1 **Title: Genetic predictors of systemic sclerosis-associated interstitial lung disease: a**

2 **review of recent literature**

3 Running title: Genetics of SSc-ILD

4

5 Carmel J.W. Stock, Elisabetta A. Renzoni

6

7 **AFFILIATIONS:**

8 Interstitial Lung Disease Unit, Royal Brompton Hospital and Imperial College London

9

10 **Corresponding author:**

11 Dr Carmel Stock, Royal Brompton and Harefield NHS Foundation Trust, Sydney Street,
12 London, SW3 6NP, United Kingdom
13 Tel: +44 207 3518456

14 e-mail: c.stock@imperial.ac.uk

15

21

22 **CONFLICT OF INTEREST STATEMENT**

23 The authors declare no conflict of interest.

**ABSTRACT**

The interplay between genetic and environmental factors is likely involved in the pathogenesis of systemic sclerosis (SSc). Interstitial lung disease associated in the context of SSc (SSc-ILD) is associated with significant morbidity, and is the leading cause of death in SSc. The spectrum of SSc-ILD severity is wide, ranging from patients with only limited and inherently stable pulmonary involvement, to those with extensive and progressive pulmonary fibrosis. In order to provide accurate prognostic information for patients, and to initiate appropriate monitoring and treatment regimens, the ability to identify patients at risk of developing severe ILD early in the disease course is crucial. Identification of genetic variants involved in disease pathogenesis can not only potentially provide diagnostic/prognostic markers, but can also highlight dysregulated molecular pathways for therapeutic targeting. A number of genetic associations have been established for susceptibility to SSc, but far fewer studies have investigated genetic susceptibility to SSc-ILD specifically. In this review we present a summary of the studies assessing genetic associations with SSc-ILD.

**KEYWORDS:** Systemic sclerosis, SSc-ILD, pulmonary fibrosis, genetics, polymorphisms

**INTRODUCTION**

Systemic sclerosis (SSc) is a connective tissue disease characterised by immune activation, fibrosis of the skin and internal organs, and widespread vasculopathy. The pattern of internal organ involvement and the natural history of the disease are highly variable. The reported frequency of interstitial lung disease in SSc (SSc-ILD) varies from 25% to 90%, depending on the detection method and disease definition.[1,2] SSc-ILD is more common in patients with the diffuse form of skin involvement, and with anti topo-isomerase autoantibodies (ATA),[3]

47     although at least half of patients with SSc-ILD do not have ATA antibodies.[4] The prominent

48     pathological ILD pattern is non-specific interstitial pneumonia (NSIP).[5] The progression of

49     SSc-ILD is highly variable, with stable and limited disease observed in the majority of

50     patients, and severe progressive disease in a substantial minority.[6]

51     Evidence for a genetic predisposition to SSc includes the observation that disease prevalence

52     in relatives of patients with SSc is significantly higher than in the general population, with a

53     reported relative risk of disease of 13 in first degree relatives, and of 15 in siblings.[7]

54     Prevalence also varies according to ethnicity. In a large US population study, the prevalence

55     of SSc was higher in individuals of African descent compared to European descent, with an

56     adjusted prevalence ratio of 1.15.[8] Choctaw native Americans have the highest reported

57     prevalence in any population (66/100 000).[9] Compared to patients of African, Japanese, and

58     Choctaw descent, the frequency of ILD is lower in SSc patients of European descent, who

59     also seem to have slower decline in lung function and better survival rates.[10]

60     Specific non-overlapping antinuclear antibodies (ANAs), including anti-centromere

61     antibodies (ACA) and ATA, also known as Scl-70, are associated with different subsets of

62     SSc. ATA autoantibodies are strongly associated with the development of SSc-ILD, while

63     ACA are protective for ILD.[11] Twin studies have shown a high concordance for ANA

64     specificity, with 90% concordance in monozygotic twins compared to 40% concordance in

65     dizygotic twins, demonstrating a strong genetic influence on ANA status.[12]

66     Genetic associations with SSc as a whole have been recently extensively reviewed

67     elsewhere.[13,14] Similarly to autoimmune diseases, a predominant genetic effect is observed

68     within the human leukocyte antigen (HLA) region. However, HLA region associations are

69     mainly confined to subgroups of patients possessing specific autoantibodies. Non-HLA genes

70     consistently associated with SSc comprise genes involved in innate immunity as well as B

71    and T cell activation, including the highly repeatable associations with interferon regulatory

72    factor 5 (*IRF5*), signal transducer and activator of transcription 4 (*STAT4*), and cell receptor

73    CD3ζ (*CD247*). [15,13,14]

74

75    **GENETIC ASSOCIATION STUDIES WITH SSC-ILD**

76    Since the discovery in the 80s that ATA autoantibodies are strongly associated with SSc-ILD,

77    there has been limited progress in enabling prediction of which SSc patients will develop

78    significant ILD. A staging system, based on the extent of fibrosis on HRCT, integrated with

79    pulmonary function as needed, provides accurate prognostic information on the clinical

80    course of SSc-ILD.[6] However, this tool can only be utilised once interstitial lung disease has

81    developed. Identification of biological or genetic markers to enable, at the time of SSc

82    diagnosis, the discrimination of patients at higher risk of developing ILD, and prediction of

83    disease progression, would result in improved clinical management of these patients.

84

85    **Major histocompatibility complex**

86    A number of HLA alleles have been associated with SSc-ILD, summarised in Table 1.

87    However, many of these studies include only small numbers of patients with SSc-ILD.

88    Selected studies, including some of the larger ones, are discussed below.

89    Fanning *et al.* reported that the strongest risk factor for SSc-ILD in a UK population (47 SSc-

90    ILD/83 non-ILD) was a combination of ATA positivity, dcSSc, and HLA DRB1*11

91    (RR=21.9, p=0.0002). In the absence of these three risk factors, DRB1*301 was a risk marker

92    for SSc-ILD, with the highest relative risk seen in ATA negative patients (RR=7.5,

93  p=0.0001).[16] The HLA-DRB1*11 association with SSc-ILD has also been demonstrated in a

94  number of different populations including Spanish,[17] and Black South African.[18] In both an

95  initial and a separate Japanese replication cohort (1st cohort - 41 SSc-ILD/147 controls, 2nd

96  cohort - 40 SSc-ILD/83 controls), the DRB5*0105 allele was significantly more common in

97  SSc-ILD patients compared to healthy controls (OR=8.07, p<0.001 and OR=17.39, p=0.009

98  respectively).[19] A number of studies of HLA alleles in Han Chinese patients have recently

99  been published. The DQB1*0501 allele was significantly more frequent in SSc-ILD

100  (OR=5.03, $p=6 \times 10^{-7}$) compared to healthy controls in the study by Zhou et al. (134 SSc-

101  ILD/239 controls). However, DQB1*0501 was also found to be associated with SSc as a

102  whole, and there was no frequency difference between the patients with and without ILD

103  (p=0.9), indicating that this association may not be subtype specific. In a study of the DPB1

104  locus by Wang et al., (199 SSc-ILD/ 78 SSc no-ILD/480 controls), DPB1*0301 was

105  associated specifically with SSc-ILD (OR=3.86, $p<10^{-7}$), with no difference in allele

106  frequency between patients without ILD and healthy controls (p=0.79), and a significant

107  difference when the two patient groups were directly compared (OR=3.56, p=0.0069).

108  DPB1*1301 was also more common in the patient group with ILD than the controls

109  (OR=2.25, $p<3.3 \times 10^{-4}$), but not in patients without ILD (p=0.17).[20] In a study of the DRB1

110  locus (295 SSc-ILD/ 138 SSc no-ILD/ 458 controls), three alleles were all significantly more

111  common in SSc-ILD compared to controls, but only DRB1*0301 was not also significantly

112  more common in the patients without lung involvement compared to controls (OR=2.47,

113  p=0.0026).[21]

114

115  **Genome-wide association studies (GWAS)**

116  Although a number of genome-wide association studies (GWAS)[22,23,24,25] and Immunochip

117  studies[26,27] have targeted SSc as a whole, to date none have been specifically designed to

118  assess genetic determinants of SSc-ILD, possibly due to the limitations on achievable cohort

119  sizes. However, post-hoc analyses of data from one of the GWAS studies was performed to

120  investigate the impact of SSc-associated single nucleotide polymorphisms (SNPs) on survival

121  and severity of ILD,[23,28] discussed below in the section on *IRF5*.

122

123  **CANDIDATE GENE STUDIES**

124  The details of the candidate gene studies discussed in this review are summarised in Table 2.

125  *IRF5*

126  The transcription factor interferon regulatory factor 5 (IRF5) induces expression of interferon

127  A and B genes and pro-inflammatory cytokines, and is critical for antiviral immunity.[29] In a

128  French population (280 SSc-ILD/760 controls), *IRF5* rs2004640 was significantly associated

129  with SSc-ILD, even after adjusting for disease duration, cutaneous involvement, and ANA on

130  multivariate analysis (OR=1.38, p=0.016).[30] A similar association was observed in a Han

131  Chinese population (227 SSc-ILD/502 controls, OR=1.38, p=0.028).[31] A three SNP haplotype

132  containing rs2004640, as well as rs3757385 and rs10954213, is a marker for a five base-pair

133  insertion/deletion polymorphism in intron 1 of *IRF5*. Analysis of the individual SNPs of this

134  haplotype showed that rs3757385 (OR=1.42, p=$5.5 \times 10^{-3}$) and rs2004640 (OR=1.54,

135  p=$9.2 \times 10^{-5}$) were significantly associated with SSc-ILD (292 SSc-ILD/989 controls),

136  although only rs2004640 remained significant following conditional regression analysis.

137  Haplotype analysis of the three SNPs showed the haplotype comprising the protective allele

138  of each SNP was significantly less common in SSc-ILD compared to controls (OR=0.64,

139    p=3.7x10$^{-4}$), and compared to non-ILD SSc patients (n=397, p=0.018).[32] However, analysis

140    of data from the 2010 GWAS study[23] to investigate the impact of SSc-associated SNPs on

141    survival and severity of ILD, using % predicted FVC as a surrogate marker of ILD severity (1

142    443 SSc in survival analysis, 914 SSc in FVC% linear regression analysis), did not find

143    rs2004640, or the three SNP haplotype, to be associated with survival or ILD severity.

144    However, the minor allele of *IRF5* rs4728142 was associated with improved survival

145    (HR=0.75, p=0.002), independent of age of onset, gender, cutaneous involvement, and

146    ANA.[33] The minor allele was also associated with less severe ILD after taking disease

147    duration into account (mean difference=2.64, p=0.019). In addition, the number of rs4728142

148    minor alleles was associated with lower expression of *IRF5* in monocytes from both patients

149    and controls.[33] Meta-analysis of data from five European populations (total of 883 SSc-ILD/4

150    012 controls), tested the above mentioned *IRF5* SNPs rs2004640 and rs4728142, plus an

151    additional SNP, rs10488631, and found all three to be significantly associated with SSc-ILD

152    compared to controls. However, all three SNPs were also significantly associated with each

153    of the other subtypes tested (lcSSc, dcSSc, ATA, ACA, no ILD), and there was no difference

154    in allele frequencies when the patients with and without each phenotype, including with and

155    without ILD (883 SSc-ILD/1 797 SSc no-ILD), were compared directly, suggesting that these

156    *IRF5* polymorphisms may be associated with SSc as a whole rather than with any specific

157    subtype.[34]

158

159    ***STAT4***

160    Signal transducer and activator of transcription 4 (STAT4) is a transcription factor associated

161    with expression of type 1 interferons, IL-12, and IL-23. *STAT4* rs7574865 is associated with

162    systemic lupus erythematosus (SLE) and rheumatoid arthritis (RA).[35] This polymorphism has

163 also been associated with SSc-ILD (316 SSc-ILD/964 controls, OR=1.42, p=0.006), with an

164 additive effect of the *IRF5* SNP rs2004640, where carriage of at least three risk alleles of

165 these two SNPs is strongly associated with SSc-ILD (OR=1.79, p=0.002), with dcSSc and

166 ATA autoantibody being independent risk factors.[36] In a study of three *STAT4* SNPs in a Han

167 Chinese population (237 SSc-ILD/534 controls), rs7574865 and rs10168266 were both

168 significantly associated with SSc-ILD compared to controls (OR=1.86, p=$1.2 \times 10^{-4}$ and

169 OR=1.73, p=$7.7 \times 10^{-4}$ respectively). The third SNP tested, rs3821236, was also associated

170 with SSc-ILD, but significance was lost following Bonferroni correction (p=0.015, OR

171 =1.54).[37] However, in a study of six populations of European ancestry (total of 450 SSc-

172 ILD/3 113 controls), rs7574865 was not associated with SSc-ILD in any of the populations

173 individually or in a meta-analysis.[38]

174

175 *CD226*

176 *CD226* encodes DNAX accessory molecule 1, involved in cell-mediated cytotoxicity of T

177 and NK cells. The non-synonymous *CD226* SNP, rs763361, has been associated with a

178 number of autoimmune diseases including type 1 diabetes mellitus, multiple sclerosis, and

179 RA.[39] A meta-analysis of three European populations (total of 662 SSc-ILD/1 642 controls)

180 found this SNP to be associated with SSc-ILD (OR=1.27, p=$2.98 \times 10^{-4}$). A trend towards a

181 significant association with SSc-ILD was also seen when the populations were analysed

182 separately.[40] A haplotype of three SNPs in *CD226*, rs763361, rs34794968, and rs727088, has

183 been significantly associated with SLE and correlated with expression levels in T cells.[41]

184 Meta-analysis testing of this haplotype in seven European populations (729 SSc-ILD/3 966

185 controls) found none of the individual SNPs to be associated with SSc-ILD, but did find that

186 one of the haplotypes containing the previously associated allele of rs763361, was over-

187    represented in the SSc-ILD subgroup compared to controls (OR=1.27, p=0.032). A trend

188    towards a significant difference in frequency of this haplotype between SSc patients with and

189    without ILD was also seen (p=0.069).[42]

190

191    *NLRP1*

192    NLR family, pyrin domain containing 1 (*NLRP1*) is the activating platform required for

193    formation of the NALP1 inflammasome, involved in activation of inflammatory processes. In

194    a three-population meta-analysis study investigating five *NRLP1* SNPs (674 SSc-ILD/1 587

195    controls), rs8182352 was significantly associated with SSc-ILD compared to controls

196    (OR=1.19, p=0.0065), and compared to the non-ILD subgroup (n=1 255, OR not stated,

197    p=0.046). An additive effect of *NRLP1* rs8182352 with the *IRF5* rs2004640 and *STAT4*

198    rs7574865 risk alleles was identified, resulting in a 1.33-fold increase in OR for SSc-ILD

199    with each additional risk allele.[43]

200

201    *IRAK1*

202    Like many autoimmune diseases, SSc is characterised by female predominance,

203    approximately 4.6:1.[44] Interleukin-1 receptor-associated kinase 1 (*IRAK1*), a protein kinase

204    involved in signalling through the Toll-like receptors/IL-1R is located on the X chromosome.

205    Two non-synonymous SNPs, rs1059702 (Phe196Ser) and rs1059703 (Leu532Ser) are in

206    complete linkage disequilibrium, and the variant forms result in increased NFκ-B activity in

207    inflammatory responses.[45] The *IRAK1* variant rs1059702, was investigated in a large study of

208    SSc in three European populations. In the Italian cohort (167 SSc-ILD/ 509 controls) both the

209    T allele and TT genotype were significantly associated with SSc-ILD (OR=2.19, p=0.007 and

210 OR=2.19, p=0.039 respectively). Only the allelic association reached statistical significance

211 (OR=1.11, p=0.047) in the German cohort (167 SSc-ILD/1 083 controls), although the TT

212 genotype frequency was also non-significantly increased in the SSc-ILD group. In the French

213 cohort (334 SSc-ILD/625 controls), the frequency of both the rs1059702 T allele and the TT

214 genotype of were increased in SSc-ILD compared to controls, but neither reached statistical

215 significance (p=0.14 for allele, p-value for genotype not stated). When the three cohorts were

216 analysed together in a meta-analysis, both the T allele and the TT genotype were significantly

217 associated with SSc-ILD (OR=1.37, $1.99 \times 10^{-4}$ and OR=2.09, $9.05 \times 10^{-4}$ respectively).[46] The

218 findings of this study have been replicated in a subsequent study of women from four

219 European cohorts (461 SSc-ILD/2 043 controls, only meta-analysis of the cohorts reported),

220 which also found rs1059702 to be significantly associated with SSc-ILD when compared to

221 both controls (OR=1.30, $p=8.46 \times 10^{-3}$) and patients without ILD (OR=1.26, p=0.025).[47]

222

223 *CTGF*

224 Connective tissue growth factor (CTGF) induces myofibroblast differentiation and increased

225 extracellular matrix (ECM) production. Serum levels of CTGF correlate with the extent of

226 pulmonary fibrosis SSc-ILD.[48] In the study by Fonseca and Lindahl *et al*, the GG genotype of

227 *CTGF* rs6918698 was significantly associated with SSc-ILD compared to controls (207 SSc-

228 ILD/500 controls), even after adjusting for gender and ANA (OR=2.0, p<0.05). The disease

229 associated G allele results in significantly higher transcriptional activity, with allele specific

230 differential binding of the transcription factors Sp1 and Sp3 to this locus.[49] This association

231 was confirmed in a Japanese cohort (188 SSc-ILD/269 controls, OR=2.0, p<0.001).[50]

232 However, in a study of seven populations of European ancestry, no significant association

233 was detected in any of the populations whether tested separately, or together in a meta-

234  analysis (total of 1 180 SSc/1 784 controls), although no further information, including

235  patient numbers, is provided with regards to the subtype analyses.[51] The most recently

236  published study of this polymorphism was performed in a small Thai cohort (34 SSc-ILD/99

237  controls) with no association identified with SSc-ILD compared to controls.[52]

238

239  *CD247*

240  The *CD247* gene encodes the T-cell surface glycoprotein zeta chain (CD3ζ), a signalling

241  component of the T cell receptor (TCR)/CD3 complex. In a French population, *CD247*

242  rs2056626 was found to be associated with SSc-ILD compared to controls (346 SSc-ILD/990

243  controls, OR=0.65, p=6.8x10$^{-3}$), and not as strongly associated in patients with no lung

244  disease compared to controls (n=554, p=0.01).[53] This finding was however not replicated in a

245  study in a Han Chinese population (198 SSc-ILD/523 controls, p=0.83).[54]

246

247  **UNREPLICATED STUDIES WITH SMALL COHORT SIZES**

248  There are a number of additional studies identifying genetic associations with SSc-ILD, but

249  in cohorts which are too small to allow meaningful conclusions, and which have not been

250  repeated in additional cohorts. These studies have been included in Table 2 for completeness,

251  but the small number of patients and lack of replication must been borne in mind while

252  interpreting these associations.

253

254  **DISCUSSION**

255    For many of the associations presented in this review there have either been conflicting

256    results published from replication studies, or, following the initial association, there have

257    been no further studies published in independent cohorts. However, in recent years there has

258    been a move towards published association studies including both discovery and internal

259    replication cohorts with meta-analysis performed on the combined cohorts, allowing greater

260    confidence in the results compared to those from small, single cohort studies. SSc-ILD is a

261    complex disease with a number of genetic factors expected to be involved in susceptibility,

262    each with only relatively modest effects. As SSc-ILD is relatively rare, most of the published

263    studies are hampered by insufficient power to detect associations when SSc phenotypic

264    subgroups are analysed separately. This must be taken into account when interpreting

265    negative association results. The majority of published studies have been performed in

266    populations of European descent. However, the prevalence of ILD is lower in SSc patients of

267    European descent than in patients of African or Japanese descent. More studies in these non-

268    European populations may aid discovery of SSc-ILD associated genes. A large collaborative

269    project entitled 'Genome Research in African American Scleroderma Patients', led by the

270    National Human Genome Institute, is currently ongoing, with the aim of discovering common

271    and low-frequency variants associated with SSc susceptibility in African Americans.[55]

272    When studying clinical subgroups, the careful definition of phenotypes is crucial to allow

273    appropriate comparisons between patients with and without a phenotype, as well as between

274    different studies. In the field of SSc-ILD genetics this has so far been hampered by the lack of

275    a standardised definition of SSc-ILD, with studies using variable definitions for the presence

276    of ILD, including the presence of ground glass or reticular shadowing on HRCT, evidence of

277    fibrosis on chest radiograph, or impaired lung function.

The disease course of SSc-ILD is highly variable. Identification of specific genetic predictors of severe/progressive SSc-ILD is crucial, both from a pathogenesis and a clinical management perspective. Use of longitudinal clinical data to further define the SSc-ILD phenotype in terms of severity or rate of progression would enable investigation of genetic variants in relation to likelihood of ILD progression and severity. The recent staging system proposed by Goh et al.,[6] which subgroups SSc-ILD as limited or extensive based on rapid estimation of CT extent, supplemented, if necessary, with FVC levels, has been shown to provide accurate prospective prognostic separation. This system could be used to provide prognostic information, even when only limited clinical data is available. The ability of the Goh staging system to predict mortality is further increased when combined with short term pulmonary function trends.[56] Use of this surrogate of disease mortality means that long term follow-up data may not be required to investigate association of genetic variants with SSc-ILD outcome.

Finally, in most studies published so far, it is difficult to disentangle the association with autoantibodies linked with SSc-ILD, such as ATA, and associations with SSc-ILD per se. Although ATA autoantibodies have a high degree of specificity for the development of ILD in SSc, they are not a sensitive marker, as more than half of SSc-ILD patients are ATA autoantibody negative.[4] Therefore, subgroup analysis of SSc-ILD cohorts according to ANA status is required to allow separation of genetic variants associated with ATA or other antibodies and those associated specifically with development of lung fibrosis.

In SSc as a whole, the genetic risk appears to be mainly linked to immune pathway genes. Whether this is the same for the genetic risks for severe or progressive SSc-ILD remains to be determined. The genetic basis for SSc-ILD would seem to be different from that of the idiopathic interstitial pneumonias, as no association is observed with the *MUC5B* variant

13

302  strongly associated with IPF.[57] The fact that immunosuppressants are observed to stabilise

303  disease in the majority of patients with progressive lung fibrosis in the context of SSc

304  suggests that immune mediated pathways are key in driving the fibrotic process, but how this

305  translates into genetic predisposition will require further study.

306  Considering the expected small effect size from each individual genetic loci, and the need to

307  analyse SSc-ILD subgroups according to clinical and serological phenotypes, the requirement

308  for sufficiently large sample sizes with well characterised phenotypes is clear. National and

309  international collaborations will be indispensable to study genetic associations specific to

310  SSc-ILD, in order to enable collection of sufficiently large patient cohorts. It is also important

311  that replication of association studies is followed by functional work to determine the

312  biological significance of disease-associated genetic variants.

313

314  **CONCLUSIONS**

315  From the published literature presented in this review, genetic variation seems to be involved

316  in susceptibility to SSc-ILD. However, to date, no specific genetic variant has been

317  unequivocally associated with SSc-ILD and/or likelihood of ILD progression. By studying

318  sufficiently large cohorts of SSc with and without ILD, carefully staged, with reliable

319  longitudinal data, we should place ourselves in a better position to identify genes associated

320  with the development and rate of progression of SSc-ILD. Knowledge of the genetic

321  susceptibility to SSc-ILD should represent a stepping stone towards a better understanding of

322  the pathobiology of severe/progressive SSc-ILD, and should enable the identification of

323  prognostic and therapeutic targets in this debilitating and potentially fatal disease.

324

331    **CONFLICT OF INTEREST STATEMENT**

332    The authors declare no conflict of interest.

333

334    **REFERENCES**

335

336    1 Schurawitzki H, Stiglbauer R, Graninger W *et al.* Interstitial lung disease in progressive

337       systemic sclerosis: high-resolution CT versus radiography. *Radiology* 1990; **176**: 755-759.

338    2 Steen VD, Conte C, Owens GR, Medsger TA, Jr.. Severe restrictive lung disease in

339       systemic sclerosis. *Arthritis Rheum.* 1994; **37**: 1283-1289.

340    3 Morelli S, Barbieri C, Sgreccia A *et al.* Relationship between cutaneous and pulmonary

341       involvement in systemic sclerosis. *J.Rheumatol.* 1997; **24**: 81-85.

342    4 Gilchrist FC, Bunn C, Foley PJ *et al.* Class II HLA associations with autoantibodies in

343       scleroderma: a highly significant role for HLA-DP. *Genes Immun.* 2001; **2**: 76-81.

344    5 Bouros D, Wells AU, Nicholson AG *et al.* Histopathologic subsets of fibrosing alveolitis in

345       patients with systemic sclerosis and their relationship to outcome. *Am.J.Respir.Crit Care*

346       *Med.* 2002; **165**: 1581-1586.

347 6 Goh NS, Desai SR, Veeraraghavan S *et al.* Interstitial lung disease in systemic sclerosis: a

348     simple staging system. *Am.J.Respir.Crit Care Med.* 2008; **177**: 1248-1254.

349 7 Arnett FC, Cho M, Chatterjee S, Aguilar MB, Reveille JD, Mayes MD. Familial occurrence

350     frequencies and relative risks for systemic sclerosis (scleroderma) in three United States

351     cohorts. *Arthritis Rheum.* 2001; **44**: 1359-1362.

352 8 Mayes MD, Lacey JV, Jr., Beebe-Dimmer J *et al.* Prevalence, incidence, survival, and

353     disease characteristics of systemic sclerosis in a large US population. *Arthritis Rheum.*

354     2003; **48**: 2246-2255.

355 9 Arnett FC, Howard RF, Tan F *et al.* Increased prevalence of systemic sclerosis in a Native

356     American tribe in Oklahoma. Association with an Amerindian HLA haplotype. *Arthritis*

357     *Rheum.* 1996; **39**: 1362-1370.

358 10 Kuwana M, Kaburaki J, Arnett FC, Howard RF, Medsger TA, Jr., Wright TM. Influence of

359     ethnic background on clinical and serologic features in patients with systemic sclerosis and

360     anti-DNA topoisomerase I antibody. *Arthritis Rheum.* 1999; **42**: 465-474.

361 11 Hesselstrand R, Scheja A, Shen GQ, Wiik A, Akesson A. The association of antinuclear

362     antibodies with organ involvement and survival in systemic sclerosis.

363     *Rheumatology.(Oxford)* 2003; **42**: 534-540.

364 12 Feghali-Bostwick C, Medsger TA, Jr., Wright TM. Analysis of systemic sclerosis in twins

365     reveals low concordance for disease and high concordance for the presence of antinuclear

366     antibodies. *Arthritis Rheum.* 2003; **48**: 1956-1963.

367 13 Murdaca G, Contatore M, Gulli R, Mandich P, Puppo F. Genetic factors and systemic

368     sclerosis. *Autoimmun.Rev.* 2016

369 14 Chairta P, Nicolaou P, Christodoulou K. Genomic and genetic studies of systemic

370     sclerosis: A systematic review. *Hum.Immunol.* 2017; **78**: 153-165.

371 15 Ramos PS, Silver RM, Feghali-Bostwick CA. Genetics of systemic sclerosis: recent

372 advances. *Curr.Opin.Rheumatol.* 2015; **27**: 521-529.

373 16 Fanning GC, Welsh KI, Bunn C, Du BR, Black CM. HLA associations in three mutually

374 exclusive autoantibody subgroups in UK systemic sclerosis patients. *Br.J.Rheumatol.* 1998;

375 **37**: 201-207.

376 17 Simeon CP, Fonollosa V, Tolosa C *et al.* Association of HLA class II genes with systemic

377 sclerosis in Spanish patients. *J.Rheumatol.* 2009; **36**: 2733-2736.

378 18 Tikly M, Rands A, McHugh N, Wordsworth P, Welsh K. Human leukocyte antigen class II

379 associations with systemic sclerosis in South Africans. *Tissue Antigens* 2004; **63**: 487-490.

380 19 Odani T, Yasuda S, Ota Y *et al.* Up-regulated expression of HLA-DRB5 transcripts and

381 high frequency of the HLA-DRB5*01:05 allele in scleroderma patients with interstitial lung

382 disease. *Rheumatology.(Oxford)* 2012; **51**: 1765-1774.

383 20 Wang J, Guo X, Yi L *et al.* Association of HLA-DPB1 with scleroderma and its clinical

384 features in Chinese population. *PLoS.One.* 2014; **9**: e87363-

385 21 He D, Wang J, Yi L *et al.* Association of the HLA-DRB1 with scleroderma in Chinese

386 population. *PLoS.One.* 2014; **9**: e106939-

387 22 Zhou X, Tan FK, Wang N *et al.* Genome-wide association study for regions of systemic

388 sclerosis susceptibility in a Choctaw Indian population with high disease prevalence.

389 *Arthritis Rheum.* 2003; **48**: 2585-2592.

390 23 Radstake TR, Gorlova O, Rueda B *et al.* Genome-wide association study of systemic

391 sclerosis identifies CD247 as a new susceptibility locus. *Nat.Genet.* 2010; **42**: 426-429.

392 24 Allanore Y, Saad M, Dieude P *et al.* Genome-wide scan identifies TNIP1, PSORS1C1,

393 and RHOB as novel risk loci for systemic sclerosis. *PLoS.Genet.* 2011; **7**: e1002091-

394    25 Zhou X, Lee JE, Arnett FC *et al.* HLA-DPB1 and DPB2 are genetic loci for systemic

395        sclerosis: a genome-wide association study in Koreans with replication in North Americans.

396        *Arthritis Rheum.* 2009; **60**: 3807-3814.

397    26 Mayes MD, Bossini-Castillo L, Gorlova O *et al.* Immunochip analysis identifies multiple

398        susceptibility loci for systemic sclerosis. *Am.J.Hum.Genet.* 2014; **94**: 47-61.

399    27 Zochling J, Newell F, Charlesworth JC *et al.* An Immunochip based interrogation of

400        scleroderma susceptibility variants identifies a novel association at DNASE1L3. *Arthritis*

401        *Res.Ther.* 2014; **16**: 438

402    28 Gorlova O, Martin JE, Rueda B *et al.* Identification of novel genetic markers associated

403        with clinical phenotypes of systemic sclerosis through a genome-wide association strategy.

404        *PLoS.Genet.* 2011; **7**: e1002178

405    29 Yanai H, Chen HM, Inuzuka T *et al.* Role of IFN regulatory factor 5 transcription factor in

406        antiviral immunity and tumor suppression. *Proc.Natl.Acad.Sci.U.S.A* 2007; **104**: 3402-

407        3407.

408    30 Dieude P, Guedj M, Wipff J *et al.* Association between the IRF5 rs2004640 functional

409        polymorphism and systemic sclerosis: a new perspective for pulmonary fibrosis. *Arthritis*

410        *Rheum.* 2009; **60**: 225-233.

411    31 Wang J, Yi L, Guo X *et al.* Association of the IRF5 SNP rs2004640 with systemic

412        sclerosis in Han Chinese. *Int.J.Immunopathol.Pharmacol.* 2014; **27**: 635-638.

413    32 Dieude P, Dawidowicz K, Guedj M *et al.* Phenotype-haplotype correlation of IRF5 in

414        systemic sclerosis: role of 2 haplotypes in disease severity. *J.Rheumatol.* 2010; **37**: 987-

415        992.

416    33 Sharif R, Mayes MD, Tan FK *et al.* IRF5 polymorphism predicts prognosis in patients

417        with systemic sclerosis. *Ann.Rheum.Dis.* 2012; **71**: 1197-1202.

418  34 Carmona FD, Martin JE, Beretta L *et al.* The systemic lupus erythematosus IRF5 risk

419  haplotype is associated with systemic sclerosis. *PLoS.One.* 2013; **8**: e54419

420  35 Remmers EF, Plenge RM, Lee AT *et al.* STAT4 and the risk of rheumatoid arthritis and

421  systemic lupus erythematosus. *N.Engl.J.Med.* 2007; **357**: 977-986.

422  36 Dieude P, Guedj M, Wipff J *et al.* STAT4 is a genetic risk factor for systemic sclerosis

423  having additive effects with IRF5 on disease susceptibility and related pulmonary fibrosis.

424  *Arthritis Rheum.* 2009; **60**: 2472-2479.

425  37 Yi L, Wang JC, Guo XJ *et al.* STAT4 is a genetic risk factor for systemic sclerosis in a

426  Chinese population. *Int.J.Immunopathol.Pharmacol.* 2013; **26**: 473-478.

427  38 Rueda B, Broen J, Simeon C *et al.* The STAT4 gene influences the genetic predisposition

428  to systemic sclerosis phenotype. *Hum.Mol.Genet.* 2009; **18**: 2071-2077.

429  39 Hafler JP, Maier LM, Cooper JD *et al.* CD226 Gly307Ser association with multiple

430  autoimmune diseases. *Genes Immun.* 2009; **10**: 5-10.

431  40 Dieude P, Guedj M, Truchetet ME *et al.* Association of the CD226 Ser(307) variant with

432  systemic sclerosis: evidence of a contribution of costimulation pathways in systemic

433  sclerosis pathogenesis. *Arthritis Rheum.* 2011; **63**: 1097-1105.

434  41 Lofgren SE, Delgado-Vega AM, Gallant CJ *et al.* A 3'-untranslated region variant is

435  associated with impaired expression of CD226 in T and natural killer T cells and is

436  associated with susceptibility to systemic lupus erythematosus. *Arthritis Rheum.* 2010; **62**:

437  3404-3414.

438  42 Bossini-Castillo L, Simeon CP, Beretta L *et al.* A multicenter study confirms CD226 gene

439  association with systemic sclerosis-related pulmonary fibrosis. *Arthritis Res.Ther.* 2012; **14**:

440  R85

441  43 Dieude P, Guedj M, Wipff J *et al.* NLRP1 influences the systemic sclerosis phenotype: a

442     new clue for the contribution of innate immunity in systemic sclerosis-related fibrosing

443     alveolitis pathogenesis. *Ann.Rheum.Dis.* 2011; **70**: 668-674.

444  44 Arora-Singh RK, Assassi S, del Junco DJ *et al.* Autoimmune diseases and autoantibodies

445     in the first degree relatives of patients with systemic sclerosis. *J.Autoimmun.* 2010; **35**: 52-

446     57.

447  45 Liu G, Tsuruta Y, Gao Z, Park YJ, Abraham E. Variant IL-1 receptor-associated kinase-1

448     mediates increased NF-kappa B activity. *J.Immunol.* 2007; **179**: 4125-4134.

449  46 Dieude P, Bouaziz M, Guedj M *et al.* Evidence of the contribution of the X chromosome

450     to systemic sclerosis susceptibility: association with the functional IRAK1 196Phe/532Ser

451     haplotype. *Arthritis Rheum.* 2011; **63**: 3979-3987.

452  47 Carmona FD, Cenit MC, Diaz-Gallo LM *et al.* New insight on the Xq28 association with

453     systemic sclerosis. *Ann.Rheum.Dis.* 2013

454  48 Sato S, Nagaoka T, Hasegawa M *et al.* Serum levels of connective tissue growth factor are

455     elevated in patients with systemic sclerosis: association with extent of skin sclerosis and

456     severity of pulmonary fibrosis. *J.Rheumatol.* 2000; **27**: 149-154.

457  49 Fonseca C, Lindahl GE, Ponticos M *et al.* A polymorphism in the CTGF promoter region

458     associated with systemic sclerosis. *N.Engl.J.Med.* 2007; **357**: 1210-1220.

459  50 Kawaguchi Y, Ota Y, Kawamoto M *et al.* Association study of a polymorphism of the

460     CTGF gene and susceptibility to systemic sclerosis in the Japanese population.

461     *Ann.Rheum.Dis.* 2009; **68**: 1921-1924.

462  51 Rueda B, Simeon C, Hesselstrand R *et al.* A large multicentre analysis of CTGF -945

463     promoter polymorphism does not confirm association with systemic sclerosis susceptibility

464     or phenotype. *Ann.Rheum.Dis.* 2009; **68**: 1618-1620.

465    52 Louthrenoo W, Kasitanon N, Wichainun R *et al.* Lack of CTGF*-945C/G Dimorphism in

466        Thai Patients with Systemic Sclerosis. *Open.Rheumatol.J.* 2011; **5**: 59-63.

467    53 Dieude P, Boileau C, Guedj M *et al.* Independent replication establishes the CD247 gene

468        as a genetic systemic sclerosis susceptibility factor. *Ann.Rheum.Dis.* 2011; **70**: 1695-1696.

469    54 Wang J, Yi L, Guo X *et al.* Lack of Association of the CD247 SNP rs2056626 with

470        Systemic Sclerosis in Han Chinese. *Open.Rheumatol.J.* 2014; **8**: 43-45.

471    55 Genome Research in African American Scleroderma Patients

472        http://www.srfcure.org/research/featured-research

473    56 Goh NS, Hoyles RK, Denton CP *et al.* Short term pulmonary function trends are

474        predictive of mortality in interstitial lung disease associated with systemic sclerosis.

475        *Arthritis Rheumatol.* 2017 [Epub ahead of print]

476    57 Stock CJ, Sato H, Fonseca C *et al.* Mucin 5B promoter polymorphism is associated with

477        idiopathic pulmonary fibrosis but not with development of lung fibrosis in systemic

478        sclerosis or sarcoidosis. *Thorax* 2013; **68**: 436-441.

479    58 Gladman DD, Kung TN, Siannis F, Pellett F, Farewell VT, Lee P. HLA markers for

480        susceptibility and expression in scleroderma. *J.Rheumatol.* 2005; **32**: 1481-1487.

481    59 Zhou XD, Yi L, Guo XJ *et al.* Association of HLA-DQB1*0501 with scleroderma and its

482        clinical features in Chinese population. *Int.J.Immunopathol.Pharmacol.* 2013; **26**: 747-751.

483    60 Briggs DC, Vaughan RW, Welsh KI, Myers A, duBois RM, Black CM. Immunogenetic

484        prediction of pulmonary fibrosis in systemic sclerosis. *Lancet* 1991; **338**: 661-662.

485    61 Zhao W, Yue X, Liu K *et al.* The status of pulmonary fibrosis in systemic sclerosis is

486        associated with IRF5, STAT4, IRAK1, and CTGF polymorphisms. *Rheumatol.Int.* 2017

487    62 Kowal-Bielecka O, Chwiesko-Minarowska S, Bernatowicz PL *et al.* The arachidonate 5-

488        lipoxygenase activating protein gene polymorphism is associated with the risk of

489 scleroderma-related interstitial lung disease: a multicentre European Scleroderma Trials and

490 Research group (EUSTAR) study. *Rheumatology.(Oxford)* 2017; **56**: 844-852.

491 63 Hoshino K, Satoh T, Kawaguchi Y, Kuwana M. Association of hepatocyte growth factor

492 promoter polymorphism with severity of interstitial lung disease in Japanese patients with

493 systemic sclerosis. *Arthritis Rheum.* 2011; **63**: 2465-2472.

494 64 Kawaguchi Y, Tochimoto A, Ichikawa N *et al.* Association of IL1A gene polymorphisms

495 with susceptibility to and severity of systemic sclerosis in the Japanese population.

496 *Arthritis Rheum.* 2003; **48**: 186-192.

497 65 Beretta L, Bertolotti F, Cappiello F *et al.* Interleukin-1 gene complex polymorphisms in

498 systemic sclerosis patients with severe restrictive lung physiology. *Hum.Immunol.* 2007;

499 **68**: 603-609.

500 66 Rech TF, Moraes SB, Bredemeier M *et al.* Matrix metalloproteinase gene polymorphisms

501 and susceptibility to systemic sclerosis. *Genet.Mol.Res.* 2016; **15**:

502 67 Manetti M, Ibba-Manneschi L, Fatini C *et al.* Association of a functional polymorphism in

503 the matrix metalloproteinase-12 promoter region with systemic sclerosis in an Italian

504 population. *J.Rheumatol.* 2010; **37**: 1852-1857.

505 68 Sumita Y, Sugiura T, Kawaguchi Y *et al.* Genetic polymorphisms in the surfactant proteins

506 in systemic sclerosis in Japanese: T/T genotype at 1580 C/T (Thr131Ile) in the SP-B gene

507 reduces the risk of interstitial lung disease. *Rheumatology.(Oxford)* 2008; **47**: 289-291.

508 69 Ates O, Musellim B, Ongen G, Topal-Sarikaya A. NRAMP1 (SLC11A1): a plausible

509 candidate gene for systemic sclerosis (SSc) with interstitial lung involvement.

510 *J.Clin.Immunol.* 2008; **28**: 73-77.

511

512    **Table 1. HLA associations with SSc-ILD**

513    Corrected p values given where available. ORs are shown as OR (95% confidence interval),

514    where available.

515

516    **Table 2. Non-HLA associations with SSc-ILD**

517    Corrected p values given where available. ORs are shown as OR (95% confidence interval),

518    where available. $^{†}$= meta-analysis or previously published studies. $^{§}$= total number of SSc

519    patients, when SSc-ILD number not given. $^{¶}$= meta-analysis of the different populations

520    included.