment (Supplementary Fig. 2). First, a common fragment to the three of them (from nucleotide 1 to 193 of *hyg*) was PCR-amplified with rechyg5F and rechyg5intR primers and cloned using PvuII/BclI. Then, each synthetic fragment was PCR amplified and cloned using BclI/NheI with rec95-BclIF and rec95-NheIR primers, in the case of pRhomyco 95%. For pRhomyco 90% and

85%, the variable fragment was cloned using BclI/PciI-BspHI with rec90-85-BclIF and rec90-85-PciIR primers.

### **Computational analyses.**

A summary of all the computational approaches conducted is depicted in Supplementary Fig. 4. For NucS, we used the structure-based alignment of bacterial and archaeal NucS (Supplementary Fig. 3). Then we followed the procedure depicted in Supplementary Fig. 4 to conduct domain analyses using sequences.

Using the different defined regions by structural bioinformatics, we first conducted sequence searches against the large database. We made non-redundant alignments of the N-terminal region (containing 63 sequences) and the C-terminal region (containing 39 sequences) and built profile hidden Markov models (HMMs). These profiles were searched using pHMMER<sup>3</sup> against the large References Proteomes file where additional proteins containing CT and NT regions in alternative proteins were found. While the NT region was found in alternative archaeal proteins (we selected two after checking the alignment F7PKA0\_9EURY and L0I7R5\_HALRX (DUF91, PF01939), the CT region was found not only in additional prokaryotic groups, but also in eukaryotic sequences (I1F1Q3\_AMPQE, I1F1Q5\_AMPQE, B3SFT8\_TRIAD, B9T981\_RICCO, and A0A015J6U4\_9GLOM). With the exception of A0A015J6U4, and B9T981, the matching regions are located within a DUF1016 domain (uncharacterized). We checked the alignments to confirm that the catalytic residues, as described in the structure of *P. abyssi*<sup>1</sup>, were conserved.

We next searched the PFAM database with NUCS\_MYCTU and found a domain, PF01939, which was trained in three NucS full sequences. From the PF01939 domain (539 sequences in 464 species), only 368 entries in 363 (archaeal and bacterial) species gave a match with the full protein (Supplementary data file 1).

To identify MutS and MutL in bacterial and archaeal species, we generated HMMs from different proteins, which would capture most of the bacterial and archaeal diversity (Supplementary Fig. 4). The MutS profile was trained with the following sequences (MUTS\_ECOLI, C1F256\_ACIC5, MUTS\_AQUAE, MUTS\_BACTN, MUTS\_CHLPN, MUTS\_CHLTE, MUTS\_CHLAA, WP\_027388865.1, MUTS\_PROM3, B5YE41\_DICT6, MUTS\_BACSU, MUTS\_STAAM, MUTS\_CLOTE, MUTS\_FUSNN, C1AEM6\_GEMAT, A6DUF8\_9BACT, MUTS\_RHOBA, MUTS\_BRUME, MUTS\_NEIMA, MUTS\_SALTY, MUTS\_VIBC3, MUTS\_PSEAE, MUTS\_DESVH, MUTS\_TREPA, A0A075WU36\_9BACT, MUTS\_THEMA, WP\_009958798.1), while the MutL profile was trained with the following sequences (MUTL\_ACIC5, MUTL\_AQUAE, MUTL\_BACTN, MUTL\_CHLPN, MUTL\_CHLTE, A9WJ86\_CHLAA,

MUTL\_DEIRA, MUTL\_DICT6, D9S9H6\_FIBSS, MUTL\_BACSU, MUTL\_STAAM, MUTL\_CLOTE, Q8RG56\_FUSNN, MUTL\_GEMAT, A6DHB3\_9BACT, Q7UMZ3\_RHOBA, MUTL\_BRUME, MUTL\_NEIMA, A0A0C5TQ33\_SALTM, MUTL\_VIBCH, MUTL\_PSEAE, Q72ET5\_DESVH, MUTL\_TREPA, A0A075WSF3\_9BACT, MUTL\_THEMA, and MUTL\_ECOLI). The models were searched against each reference proteome where bit scores <50 were discarded to filter and analyze the results. Only full proteins were retrieved by forcing a length > 75%, so partial hits would be excluded.

For actinobacterial and archaeal species not showing NucS, searches in all the particular genomes (only for complete genomes) were done using translated searches against their genomes using tblastn with NucS\_MYCTU, MutS\_ECOLI, and MUTL\_ECOLI for actinobacterial, and NucS\_PYRAB, MutS\_HALMA, and MUTL\_HALSA for archaeal proteins.

## References for Supplementary Information

- 1 Ren, B. *et al.* Structure and function of a novel endonuclease acting on branched DNA substrates. *The EMBO journal* **28**, 2479-2489, doi:emboj2009192 [pii] 10.1038/emboj.2009.192 (2009).
- 2 Bierman, M. *et al.* Plasmid cloning vectors for the conjugal transfer of DNA from *Escherichia coli* to *Streptomyces* spp. *Gene* **116**, 43-49 (1992).
- 3 Eddy, S. R. A new generation of homology search tools based on probabilistic inference. *Genome Inform* **23**, 205-211 (2009).