

# **Distinct Biological Potential of *Streptococcus gordonii* and *Streptococcus sanguinis* Revealed by Comparative Genome Analysis**

**Wenning Zheng<sup>1,2</sup>; Mui Fern Tan<sup>1,2</sup>; Lesley A. Old<sup>4</sup>; Ian C. Paterson<sup>2,3</sup>; Nicholas S Jakubovics<sup>4,5\*</sup>; Siew Woh Choo<sup>1,2\*</sup>**

<sup>1</sup>Genome Informatics Research Laboratory, High Impact Research Building (HIR) Building, University of Malaya, 50603 Kuala Lumpur, Malaysia

<sup>2</sup>Department of Oral and Craniofacial Sciences, Faculty of Dentistry, University of Malaya, 50603 Kuala Lumpur, Malaysia

<sup>3</sup>Oral Cancer Research and Coordinating Centre, Faculty of Dentistry, University of Malaya, 50603 Kuala Lumpur, Malaysia

<sup>4</sup>Centre for Oral Health Research, School of Dental Sciences, Newcastle University, Newcastle upon Tyne, NE2 4BW, United Kingdom

<sup>5</sup>Genome Solutions Sdn Bhd, Suite 8, Innovation Incubator UM, Level 5, Research Management & Innovation Complex, University of Malaya, 50603 Kuala Lumpur, Malaysia

**\*=Corresponding authors**

**Siew Woh Choo**

**Email: [lawrence.choo@xjtlu.edu.cn](mailto:lawrence.choo@xjtlu.edu.cn)**

**Nicholas S. Jakubovics**

**Tel: +44 (0) 191 208 6796 ; Email: [nick.jakubovics@ncl.ac.uk](mailto:nick.jakubovics@ncl.ac.uk)**

**S1 Table: Functional enrichment analyses shows Ss unique core genes enriched in compound biosynthetic process (8) and cobalamin biosynthesis process (20).**

Enriched biological process	Rast_ID	Enriched <i>Streptococcus sanguinis</i> unique core genes
Porphyrin-containing compound biosynthetic process	SK36.peg.479 SK36.peg.476 SK36.peg.475 SK36.peg.474 SK36.peg.465 SK36.peg.461 SK36.peg.2040 SK36.peg.2038	Glutamate-1-semialdehyde aminotransferase (EC 5.4.3.8) Porphobilinogen deaminase (EC 2.5.1.61) Glutamyl-tRNA reductase (EC 1.2.1.70) Siroheme synthase / Precorrin-2 oxidase (EC 1.3.1.76) Uroporphyrinogen-III methyltransferase (EC 2.1.1.107) / Uroporphyrinogen-III synthase (EC 4.2.1.75) Cobalt-precorrin-4 C11-methyltransferase (EC 2.1.1.133) /cobM FIG01117915: hypothetical protein FIG01118726: hypothetical protein
Cobalamin biosynthetic process	SK36.peg.481 SK36.peg.480 SK36.peg.472 SK36.peg.502 SK36.peg.469 SK36.peg.500 SK36.peg.468 SK36.peg.467 SK36.peg.466 SK36.peg.464 SK36.peg.463 SK36.peg.462 SK36.peg.461 SK36.peg.460 SK36.peg.459 SK36.peg.458 SK36.peg.457 SK36.peg.456 SK36.peg.455 SK36.peg.454	Cobalamin synthase Adenosylcobinamide-phosphate guanylyltransferase (EC 2.7.7.62) / cobU Cobyric acid synthase Nicotinate-nucleotide--dimethylbenzimidazole phosphoribosyltransferase (EC 2.4.2.21) / cobT Additional substrate-specific component CbiN of cobalt ECF transporter L-threonine 3-O-phosphate decarboxylase (EC 4.1.1.81) Substrate-specific component CbiM of cobalt ECF transporter Cobalt-precorrin-2 C20-methyltransferase (EC 2.1.1.130) Sirohydrochlorin cobaltochelataase CbiK (EC 4.99.1.3) Cobalt-precorrin-6x reductase (EC 1.3.1.54) Cobalt-precorrin-3b C17-methyltransferase /cbiH Cobalamin biosynthesis protein CbiG Cobalt-precorrin-4 C11-methyltransferase (EC 2.1.1.133) Cobalt-precorrin-6y C15-methyltransferase [decarboxylating] (EC 2.1.1.-) Cobalt-precorrin-6y C5-methyltransferase (EC 2.1.1.-) Cobalt-precorrin-6 synthase, anaerobic Cobalt-precorrin-8x methylmutase (EC 5.4.1.2) Cobalt-precorrin-8x methylmutase (EC 5.4.1.2) Adenosylcobinamide-phosphate synthase /cbiP Cobyric acid A, C-diamide synthase /cobB/cbiA

**S2 Table: Overview of putative prophages including the size of the prophage, the number of CDS, ATT-site status and GC content.**

<b>Prophages</b>	<b>Length (kb)</b>	<b>CDS</b>	<b>ATT-site identified</b>	<b>GC content</b>
7863_1	5.8kb	6	No	40.99%
FSS8_1	43.2kb	58	Yes	38.80%
SK12_1	36.5kb	54	No	41.25%
SK184_1	59.2kb	57	Yes	41.08%
SK184_3	48.7kb	75	Yes	40.74%
Channon_2	39.4Kb	62	Yes	38.76%
FSS4_1	30.9Kb	25	Yes	43.52%
FSS4_2	16.4Kb	29	Yes	40.31%
MB451_1	23.3Kb	26	Yes	43.51%
SK184_2	36kb	21	Yes	37.94%
SK184_4	6.9kb	11	No	43.49%
MB666_1	47Kb	32	Yes	40.29%

**S3 Table: Overview of the putative genomic island (GI) details including the size of the GI, the number of CDS, GC contents and key genes incorporated in each GI.**

GI	Size (bp)	Number of CDSs	GC content	Key Genes
GI_5	5253	5	33.70%	DNA recombination and repair protein RecF; FIG001621: Zinc protease; FIG009210: peptidase, M16 family and Transcriptional regulator in cluster with unspecified monosaccharide ABC transport system
GI_14	10312	6	20.50%	hypothetical proteins
GI_16	5085	6	38.90%	FIG007079: UPF0348 protein family; FIG145533: Methyltransferase (EC 2.1.1.-); lojap protein;Hydrolase (HAD superfamily), YqeK and Nicotinate-nucleotide adenylyltransferase (EC 2.7.7.18)
GI_31	5557	6	33.60%	Permease of the drug/metabolite transporter (DMT) superfamily; TetR/AcrR family transcriptional regulator
GI_43	7035	10	42.20%	Integrase; Chromosome segregation helicase and MutT/nudix family protein; 7,8-dihydro-8-oxoguanine-triphosphatase
GI_45	5556	8	33.60%	Chromosome (plasmid) partitioning protein ParB / Stage 0 sporulation protein J; Serine protease, DegP/HtrA, do-like (EC 3.4.21.-); LSU m3Psi1915 methyltransferase RlmH and Competence pheromone precursor
GI_47	7627	12	43.20%	Integrase; Chromosome segregation helicase; MutT/nudix family protein; 7,8-dihydro-8-oxoguanine-triphosphatase; acetyltransferase,GNAT family;Ribosomal protein L11 methyltransferase (EC 2.1.1.-); Ribosomal RNA small subunit methyltransferase E (EC 2.1.1.-), and Mobile element protein (2 units)
GI_51	7355	9	31.90%	Chromosome (plasmid) partitioning protein ParB / Stage 0 sporulation protein J; Serine protease, DegP/HtrA, do-like (EC 3.4.21.-);LSU m3Psi1915 methyltransferase RlmH; Competence pheromone precursor; Histidine kinase of the competence regulon ComD; Response regulator of the competence regulon ComE; GTP-binding and nucleic acid-binding protein YchF; Peptidyl-tRNA hydrolase (EC 3.1.1.29) and Transcription-repair coupling factor
GI_53	4194	5	44.40%	CAAX amino terminal protease family family and Transcriptional regulator, TetR family
GI_55	5516	7	48.40%	V-type ATP synthase subunit C, E, F,G, I and K (EC 3.6.3.14) and Acetyltransferase, GNAT family
GI_58	7364	8	32.10%	Chromosome (plasmid) partitioning protein ParB / Stage 0 sporulation protein J; Serine protease, DegP/HtrA, do-like (EC 3.4.21.-); LSU m3Psi1915 methyltransferase RlmH; Competence pheromone precursor; Histidine kinase of the competence regulon ComD; Response regulator of the competence regulon ComE; GTP-binding and nucleic acid-binding protein YchF and Peptidyl-tRNA hydrolase (EC 3.1.1.29)
GI_67	4094	5	41.90%	Topoisomerase IV subunit B (EC 5.99.1.-) and lipoprotein, putative
GI_75	4183	5	44.60%	CAAX amino protease and Transcriptional regulator, TetR family

**S4 Table: The genome sequencing results of 19 isolated *Streptococcus* strains using Next Generation Sequencing Illumina Hiseq 2000 platform.**

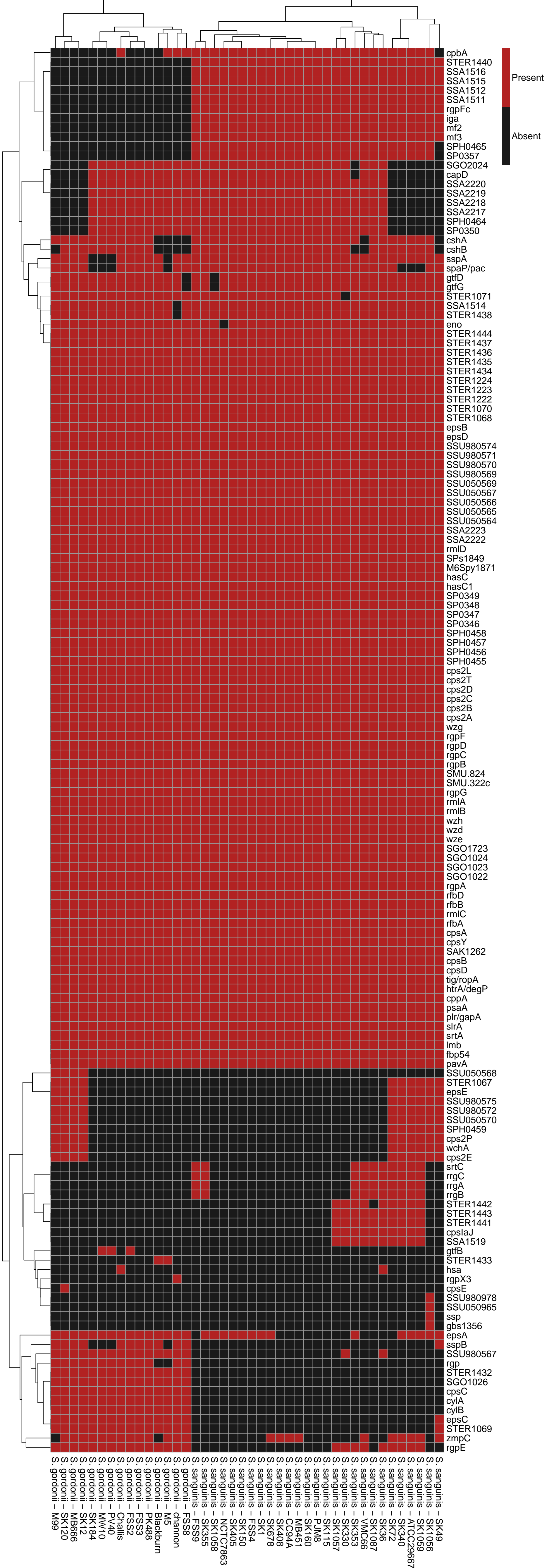
<b>Strain Name</b>	<b>Yield (Mbases)</b>	<b>Number of Reads</b>	<b>Mean Quality Score (PF)</b>
<b>PV40</b>	826	8264126	36.24
<b>NCTC 7863</b>	822	8222024	35.69
<b>Blackburn</b>	1180	11797640	36.55
<b>Channon</b>	695	6949680	36.49
<b>FSS2</b>	882	8823944	36.67
<b>FSS3</b>	944	9442900	37.11
<b>FSS4</b>	1294	12943882	36.6
<b>FSS8</b>	1010	10102148	36.69
<b>FSS9</b>	988	9877224	36.61
<b>M5</b>	678	6784012	36.43
<b>M99</b>	666	6657462	36.19
<b>MB451</b>	1127	11271462	36.59
<b>MB666</b>	1095	10949508	36.81
<b>MW10</b>	1069	10693318	36.94
<b>PJM8</b>	1054	10543052	36.63
<b>PK488</b>	624	6240768	36.14
<b>SK12</b>	878	8782388	36.81
<b>SK120</b>	732	7324194	36.96
<b>SK184</b>	680	6795252	36.36

**S5 Table: The sequencing coverage of 19 isolated *Streptococcus* strains based on the reference genome sizes of Sg Challis (2.2Mb) and Ss SK36 (2.39Mb).**

<b>Strain Name</b>	<b>Total Read Length (bp)</b>	<b>Genome Size (bp)</b>	<b>Sequencing Coverage</b>
<b>PV40</b>	826412600	2200000	375.64
<b>Blackburn</b>	721326800	2200000	327.88
<b>Channon</b>	694968000	2200000	315.89
<b>FSS2</b>	882394400	2200000	401.09
<b>FSS3</b>	944290000	2200000	429.22
<b>FSS8</b>	1010214800	2200000	459.19
<b>M5</b>	678401200	2200000	308.36
<b>M99</b>	665746200	2200000	302.61
<b>MB666</b>	1094950800	2200000	497.70
<b>MW10</b>	1069331800	2200000	486.06
<b>NCTC7863</b>	822202400	2390000	344.02
<b>FSS4</b>	1294388200	2390000	541.59
<b>FSS9</b>	987722400	2390000	413.27
<b>MB451</b>	1127146200	2390000	471.61
<b>PJM8</b>	1054305200	2390000	441.13
<b>PK488</b>	624076800	2200000	283.67
<b>SK12</b>	878238800	2200000	399.20
<b>SK120</b>	732419400	2200000	332.92
<b>SK184</b>	679525200	2200000	308.88

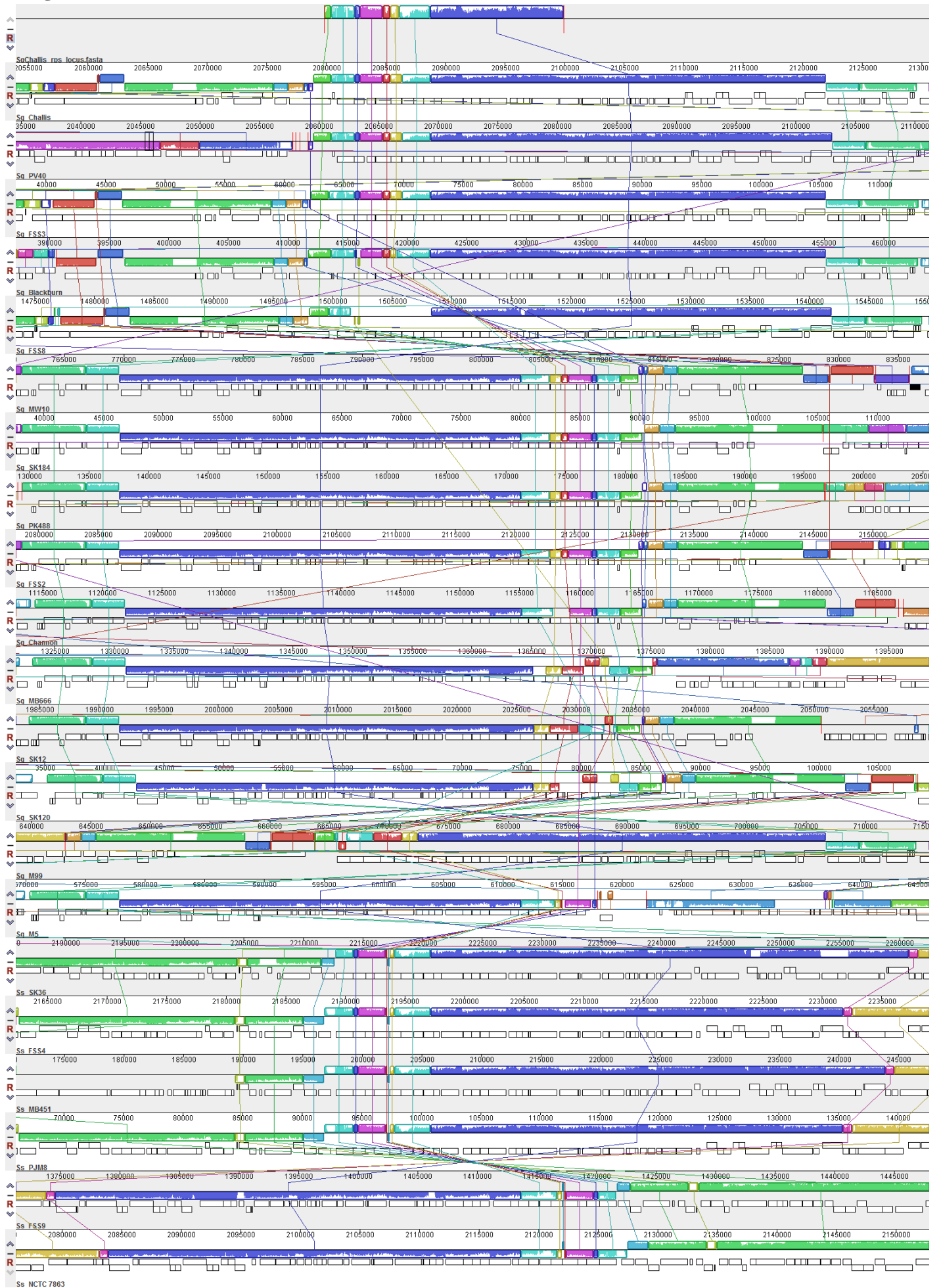
**S1 Figure. The heatmap generated for comparative pathogenomics analysis between 15 Sg strains and 27 Ss strains using a threshold of 50% sequence identity and 50% sequence coverage.**

PathoProT - SI:50%, SC:50%

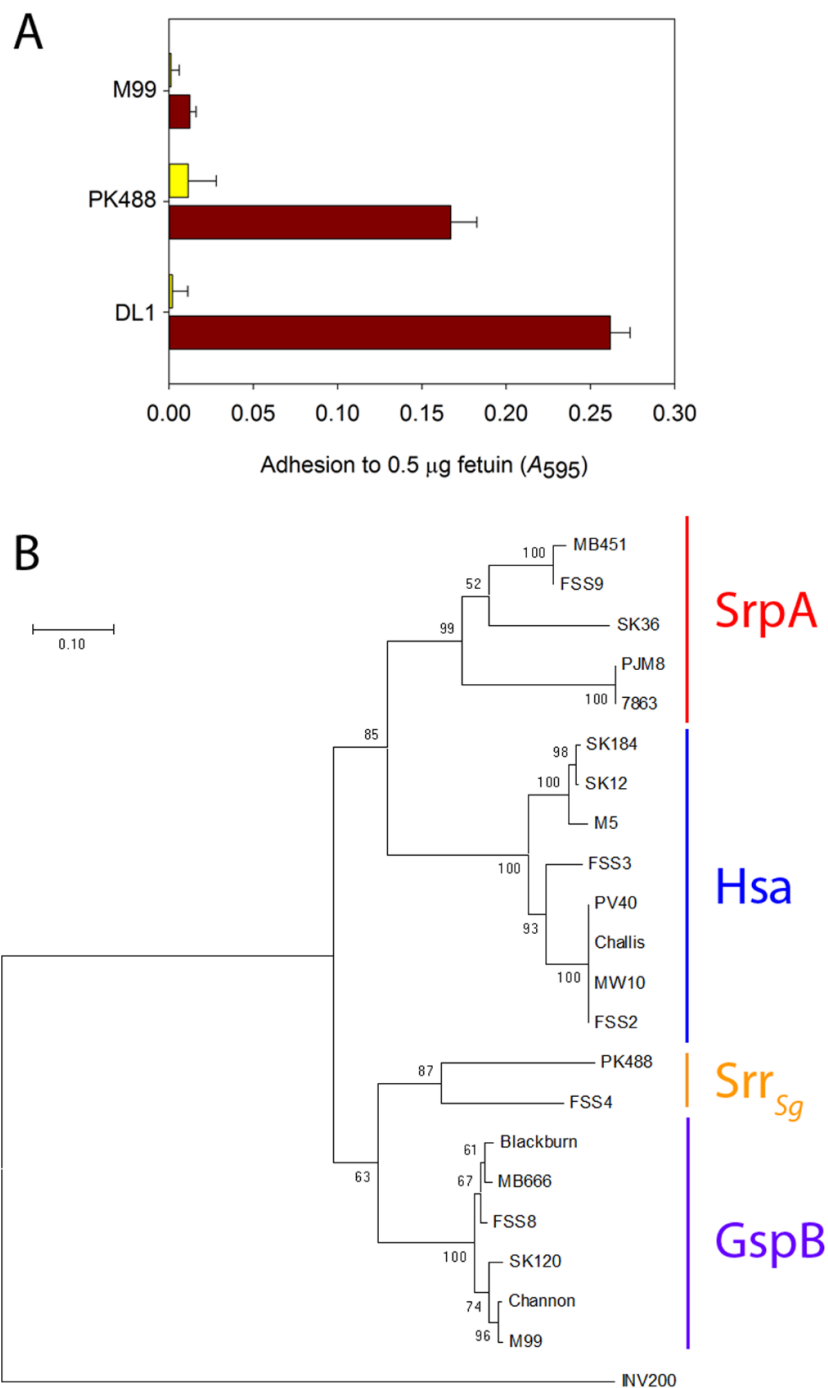




**S2 Figure. The visualization of Sg Challis-type polysaccharide gene cluster structure in Sg and Ss using Mauve software.**



**S3 Figure.** A. Adhesion of Sg M99, PK488 and DL1 to fetuin (red bars) or sialidase-treated fetuin (yellow bars). The substrate was immobilised on a plastic surface and exposed to bacteria for 2 h at 37°C. After washing, bound bacterial cells were stained with crystal violet and quantified by measuring  $A_{595}$  as described previously<sup>1</sup>. Bars represent means from three independent experiments and standard errors are shown. B. Phylogenetic analysis of N-terminal binding regions of serine-rich region proteins from Ss and Sg. The evolutionary history was inferred using the Neighbor-Joining method using *S. pneumoniae* INV200 PsrP as an outgroup<sup>2</sup>. The percentage of replicate trees in which the associated taxa clustered together in the bootstrap test (500 replicates) are shown next to the branches<sup>3</sup>. The tree is drawn to scale, with branch lengths in the same units as those of the evolutionary distances used to infer the phylogenetic tree. The evolutionary distances were computed using the Poisson correction method<sup>4</sup> and are shown as the number of amino acid substitutions per site. Evolutionary analyses were conducted in MEGA7<sup>5</sup>. The binding region sequences of Sg PK488 and FSS4 clustered separately from the relatively tight GspB and Hsa variant groups (Sg strains) or the SrpA-type proteins, which were found only in Ss strains.



## Supplemental references

1. Jakubovics, N.S., Brittan, J. L., Dutton, L.C. & Jenkinson, H.F. Multiple adhesin proteins on the cell surface of *Streptococcus gordonii* are involved in adhesion to human fibronectin. *Microbiology* **155**, 3572-3580 (2009).
2. Saitou N. & Nei M. The neighbor-joining method: A new method for reconstructing phylogenetic trees. *Molecular Biology and Evolution* **4**, 406-425 (1987).
3. Felsenstein J. Confidence limits on phylogenies: An approach using the bootstrap. *Evolution* **39**, 783-791 (1985).
4. Zuckerkandl E. & Pauling L. Evolutionary divergence and convergence in proteins. Edited in *Evolving Genes and Proteins* by V. Bryson and H.J. Vogel, pp. 97-166. Academic Press, New York (1965).
5. Kumar S., Stecher G., & Tamura K. MEGA7: Molecular Evolutionary Genetics Analysis version 7.0 for bigger datasets. *Molecular Biology and Evolution* **33**:1870-1874 (2016).