



**HELMHOLTZ  
ZENTRUM FÜR  
INFEKTIONSFORSCHUNG**

**This is a pre- or post-print of an article published in  
Sagar, V., Bergmann, R., Nerlich, A., McMillan, D.J.,  
Schmitz, D.P.N., Chhatwal, G.S.**

**Variability in the distribution of genes encoding  
virulence factors and putative extracellular proteins of  
Streptococcus pyogenes in India, a region with high  
streptococcal disease burden, and implication for  
development of a regional multisubunit vaccine  
(2012) Clinical and Vaccine Immunology, 19 (11), pp.  
1818-1825.**



1 **Variability in the distribution of genes encoding virulence factors and putative**  
2 **extracellular proteins of *Streptococcus pyogenes* in India – a region with high**  
3 **streptococcal disease burden: implication for the development of a regional**  
4 **multisubunit vaccine**

5

6 Vivek Sagar<sup>1\*</sup>, René Bergmann<sup>1\*</sup>, Andreas Nerlich<sup>1\*§</sup>, David J. McMillan<sup>2</sup>, D. Patric  
7 Nitsche-Schmitz<sup>1</sup> and Gursharan S. Chhatwal<sup>1#</sup>

8

9 <sup>1</sup>Department of Medical Microbiology, Helmholtz Centre for Infection Research,  
10 Inhoffenstrasse 7, 38124 Braunschweig,

11 <sup>2</sup>Queensland Institute for Medical Research, Brisbane, Australia

12

13\*equal contribution

14

15#correspondence: Gursharan S. Chhatwal, PhD

16 E-mail: [gsc@helmholtz-hzi.de](mailto:gsc@helmholtz-hzi.de)

17

18Running Title: Genotyping of Indian GAS isolates

19

20

21

22

23<sup>§</sup>present address: University of Veterinary Medicine Hannover Foundation, Institute of  
24Microbiology, Centre for Infection Medicine, Bischofsholer Damm15, 30173

25Hannover

## 1Abstract

2*Streptococcus pyogenes* causes a wide variety of human diseases and is a  
3significant cause of morbidity and mortality. Attempts to develop a vaccine were  
4hampered by genetic diversity of *S. pyogenes* across different regions of the world.  
5This study aims to identify streptococcal antigens suitable for a region-specific  
6vaccine in India. We used a two step approach, firstly performing epidemiological  
7analysis to identify the conserved antigens among Indian isolates. The second step  
8consisted of validating those identified antigens by serological analysis. The 201  
9streptococcal clinical isolates from India used in this study represented 69 different  
10*emm* types with *emm12* as the most prevalent. Virulence profiling of the North and  
11South Indian *S. pyogenes* isolates with a custom-designed streptococcal virulence  
12microarray identified 7 conserved putative vaccine candidates. Collagen-like surface  
13protein (SCI), putative secreted 5 prime nucleotidase (PSNT) and C5a peptidase  
14were found in 100 percent of the isolates, while R28, a putative surface antigen  
15(PSA) and a hypothetical protein (HYP) were found in 90 percent of the isolates. A  
16fibronectin binding protein, SfbI, was present in only 78 percent of the isolates. In  
17order to validate the identified potential vaccine candidates, 185 serum samples  
18obtained from patients with different clinical manifestations were tested for  
19antibodies. Irrespective of clinical manifestations, serum samples showed high  
20antibody titers to all proteins except for SCI and R28. Thus, the data indicates that  
21PSNT, C5a peptidase, PSA, HYP and SfbI are promising candidates for a region-  
22specific streptococcal vaccine for the different parts of India.

23

## 1Introduction

2*Streptococcus pyogenes* (group A streptococcus, GAS) is exclusively a human  
3pathogen and the etiological agent of a wide variety of diseases that vary in clinical  
4severity, while also being a significant cause of morbidity and mortality . These  
5diseases include pharyngitis, impetigo, scarlet fever, post-streptococcal  
6glomerulonephritis, invasive diseases, rheumatic fever (RF) and rheumatic heart  
7disease (RHD) . While rheumatic fever and rheumatic heart disease are the greatest  
8cause of mortality in developing nations, deaths in developed nations are mainly  
9attributable to invasive diseases .

10The differences in the prevalence and molecular epidemiology of GAS isolates reflect  
11the differences in the importance of RF/RHD and invasive diseases in these  
12populations. GAS carriage and infection are prevalent in many developing nations,  
13with a large number of different *emm* types circulating at one time, with no one type  
14being dominant. In contrast, a limited number of specific *emm* types are predominant  
15in developed nations, and are often associated with specific clinical manifestations .

16

17Despite the majority of GAS associated deaths occurring in developing nations, the  
18majority of comparative genetic and genomic studies have focussed on isolates from  
19developed nations . These studies report a high degree of genetic diversity between  
20isolates with divergent *emm* types, as well as diversity within an *emm* type. In regions  
21with high streptococcal disease burden, where GAS isolates are more likely to come  
22into direct contact, the probability of lateral gene transfer (LGT) is increased. Greater  
23inter- and intra-*emm*-type genetic diversity is likely to exist. Such diversity and  
24potential changes in population structure are important considerations when  
25designing vaccine candidates that should provide coverage against the entire GAS

1population.

2

3In this study, we used a two step approach in order to identify promising candidates  
4for a region-specific vaccine. In the first step we conducted epidemiological analysis  
5in order to identify antigens conserved in different parts of India, in the second step  
6we purified these antigens and performed serological analysis using convalescent  
7serum samples. We used virulence gene profiling to assess genetic diversity and  
8population structure of GAS in India, a country where streptococcal disease burden is  
9high . We found a high degree of genetic diversity between isolates of different *emm*  
10types, but relatively conserved genotypes within an *emm* type. We assessed the  
11distribution of seven genes conserved in Indian isolates encoding current vaccine  
12targets in this population, and report on serological responses against each of these  
13antigens.

14

## 1Experimental procedures

### 2Bacterial strains and human sera

3Bacterial isolates and sera used in this study were collected as part of ASSIST  
4program, funded by the European Commission, with three Indian and three European  
5partners ([http://www.helmholtz-  
6hzi.de/en/research/research\\_projects/view/projekt/projekt/assist/](http://www.helmholtz-hzi.de/en/research/research_projects/view/projekt/projekt/assist/)). The samples were  
7collected during 2007 - 2010 in two defined areas in Chandigarh (North India) and  
8Vellore (South India), which are about 3000 km apart, have high streptococcal  
9disease burden and different climatic conditions. Besides hospital patients, 3000 and  
102400 school children were screened in Chandigarh and Vellore, respectively. Sixty-  
11five isolates from Chandigarh and 136 from Vellore were included in this study. The  
12isolates were also classified on the basis of recovery from the throat of asymptomatic  
13carriers (n=44), or patients presenting with pharyngitis (n=20). Another 32 isolates  
14were collected from the skin in these surveys. Thirty-four isolates were collected from  
15patients presenting with invasive disease at clinics in Chandigarh and Vellore. The  
16isolation site of 71 isolates was unknown. Serum samples were collected from  
17individuals in North (n=110) and South India (n=75). Eighty-two of these samples  
18were obtained from patients with RHD, 24 collected from patients with RF and 9  
19sourced from patients with invasive disease. Another 15 were obtained from patients  
20presenting with symptoms of pharyngitis, 30 were asymptomatic patients positive for  
21GAS and 25 sourced from healthy people from the survey areas with no current signs  
22of streptococcal infection. Because of the possibility of subclinical exposure in the  
23survey areas due to high disease burden, we included 10 serum samples collected  
24from healthy people from Germany as controls. The serum samples and  
25streptococcal isolates were from the same patients in case of invasive diseases. For

1all other clinical manifestations the isolates and serum samples were from the  
2defined survey areas, but not from the same patients.

3

#### 4**DNA extraction and *emm* typing**

5Genomic DNA was isolated using zirconium beads in combination with the Qiagen  
6DNeasy kit (Qiagen, Hilden, Germany). The *emm* type of individual strains was  
7determined by amplification and sequencing the 5'-end region of the *emm* gene  
8described by Beall and Facklam . The resulting nucleotide sequences were blasted  
9against the *emm* gene database  
10(<http://www.cdc.gov/ncidod/biotech/strep/strepblast.htm>) to determine the *emm* type  
11of individual isolates.

12

#### 13**Microarray, hybridization and data analysis**

14The microarray used in this study has been described in detail in a previous study .  
15Sample preparation, array hybridization and data processing were performed as  
16described previously . Briefly, isolated genomic DNA was digested with *A**lu**I*, labelled  
17with biotin/streptavidin-Cy5 and fluorescence signals were quantified after  
18hybridization using ImaGene software (BioDiscovery). A two component mixture  
19model was fitted to the background-corrected and log<sub>2</sub>-transformed data by a  
20maximum likelihood method. A discriminant function was used to represent the  
21propensity of a gene for being present or absent in a particular isolate. Discriminant  
22values were stored in a signal probability matrix and used to construct dendrograms  
23using Bayesian agglomerative hierarchical clustering . All routines for statistical  
24calculations were implemented in the R statistics package ([www.r-project.org](http://www.r-project.org)).

25

26

2  
3

## 1**Cloning and expression of cell surface proteins**

2Out of 219 genes on the array, 7 genes were selected which showed 75% or more  
3frequency in 201 strains and encode for cell wall proteins. These proteins, which  
4included collagen like-surface protein (SCI), putative secreted 5'-nucleotidase  
5(PSNT), C5a peptidase (C5a), R28 (R28), putative surface antigen (PSA) and a  
6hypothetical transposase (HYP) were amplified (Table 1) and cloned into the pGEX-  
76P-1 vector (GE health care). Cloning of streptococcal fibronectin binding protein I  
8(Sfbl) was described by Talay *et al.* . The expression of the recombinant GST-fusion  
9proteins in *E. coli* BL21(DE3) was induced by addition of IPTG to 0.5 mM. Cultures  
10were harvested and lysed using a French Press (SLM instruments inc.). Cellular  
11debris was removed by centrifugation and the presence of recombinant proteins in  
12supernatants confirmed by SDS-PAGE.

13

## 14**Antibody responses to putative vaccine candidates**

15ELISA was performed using glutathione coated microtitre plates (Thermo Fisher  
16Scientific) with *E. coli* lysates (0.25 µg/µl) containing GST-tagged recombinant  
17proteins following the method of Sehr *et al.* . PBS and *E. coli* lysate containing GST  
18alone served as controls. A peptide representing the IgG binding motif from FOG, a  
19streptococcal protein with IgG-Fc fragment binding capacity was used as the  
20normalization control. HRP conjugated goat antihuman IgG antibody and ABTS (2,2'-  
21azino-bis(3-ethylbenzthiazoline-6-sulphonic acid) were used for color development  
22that was measured at 405 nm. Absorbance values were subsequently normalized  
23against values obtained from wells containing PBS, GST and the FOG IgG binding  
24peptide. Differences in absorbance between groups were examined using one-way  
25ANOVA. Because of the high number of serum samples and 8 proteins, the  
26experiment was designed in such a way that maximum number of serum samples

1could be covered by one experiment in order to avoid experiment to experiment  
2variations. Therefore we chose to monitor the immunogenicity of proteins in all serum  
3samples at the dilution of 1:200 and expressed the result as OD.

4

## 5**Results**

### 6***emm* sequence type distribution of GAS isolates from North and South India**

7The 201 isolates collected represent 69 different *emm* sequence types  
8(Supplementary Table S1) and 83 different *emm* sequence subtypes (data not  
9shown). The 20 most prevalent *emm* types (Figure 1) were *emm12* (6.5%), *emm11*  
10(5.5%), *emm49* (4.5%), *emm28*, *emm80*, st1389 (3.5% each), *emm3* (3.0%), *emm4*,  
11*emm44*, *emm75* and *emm112* (2.5% each). The types *emm2*, *emm22*, *emm69*,  
12*emm74*, *emm77*, *emm93*, *emm110*, *emm104* and *emm108* accounted for 2% of all  
13strains recovered. The 20 prevalent *emm* types represented 57.7% of all GAS  
14isolated in both regions.

15

16Of these twenty prevalent *emm* types, ten (*emm4*, *emm11*, *emm12*, *emm44*, *emm49*,  
17*emm74*, *emm80*, *emm110*, *emm112*, and st1389) were found both in Northern and  
18Southern India. *Emm3*, *emm28*, *emm69*, *emm77*, *emm93*, *emm104*, and *emm108*  
19were only found in Southern India, whereas *emm2*, *emm22* *emm75* were specific for  
20Northern India. With 55 different *emm* sequence types the South Indian strains  
21showed a greater heterogeneity as compared to 32 different *emm* sequence types  
22found within the North Indian strains. This is reflected in the Simpsons Index of  
23Diversity (*D*) for *emm* types in Southern India (*D*= 0.981; 95% Confidence interval,  
240.976-0.986) compared to Northern India (*D*= 0.957; 95% Confidence interval, 0.936-  
250.977).

1

## 2**Overview of gene distribution in Indian Isolates**

3Of the 219 streptococcal virulence factors and extracellular surface proteins  
4(VF/ECP) represented on the array, 91 (41.5%) were found in all isolates  
5(Supplementary Table S2), and 150 (68.5%) were found in more than 80% of the  
6isolates. When we analyzed the distribution of the VF/ECP genes with respect to  
7geographic region, 102 genes (46.6%) and 100 genes (45.7%) were present in all  
8North Indian and South Indian isolates, respectively. In order to identify region-  
9specific virulence genes we statistically analyzed the distribution of VF/ECP genes in  
10both regions. We identified one gene in the North Indian isolates and 15 genes in the  
11South Indian isolates with a statistically significant difference in distribution (Fisher  
12exact test,  $P < 0.05$ , Table 2), reflecting the higher *emm* type heterogeneity in South  
13India. However, none of the 16 genes was found exclusively in one of the regions.

14

## 15**Genetic relationships amongst isolates**

16*Emm* typing is still largely accepted as sufficient for defining related strains in GAS  
17research. However, epidemiological and molecular observations also demonstrate  
18that LGT involving the *emm* gene occurs, suggesting that isolates of the same *emm*  
19type that are temporally or geographically displaced may be genetically diverse, and  
20be more closely related to isolates of a different *emm* type . To investigate the  
21relationship between all isolates we used a Bayesian agglomerative hierarchical  
22cluster algorithm (BHC) to construct a dendrogram of all isolates that were  
23represented by  $\geq 3$  isolates per *emm* type (n = 154 isolates), based on the presence  
24or absence of all genes represented on the array. Using this approach we found  
25isolates of the same *emm* type to predominantly cluster together (Figure 2).  
26However, there were instances where this did not occur, suggesting that LGT of the

2

3

1 *emm* gene may have occurred in these instances. We next analyzed the association  
2 of the clusters with regard to disease type. As shown in Fig 2, there was only one  
3 distinct sub-cluster solely composed of invasive isolates consisting of one *emm4*, one  
4 *emm28* and three *emm49* strains. Using BHC clustering we did not find any  
5 significant clustering of virulence factors with invasive and non-invasive isolates.  
6 However, when we analyzed the association of virulence factors with invasiveness  
7 using the Fischer's Exact Test ( $P < 0.05$ ), we found 13 genes positively associated  
8 with invasiveness (Supplementary Table S3). Six genes were negatively associated  
9 with invasiveness (odds ratio  $< 1$ ).

10

### 11 **Serological response to vaccine antigens**

12 We examined the distribution of the shortlisted seven GAS vaccine candidates in the  
13 India population, and found them to be highly conserved across the population (Table  
14 3). Collagen-like surface protein (SCI), putative secreted 5' nucleotidase (PSNT) and  
15 C5a peptidase were found in 100% of isolates. R28, the putative surface antigen  
16 (PSA) and the hypothetical protein (HYP) were found in more than 90% of isolates.  
17 Streptococcal fibronectin-binding protein I (Sfbl) was found in 78.5% of isolates.

18

19 Serum antibodies indicate that epitopes of a particular protein are visible to the host  
20 immune system, and these proteins are likely to be recognized by antibodies  
21 generated against corresponding vaccine candidates. Therefore, we next  
22 investigated serological responses to these seven identified GAS vaccine  
23 candidates. We expressed these proteins as fusion with GST and used them in  
24 ELISA with human sera (Figure 3).

25

1Control healthy sera had very low immune response against all of the proteins. The  
2majority of serum samples from India had elevated responses to PSNT, C5a  
3peptidase and PSA, demonstrating the high incidence of GAS infection in these  
4populations. PSA-specific IgG responses were statistically more frequent in the  
5pharyngitis group than in healthy and RHD groups ( $P < 0.05$ ). Lower antigen specific  
6antibody responses were observed with HYP and SfbI. For these two proteins the  
7highest responses were observed with sera collected from pharyngitis patients and  
8patients with GAS throat colonization at the time of sera collection. For HYP, IgG  
9responses in the pharyngitis group were statistically significantly higher than the  
10healthy from non-survey and survey areas and RF groups ( $P < 0.05$ ). The SfbI  
11response for both the carrier group and pharyngitis group was statistically  
12significantly higher than control and RHD groups ( $P < 0.05$ ) (Supplementary Table  
13S4). The majority of sera did not respond, or responded weakly to R28 and SCI. No  
14differences in responses were observed in sera collected from North and South India.  
15For 4/7 of the proteins used here the response in pharyngitis sera was significantly  
16higher than non-survey area healthy sera (Supplementary Table S4) indicating the  
17possibility that these proteins are most likely expressed in the early phase of  
18infection.

19

20

## 21 **Discussion**

22This study is the most comprehensive assessment of the genetic repertoire of GAS  
23carried out in a country with high streptococcal disease burden. The twenty most  
24common *emm* types accounted for 57.7% of the isolates collected, reinforcing the  
25high degree of diversity in the population based on *emm* typing and is consistent with

1results of other studies of GAS population structure carried out in developing  
2nations . Within the Indian population, 91 genes were conserved in 100% of isolates.  
3Using the same array, we previously reported that 129 and 125 genes were 100%  
4conserved in GAS isolates from the Netherlands and United States respectively. The  
5majority (80 genes) of conserved genes in the isolates from Europe and India were  
6similar, indicating that they are part of the GAS core genome and encode proteins  
7that are critical for virulence or confer biological functions essential to the fitness of  
8the organism, and are therefore promising vaccine candidates.

9

10The M-protein remains a favored GAS vaccine candidate. The predominant  
11bactericidal antibody response raised after GAS infection targets the amino-terminal  
12of the M-protein, the same region that is used for *emm*-typing. However antibodies  
13raised against this region have traditionally been thought to be type specific i.e.  
14antibodies raised against the amino terminus of one M-protein will not recognise the  
15amino terminus of other M-proteins. Candidate vaccines that target this region  
16therefore must include amino-termini from multiple M-proteins. Of the 69 *emm*  
17sequences reported in this study, 14 are represented in the 26-valent amino-terminal  
18GAS vaccine candidate and represent 31.8% of isolates recovered. The highest  
19proportion of isolates included in the vaccine was found in South Indian GAS isolates  
20(17.4%) and included 10 vaccine related *emm* types. Of the North Indian Isolates  
2114.4% were covered by the vaccine, and included 6 vaccine related *emm* types.  
22Recently a similar 30-valent vaccine has also been reported . Nineteen of the *emm*  
23types present in the 30mer vaccine were present in isolates in this study. The same  
24study also reported that cross-reactivity was observed with another 24 M-proteins not  
25present in the vaccine. Fourteen of these cross-reactive *emm* types were present in  
26this study. In theory, the 30-valent vaccine would induce antibodies that are effective

1against 58% of the isolates recovered in this study. The *emm* types present in the 26-  
2valent vaccine, and 30-valent vaccine were chosen based on their importance in  
3North American and European contexts. A re-tailored vaccine, containing amino-  
4termini from M-types common in India or other regions where streptococcal disease  
5is prevalent, may increase vaccine coverage by this approach. However, given the  
6large numbers of *emm* types circulating, and the differences in M-types present in  
7North and South India, such a vaccine may prove difficult to design.

8

9Both traditional and reverse vaccinology approaches have been used to identify  
10alternatives to the M-protein . Rodriguez-Ortega *et al.* used a proteomic approach for  
11identifying streptococcal surface exposed proteins for their use as vaccine  
12candidates. They however used only three *emm* types that were not prevalent in  
13India. Reverse vaccinology, a genome-based approach to vaccine development, has  
14also been used to identify streptococcal vaccine candidates . This approach,  
15however, requires the whole genome sequence of the isolates and to date no  
16prevalent *emm type* strain in India has been sequenced. Proteins that confer  
17essential virulence and biological functions, and are therefore encoded in the core  
18genome, are attractive vaccine targets. In addition to functionality through  
19bactericidal or neutralizing activity, antibodies raised against critical proteins may  
20contribute to prevention of infection by inhibiting protein function. Multi-subunit  
21vaccines containing several proteins have added appeal, as abolition of multiple  
22biological activities may attenuate virulence of the organism further. The targeting of  
23several proteins also potentially has an additive effect with respect to opsonic or  
24neutralizing antibodies. A multi-subunit approach to vaccine can also reduce the  
25probability of vaccine escape that may occur after loss of epitopes targeted by the  
26vaccine that may occur through mutation or lateral gene transfer .

1The seven proteins we examined serologically have all been proposed as potential  
2vaccine candidates, and were chosen here because they represent both well  
3characterized and relatively new targets. All these proteins are also predicted to be  
4surface associated. Our array data support this selection because six of these  
5proteins were conserved in all of the invasive isolates tested in this study. R28 is a  
6cell surface virulence protein which has many repetitive sequence and has homology  
7with a protein found in group B streptococcus . SCI is a collagen like protein which is  
8known to have collagen like sequence . HYP has homology with FbaA proteins ,  
9suggesting it to have fibronectin binding ability. PSNT is a putative surface  
10nucleotidase protein likely to be involved in nutrient acquisition . The corresponding  
11PNST of *Haemophilus influenzae* has been reported to have protective efficacy in a  
12rat infection model . PSA is a homologue to spy0843 of *S. pyogenes* MGAS5005  
13which contributed to a protective host immune response in a mouse infection model .

14

15The identified antigens were validated for their vaccine potential by serological  
16analysis. For invasive disease we had the matching isolates and serum samples,  
17because these were acute manifestations that required hospitalization. For other  
18manifestations, which were not acute, it was not possible to have the matching  
19isolate serum samples, so we had to rely on the samples from defined survey areas.  
20The immunogenicity of the identified proteins were tested using ELISA in order to  
21determine their suitability as vaccine candidates. For simplicity, antibodies were  
22determined using a single dilution and results were expressed as OD values. For  
23further studies on these antigens, however, the end point titers or the use of a  
24reference standard serum will be required . Our data demonstrates that antibodies to  
25five of the seven GAS proteins are present in sera in individuals living in a region

1where streptococcal disease is common. In general, there was no difference between  
2antibody levels between individuals, irrespective of disease status of individuals.  
3However two proteins, HYP and SfbI, had higher antibody responses in sera from  
4pharyngitis and asymptomatic throat carriers when compared to sera from RF/RHD  
5and asymptomatic non-carriers. Fibronectin binding proteins have been shown to be  
6involved in adherence to epithelial cells and may therefore be important in  
7colonization and early infection. Responses to PSA, which binds to the surface of  
8epithelial cells were also significantly higher in sera from pharyngitis patients. Given  
9the increased antibody response to these three antigens (HYP, SfbI and PSA) in  
10carriers and patients with pharyngitis, they may be excellent candidates for inclusion  
11in a vaccine designed to prevent initial colonization. In this study, generally a lower  
12immune response was detected against all the proteins with serum samples from  
13invasive diseases. This could be due to the fact that the invasive disease is an acute  
14manifestation where the time is too short for the generation of sufficient antibodies to  
15give a high response.

16

17In conclusion, PSNT, C5a peptidase, PSA, HYP and SfbI are promising candidates  
18for a region-specific streptococcal vaccine for the Indian continent. Among these five  
19candidates PSA and PSNT are not well studied and can be further explored in future  
20studies. The two-step approach used in this study to identify vaccine candidates with  
21regional specificity seems to be very promising. However, more isolates and more  
22defined serum samples should be analyzed in further work to support this study.

23

24

## 25**Acknowledgments**

26The authors would like to thank N. Janze for excellent technical assistance. This work

1was supported by the European Community's Sixth Framework Programme ASSIST  
2under contract number 032390. We are grateful to ASSIST epidemiology team  
3members, especially A. Kumar, K.N. Brahmadathan, V. Abraham and Y. Sharma for  
4providing streptococcal isolates and serum samples. We thank A.P. Oxley for  
5carefully reading the manuscript.

## 2References

31. **Abdissa, A., D. Asrat, G. Kronvall, B. Shittu, D. Achiko, M. Zeidan, L. K. Yamuah, and A. Aseffa.** 2006. High diversity of group A streptococcal emm types among healthy schoolchildren in Ethiopia. *Clin Infect Dis* **42**:1362-1367.
62. **Agarwal, A. K., M. Yunus, J. Ahmad, and A. Khan.** 1995. Rheumatic heart disease in India. *J R Soc Health* **115**:303-304, 309.
83. **Banks, D. J., S. F. Porcella, K. D. Barbian, S. B. Beres, L. E. Philips, J. M. Voyich, F. R. DeLeo, J. M. Martin, G. A. Somerville, and J. M. Musser.** 2004. Progress toward characterization of the group A Streptococcus metagenome: complete genome sequence of a macrolide-resistant serotype M6 strain. *J Infect Dis* **190**:727-738.
134. **Beall, B., R. Facklam, and T. Thompson.** 1996. Sequencing emm-specific PCR products for routine and accurate typing of group A streptococci. *J Clin Microbiol* **34**:953-958.
165. **Beres, S. B., E. W. Richter, M. J. Nagiec, P. Sumby, S. F. Porcella, F. R. DeLeo, and J. M. Musser.** 2006. Molecular genetic anatomy of inter- and intraserotype variation in the human bacterial pathogen group A Streptococcus. *Proc Natl Acad Sci U S A* **103**:7059-7064.
206. **Bessen, D. E., and S. K. Hollingshead.** 1995. Horizontal transfer and mosaic-like emm gene structures in group A streptococci. *Dev Biol Stand* **85**:169-173.
227. **Carapetis, J. R., A. C. Steer, E. K. Mulholland, and M. Weber.** 2005. The global burden of group A streptococcal diseases. *Lancet Infect Dis* **5**:685-694.
248. **Cole, J. N., A. Henningham, C. M. Gillen, V. Ramachandran, and M. J. Walker.** 2008. Human pathogenic streptococcal proteomics and vaccine development. *Proteomics Clin Appl* **2**:387-410.
279. **Cunningham, M. W.** 2000. Pathogenesis of group A streptococcal infections. *Clin Microbiol Rev* **13**:470-511.
2910. **Dale, J. B., T. A. Penfound, E. Y. Chiang, and W. J. Walton.** 2011. New 30-valent M protein-based vaccine evokes cross-opsonic antibodies against non-vaccine serotypes of group A streptococci. *Vaccine* **29**:8175-8178.
3211. **Dey, N., D. J. McMillan, P. J. Yarwood, R. M. Joshi, R. Kumar, M. F. Good, K. S. Sriprakash, and H. Vohra.** 2005. High diversity of group A Streptococcal emm types in an Indian community: the need to tailor multivalent vaccines. *Clin Infect Dis* **40**:46-51.
3612. **Eshaghi, M., A. M. Ali, F. Jamal, and K. Yusoff.** 2002. Existence of two emm-like "mrp" and "emm" genes in the mga regulon of the Streptococcus pyogenes strain ST4547. *J Biochem Mol Biol Biophys* **6**:23-28.
3913. **Ferretti, J. J., W. M. McShan, D. Ajdic, D. J. Savic, G. Savic, K. Lyon, C. Primeaux, S. Sezate, A. N. Suvorov, S. Kenton, H. S. Lai, S. P. Lin, Y. Qian, H. G. Jia, F. Z. Najar, Q. Ren, H. Zhu, L. Song, J. White, X. Yuan, S. W. Clifton, B. A. Roe, and R. McLaughlin.** 2001. Complete genome sequence of an M1 strain of Streptococcus pyogenes. *Proc Natl Acad Sci U S A* **98**:4658-4663.
4414. **Hu, M. C., M. A. Walls, S. D. Stroop, M. A. Reddish, B. Beall, and J. B. Dale.** 2002. Immunogenicity of a 26-valent group A streptococcal vaccine. *Infect Immun* **70**:2171-2177.
4715. **Loimaranta, V., J. Hytonen, A. T. Pulliainen, A. Sharma, J. Tenovuo, N. Stromberg, and J. Finne.** 2009. Leucine-rich repeats of bacterial surface proteins serve as common pattern recognition motifs of human scavenger receptor gp340. *J*

- 1 Biol Chem **284**:18614-18623.
216. **Lukomski, S., K. Nakashima, I. Abdi, V. J. Cipriano, R. M. Ireland, S. D. Reid,**  
3 **G. G. Adams, and J. M. Musser.** 2000. Identification and characterization of the scl  
4 gene encoding a group A Streptococcus extracellular protein virulence factor with  
5 similarity to human collagen. *Infect Immun* **68**:6542-6553.
617. **Maione, D., I. Margarit, C. D. Rinaudo, V. Masignani, M. Mora, M. Scarselli, H.**  
7 **Tettelin, C. Brettoni, E. T. Iacobini, R. Rosini, N. D'Agostino, L. Miorin, S.**  
8 **Buccato, M. Mariani, G. Galli, R. Nogarotto, V. Nardi Dei, F. Vegni, C. Fraser,**  
9 **G. Mancuso, G. Teti, L. C. Madoff, L. C. Paoletti, R. Rappuoli, D. L. Kasper, J.**  
10 **L. Telford, and G. Grandi.** 2005. Identification of a universal Group B streptococcus  
11 vaccine by multiple genome screen. *Science* **309**:148-150.
1218. **McMillan, D. J., R. G. Beiko, R. Geffers, J. Buer, L. M. Schouls, B. J. Vlamincx,**  
13 **W. J. Wannet, K. S. Sriprakash, and G. S. Chhatwal.** 2006. Genes for the majority  
14 of group a streptococcal virulence factors and extracellular surface proteins do not  
15 confer an increased propensity to cause invasive disease. *Clin Infect Dis* **43**:884-891.
1619. **McMillan, D. J., R. Geffers, J. Buer, B. J. Vlamincx, K. S. Sriprakash, and G. S.**  
17 **Chhatwal.** 2007. Variations in the distribution of genes encoding virulence and  
18 extracellular proteins in group A streptococcus are largely restricted to 11 genomic  
19 loci. *Microbes Infect* **9**:259-270.
2020. **McShan, W. M., J. J. Ferretti, T. Karasawa, A. N. Suvorov, S. Lin, B. Qin, H. Jia,**  
21 **S. Kenton, F. Najjar, H. Wu, J. Scott, B. A. Roe, and D. J. Savic.** 2008. Genome  
22 sequence of a nephritogenic and highly transformable M49 strain of *Streptococcus*  
23 *pyogenes*. *J Bacteriol* **190**:7773-7785.
2421. **Miura, K., A. C. Orcutt, O. V. Muratova, L. H. Miller, A. Saul, and C. A. Long.**  
25 2008. Development and characterization of a standardized ELISA including a  
26 reference serum on each plate to detect antibodies induced by experimental malaria  
27 vaccines. *Vaccine* **26**:193-200.
2822. **Musser, J. M., V. Kapur, J. Szeto, X. Pan, D. S. Swanson, and D. R. Martin.** 1995.  
29 Genetic diversity and relationships among *Streptococcus pyogenes* strains expressing  
30 serotype M1 protein: recent intercontinental spread of a subclone causing episodes of  
31 invasive disease. *Infect Immun* **63**:994-1003.
3223. **Nitsche-Schmitz, D. P., H. M. Johansson, I. Sastalla, S. Reissmann, I. M. Frick,**  
33 **and G. S. Chhatwal.** 2007. Group G streptococcal IgG binding molecules FOG and  
34 protein G have different impacts on opsonization by C1q. *J Biol Chem* **282**:17530-  
35 17536.
3624. **Pai, R., M. R. Moore, T. Pilishvili, R. E. Gertz, C. G. Whitney, and B. Beall.** 2005.  
37 Postvaccine genetic structure of *Streptococcus pneumoniae* serotype 19A from  
38 children in the United States. *J Infect Dis* **192**:1988-1995.
3925. **Panchaud, A., L. Guy, F. Collyn, M. Haenni, M. Nakata, A. Podbielski, P.**  
40 **Moreillon, and C. A. Roten.** 2009. M-protein and other intrinsic virulence factors of  
41 *Streptococcus pyogenes* are encoded on an ancient pathogenicity island. *BMC*  
42 *Genomics* **10**:198.
4326. **Pfoh, E., M. R. Wessels, D. Goldmann, and G. M. Lee.** 2008. Burden and economic  
44 cost of group A streptococcal pharyngitis. *Pediatrics* **121**:229-234.
4527. **Rato, M. G., A. Nerlich, R. Bergmann, R. Bexiga, S. F. Nunes, C. L. Vilela, I.**  
46 **Santos-Sanches, and G. S. Chhatwal.** 2011. Virulence gene pool detected in bovine  
47 group C *Streptococcus dysgalactiae* subsp. *dysgalactiae* isolates by use of a group A S.  
48 *pyogenes* virulence microarray. *J Clin Microbiol* **49**:2470-2479.
4928. **Reid, S. D., N. M. Green, G. L. Sylva, J. M. Voyich, E. T. Stenseth, F. R. DeLeo,**  
50 **T. Palzkill, D. E. Low, H. R. Hill, and J. M. Musser.** 2002. Postgenomic analysis of  
51 four novel antigens of group a streptococcus: growth phase-dependent gene

- transcription and human serologic response. *J Bacteriol* **184**:6316-6324.
229. **Rodriguez-Ortega, M. J., N. Norais, G. Bensì, S. Liberatori, S. Capo, M. Mora, M. Scarselli, F. Doro, G. Ferrari, I. Garaguso, T. Maggi, A. Neumann, A. Covre, J. L. Telford, and G. Grandi.** 2006. Characterization and identification of vaccine candidate proteins through analysis of the group A *Streptococcus* surface proteome. *Nat Biotechnol* **24**:191-197.
730. **Rohde, M., R. M. Graham, K. Branitzki-Heinemann, P. Borchers, C. Preuss, I. Schleicher, D. Zahner, S. R. Talay, M. Fulde, K. Dinkla, and G. S. Chhatwal.** 2011. Differences in the aromatic domain of homologous streptococcal fibronectin-binding proteins trigger different cell invasion mechanisms and survival rates. *Cell Microbiol* **13**:450-468.
1231. **Sagar, V., D. K. Bakshi, S. Nandi, N. K. Ganguly, R. Kumar, and A. Chakraborti.** 2004. Molecular heterogeneity among north Indian isolates of Group A *Streptococcus*. *Lett Appl Microbiol* **39**:84-88.
1532. **Sagar, V., R. Kumar, N. K. Ganguly, and A. Chakraborti.** 2008. Comparative analysis of emm type pattern of Group A *Streptococcus* throat and skin isolates from India and their association with closely related SIC, a streptococcal virulence factor. *BMC Microbiol* **8**:150.
1933. **Sakota, V., A. M. Fry, T. M. Lietman, R. R. Facklam, Z. Li, and B. Beall.** 2006. Genetically diverse group A streptococci from children in far-western Nepal share high genetic relatedness with isolates from other countries. *J Clin Microbiol* **44**:2160-2166.
2334. **Savage, R. S., K. Heller, Y. Xu, Z. Ghahramani, W. M. Truman, M. Grant, K. J. Denby, and D. L. Wild.** 2009. R/BHC: fast Bayesian hierarchical clustering for microarray data. *BMC Bioinformatics* **10**:242.
2635. **Sehr, P., K. Zumbach, and M. Pawlita.** 2001. A generic capture ELISA for recombinant proteins fused to glutathione S-transferase: validation for HPV serology. *J Immunol Methods* **253**:153-162.
2936. **Sette, A., and R. Rappuoli.** 2010. Reverse vaccinology: developing vaccines in the era of genomics. *Immunity* **33**:530-541.
3137. **Smeesters, P. R., A. Vergison, D. Campos, E. de Aguiar, V. Y. Miendje Deyi, and L. Van Melderen.** 2006. Differences between Belgian and Brazilian group A *Streptococcus* epidemiologic landscape. *PLoS One* **1**:e10.
3438. **Smoot, J. C., K. D. Barbian, J. J. Van Gompel, L. M. Smoot, M. S. Chaussee, G. L. Sylva, D. E. Sturdevant, S. M. Ricklefs, S. F. Porcella, L. D. Parkins, S. B. Beres, D. S. Campbell, T. M. Smith, Q. Zhang, V. Kapur, J. A. Daly, L. G. Veasy, and J. M. Musser.** 2002. Genome sequence and comparative microarray analysis of serotype M18 group A *Streptococcus* strains associated with acute rheumatic fever outbreaks. *Proc Natl Acad Sci U S A* **99**:4668-4673.
4039. **Stalhammar-Carlemalm, M., T. Areschoug, C. Larsson, and G. Lindahl.** 2000. Cross-protection between group A and group B streptococci due to cross-reacting surface proteins. *J Infect Dis* **182**:142-149.
4340. **Stalhammar-Carlemalm, M., T. Areschoug, C. Larsson, and G. Lindahl.** 1999. The R28 protein of *Streptococcus pyogenes* is related to several group B streptococcal surface proteins, confers protective immunity and promotes binding to human epithelial cells. *Mol Microbiol* **33**:208-219.
4741. **Steer, A. C., I. Law, L. Matatolu, B. W. Beall, and J. R. Carapetis.** 2009. Global emm type distribution of group A streptococci: systematic review and implications for vaccine development. *Lancet Infect Dis* **9**:611-616.
5042. **Talay, S. R., A. Zock, M. Rohde, G. Molinari, M. Oggioni, G. Pozzi, C. A. Guzman, and G. S. Chhatwal.** 2000. Co-operative binding of human fibronectin to

- 1 Sfbl protein triggers streptococcal invasion into respiratory epithelial cells. *Cell*  
2 *Microbiol* **2**:521-535.
343. **Terao, Y., S. Kawabata, E. Kunitomo, J. Murakami, I. Nakagawa, and S.**  
4 **Hamada.** 2001. Fba, a novel fibronectin-binding protein from *Streptococcus*  
5 *pyogenes*, promotes bacterial entry into epithelial cells, and the fba gene is positively  
6 transcribed under the Mga regulator. *Mol Microbiol* **42**:75-86.
744. **Zagursky, R. J., P. Ooi, K. F. Jones, M. J. Fiske, R. P. Smith, and B. A. Green.**  
8 2000. Identification of a *Haemophilus influenzae* 5'-nucleotidase protein: cloning of  
9 the nucA gene and immunogenicity and characterization of the NucA protein. *Infect*  
10 *Immun* **68**:2525-2534.
- 11  
12

1**Tables**2**Table 1** Oligonucleotides used in this study.

primer name	sequence 5' - 3'	Description	Annealing temperature (°C)
Scl F	G <b>CGAATTC</b> GAGGTTTCTTCTACGACTATGA	<b>EcoRI</b> restriction site	65
Scl R	G <b>CGTCGAC</b> ACGTCTGTGGTTGTTGGCTA	<b>SalI</b> restriction site	
PSNT F	G <b>CGAATTC</b> GATCAAGTTGATGTGCAATTCC	<b>EcoRI</b> restriction site	65
PSNT R	G <b>CGTCGAC</b> AGTGGAAGTAGAGATAGTATTT	<b>SalI</b> restriction site	
HYP F	G <b>CGAATTC</b> GTAGATGGCATCCCTCCAAT	<b>EcoRI</b> restriction site	65
HYP R	G <b>CGTCGAC</b> CTGACTCATGGGCCCTAA	<b>SalI</b> restriction site	
PSA F	G <b>CGAATTC</b> GTCAAAGAGCCGATTCTTAAACA	<b>EcoRI</b> restriction site	65
PSA R	G <b>CGTCGAC</b> TATTGCAGAGTGTCGTCCT	<b>SalI</b> restriction site	
R28 F	G <b>CGAATTC</b> TCTACAATTCCAGGGAGTGC	<b>EcoRI</b> restriction site	65
R28 R	G <b>CGTCGAC</b> CCCTTTGACTTGCTGATTTTTACC	<b>SalI</b> restriction site	
Scp F	G <b>CGGATCCA</b> ATACTGTGACAGAAGACACTCC	<b>BamHI</b> restriction site	55
Scp R	GCT <b>GTCGAC</b> TTATTACGCTCCTGCTCCTTGTTGGC G	<b>SalI</b> restriction site	

3

4

2

3

1**Table 2** Genes differentially conserved in the North and South Indian GAS  
 2populations (Fishers exact test,  $P < 0.05$ ).

GeneID	Gene name	Conservation (%)	
		North India	South India
SPy0116	hypothetical protein	56.9	91.9
SPy0317	conserved hypothetical protein	55.4	92.6
SPy1006	put. lysin - phage associated	96.9	85.3
SPy2009	hypothetical protein (transposase)	83.1	94.9
SpyM3_0130	streptolysin O	63.1	94.1
SpyM3_0304	conserved hypothetical protein	78.5	94.9
SpyM3_0343	hypothetical protein	93.8	100.0
SpyM3_0653	put. ABC transporter substrate-binding protein	87.7	98.5
SpyM3_0815	put. hemolysin III	70.8	90.4
SpyM3_0823	hypothetical protein	83.1	97.1
SpyM3_0833	put. citrate lyase beta subunit	86.2	97.1
SpyM3_0862	put. DNA/pantothenate metabolism flavoprotein	40.0	91.9
SpyM3_0999	hypothetical protein	93.8	100.0
SpyM3_1390	put. penicillin-binding protein 1A	92.3	99.3
SpyM3_1718	surface lipoprotein DppA	67.7	91.2
SpyM3_1762	hypothetical protein	93.8	99.3

1**Table 3** Distribution of genes encoding putative vaccine antigens in the Indian GAS  
 2population and comparison to isolates from the Netherlands (McMillan, 2006).

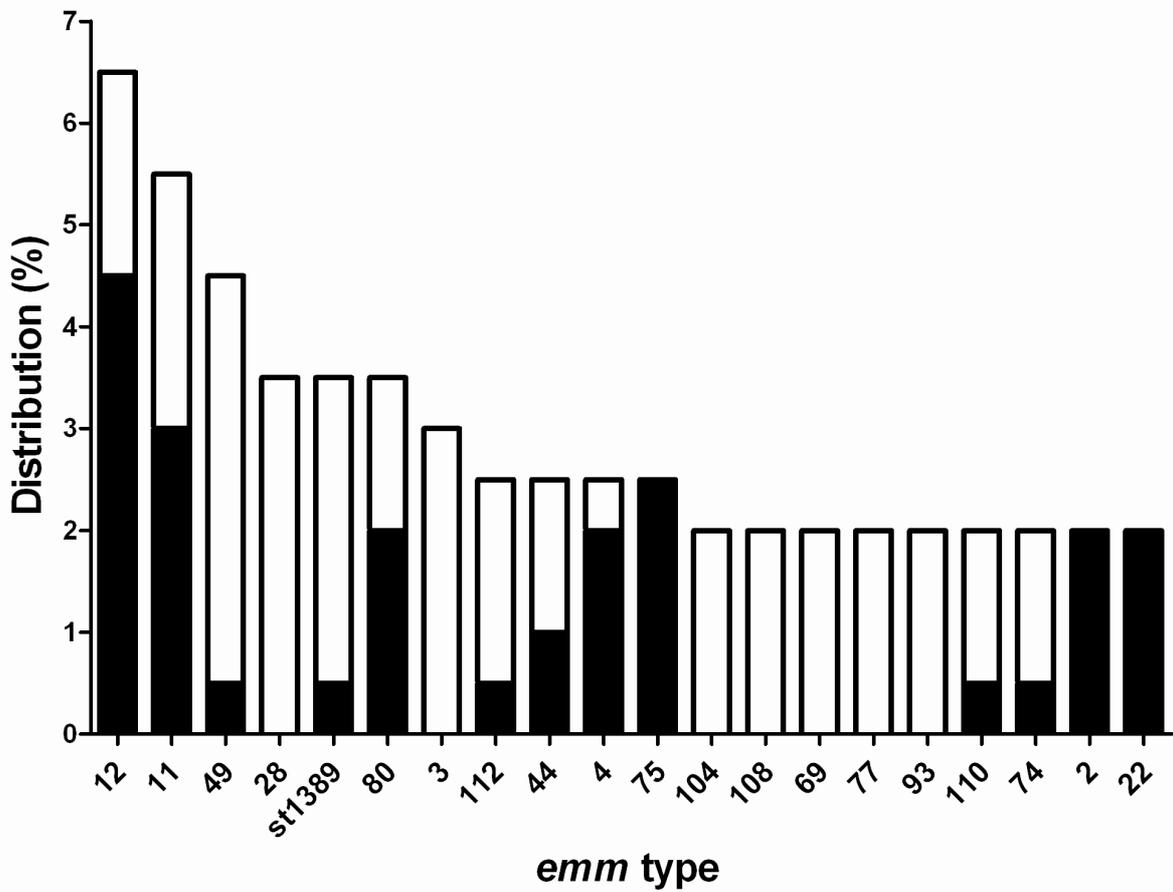
<b>Protein</b>	<b>GeneID</b>	<b>Distribution India %</b>	<b>Distribution The Netherlands %</b>
Collagen like surface protein (SCI)	spy1983	100	100
Putative secreted 5'-nucleotidase (PSNT)	spyM3_0591	100	100
C5a peptidase	SpyM3_1726	100	100
R28	AF091393	97.6	55
Putative surface antigen (PSA)	spyM3_0569	99.5	100
Hypothetical protein (HyP)	Spy2009	91.2	70
Sfbl	X67947	78.5	68

3

4

## Figures

2

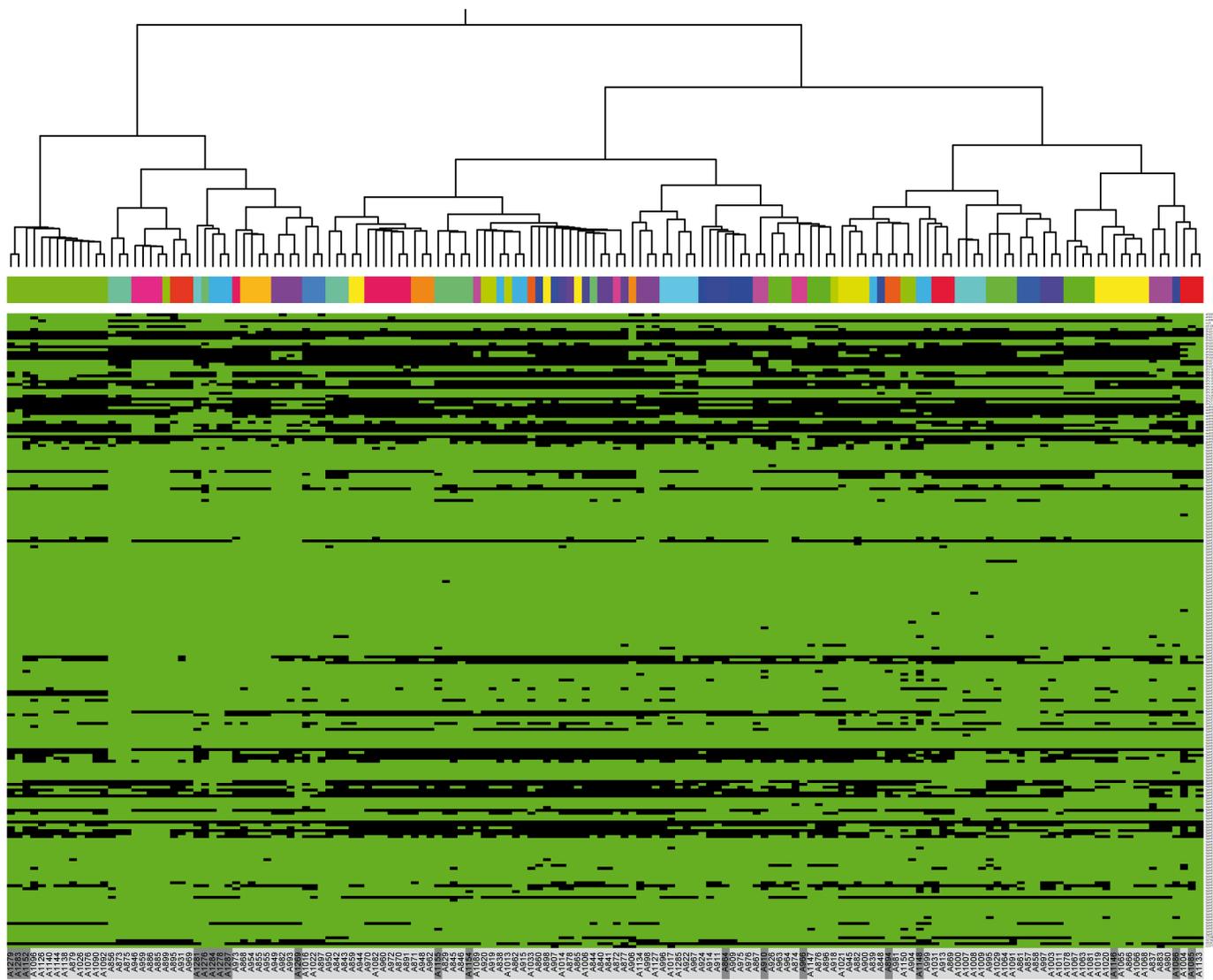


3

**Figure 1 Distribution of the 20 most common GAS emm types.** Black bars represent emm types found in North India, while open bars represent isolates found in South India. The 20 prevalent emm types only represented 58% of all GAS isolated in both regions.

8

- emm-types
- 1-2
  - 100
  - 102
  - 104
  - 108
  - 11
  - 110
  - 112
  - 113
  - 12
  - 15
  - 18
  - 2
  - 22
  - 28
  - 3
  - 4
  - 44
  - 49
  - 53
  - 60
  - 63
  - 69
  - 74
  - 75
  - 77
  - 80
  - 82
  - 85
  - 92
  - 93
  - st1389
  - st2147



Genes

1

\*

2

3

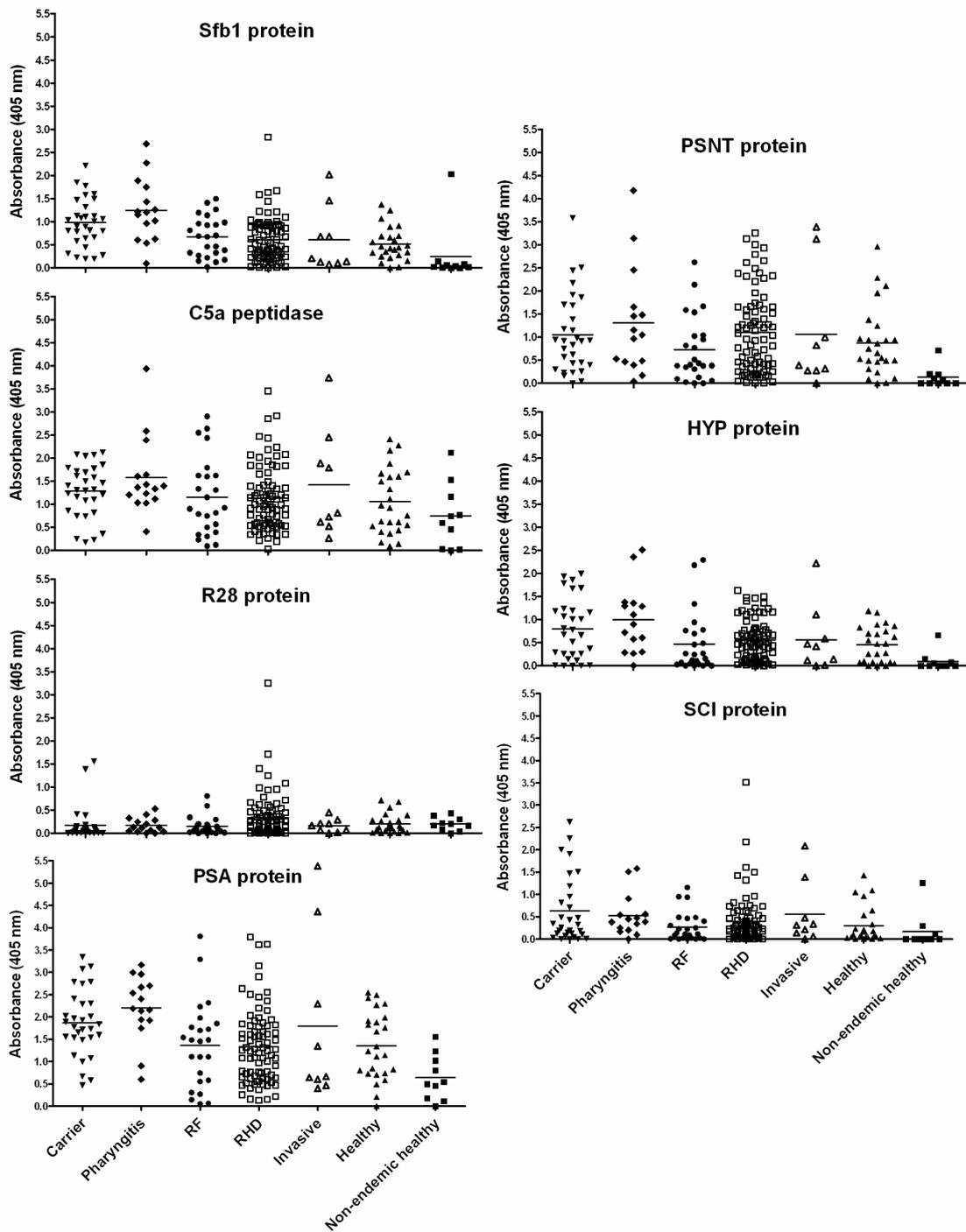
1

2

3**Figure 2** Bayesian agglomerative hierarchical clustering of 154 GAS isolates collected in North and South India. The dendrogram shows  
4the clustering of the isolates into groups with similar genetic profiles. Rows represent VF/ECP genes. Genes that are present in specific  
5isolates are shown in green and genes that are absent or highly divergent are shown in black. Columns represent isolates and the colour  
6bar below the dendrogram represents the different emm types. Strain numbers of invasive isolates are shaded in dark grey whereas non-  
7invasive isolates are shaded in light grey. The asterisk indicates a sub-cluster consisting solely of invasive isolates.

2

3



1

2**Figure 3 Serological responses to streptococcal antigens in sera.** Samples  
 3were collected from healthy individuals, asymptomatic carriers and patients  
 4presenting with pharyngitis, RF, RHD and invasive diseases. Control sera were  
 5collected from patients living in Germany. Statistical differences between the mean  
 6absorbance for each group were determined using one-way ANOVA and presented  
 7in supplementary table S4.

2

3