# The Evolutionary Moulding in plant-microbial symbiosis: matching population diversity of rhizobial *nod*A and legume *NFR5* genes

*Anna A. Igolkina[1, 2*], Georgii A. Bazykin[3, 4], Elena P. Chizhevskaya[1], Nikolai A. Provorov[1],*

*Evgeny E. Andronov[1, 5, 6]*

[1]*All-Russian Research Institute of Agricultural Microbiology of the Russian Academy of Agricultural Sciences, Russia,*

[2]*Saint Petersburg State Polytechnic University, Russia,*

[3]*Institute for Information Transmission Problems (Kharkevich Institute) of the Russian Academy of Sciences, Russia,*

[4]*Skolkovo Institute of Science and Technology, Russia,*

[5]*V V Dokuchaeva Soil Institute, Russian Academy of Agricultural Sciences, Russia,*

[6]*Saint Petersburg State University, Russia*

[*]*igolkinaanna11@gmail.com*

## Abstract

We propose the Evolutionary Moulding hypothesis that population diversities of partners in nitrogen-fixing rhizobium-legume symbiosis are matched, and tested it in nucleotide polymorphism of symbiotic genes encoding two components of the plant-bacteria signalling system. The first component is the rhizobial *nod*A acyltransferase involved in the fatty acid tail decoration of Nod factor (rhizobia signalling molecule). The second component is the plant *NFR5* receptor, putatively required for Nod-factor binding.

We collected three wild growing legume species together with soil samples adjacent to the roots (soil pool) from one large 25-year fallow: *Vicia sativa*, *Lathyrus pratensis* and *Trifolium hybridum* nodulated by one of the two *Rhizobium leguminosarum* biovars (*viciae* and *trifolii*). For each plant species we prepared three pools for DNA extraction: the plant pool (30 plant indiv.), the nodule pool (90 nodules) and the soil pool (30 samples). *NFR5* gene libraries from the plant pool and *nod*A gene libraries from nodule and soil pools were sequenced by Sanger technology and High-throughput pyrosequencing, respectively. Analysis of the data demonstrated concordance in population diversities of one symbiotic partner (rhizobia) the second partner (legume host), in line with the Evolutionary Moulding hypothesis. This effect

35    was evinced by the following observations for each plant species: (1) significantly increased

36    diversity in the nodule *nod*A popset (set of gene sequences derived from the nodule population)

37    compared to the soil popset; (2) a monotonic relationship between the diversity in the plant

38    *NFR5* gene popset and the nodule rhizobial *nod*A gene popset; and (3) higher topological

39    similarity of the *NFR5* gene tree with the *nod*A gene tree of the nodule popset, than with the

40    *nod*A gene tree of the soil popset. Both nonsynonymous diversity and Tajima's D were increased

41    in the nodule popsets compared to the soil popsets, consistent with relaxation of negative

42    selection and/or admixture of balancing selection underlying the Evolutionary Moulding effect.

43    We propose that the observed genetic concordance arises from the selection of particular

44    characteristics of the nodule *nod*A genes by the host plant.

45

46

**Introduction**

47

48

49      One of the earliest studied types of symbiosis, the host–parasite interaction, was described

50      by Flor's Gene-for-Gene concept (Flor 1971; Flor 1942) and, in fact, the first mathematical

51      model of coevolution was explicitly based on the assumption of a Gene-for-Gene (GFG)

52      interaction (Thompson and Burdin 1992; Mode 1958). Further analyses of host-parasite

53      interactions revealed alternative concepts, namely matching-allele(MA) (Frank 1994), inverse-

54      matching-allele (IMA) (Otto and Michalakis 1998) and inverse-gene-for-gene (IGFG) (Fenton,

55      Antonovics, and Brockhurst 2009), which together with the GFG represent the opposite end of

56      the same continuum of host–parasite specificity (Agrawal and Lively 2002). These theoretical

57      concepts, developed for antagonistic systems, found their reflection in mutualistic symbiotic

58      systems (Lewis-Henderson and Djordjevic 1991; Sadowsky et al. 1991; Cregan, Sadowsky, and

59      Keyser 1991; Parker 1999; Sachs, Essenberg, and Turcotte 2011). Later, various mathematical

60      and theoretical models were developed to describe different types of symbiotic interactions

61      (Yoder 2016; Kwiatkowski, Engelstadter, and Vorburger 2012). One aspect of symbiotic

62      interactions that previously had much attention in models is the difference in evolutionary rates

63      between interacting species. For example, the so-called Red Queen model (Van Valen 1973;

64      Paterson et al. 2010; Pal et al. 2007) suggests the metaphor of an evolutionary arms race,

65      whereby a more rapidly evolving species in a host–parasite interaction has an adaptive advantage

66      over a more slowly evolving species. For mutualistic symbiosis, the alternative Red King model

67      proposes that the slower-evolving species are favoured (Bergstrom and Lachmann 2003).

68      However, the recent study showed that the preference of the particular model is not uniquely

69      determined by the type of symbiosis (Rubin and Moreau 2016). Specifically, the Red Queen

70      model initially introduced for antagonism was detected in the ant-plant mutualism. Here, we

71      suggested another aspect of symbiotic interactions that can potentially provide the basis for

72      models describing both mutualism and antagonism.

73

74      Our hypothesis, which we refer to as Evolutionary Moulding, is that the population

75      diversity levels and structures are matched between symbionts. We proposed that this hypothesis

76      can be a common aspect of symbiotic systems and tested it in the nitrogen-fixing rhizobium-

77      legume symbiosis. Previous analysis of the symbiosis between *Rhizobium galegae* and its host

78      plant *Galega* indicated correspondence of population diversity levels between microsymbionts

79      and the host *Galega* species (Andronov et al. 2003; Österman et al. 2011). In particular, a more

80      genetically diverse *Galega orientalis* population harbours a more diverse root-nodule rhizobial

81      population, while its less diverse sympatric counterpart *Galega officinalis* forms symbiosis with

82    a less diverse rhizobial population. This observation is related to the well-studied phenomenon of

83    shaping the genetic structure of the rhizobial population through the selection of specific

84    rhizobial genotypes by the host plant (Paffetti et al. 1998; Laguerre et al. 2003; Heath and Tiffin

85    2007; Depret and Laguerre 2008). The reverse effect, i.e. plant evolution driven by rhizobia, has

86    also been discussed (Martínez-Romero 2009), and is consistent with the recently detected

87    evidence for balancing selection on the symbiotic legume genes in a large legume genome

88    dataset (Yoder 2016). Therefore, we expected that the interplay of symbiotic populations leads to

89    concordance between the diversity levels in their symbiotic genes.

90

91    In our study, we focused on two genes encoding the essential components of rhizobium-

92    legume signalling system which are associated with each other through a lipochito-

93    oligosaccharide called Nod factor (Figure 1). The first component is the rhizobial *nod*A gene

94    which encodes an acyltransferase enzyme essential in Nod factor (NF) biosynthesis, specifically

95    in the attachment of the long-chain fatty acid tail to the oligosaccharide backbone (Dénarié,

96    Debellé, and Promé 1996; Esseling and Emons 2004; Oldroyd 2013). The second component is

97    one of the plant symbiotic receptor genes, *NFR5*, which is a homologue of *LjNFR5*, *MtNFP*,

98    *PsSym10* genes. Its product recognises NFs (signalling molecules) by three extracellular LysM

99    domains and triggers the formation of root nodule primordia giving the green light to the process

100    of bacterial infection (Oldroyd 2013). The NFs are major determinants of host specificity:

101    rhizobia produce NFs with different structures and host plants percept only those Nod factors

102    that have a certain composition (Mergaert and Montagu 1997). The variation of NFs structure is

103    observed not only between rhizobia species but also at the intra-species level (Spaink 2000), one

104    rhizobia species produces a mixture of NFs that vary in the fatty tail modifications. As proposed,

105    the *nod*A product can vary in its fatty acid specificity, thus contributing to the bacterial host

106    range (Moulin, Béna, and Stępkowski 2004; Dénarié, Debellé, and Promé 1996; Ritsema et al.

107    1996; Roche et al. 1996). It is logical to assume that the *nod*A gene diversity in a rhizobial

108    population can reflect the structural variation of NFs produces by this population. Indeed, it has

109    been shown that minor differences in the structure of fatty acids tail can affect intra-species host

110    specificity (Li et al. 2011). On the host-plant side, NFs are recognized by high-affinity legume

111    receptors (Broghammer et al. 2012; Moulin, Béna, and Stępkowski 2004). Studies on the model

112    legumes revealed *NFR5* as one of the major receptors to percept NFs (Radutoiu et al. 2007).

113    Mutant analysis showed that single amino acid differences in one domain of the *NFR5* receptor

114    change recognition of Nod factor variants (Broghammer et al. 2012; Radutoiu et al. 2007). Such

115    mediation of NFs between rhizobial *nod*A and legume *NFR5* genes make them good candidates

116    for Evolutionary Moulding.

117

118      Moreover, it has been shown that the topology of the *nod*A gene tree follows the
119    corresponding host plant tree more strictly than the 16S rRNA-based rhizobial phylogeny
120    (Dobert, Breil, and Triplett 1994; Suominen et al. 2016). Thus, we proposed that the
121    Evolutionary Moulding can be observed by comparing the genetic diversity levels of the
122    symbiotic genes of both partners.

123

124      We investigated the *nod*A-*NFR5* symbiotic systems of three wild growing legume
125    species (*Vicia sativa*, *Lathyrus pratensis* and *Trifolium hybridum*) and their rhizobial
126    microsymbionts. Sampling of the experimental material was performed uniformly on the one
127    large natural fallow (more than 25 years) field in order to avoid the influence of ecological
128    factors. *Vicia* and *Lathyrus* species represent the same cross-inoculation group nodulated with
129    *Rhizobium leguminosarum* bv. *viciae* strains, while *Trifolium* belongs to a separate group
130    nodulated with *R.leguminosarum* bv. *trifolii*. One of the important traits of the rhizobium-legume
131    symbiosis is the annual cycle of rhizobia, consisting of nodule formation with consequent
132    amplification of rhizobia inside of the nodule, followed by a release of the nodule rhizobia back
133    into the soil after nodule degradation probably leading to an increase of the frequency of this
134    rhizobial genotype in the soil. Therefore, we analysed both soil and nodule populations of
135    rhizobia, which affect each other.

136

137      Testing the Evolutionary Moulding hypothesis required comparison of structural
138    (topological) characteristics of plant and rhizobia populations. Traditionally, the topological
139    similarity between two populations is estimated as the congruence of two respective labelled
140    trees (Leigh et al. 2011). Here, we proposed a novel method to compare topologies of two gene
141    trees with unlabelled leaves. The method was based on the gCEED approach (Choi and Gomez
142    2009) that translates each population to the Gaussian mixture model in a K-dimensional space.
143    This method can be classified as a kind of beta-diversity metrics, which, by analogy with
144    taxonomic (Jost, Chao, and Chazdon 2011) and phylogenetic (e.g. UniFrac (Lozupone et al.
145    2011)), could be denoted as "topological beta-diversity". We applied it to show that the tree
146    structures are concordant between the two symbiont species.

147

148      The main aim of this study was to estimate and to compare the population diversity of the
149    symbiotic genes between the soil and nodule rhizobial populations and to link the observed
150    differences to those between the symbiotic determinants of the host plants. Using common

151     diversity statistics and a novel approach for comparing the structures of partner's populations

152     from pooled sequencing data we confirmed the Evolutionary Moulding hypothesis.

153

154                           **Materials and Methods**

155

156     **Sampling**

157

158        Three wild growing legume species (30 samples per species) together with rhizosphere soil

159     - the common vetch *Vicia sativa*, the meadow vetchling *Lathyrus pratensis* and the alsike clover

160     *Trifolium hybridum* – were uniformly collected from the large natural fallow (more than 25

161     years) field near the town Vyritsa (Gatchinskii region of Leningradskaya oblast, Russia,

162     59°24′7.74″ N; 30°15′28.74″ E). All sampled plants had formed nitrogen-fixing symbiotic root

163     nodules which were selected and thoroughly washed. Soil samples were collected from the close

164     proximity to the plant roots (1-5 cm). Here and below we named these samples as "soil

165     samples". Three nodules from each plant sample were picked from the main or the closest to the

166     main lateral roots. For each legume species we prepared three pools for DNA extraction: the

167     plant pool (30 leaf pieces, 0.1g each), the nodule pool (90 nodules, 3 nodules per plant

168     individual) and the soil pool (30 soil samples , 0.2g each).

169

170     ***nod*A amplification and sequencing**

171

172        DNA was isolated from soil and nodule pools by bead beating homogenization (Precellys

173     24) and purification (PowerSoil DNA Isolation Kit, MoBio, United States). Two pairs of nested

174     degenerate oligonucleotide primers were designed for *nod*A gene of *R.leguminosarum* bv. *viciae*

175     and bv. *trifolii*. The first round of nested PCR with external primers – forward (5'-

176     DGGHYTGTAYGGAGTGC-3') and reverse (5'- AGYTCSSACCCRTTT -3') – produces a

177     324bp amplicon product; the second round with inner primers – forward (5'-

178     YTDGGMATCGCHCACT-3') and reverse (5'- RDACGAGBACRTCTTCRGT-3') – produces

179     a 217 bp amplicon product. The reaction conditions in the first and the second round of the PCR

180     consisted of the initial denaturation step at 94°C for 3 min followed by 35 cycles with

181     denaturation at 94°C for 30 s, primer annealing at 50°C for 30 s and extension at 72°C for 1 min.

182     The bar-coded PCR products from six *nod*A libraries (soil and nodule libraries from three plants

183     species) were sequenced with a Roche 454 GS Junior (following the manufacturer's protocols)

184     generating an average of 3000-4000 reads per library. All obtained sequences were subjected to

185     filtration by quality (quality score higher than 25), length (longer than 170 bp) and separating

186  into libraries according to barcodes in QIIME. The sequencing data were deposited at the NCBI

187  short read archive under the bioproject number PRJNA297503. We introduced the term "popset"

188  to designate a set of *nod*A gene sequences from nodule or soil rhizobia population.

189

190  The multiple alignments of the remained sequences within each popset were performed

191  with ClustalW as implemented in MEGA (Tamura et al. 2013). Sequences with frameshift errors

192  were removed. The resultant multiple alignment for each popset did not contain gaps.

193

194  **NFR5 amplification and sequencing**

195

196  DNA from plant leaf pools was isolated by AxioPrep kit (Axigen) and was used as the

197  template DNA for PCR amplifications. Approximately 0.9 kb DNA fragments encoding all three

198  LysM domains of the plant receptor gene, *NFR5*, were amplified with the following pairs of

199  primers: forward "NFR5-for4" (5'AAGTCTTGGTTGTTACTTGCC-3') and reverse "NFR5-

200  Grev3" (5'-CACCTGAAAGTAACTTATCYGCA-3') for *Vicia sativa*; forward "NFR5-for4"

201  and reverse "NFR5-Grev3" (5'-TGCAGTCTCAGCTAATGAAGTAC-3') for *Lathyrus*

202  *pratensis*; forward "NFR5-for4" and reverse "NFR5-Grev6" (5'-

203  CATACATTGTTGGCTTGCTTAC-3') for *Trifolium hybridum*. The standard PCR protocol was

204  used: initial denaturation at 95°C for 3 min, 30 cycles with denaturation at 94°C for 30 s, primer

205  annealing at 48°C for 30 s, extension at 72°C for 1 min and final extension for 4 min. PCR

206  fragments were extracted from agarose gel ((Onishchuk et al. 2015)) and cloned into the plasmid

207  pTZ57R/T (Thermo Scientific, Lithuania). For each plant species 100 randomly selected cloned

208  fragments of *NFR5* genes were sequenced by Sanger method in an automated ABI 3500xL

209  sequencer (Applied Biosystems) using standard M13 (-20) and (-26) primers. Sequences were

210  deposited in the GenBank database under the PopSet accession number 1041522217. The

211  multiple alignment of 100 sequences was performed with ClustalW as implemented in

212  MEGA6(Tamura et al. 2013).

213

214  **Gene trees**

215

216  The total *nod*A gene sequences from nodule and soil popsets aligned with ClustalW were

217  clustered into operational taxonomic units (OTUs) at 95% nucleotide sequence identity threshold

218  using the UCLUST algorithm implemented in QIIME 1.9. Neighbor-Joining (NJ) dendrogram

219  based on the numbers of differences between representatives of each OTU was constructed in

220    MEGA6 and rooted using the outgroup *nod*A gene sequence of *Sinorhizobium meliloti*
221    (GenBank ID AZNW01000092.1).

222

223    **Diversity analysis**

224

225    We quantified the diversity of haplotypes among rhizobia popsets using three indices: bias-
226    corrected Chao1 index of the haplotype richness (Chao 1984), Simpson 1-D index of haplotype
227    evenness (Simpson 1949) and Shannon H entropy index (Shannon 1948). Nucleotide diversity $\pi$
228    was calculated as the mean difference between two randomly picked sequences in the popset.

229

230    In order to provide statistical support for differences between nodule and soil popsets,
231    distributions of test statistics were obtained by independently subsampling with replacement of
232    2500 sequences from each popset in 2000 trials and bootstrapping of nucleotide positions 1000
233    times. All of the distributions passed the normality test (Kolmogorov-Smirnov test, p-value >
234    0.01). We tested the following null hypothesis for each diversity index: the diversity index value
235    for a nodule popset was lower or equal than the diversity index value for the corresponding
236    nodule popset. Welch's t-test was used to evaluate the null hypothesis at 0.01 significance level.

237

238    To compare the levels of nucleotide diversity $\pi$ between plant and nodule rhizobia popsets,
239    we additionally constructed the distributions of $\pi$ statistics for each plant population
240    subsampling with replacement of 70 sequences in 2000 trials from each plant *NFR5* popset
241    which initially contained 100 sequences. After that, we randomly formed pairs of $\pi$ diversity
242    values from the obtained distributions for the plants popsets and the respective nodule popsets.
243    The total set of 6000 pairs (2000 per plant sp.) was taken as a paired sample data. The
244    relationship between host plants (*NFR5* gene) and nodule rhizobial (*nod*A gene) population
245    diversity was assessed as the value of the Spearman's rank correlation coefficient for the paired
246    sample data. The values higher than 0.8 were taken to imply monotonic relationship in $\pi$
247    between plant and nodule rhizobial posets.

248

249    All calculations were implemented in MATLAB. The link to the GitHub repository
250    containing MATLAB scripts for the diversity analysis is provided in the Supplementary file.

251

252    **Statistical methods for detecting selection**

253

254    The dN/dS ratio is a widely used measure to quantify the selection pressure acting on a set
255    of homologous protein-coding gene regions, where dN and dS are two measures of divergence
256    between species, with dN corresponding to the number of nonsynonymous substitutions per non-
257    synonymous site and dS to the number of synonymous substitutions per synonymous site. pN
258    and pS statistics are analogous to dN and dS statistics but are used for levels of polymorphism
259    within a population rather than divergence (Kryazhimskiy and Plotkin 2008).
260
261    Before calculating the pN/pS ratio, we analysed the independency of the pN and pS
262    statistics. The analysis of linkage between nucleotide positions within the sequenced region of
263    the *nod*A gene was performed within each joint popset (nodule + soil) as follows. First, for each
264    pair of polymorphic nucleotide positions, we applied the chi-square test (with pooling) and
265    extracted pairs of positions that were significantly non-independent (FDR p-value < 0.001).
266    Second, for each position in the obtained pairs, we determined whether it was mostly a
267    synonymous or non-synonymous polymorphic site (See Supplementary text). Third, we counted
268    the number of sites that formed pairs with both synonymous and non-synonymous polymorphic
269    sites and the total number of sites in linked pairs. We considered pN and pS to be independent of
270    each other if the former number was twice as low as the latter number. Under the non-
271    independence of pN and pS, the pS statistic does not reflect neutral mutations and the pN/pS
272    ratio became inconsistent. Thus, we analysed pN and pS separately, mostly focusing on pN as
273    natural selection acts on nonsynonymous changes. Comparing the nodule and soil popsets, we
274    assumed that the increase of pN in one of them indicates the potential presence of stronger
275    positive selection in it or the presence of stronger negative selection in the other popset. We
276    hypothesized that the pN value in a soil rhizobial popset is higher of equal than the pN value in
277    the respective nodule popset. We tested this hypothesis for each legume species separately
278    performing Welch's t-test considering 0.01 level of significance as described above.
279
280    Tajima's D statistics represents the difference between the observed and the theoretically
281    expected nucleotide diversity (Tajima 1989). If the mutations are neutral and the population
282    adheres to the Wright-Fisher assumptions (Hartl and Clark 2007), Tajima's D equals zero.
283    Values significantly higher than zero indicate the deficit of rare alleles (e.g. due to a recent
284    decrease in population size or balancing selection), and values significantly lower than zero
285    indicate the excess of rare alleles (e.g. due to a recent population size expansion or purifying
286    selection). In order to identify the selection type underlying the transformation of a soil popset to
287    the respective nodule popset, we compared the values of Tajima's D between the popsets. We

288  considered the hypothesis that the Tajima's D in the soil popset is higher or equal to that in

289  nodule popset. We tested this hypothesis using Welch's t-test at 0.01 level of significance.

290

291  The calculations were performed using MATLAB PGEToolbox (Population Genetics

292  Evolution Toolbox). The link to the GitHub repository containing MATLAB scripts for the

293  analysis of selection is provided in the Supplementary file.

294

295

296

297  **Topological organization of diversity in plant and rhizobia popsets**

298

299  Let two populations of different sizes are represented by a set of aligned sequences. Let $p$

300  be a population index, $p \in \{1,2\}$. At the first step, the method identifies the unique haplotypes

301  and their frequencies in each population, denoted as $\{(h_i^p, f_i^p)\}$ , $i = \overline{1, n_p}$, where $n_p$ is the

302  number of unique haplotypes in the population $p$. Let $D_{i,j}^p$ be the symmetric distance matrix

303  between each pair of haplotypes in a population $p$, $i = \overline{1, n_p}$. At the second step, the hierarchical

304  agglomerative clustering method merges haplotypes into clusters until the number of clusters

305  equals to predefined number $m$. Let fix a population $p$ and omit this index. The clustering

306  algorithms starts with placing each haplotype to its own cluster which is described by two

307  parameters: frequency $f_i$ and mean difference $\sigma_i$ (initialized to zero). Then, the following

308  procedure is repeated until $m$ clusters is achieved. Two clusters, $i'$-th and $j'$-th with the smallest

309  pairwise distance ( $i', j' : D_{i',j'} = \min\limits_{i,j=\overline{1,n_p}} D_{i,j}$ ) are merged into one new cluster. A distance from

310  this cluster to a $k$-th cluster is calculated as follows: $\frac{D(i',k)f(i')+D(j',k)f(j')}{f(i')+f(j')}$ . The frequency of the

311  new cluster is $f(i) + f(j)$ and the mean difference of the new cluster is $D_{i',j'}$.

312

313  The described hierarchical clustering is applied to each populations and yields the $m$

314  clusters of haplotypes with the reduced distance matrix between them $D_{i,j}^p$, the frequencies $f_i^p$

315  and the mean differences within clusters $\sigma_i^p$, $i, j = \overline{1, m}$. In order to normalize $D_{i,j}^p$ and $\sigma_i^p$ value

316  between two populations we divided these values by the median across $D_{i,j}^p$. If a cluster contains

317  only one haplotype and its $\sigma_i^p$ equals to zero, we set $\sigma_i^p = \min\limits_{i=\overline{1,m}; \sigma_i^p \neq 0} \frac{\sigma_i^p}{f_i^p}$.

318

319  At the third step, the set of clusters for each population was translated into $K$-dimensional

320  Euclidian space by Metric multidimensional scaling (Metric MDS) that transforms the distance

321  matrix $\left\{D_{i,j}^{p}\right\}_{i,j=\overline{1,m}}$ into a set of coordinates $\left\{x_{i}^{p}\right\}_{i,=\overline{1,m}}$. Then, the Gaussian mixture model

322  (GMM) is introduced for the population $p$ as follows:

323
$$G_p(x) = \sum_{i=1}^{m} f_i\, N\!\left(x\middle|x_i^p, \sigma_i\right),$$

324

325  In the current project we took $K = 3$ as it was enough for low values of the stress function in

326  MDS procedure.

327

328  The adjustment of two GMMs is carried out by Procrustes superimposition: the

329  minimization of the dissimilarity between mixtures ($\Delta G$) via translating, rotating and mirror

330  reflection. The lower $\Delta G$ value is, the more similar two GMMs are and, consequently, more

331  similar topologies of two population structures are.

332

333  For each plant species we performed two comparisons: "plant population vs. nodule

334  popset" and "plant population vs. soil popset". We joint rhizobial nodule and soil populations

335  before clustering and MDS, and separated them before the Procrustes analysis. This

336  manipulation ensured that the same haplotypes in both comparisons were taken into account in

337  the same way. We tested the null hypothesis that the $\Delta G$ value in the first comparison is higher

338  or equal to the $\Delta G$ value in the second comparison. In other words, the similarity between nodule

339  rhizobia and plant popset topologies is not higher than the similarity between soil rhizobia and

340  plant popset topologies. For each of the two comparisons, we obtained the set of $\Delta G$ values

341  bootstrapping of sequences in rhizobial popsets. To test the hypothesis, we compared two

342  obtained sets of $\Delta G$ values by one-sided Mann-Whitney U test with 0.01 level of significance.

343

344  For visual comparison of plant and rhizobial populations we constructed tanglegrams

345  based on adjusted GMMs after Procrustes superimposition. A tanglegram is a diagram with a

346  pair of two binary trees with matching leaves connected by edges. To construct it we built two

347  NJ gene trees for $m$ plant *NFR5* clusters and $m$ rhizobium *nod*A clusters based on the between-

348  cluster distance matrices and plotted two trees face to face. A pair of leaves from two trees was

349  connected by an edge if a 3D point corresponded to one leave was located within the five closest

350  points to a point of another leave and vice versa.

351

352    The link to the GitHub repository containing MATLAB scripts for topological beta-
353    diversity analysis is provided in the Supplementary file.

354

355                                        **Results**

356

357    ***nod*A gene: OTUs and cluster analysis**

358

359    *Nod*A gene libraries (for root nodules and for soil samples) were sequenced by NGS
360    technology. Sequencing and filtration of *nod*A gene libraries produced a total of 22463
361    sequences, for an average of 3750 sequences per popset. Clustering by 95% sequence identity
362    produced 15 OTUs, which frequencies differed between the popsets. The NJ tree constructed for
363    OTU representatives showed that the rhizobium population consisted of two clusters (orange and
364    blue in Figure 2). Almost all sequences from the first cluster (9 OTUs) belonged to Vicia and
365    Lathyrus popsets, while sequences from the second cluster (6 OTUs) were mostly detected in
366    Trifolium popsets. The converse was also true: most of the Vicia-Lathyrus popset belonged to
367    the first cluster (90% average), while most of the Trifolium popset belonged to the second
368    cluster. This distribution of OTUs was in agreement with the "cross-inoculation groups" concept,
369    which refers to the fact that separate groups of legume species can be successfully inoculated by
370    only the specific groups of rhizobia. We attributed the first cluster to *R.leguminosarum* bv. *viciae*
371    and the second to *R.leguminosarum* bv. *trifolii.* For further analysis, we kept in Vicia-Lathyrus
372    popsets only sequences from the first cluster, and in Trifolium popsets, only sequences from the
373    second cluster.

374

375    **Differences between soil and nodule populations of rhizobia**

376

377    For each of the six *nod*A gene popsets, we examined the following measures of population
378    diversity: number of haplotypes (Supplementary Figure S1), Chao1 index of richness, Simpson
379    1-D index of evenness, Shannon H index of entropy, and nucleotide diversity $\pi$. All diversity
380    indexes are significantly higher in nodule popsets compared to soil popsets (Figure 3.A-D;
381    Welch's t-test, p<0.01).

382

383

384    **Analysis of selection**

385

386    Analysis of population structures revealed the significant linkage between synonymous
387    and non-synonymous polymorphic sites, suggesting that the pN and pS statistics were not
388    independent (see Supplementary Text). Further comparison of pN/pS values between nodule and
389    soil popsets did not establish the significant difference (Supplementary Fig. S7). Despite the
390    inconsistency of pN/pS in our case, we observed that the values of both statistics, pN and pS,
391    were significantly increased (Welch's t-test, p-values < 0.01) in the nodule popsets, indicating
392    that both non-synonymous and synonymous diversity was elevated there (Table 1). The values of
393    Tajima's D were significantly lower than 0 in all of the rhizobial *nod*A popsets, indicating the
394    presence of negative selection (Table 1). However, within each plant, they were significantly
395    higher in nodule popsets than in the soil popsets (p-values < 0.01), suggestive of relaxation of
396    negative selection or admixture of balancing selection in the former. Trifolium popsets displayed
397    the highest values of this statistic, consistently with the strongest admixture of balancing
398    selection, while the Lathyrus popsets, the lowest.

399

400    **Table 1.** Values of pN, pS, pN/pS and Tajima's D statistics for the popsets of different origin.
401    The difference in pN/pS values between nodule and soil popsets for each plant was not
402    significant. The significance (p-value < 0.01) of the difference in values between nodule and soil
403    population for each plant are marked with "**".

| Popset | | pN | pS | pN/pS | Tajima's D |
|---|---|---|---|---|---|
| Vicia | nodule | 0.0170 | 0.0228 | 0.7450 | -2.2644 |
| | soil | 0.0075 ** | 0.0099 | 0.7600 | -2.5219 ** |
| Lathyrus | nodule | 0.0065 | 0.0086 | 0.7570 | -2.5609 |
| | soil | 0.0051 ** | 0.0072 | 0.7186 | -2.5977 * |
| Trifolium | nodule | 0.0245 | 0.0335 | 0.7301 | -2.0294 |
| | soil | 0.0167 ** | 0.0232 | 0.7184 | -2.2496 ** |

404

405

406    **Relationship between the bacteria and host plant diversities**

407

408    We calculated the diversity levels $\pi$ in each plant population and found that the ranking of
409    the three species was the same as the ranking based on the $\pi$ nucleotide diversity values in
410    corresponding nodule popsets. By bootstrapping the plant and rhizobial nodule popsets we
411    estimated the Spearman correlation between these popsets's nucleotide diversities. The 0.89
412    value, that was higher than the predefined threshold indicated that the monotonic relationship

413    between diversities in popsets of plant *NFR5* gene sequences and in bacterial *nod*A nodule

414    popsets is statistically significant (Figure 4).

415

416    **Concordance of gene trees**

417

418    A visual comparison of topologies of plant *NFR5* trees with those of corresponding

419    rhizobial *nod*A trees for nodule and soil popsets revealed that the topology of clades in the plant

420    trees was more similar to the topology of clades in the nodule trees than in the soil trees

421    (Supplementary Figure S3). Based on this observation, we proposed that the gene tree of the

422    nodule rhizobia is more similar to that of the host plant than the gene tree of the soil rhizobia.

423

424    To formally test this, we developed a method for comparing structures (topologies)

425    between two popsets. We tested the null hypothesis that the topological similarity between a

426    nodule rhizobia popset and a plant popset is not higher than the topological similarity between

427    the soil rhizobia popset and the plant popset. This hypothesis was rejected at the 0.01 level of

428    significance. We constructed tanglegrams that also illustrated a higher similarity of the topology

429    of the nodule rhizobial *nod*A gene trees (Figure 5, left tanglegram) than the soil rhizobial *nod*A

430    gene trees (Figure 5, right tanglegram) to that of plant NFR5 gene tree.

431

432                                    **Discussion**

433

434    Symbiotic interactions represent a special case of ecological interactions when one of the

435    partners provides an "environment" for another. In "On the Origin of Species" (Darwin 1872),

436    Charles Darwin proposed that "the life of each species depends in a more important manner on

437    the presence of other already defined organic forms, than on climate". This is particularly true

438    for organisms in deeply integrated symbiotic systems.

439

440    Here, we traced the coordination in levels of population diversities between partners

441    within the essential components of the rhizobium-legume signalling system: plant symbiotic

442    receptor gene *NFR5* and *Rhizobium* symbiotic gene *nod*A involved in the synthesis of signalling

443    molecules Nod factors, ensuring the first stage of partner recognition (Oldroyd 2013). The

444    matching was detected in three phenomena. The first is the increase of population diversity

445    levels in nodule rhizobial population (mutualist phase) in a comparison with soil rhizobial

446    popsets (free-living phase). The second is the monotonic relationship between the host plant

447    diversity and the nodule rhizobial diversity. The third is the similarity in topological diversity

448    between host plant population and nodule rhizobial population. All of these observations

449    illustrated the aspect of symbiotic interactions, which we refer to as the Evolutionary Moulding

450    effect. Under this effect the soil rhizobial population *nod*A gene pool "is moulded into rigid

451    matrix" of host-plant population *NFR5* gene pool forming nodule pool of *nod*A and ensuring (1)

452    the sufficient increase of the diversity level in nodule pool to meet the host plant needs (2) the

453    same trend in the matching diversity levels between rhizobial and host-plant across legume

454    species (3) the "topological" matching of *NFR5* with *nod*A gene trees between host-plant popset

455    with nodule rhizobial popset, but not with soil rhizobial popset.

456

457    The increased diversity (evenness, richness and nucleotide polymorphism) in nodule

458    popsets in comparison with respective soil popsets seem unexpected while only 90 nodules

459    formed each nodule popset. Indeed, as the formation of each nodule popset involves a bottleneck

460    of only ~90 nodules the nodule population can be expected to have a reduced population size and

461    therefore, reduced diversity. However, if the nodule population was formed non randomly, this

462    increase is not surprising. There are several possible explanations to the observed trends: (i)

463    selection imposed by plants on soil populations which leads to increased diversity in nodule

464    populations; (ii) increased mutation rate in bacterial populations under stress conditions inside

465    root nodules (Krasinikov and Melkumova 1963; Roumiantseva et al. 2004); (iii) negative

466    selection in the soil which reduces diversity there. The last explanation is unlikely, as the

467    selection pressure in the soil hardly affects the symbiotic genes. The second explanation

468    remained only hypothetical as it was not well described and had no much attention in recent

469    studies (Krasinikov and Melkumova 1963; Roumiantseva et al. 2004). By exclusion, the first

470    explanation is the most likely. The analysis of the selection imposed by plants revealed

471    significantly increased nonsynonymous diversity (pN) and Tajima's D values in the nodule

472    popsets. This may be indicative of weaker negative selection in a nodule popset in a comparison

473    with the respective soil popset but is also consistent with a contribution of balancing selection or

474    presence of stronger population structure in the former. Earlier it was shown that in symbiotic

475    systems, besides the above-mentioned types of selection, negative frequency-dependent selection

476    in favour of rare genotypes during the competition of rhizospheric bacteria for root

477    nodulation(Amarger and Lobreau 1982) may also play an important role (Andronov et al. 2015;

478    Provorov and Vorobyov 2000; Provorov and Vorobyov 2006).

479

480    The pronounced influence of plant on the formation of nodule population was

481    demonstrated in the monotonic correspondence between diversity in plant popsets and diversity

482    in respective nodule rhizobial popsets (Spearman correlation = 0.89). Another convincing

483     evidence of the plant-imposed effect was obtained with the specially developed method to

484     compute "topological beta-diversity" – difference in topological structures of two population sets

485     (plant and rhizobia) of sequences. The similarity between *NFR5* and *nod*A gene trees of plant

486     and nodule popsets was significantly higher than the same between plant and soil popsets. This

487     result demonstrated the transformation of the initial soil *nod*A pool by the template of the host

488     plant receptor pool. Our conclusion is in line with the numerous works studying the interplay

489     between diversities of host plant and rhizobia (Paffetti et al. 1998; Andronov et al. 2003; Bena et

490     al. 2005; Bailly et al. 2006; Depret and Laguerre 2008; Rangin et al. 2008; Barrett et al. 2016;

491     Vuong, Thrall, and Barrett 2016; Österman et al. 2011), however, in most cases we cannot

492     directly compare these studies due to the differences in the experimental design.

493

494     The observed similarity between nodule rhizobial *nod*A gene popsets and plant *NFR5*

495     receptor gene popsets revealed the hierarchical organization of effective interaction: two

496     symbionts should be genetically compatible at the single organism level and also at the

497     population level. The process of forming this interaction could be explained metaphorically as

498     the Evolutionary Moulding: shaping the population structure of one symbiont using the

499     population structure of another symbiont as a "matrix". The important point in this shaping is the

500     difference between evolutionary rates in plants and bacteria. The bacteria have a significantly

501     higher evolutionary rate than plants, therefore the diversity of *nod*A gene in bacterial

502     populations, like the flexible genetic material in the Evolutionary Moulding, reflected the shape

503     of more "rigid" diversity of *NFR5* receptor gene in plant populations. We hypothesise that under

504     the Evolutionary Moulding effect two symbiotic populations tend to relax the incoordination of

505     genetic diversities between two parts of the symbiont-host signalling system, that is mostly

506     achieved by a faster evolving partner. In the current project the faster evolving partner is

507     rhizobia, however, we assume that if plants hypothetically had evolved with higher rates than

508     bacteria, we would have expected that the rhizobial component would have played a role of the

509     matrix (or mould) in terms of the Evolutionary Moulding. It should be also highlighted, that

510     under the Evolutionary Moulding, the transformation of the "more flexible" symbiont population

511     does not necessarily lead to its increased population diversity. When a population of a "more

512     rigid" symbiont is, for example, homogeneous (conservative in polymorphism), the population

513     diversity of the flexible symbiont could decrease.

514

515     According to the Evolutionary Moulding effect, the relationship in population diversity

516     between rhizobia and host-plant may be observed not only within the pair of *nod*A-*NFR5* genes

517     (which are related through the Nod factor) but also within any pair of interplaying genes from

518 plant and bacterial sides. We propose that genome-wide searching of "matching" genes under the
519 Evolutionary Moulding can be an extension to the traditional methods of functional analysis of
520 genes.

521

522 At present, we can only hypothesize the molecular mechanism of the Evolutionary
523 Moulding in rhizobium-legume symbiosis. We suppose that polymorphism in symbiotic genes of
524 a rhizobial population is probably associated with the structural diversity of Nod factors
525 produced by this population, in particular with variations in the unsaturated fatty acid tail.
526 Nowadays, an impressive progress is achieved in the resolution and accuracy of the methods to
527 analyse the Nod factor structure (Poinsot et al. 2016) and its docking to the receptor proteins
528 (Malkov et al. 2016). We believe that such approaches will facilitate the analysis of Nod factor
529 structural variation produced by rhizobial populations and, as a result, the molecular mechanisms
530 in the Evolutionary Moulding.

531

532

533 **Acknowledgements**

535

536 **References**

537 Agrawal, A, and C M Lively. 2002. 'Infection Genetics: Gen-for-Gene versus Matching Alleles
538 Models and All Points in between'. *Evolutionary Ecology Research* 4: 79–90.
539 Amarger, N., and J. P. Lobreau. 1982. 'Quntitative Study of Nodulation Competitiveness in
540 Rhizobium Strains'. *Applied and Environmental Microbiology* 44 (3): 583–88.
541 Andronov, E.E., A.A. Igolkina, A.K. Kimeklis, N.I. Vorobyov, and N.A. Provorov. 2015.
542 'Characteristics of Natural Selection in Populations of Nodule Bacteria (Rhizobium
543 Leguminosarum) Interacting With Different Host Plants'. *Genetika* 51 (10).
544 Andronov, E.E., Z Terefework, M L Roumiantseva, N I Dzyubenko, O P Onichtchouk, O N
545 Kurchak, J P W Young, B V Simarov, and K Lindstro. 2003. 'Symbiotic and Genetic
546 Diversity of Rhizobium Galegae Isolates Collected from the Galega Orientalis Gene Center
547 in the Caucasus'. *Applied and Environmental Microbiology* 69 (2): 1067–74.
548 doi:10.1128/AEM.69.2.1067.
549 Bailly, Xavier, Isabelle Olivieri, Stéphane De Mita, Jean-Claude Cleyt-Marel, and Gilles Bena.
550 2006. 'Recombination and Selection Shape the Molecular Diversity Pattern of Nitrogen-
551 Fixing Sinorhizobium Sp . Associated to Medicago'. *Molecular Ecology* 15: 2719–34.
552 doi:10.1111/j.1365-294X.2006.02969.x.

553  Barrett, Luke G, Peter C Zee, James D Bever, Joseph T Miller, and Peter H Thrall. 2016.
554       'Evolutionary History Shapes Patterns of Mutualistic Benefit in Acacia – Rhizobial
555       Interactions'. *Evolution*, 1–13. doi:10.1111/evo.12966.
556  Bena, G., A Lyet, T Huguet, and I Olivieri. 2005. 'Medicago – Sinorhizobium Symbiotic
557       Specificity Evolution and the Geographic Expansion of Medicago'. *European Society For*
558       *Evolutionary Biology* 18: 1547–58. doi:10.1111/j.1420-9101.2005.00952.x.
559  Bergstrom, Carl T, and Michael Lachmann. 2003. 'The Red King Effect : When the Slowest
560       Runner Wins the Coevolutionary Race'. *PNAS* 100 (2): 593–98.
561  Broghammer, Angelique, Lene Krusell, Mickaël Blaise, Jørgen Sauer, John T Sullivan, and
562       Nicolai Maolanon. 2012. 'Legume Receptors Perceive the Rhizobial Lipochitin
563       Oligosaccharide Signal Molecules by Direct Binding'. *PNAS* 109 (34): 13859–64.
564       doi:10.1073/pnas.1205171109.
565  Chao, A. 1984. 'Non-Parametric Estimation of the Number of Classes in a Population'.
566       *Scandinavian Journal of Statistics* 11: 265–70.
567  Choi, Kwangbom, and Shawn M Gomez. 2009. 'Comparison of Phylogenetic Trees through
568       Alignment of Embedded Evolutionary Distances'. *BMC Bioinformatics* 15: 1–15.
569       doi:10.1186/1471-2105-10-423.
570  Cregan, P B, M J Sadowsky, and H H Keyser. 1991. 'Gene-for-Gene Interaction in the Legume-
571       Rhizobium Symbiosis'. In *The Rhizosphere and Plant Growth*, 2:163–71. Dordrecht:
572       Springer Netherlands. doi:10.1007/978-94-011-3336-4_32.
573  Darwin, C. 1872. *On the Origin of Species by Means of Natural Selection, Or, the Preservation*
574       *of Favoured Races in the Struggle for Life*. London: J. Murray.
575  Dénarié, Jean, Frédéric Debellé, and Jean-Claude Promé. 1996. 'Rhizobium Lipo-
576       Chitooligosaccharide Nodulation Factors: Signaling Molecules Mediating Recognition and
577       Morphogenesis'. *Annual Review of Biochemistry* 65 (1): 503–35.
578       doi:10.1146/annurev.bi.65.070196.002443.
579  Depret, Géraldine, and Gisèle Laguerre. 2008. 'Plant Phenology and Genetic Variability in Root
580       and Nodule Development Strongly Influence Genetic Structuring of Rhizobium
581       Leguminosarum Biovar Viciae Populations Nodulating Pea'. *New Phytologist* 179: 224–35.
582       doi:10.1111/j.1469-8137.2008.02430.x.
583  Dobert, R.C., B.T. Breil, and E.W. Triplett. 1994. 'DNA Sequence of the Common Nodulation
584       Genes of Bradyrhizobium Elkanii and Their Phylogenetic Relationship to Those of Other
585       Nodulating Bacteria'. *Mol Plant Microbe Interact* 7 (5): 564–72.
586  Esseling, J.J., and A.M.C. Emons. 2004. 'Dissection of Nod Factor Signalling in Legumes : Cell
587       Biology , Mutants and Pharmacological Approaches' 214 (December 2003): 104–13.

588  Fenton, Andrew, Janis Antonovics, and Michael A. Brockhurst. 2009. 'Inverse-Gene-for-Gene

589  Infection Genetics and Coevolutionary Dynamics'. *The American Naturalist* 174 (6): E230–

590  42. doi:10.1086/645087.

591  Flor, H.H. 1942. 'Inheritance of Pathogenicity in Melampsora Lini'. *Phytopath* 32: 653–669.

592  Flor, H H. 1971. 'Current Status of the Gene-Fob-Gene Concept'. *Annu.Rev .Phytopathol.* 9:

593  275–96.

594  Frank, S. A. 1994. 'Recognition and Polymorphism in Host-Parasite Genetics'. *Philosophical*

595  *Transactions of the Royal Society B: Biological Sciences* 346 (1317): 283–93.

596  doi:10.1098/rstb.1994.0145.

597  Hartl, Daniel L., and Andrew G. Clark. 2007. *Principles of Population Genetics (4th Ed.)*.

598  Sunderland, MA: Sinauer Associates.

599  Heath, Katy D, and Peter Tiffin. 2007. 'Context Dependence in the Coevolution of Plant and

600  Rhizobial Mutualists'. *Proc. R. Soc. B* 274 (May): 1905–12. doi:10.1098/rspb.2007.0495.

601  Jost, Lou, Anne Chao, and Robin L. Chazdon. 2011. 'Compositional Similarity and B (Beta)

602  Diversity'. In *Biological Diversity - Frontiers in Measurement and Assessment*, 66–84.

603  Krasinikov, N.A., and T.N. Melkumova. 1963. 'Variability of Nodule Bacteria inside Nodules of

604  Leguminous Plants (Article in Russian)'. *Proceedings of the USSR Academy of Sciences,*

605  *Biol Series* 5: 693 – 706.

606  Kryazhimskiy, Sergey, and Joshua B Plotkin. 2008. 'The Population Genetics of dN / dS'. *PloS*

607  *Genetics* 4 (12). doi:10.1371/journal.pgen.1000304.

608  Kwiatkowski, M, J Engelstadter, and C Vorburger. 2012. 'On Genetic Specificity in Symbiont-

609  Mediated Host-Parasite Coevolution'. *PloS Computational Biology* 8 (8).

610  doi:10.1371/journal.pcbi.1002633.

611  Laguerre, Gisele, Philippe Louvrier, Marie-reine Allard, and Amarger Noelle. 2003.

612  'Compatibility of Rhizobial Genotypes within Natural Populations of Rhizobium

613  Leguminosarum Biovar Viciae for Nodulation of Host Legumes'. *Applied and*

614  *Environmental Microbiology* 69 (4): 2276–83. doi:10.1128/AEM.69.4.2276.

615  Leigh, Jessica W., Fraṇois Joseph Lapointe, Philippe Lopez, and Eric Bapteste. 2011.

616  'Evaluating Phylogenetic Congruence in the Post-Genomic Era'. *Genome Biology and*

617  *Evolution* 3 (1): 571–87. doi:10.1093/gbe/evr050.

618  Lewis-Henderson, W. R., and M. A. Djordjevic. 1991. 'A Cultivar-Specific Interaction between

619  Rhizobium Leguminosarum Bv. Trifolii and Subterranean Clover Is Controlled by nodM,

620  Other Bacterial Cultivar Specificity Genes, and a Single Recessive Host Gene'. *Journal of*

621  *Bacteriology* 173 (9): 2791–99.

622  Li, Ronghui, Maggie R Knox, Anne Edwards, Bridget Hogg, T H Noel Ellis, Gehong Wei, and J

623      Allan Downie. 2011. 'Natural Variation in Host-Specific Nodulation of Pea Is Associated
624          with a Haplotype of the SYM37 LysM-Type Receptor-Like Kinase'. *MPMI* 24 (11): 1396–
625          1403.

626      Lozupone, Catherine, Manuel E Lladser, Dan Knights, Jesse Stombaugh, and Rob Knight. 2011.
627          'UniFrac: An Effective Distance Metric for Microbial Community Comparison.' *The ISME*
628          *Journal* 5 (2). Nature Publishing Group: 169–72. doi:10.1038/ismej.2010.133.

629      Malkov, Nikita, Judith Fliegmann, Charles Rosenberg, Virginie Gasciolli, Antonius C. Cj
630          Timmers, Alessandra Nurisso, Julie Cullimore, and J.-J. Jean-Jacques Bono. 2016.
631          'Molecular Basis of Lipo-Chitooligosaccharide Recognition by the Lysin Motif Receptor-
632          like Kinase LYR3 in Legumes'. *Biochemical Journal*, no. 2016: BCJ20160073.
633          doi:10.1042/BCJ20160073.

634      Martínez-Romero, Esperanza. 2009. 'Coevolution in Rhizobium -Legume Symbiosis?' *DNA and*
635          *Cell Biology* 28 (8): 361–70. doi:10.1089/dna.2009.0863.

636      Mergaert, Peter, and Marc Van Montagu. 1997. 'Molecular Mechanisms of Nod Factor
637          Diversity'. *Molecular Microbiology* 25 (5): 811–17.

638      Mode, Charles J. 1958. 'A Mathematical Model for the Co-Evolution of Obligate Parasites and
639          Their Hosts'. *Evolution* 12 (2): 158. doi:10.2307/2406026.

640      Moulin, Lionel, Gilles Béna, and T Stępkowski. 2004. 'Phylogenetic Analyses of Symbiotic
641          Nodulation Genes Support Vertical and Lateral Gene Co- Transfer within the
642          Bradyrhizobium Genus'. *Molecular Phylogenetics and Evolution* 30: 720–32.
643          doi:10.1016/S1055-7903(03)00255-0.

644      Oldroyd, Giles E D. 2013. 'Speak , Friend , and Enter : Signalling Systems That Promote
645          Beneficial Symbiotic Associations in Plants'. *Nature Publishing Group* 11 (4). Nature
646          Publishing Group: 252–63. doi:10.1038/nrmicro2990.

647      Onishchuk, O P, E P Chizhevskaya, O N Kurchak, E E Andronov, and B V Simarov. 2015.
648          'Identification of New Genes of Nodule Bacteria Sinorhizobium Meliloti Involved in the
649          Control of Efficiency of Symbiosis with Alfalfa Medicago Sativa'. *Russian Journal of*
650          *Genetics: Applied Research* 5 (2): 126–31. doi:10.1134/S2079059715020070.

651      Österman, J., E. P. Chizhevskaja, E. E. Andronov, D. P. Fewer, Z. Terefework, M. L.
652          Roumiantseva, O. P. Onichtchouk, et al. 2011. 'Galega Orientalis Is More Diverse than
653          Galega Officinalis in Caucasus-Whole-Genome AFLP Analysis and Phylogenetics of
654          Symbiosis-Related Genes'. *Molecular Ecology* 20 (22): 4808–21. doi:10.1111/j.1365-
655          294X.2011.05291.x.

656      Otto, Sarah P., and Yannis Michalakis. 1998. 'The Evolution of Recombination in Changing
657          Environments'. *Trends in Ecology & Evolution* 13 (4): 145–51. doi:10.1016/S0169-

658    5347(97)01260-3.

659  Paffetti, Donatella, Fabrice Daguin, Silvia Fancelli, Stefano Gnocchi, Francesca Lippi, Carla
660    Scotti, Marco Bazzicalupo, Dipartimento Ortoflorofrutticoltura, G Donizzetti, and Biologia
661    Animale. 1998. 'Influence of Plant Genotype on the Selection of Nodulating Sinorhizobium
662    Meliloti Strains by Medicago Sativa'. *Antonie van Leeuwenhoek* 73: 3–8.

663  Pal, Csaba, Mar´ıa D. Macia, Antonio Oliver, Ira Schachar, and Angus Buckling. 2007.
664    'Coevolution with Viruses Drives the Evolution of Bacterial Mutation Rates'. *Nature* 450
665    (December): 1079–81. doi:10.1038/nature06350.

666  Parker, M A. 1999. 'Mutualism in Metapopulations of Legumes and Rhizobia'. *The American*
667    *Naturalist* 153 (May 1999): S48–60. doi:10.1086/303211.

668  Paterson, Steve, Tom Vogwill, Angus Buckling, Rebecca Benmayor, Andrew J Spiers, Nicholas
669    R Thomson, Mike Quail, et al. 2010. 'Antagonistic Coevolution Accelerates Molecular
670    Evolution'. *Nature* 464 (7286). Nature Publishing Group: 275–78.
671    doi:10.1038/nature08798.

672  Poinsot, Verena, Matthew B. Crook, Stephanie Erdn, Fabienne Maillet, Adeline Bascaules, and
673    Jean-Michel Ane. 2016. 'New Insights into Nod Factor Biosynthesis : Analyses of
674    Chitooligomers and Lipo-Chitooligomers of Rhizobium Sp . IRBG74 Mutants R E'.
675    *Carbohydr Res.* 434: 83–93. doi:10.1016/j.carres.2016.08.001.

676  Provorov, N.A., and N.I. Vorobyov. 2000. 'Population Genetics of Rhizobia: Construction and
677    Analysis of an "infection and Release" Model'. *J. Theor. Biol.* 205 (1): 105–19.

678  ———. 2006. 'Interplay of Darwinian and Frequency-Dependent Selection in the Host-
679    Associated Microbial Populations.' *Theor. Population Biol* 70 (3): 262–72.

680  Radutoiu, Simona, Lene H Madsen, Esben B Madsen, Anna Jurkiewicz, Eigo Fukai, Esben M H
681    Quistgaard, Anita S Albrektsen, Euan K James, S Thirup, and Jens Stougaard. 2007. 'LysM
682    Domains Mediate Lipochitin – Oligosaccharide Recognition and Nfr Genes Extend the
683    Symbiotic Host Range'. *The EMBO Journal* 26 (17): 3923–35.
684    doi:10.1038/sj.emboj.7601826.

685  Rangin, Cecile, Brigitte Brunel, Jean-Claude Cleyt-Marel, Marie-mathilde Perrineau, and Bena
686    Gilles. 2008. 'Effects of Medicago Truncatula Genetic Diversity , Rhizobial Competition ,
687    and Strain Effectiveness on the Diversity of a Natural Sinorhizobium Species Community'.
688    *Applied and Environmental Microbiology* 74 (18): 5653–61. doi:10.1128/AEM.01107-08.

689  Ritsema, Tita, A.H.M. Wijfjes, B.J.J. Lugtenberg, and H.P. Spaink. 1996. 'Rhizobium
690    Nodulation Protein NodA Is a Host-Specific Determinant of the Transfer of Fatty Acids in
691    Nod Factor Biosynthesis'. *Mol Gen Genet* 251: 44–51.

692  Roche, Philippe, Fabienne Maillet, Claire Plazanet, and Jaen Denarie. 1996. 'The Common

693      nodABC Genes of Rhizobium Meliloti Are Host-Range Determinants'. *PNAS* 93
694      (December): 15305–10.

695  Roumiantseva, M L, E E Andronov, V V Sagulenko, O P Onishchuk, N A Provorov, and B V
696      Simarov. 2004. 'Instability of Cryptic Plasmids in Strain Sinorhizobium Meliloti P108 in
697      the Course of Symbiosis with Alfalfa Medicago Sativa'. *Genetics of Microorganisms* 40
698      (4): 356–62.

699  Rubin, Benjamin E R, and Corrie S Moreau. 2016. 'Of Evolution in Ant – Plant Mutualisms'.
700      *Nature Communications* 7. Nature Publishing Group: 1–11. doi:10.1038/ncomms12679.

701  Sachs, Joel L., Carla J. Essenberg, and Martin M. Turcotte. 2011. 'New Paradigms for the
702      Evolution of Beneficial Infections'. *Trends in Ecology and Evolution* 26 (4). Elsevier Ltd:
703      202–9. doi:10.1016/j.tree.2011.01.010.

704  Sadowsky, M J, P B Cregan, M Gottfert, a Sharma, D Gerhold, F Rodriguez-Quinones, H H
705      Keyser, H Hennecke, and G Stacey. 1991. 'The Bradyrhizobium Japonicum nolA Gene and
706      Its Involvement in the Genotype-Specific Nodulation of Soybeans.' *Proceedings of the*
707      *National Academy of Sciences of the United States of America* 88 (2): 637–41.
708      doi:10.1073/pnas.88.2.637.

709  Shannon, C. E. 1948. 'A Mathematical Theory of Communication'. *The Bell System Technical*
710      *Journal* 27 (April 1924): 379–423.

711  Simpson, E.H. 1949. 'Measurement of Diversity'. *Nature* 163: 688–688.

712  Spaink, Herman P. 2000. 'Root Nodulation and Infection Factors Produced by Rhizobial
713      Bacteria'. *Annu. Rev. Microbiol. 2000.* 54: 257–88.

714  Suominen, Leena, Christophe Roos, Lars Paulin, Leena Suominen, Christophe Roos, Gilles
715      Lortet, Lars Paulin, and Kristina Lindstro. 2016. 'Identification and Structure of the
716      Rhizobium Galegae Common Nodulation Genes : Evidence for Horizontal Gene Transfer
717      Genes : Evidence for Horizontal Gene Transfer', no. January.
718      doi:10.1093/oxfordjournals.molbev.a003891.

719  Tajima, Fumio. 1989. 'Statistical Method for Testing the Neutral Mutation Hypothesis by DNA
720      Polymorphism'. *Genetics* 123: 585–95.

721  Tamura, Koichiro, Glen Stecher, Daniel Peterson, Alan Filipski, and Sudhir Kumar. 2013.
722      'MEGA6 : Molecular Evolutionary Genetics Analysis Version 6 . 0'. *Mol. Biol. Evol.* 30
723      (12): 2725–29. doi:10.1093/molbev/mst197.

724  Thompson, J.N., and J.J. Burdin. 1992. 'Gene-for-Gene Coevolution between Plants and
725      Parasites'. *Nature* 360: 121–25.

726  Van Valen, Lee. 1973. 'A New Evolutionary Law'. *Evolutionary Theory* 1: 1–30.

727  Vuong, Holly B, Peter H Thrall, and Luke G Barrett. 2016. 'Host Species and Environmental

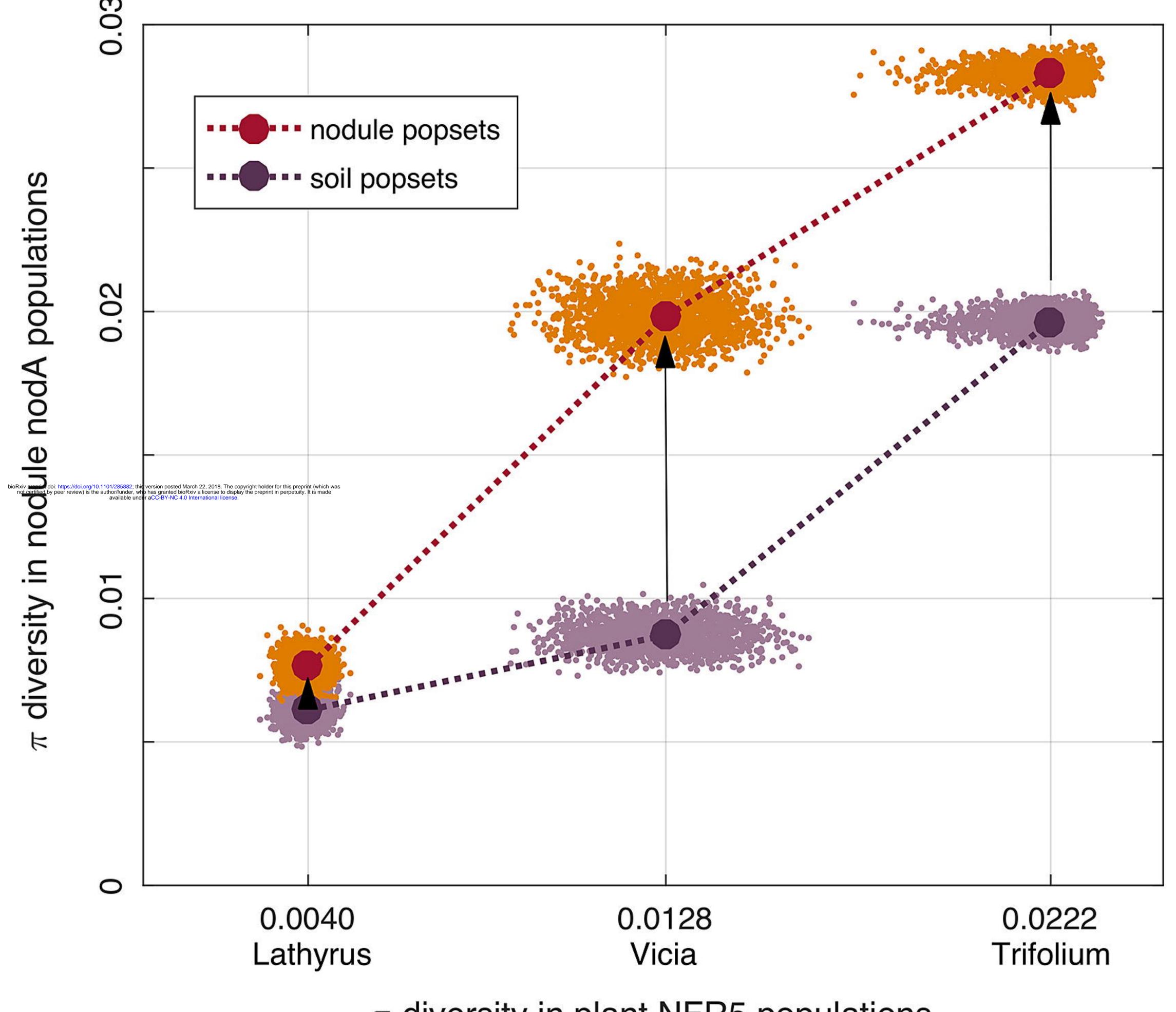728    Variation Can in Fl Uence Rhizobial Community Composition'. *Journal of Ecology*.

729    doi:10.1111/1365-2745.12687.

730  Yoder, Jeremy B. 2016. 'Understanding the Coevolutionary Dynamics of Mutualism with

731    Population Genomics'. *American Journal of Botany* 103 (10): 1742–52.

732    doi:10.3732/ajb.1600154.

733

734

735

736

737

738  **Figure 1** A part of the signal transduction system that governs the rhizobium-legume symbiosis.

739  The rhizobial *nod*A gene encodes the acyltransferase that participates in the attachment of the

740  hydrophobic long-chain fatty acid tail to the Nod factor core. Plant *NFR5* gene encodes the

741  symbiotic receptor recognising the rhizobial Nod factor followed by the symbiosis formation.

742

743  **Figure 2.** Clustering of *nod*A gene sequences after OTU-picking analysis. Columns correspond

744  to six popsets (see text). Values in cells represent the numbers of OTU sequences in a popset,

745  widths of rectangles reflect to the log of these values. The NJ tree of OTU-representatives forms

746  two clades corresponding to *Rhizobium leguminosarum* biovars from different cross-inoculation

747  groups: bv. *viciae* and bv. *trifolii*.

748

749  **Figure 3**. **(A-D)**: Diversity measures.  Letters "V", "L", "T" denote Vicia, Lathyrus, Trifolium

750  popsets respectively. Letters "Soil" and "Node" denote soil and nodule popsets. Precise formulas

751  to compute diversity measures are shown in Supplementary TableS1. Mean values of diversity

752  measures are presented in the Supplementary TableS2.

753

754  **Figure 4.** The monotonic relationship between the $\pi$ diversity levels in nodule popsets and plant

755  popsets. Dots represent the distribution of $\pi$ obtained by bootstrapping.

756

757  **Figure 5.** Comparison of plant popsets with the bacteria popsets from nodules and soil. **(A)**:

758  Projections of Gaussian mixture models for three plant *NFR5* popsets and six rhizobial *nod*A

759  popsets after the Procrustes analysis on the XoY plane (see Supplementary Figures S4-S6 for

760  other projections). The values of $\Delta G$ correspond to the difference between the GMMs for plant

761  and rhizobial (nodule or soil) popsets; differences between $\Delta G$ values in each row are significant

762  (p<0.05). A visual comparison of projection confirms this trend. For example, in the "Vicia" row

763    the rhizobium nodule popset has two peaks that remind two peaks in the plant popset and more

764    distinct than in the rhizobium soil popset **(B)**: Tanglegrams for each plant species: between

765    *NFR5* population and rhizobial *nod*A populations from nodule (left tanglegrams) and soil (right

766    tanglegrams).

767

768

**A**

| | Plant populations | Nodule popsets | Soil popsets |
|---|---|---|---|

Vicia $_{XoY}$

Nodule popsets: $\Delta G = 0.0011$
Soil popsets: $\Delta G = 0.0021$

Lathyrus $_{XoY}$

Nodule popsets: $\Delta G = 0.0008$
Soil popsets: $\Delta G = 0.0010$

Trifolium $_{XoY}$

Nodule popsets: $\Delta G = 0.0013$
Soil popsets: $\Delta G = 0.0020$

**B**

Plant — Nodule    Plant — Soil