

This is a repository copy of *Evolutionary dynamics of insertion sequences in relation to the evolutionary histories of the chromosome and symbiotic plasmid genes of Rhizobium etli populations*.

White Rose Research Online URL for this paper:

<https://eprints.whiterose.ac.uk/42645/>

Version: Submitted Version

Article:

Lozano, Luis, Hernandez-Gonzalez, Ismael, Bustos, Patricia et al. (5 more authors) (2010) Evolutionary dynamics of insertion sequences in relation to the evolutionary histories of the chromosome and symbiotic plasmid genes of Rhizobium etli populations. Applied and Environmental Microbiology. pp. 6504-6513. ISSN 0099-2240

<https://doi.org/10.1128/AEM.01001-10>

Reuse

Items deposited in White Rose Research Online are protected by copyright, with all rights reserved unless indicated otherwise. They may be downloaded and/or printed for private study, or other acts as permitted by national copyright laws. The publisher or other rights holders may allow further reproduction and re-use of the full text version. This is indicated by the licence information on the White Rose Research Online record for the item.

Takedown

If you consider content in White Rose Research Online to be in breach of UK law, please notify us by emailing eprints@whiterose.ac.uk including the URL of the record and the reason for the withdrawal request.

promoting access to White Rose research papers



Universities of Leeds, Sheffield and York
<http://eprints.whiterose.ac.uk/>

This is an author produced version of a paper published in
APPLIED AND ENVIRONMENTAL MICROBIOLOGY

White Rose Research Online URL for this paper:
<http://eprints.whiterose.ac.uk/42645>

Published paper

Lozano, L; Hernandez-Gonzalez, I; Bustos, P, et al (2010)
Evolutionary Dynamics of Insertion Sequences in Relation to the Evolutionary
Histories of the Chromosome and Symbiotic Plasmid Genes of *Rhizobium etli*
Populations
APPLIED AND ENVIRONMENTAL MICROBIOLOGY
76 (19) 6504-6513
<http://dx.doi.org/10.1128/AEM.01001-10>

1 **Evolutionary Dynamics of Insertion Sequences in Relation to the Evolutionary**
2 **Histories of the Chromosome and Symbiotic Plasmid of *Rhizobium etli* Populations**

3

4

5

6 Luis Lozano^{1*}, Ismael Hernández-González¹, Patricia Bustos¹, Rosa I. Santamaría¹,
7 Valeria Souza², J. Peter W. Young³, Guillermo Dávila¹ and Víctor González¹.

8

9 ¹Centro de Ciencias Genómicas, UNAM, Av. Universidad N/C Col. Chamilpa, Apdo.
10 Postal 565-A Cuernavaca, Morelos, México.

11 ²Departamento de Ecología Evolutiva, Instituto de Ecología, Universidad Nacional
12 Autónoma de México, AP 70-275, CU, Coyoacán 04510, Mexico DF, México.

13 ³Department of Biology, University of York, PO Box 373, YO10 5YW, York, UK.

14

15 *Corresponding Author: Luis Lozano

16 llozano@ccg.unam.mx

17 Phone: 01 777 3291690

18 Fax: 01 777 3175581

19 Centro de Ciencias Genómicas, UNAM. Av. Universidad N/C,
20 Col. Chamilpa. Apdo. Postal 565-A. Cuernavaca, Morelos,
21 Mexico

22

23 Running Head: IS and pSym Dynamics in *R. etli* Populations

24 Keywords: insertion sequence, symbiotic plasmid, genome context

1 **ABSTRACT** (235 words)

2 Insertion sequences (IS) are mobile genetic elements that are distributed in many
3 prokaryotes. In particular, in the genomes of the symbiotic nitrogen fixing bacteria
4 collectively known as rhizobia, IS are fairly abundant in plasmids or chromosomal
5 islands that carry the genes needed for symbiosis. Here, we report an analysis of the
6 distribution and genetic conservation of the IS found in the genome of *Rhizobium etli*
7 CFN42, in a collection of 87 *Rhizobium* strains belonging to populations of different
8 geographical origins. We used PCR to generate presence/absence profiles of the 39 IS
9 found in *R. etli* CFN42, and evaluated whether the IS were located in consistent
10 genomic contexts. We found that the IS from the symbiotic plasmid were frequently
11 present in the analyzed strains, whereas the chromosomal IS were observed less
12 frequently. We then examined the evolutionary dynamics of these strains based on a
13 population genetic analysis of two chromosomal housekeeping genes (*glyA* and *dnaB*)
14 and three symbiotic sequences (*nodC* and the two IS elements). Our results indicate that
15 the IS contained within the symbiotic plasmid have a higher degree of genomic context
16 conservation and a lower nucleotide diversity, genetic differentiation and fewer
17 recombination events compared with the chromosomal housekeeping genes. These
18 results suggest that the *R. etli* populations diverged recently in Mexico, that the
19 symbiotic plasmid also had a recent origin, and that the IS elements have undergone a
20 process of infection and expansion cycle.

21

22

23

24

25

1 **Introduction**

2

3 Insertion sequences (IS) are the smallest transposable elements found in
4 prokaryotes (usually less than 3 kb in size). They encode a transposase and may also
5 encode small hypothetical proteins (4, 9). IS are distinguished by their ability to move
6 within prokaryotic replicons, including both the chromosome and plasmids, and copy
7 themselves into various genomic sites. In this manner, IS elements can inactivate or
8 alter the expression of adjacent genes (4). When IS occur in two or more identical
9 copies within a genome, they can participate in various types of genetic rearrangements
10 (e.g., duplications, inversions and deletions) suggesting that IS may play an important
11 role in the evolution of their hosts by promoting genomic plasticity (34). Due to these
12 evolutionary dynamics, the diversity and distribution of IS elements differs greatly
13 between taxa and even within strains of the same species (27).

14 Various theories have been proposed to explain the evolution of IS elements in
15 laboratory model strains and environmental bacterial populations (8, 18, 25, 29). Two
16 main hypotheses seek to explain how these elements are maintained over the long term
17 in their host genomes: The first proposes that they occasionally generate beneficial
18 mutations, and therefore may represent a selective advantage to their hosts (34). The
19 second suggests that IS elements are genomic parasites that are maintained by their high
20 rate of transposition and might be disseminated among different bacterial lineages by
21 horizontal gene transfer (HGT). Data supporting the second hypothesis has shown that
22 some IS elements may transpose at high rates upon entering a new host (43). After the
23 initial infection, however, purifying selection may continuously remove these elements
24 from the genome. Thus, IS may be under an infection-expansion-extinction cycle that
25 allows them to remain in different bacterial populations within the gene pool (43).

1 These two hypotheses are not contradictory, and the evolutionary dynamics and
2 distribution of IS may differ greatly depending on several factors, including (most
3 notably) the rate of transposition and HGT, as well as selective pressures, population
4 size and the host's habitat (18, 25, 27, 29, X, Y).

5 In the nitrogen-fixing symbiotic bacteria of the genera *Rhizobium*,
6 *Sinorhizobium*, *Mesorhizobium*, *Bradyrhizobium* (of the α -proteobacteria), *Cupriavidus*,
7 and *Burkholderia* (of the β -proteobacteria), IS are particularly abundant in symbiotic
8 plasmids (pSym) and symbiotic chromosomal islands (SI) (2, 12, 14, 19, 20, 42). SI and
9 pSym encode most of the genes needed to establish symbiosis in the roots of
10 leguminous plants through nodule formation and nitrogen fixation (11). It is generally
11 believed that these elements entered the rhizobial genomes through HGT (39, 40).
12 Comparative genomic analyses have shown that both pSym and SI are highly variable
13 with the exception of a common set of genes encoding factors critical to nitrogen
14 fixation (*nif*) and nodulation (*nod*) (5, 14). SI and pSym have been found to have lower
15 GC contents and different codon usages than the corresponding chromosomal and non-
16 symbiotic plasmid sequences, suggesting that they were recently acquired by HGT.

17 Some of these symbiotic elements, such as in the cases the pSym of *R. etli*
18 CFN42 and the SI of *Mesorhizobium loti*, are conjugative and mobile (30, 32). Genomic
19 analysis of *R. etli* CFN42 revealed the presence of 39 IS belonging to different families
20 (14); these were found in the chromosome (11 IS), the 371-kb symbiotic plasmid (14
21 IS), the smaller 192-kb conjugative plasmid, p42a (13 IS) and two other plasmids, p42b
22 and p42c (2 IS). Interestingly, this particular strain shows no evidence of IS disrupting
23 ORFs or having transpositional activity. However, another 42 incomplete IS may be
24 found in the chromosome, pSym, and the conjugative plasmid; these incomplete
25 sequences are truncated or contain stop codons in their coding sequences.

1 Here, we focused on the dynamics and distribution of IS in different populations
2 of the nitrogen-fixing symbiont *Rhizobium etli*. Since the maintenance of IS in bacterial
3 species might depend on their transpositional activity and horizontal transfer rate, the
4 identification of IS in the same genomic contexts across different strains of the same
5 species could provide new insights into their persistence and divergence over short
6 evolutionary periods. To examine the evolutionary dynamics of IS in natural
7 populations of *R. etli*, we characterized the distributions, genomic contexts and
8 sequence variations of IS in isolates of *R. etli* from three populations of different
9 origins, as well as in some other *Rhizobium* species. More specifically, we used PCR to
10 generate presence/absence profiles of the 39 IS found in *R. etli* CFN42, in a collection
11 of 87 strains representing different geographical sites and a gradient of domestication of
12 the bacterial host, the common bean (*Phaseolus vulgaris*). We also evaluated whether
13 the IS were conserved in the same genomic context relative to their position in *R. etli*
14 CFN42, and determined the nucleotide sequence of two IS found in most of the isolates.
15 Several population genetic tests applied to these IS, another pSym gene (*nodC*) and two
16 chromosomal housekeeping genes (*dnaB* and *glyA*) suggest that these two IS elements
17 have been inherited vertically and represent recent components of the *R. etli* gene pool.
18 Finally, the present study strongly suggests that symbiotic plasmids have a recent origin
19 within the *R. etli* populations.

20

21 **MATERIALS AND METHODS**

22

23 **Bacterial strains, growth conditions and DNA extraction**

24 The 87 *Rhizobium* strains used in this study correspond to three different
25 collections (Table 1). Two of the strain collections are from Mexico. The first collection

1 (36) was derived from two plots in a traditional milpa system of native bean landraces
2 (San Miguel Acuexcomac, Puebla); the collection consists of 30 strains that were the
3 dominant strains for several years. The second collection (13) comes from the
4 Michoacan-Guanajuato area, which is the reported center of origin of bean
5 domestication (24); the 30 strains in this collection include a bean plant domestication
6 gradient from wild non-domesticated *Phaseolus vulgaris* to milpa landrace cultivars, as
7 well as wild bean plants that are probably the descendants of cultivated plants. The third
8 collection (33) represents *R. etli* from Spain and includes 27 strains obtained from 21
9 soil samples collected along the Guadalquivir River Valley. These are believed to
10 represent either original native rhizobia or a sub-sample of New-World rhizobia that
11 travel along the original bean seeds (Table 1). Some strains from the Puebla and
12 Spanish collections were previously characterized as *R. gallicum*, *R. giardini*, *R. fredii*
13 and *R. leguminosarum* (Table 1).

14 The various *Rhizobium* strains were grown in LY medium for 24-48 h. Genomic
15 DNA was extracted with a GenomicPrep DNA Isolation Kit (Amersham Biosciences)
16 and the DNA concentration was determined with a DyNA Quant 200 fluorometer
17 (Hoefer). Two complete sequenced *R. etli* strains, CFN42 and CIAT652, were used as
18 references strains.

19

20 **PCR amplification and DNA sequencing**

21 We first localized each of the 39 selected IS within the CFN42 genome, thereby
22 “anchoring” our examination of whether these IS conserved their genomic locations
23 (i.e., their synteny) across the wild *R. etli* isolates from different geographical regions.
24 These 39 IS elements represent 11 different families and some of these families have
25 identical or near identical copies: 6 and 4 copies for 2 different elements of the IS66

1 family, 5 copies for IS630, 2 for IS1111 and 2 for IS21. We then designed specific PCR
2 primers spanning the immediate neighborhoods of the 39 studied IS. These primers
3 allowed us to test for conservation (i.e. amplification of a fragment of similar size to the
4 one predicted in CFN42). In cases where we found two contiguous IS elements, we
5 designed primers in the neighboring genes of both IS. There were three possible results:
6 i) no PCR product, indicating that there were no identical priming sites in the wild
7 isolate (this was not an absolute positive or negative result for the presence of the IS); ii)
8 a small DNA fragment, equivalent to the distance between the two sequences near the
9 insertion site but lacking the IS (a negative result for the IS); and iii) a large DNA
10 fragment of similar size to that in CFN42 (a positive result for the IS). The PCR
11 reactions were performed in a DNA thermal cycler (Gene Amp 9700; Applied
12 Biosystems) in a 25 µl reaction mixture containing approximately 10 ng genomic
13 template DNA, 2 mM MgCl₂, 1 mM dNTPs, 5 pmol of each primer, and 2 U of TAQ
14 Polymerase (AltaEnzymes). The mixtures were subjected to 5 min denaturation at 94°C
15 followed by 30 cycles of 1 min at 94°C, 1 to 3 min at 58 to 62°C, and 3 min at 72°C.

16 For a comparison of evolutionary dynamics, we performed direct sequencing on
17 the following: two IS elements, ISRel4 and ISRel2 from pSym, which were present in
18 the highest proportion of test samples (see Results and Fig. 2); *nodC*, which is a pSym
19 gene encoding an N-acetylglucosaminyltransferase that participates in the nodulation
20 process; and two chromosomal housekeeping genes, *glyA* (serine
21 hydroxymethyltransferase) and *dnaB* (replicative DNA helicase), which have been
22 proposed to serve as predictors of genome relatedness because their sequence
23 divergence rates reflect the overall rate of genome divergence (44). Internal PCR
24 primers were designed for the latter three genes to obtain partial sequences of each
25 gene. The sequencing reaction mixtures contained approximately 10 ng genomic

1 template DNA, 1 mM MgCl₂, 1 mM dNTPs, 5 pmol of each primer, and 1U of rTth
2 Polymerase XL (Applied Biosystems). The mixtures were subjected to 5 min
3 denaturation at 92°C followed by 30 cycles of 30 sec at 92°C and 1 to 6 min at 58 to
4 62°C. The PCR products were purified with an Exo/SAP kit (Affymetrix) and
5 sequenced using a Dye-terminator cycle sequencing kit (Perkin Elmer Applied
6 Biosystems). The sequencing reactions were run on an ABI 3730 sequencer (Applied
7 Biosystems).

8

9 **Sequence analysis, alignments and phylogenetic reconstruction**

10 Sequence quality analysis, assembly and comparison were performed using the
11 SeqScape software V2.5 (Applied Biosystems). For IS characterization, we compared
12 all of the putative transposases in the complete annotated genome sequence of *R. etli*
13 CFN42, to the nr Database of the NCBI and the Insertion Sequence Database (35), with
14 BLASTp (1).

15 To define a complete and (most likely) functional IS, we used the same criteria
16 applied by the authors of the Insertion Sequence Database, as follows: (i) similar
17 organization of the genes; (ii) sequence similarity of the transposases and inverted
18 repeats; and (iii) the presence of direct repeats (14). Gene alignments for the assessed
19 genes and IS elements were done using the MUSCLE program (7). For each alignment,
20 the best substitution model was determined using FindModel (31). Phylogenetic
21 reconstruction was performed using the maximum likelihood method implemented in
22 PHYML (16). The analysis was carried out with a nonparametric bootstrap analysis of
23 100 replicates for each alignment. The complete chromosome and plasmid sequences of
24 *R. etli* CFN42 and CIAT652 were compared with the Mummer program (23). The
25 genomic contexts of the IS in the two strains were examined using Perl programs. The

1 rate of synonymous and non-synonymous substitutions, dn/ds ratio, was evaluated using
2 SNAP (22).

3

4 **Population genetic parameter estimation**

5 DNASP version 4 (26) was used to assess the following parameters: the average
6 nucleotide divergence between populations (D_{xy}); the average nucleotide diversity
7 within populations; the average nucleotide divergence per site (π); the average
8 nucleotide diversity at synonymous (π_s) and non-synonymous sites (π_{ns}); gene flow
9 estimates (F_{st} and Nm); genetic differentiation estimates (K_s , K_{st} , Z and S_{nn}); the numbers
10 of shared haplotypes, fixed and shared polymorphisms; and Tajima's D neutrality test.

11

12 **Recombination analysis**

13 Recombination tests were performed with RDP3 (28) and SplitsTree4 (17) for the five
14 genes studied in each of the collections. The RDP3 software was applied to detect and
15 analyze recombination signals using eight different methods (RDP, GENECONV,
16 Bootscan, MaxChi, Chimaera, SiScan, 3SEQ and LARD), with 100 permutation steps, a
17 Bonferroni correction for multiple tests, and a P -value threshold of 0.05. Recombination
18 events and breakpoints were accepted if they were detected by two or more methods
19 that used different approximations to detect recombination. The recombination events
20 were then confirmed independently for each recombinant strain, as suggested in the
21 literature (RDP website user manual; <http://darwin.uvigo.es/rdp/rdp.html>). Specific
22 parameter modifications were used for each gene alignment depending the number of
23 sequences in the analysis. Split decomposition and neighbor net analyses were
24 performed with the SplitsTree4 software. Both phylogenetic networks represent
25 incompatibilities within and between datasets with the use of splits. For the neighbor net

1 analysis, we applied a split filter based on the network weight and a threshold value of
2 95% (3).

3

4 **Nucleotide sequence accession numbers**

5 The relevant sequence data have been deposited in the GenBank database under
6 accession numbers GUO84443 to GUO84796.

7

8 **RESULTS**

9 **Differentiation of *R. etli* populations**

10 To study the evolutionary dynamics of the IS of *R. etli*, we first asked if the three
11 studied collections of *R. etli* strains, which were obtained from the root nodules of *P.*
12 *vulgaris* plants located in different areas, were geographically structured. The
13 phylogenies constructed based on the chromosomal housekeeping genes, *glyA* and
14 *dnaB*, were relatively similar to each other, with a few differences in the groupings of
15 some strains (data not shown). In order to obtain a phylogenetic reconstruction that
16 more accurately represent the evolutionary relationships among the 87 strains contained
17 within the three collections, we concatenated the sequence alignments (2238 bp) of both
18 genes. We obtained three large *R. etli* clusters that broadly corresponded to the three
19 different geographical areas from which the strains were collected. The Puebla
20 collection formed one group with some intermingled strains from the Michoacan-
21 Guanajuato and Spanish collections (Fig. 1). The second group consisted of the
22 Michoacan-Guanajuato strains with three intermingled sequences belonging to the
23 Spanish collection. The third *R. etli* group represented the Spanish collection with a
24 single intermingled sequence from the Puebla collection. The Michoacan-Guanajuato
25 and Spanish collections were more closely related to one another than to the Puebla

1 collection (Fig. 1). A distant fourth group, which was later used as outgroup, contained
2 several *R. gallicum* strains from Puebla and Spain. Since this phylogenetic
3 reconstruction clearly separates the three different collections according to their
4 geographic origins, we considered them to be different populations in the following
5 analyses.

6

7 **IS profiles of the *R. etli* populations**

8 We next examined which IS from *R. etli* CFN42 were also maintained in the
9 same genomic context among the *R. etli* isolates from the three different populations.
10 Our PCR-based profiling (see Methods) yielded either positive or negative results for 24
11 of the 39 IS, and showed marked differences in the conservation of the various IS
12 among the tested populations (Fig 2). ISRel4 and ISRel2 (from pSym of *R. etli* CFN42)
13 were the most common IS, present in more than 50% of the individual strains.
14 Sequencing experiments showed that ISRel9 (also from pSym) in most cases was
15 actually a hypothetical protein (hyp304) not related to any IS element and similar in size
16 to ISRel9. Only five isolates (all from Puebla) had an intact ISRel9 in the expected
17 genomic context; these elements had an average identity of 100%. ISRel5 and ISRel12
18 of pSym were relatively frequent, as they were present in 25% of the tested strains. The
19 other seven IS from pSym were not common. The IS from p42a were found in the
20 Puebla and Michoacan-Guanajuato populations, but not in the collection from Spain,
21 while the chromosomal IS were found only in the Puebla strains. Interestingly, the
22 CFN42 reference strain was found to be close to the Michoacan-Guanajuato population
23 (Fig. 1).

24

25 **DNA sequence divergence of ISRel4 and ISRel2**

1 The observation that some IS occur in the same genomic context suggests that
2 these IS may have been present since the arrival of the pSym in the populations, as it is
3 far less likely that independent transposition events could have inserted the IS into the
4 same genomic sites. It is important to mention that the ISRel4 and ISRel2 belong to IS
5 families, ISAs1 and IS481 respectively, that do not have specific target sequences for
6 transposition. To investigate the evolutionary dynamics of the two IS elements, ISRel4
7 and ISRel2, that show conservation of their genomic contexts, we looked for differences
8 in the degree of diversification and selection pressures among these IS elements, as well
9 as in other symbiotic and chromosomal housekeeping genes (Supp. Table 1). The DNA
10 sequence conservation was very high for ISRel4 and ISRel2, with an average identity of
11 98.7% and 99.4%, respectively. This is surprisingly high compared to the sequence
12 identities for the housekeeping genes *glyA* and *dnaB*, which were 75.3% and 76.4%,
13 respectively. Given that the diversification of chromosomal genes could differ from that
14 on the symbiotic plasmid, we sequenced the *nodC* gene, which encodes an N-
15 acetylglucosaminyltransferase involved in the synthesis of sugar backbones for the
16 production of lipo-oligosaccharides (critical to nodulation signaling). This was done in
17 order to have another gene to trace the diversification differences between the pSym and
18 the chromosome. The average nucleotide identity among the 44 analyzed *nodC*
19 sequences was 96.7%. The nucleotide diversity (π) values for the pSym sequences
20 (ISAs1, IS481 and *nodC*) were very low compared with those of *glyA* and *dnaB* (Table
21 2). The total average π values of *dnaB* and *glyA* were 0.066 and 0.088, respectively,
22 whereas those for ISRel4, ISRel2, and *nodC* were 0.002, 0.002, and 0.014, respectively
23 (Table 2). Since the number of polymorphic sites and the π values for ISRel4, ISRel2,
24 and *nodC* were very low, we hypothesize that these sequences may have entered the
25 gene pool of the three populations relatively recently.

1

2 **Recent origins of ISRel4, ISRel2, and pSym**

3 To evaluate the hypothesis that ISRel4 and ISRel2 may have originated recently,
4 we measured their average nucleotide diversities at synonymous sites (π_s) and non-
5 synonymous sites (π_{ns}), and compared these values with those for *nodC*, *glyA*, and
6 *dnaB*. The π_s value reflects the age of an allele in the genetic pool; genes with lower π_s
7 values are generally considered to have recent common ancestor. Our results revealed
8 that the π_s values were low for ISRel4, ISRel2, and *nodC*, but high for *dnaB* and *glyA*,
9 whereas the π_{ns} values were low for all five sequences (Table 2). Similar results were
10 obtained when these π_s and π_{ns} values were measured within each population (Supp.
11 Table 2). These findings indicate that both the IS elements and *nodC* might have entered
12 recently in the gene pool of the three populations.

13 Overall, the dn/ds ratio was lower for the chromosomal genes than for the pSym
14 genes. To test what this might mean in terms of pSym evolution, we compared the
15 shared haplotypes for chromosomal and pSym genes among the three populations. We
16 did not find any shared *dnaB* and *glyA* (i.e. chromosomal) haplotypes between the
17 populations. In contrast, we identified a number of shared haplotypes for the pSym
18 genes (*nodC*, ISRel2 and ISRel4) (data not shown). These results suggest either that
19 pSym has a relatively recent origin, or that there is high gene flow among pSym
20 sequences; regardless of either option, it seems that the pSym sequences have a different
21 evolutionary history than the chromosomal genes. To further test the possible recent
22 origin of pSym sequences in the gene pool of the three populations, we compared the
23 mean genetic divergence values between populations (D_{xy}) to the π values describing
24 the mean genetic divergence within the populations. Higher values of D_{xy} indicate
25 increasing time since population divergence. Table 3 shows that the D_{xy} values were

1 higher than the mean within-group divergence (π) for the chromosomal genes,
2 potentially reflecting phylogenetic differentiation among the three populations. In the
3 case of the pSym sequences, however, the Dxy and π values were very low and did not
4 show clear phylogenetic differentiation between the populations. This further supports
5 the idea that the pSym sequences have a relatively recent origin in the gene pool of the
6 three populations, and have followed a different evolutionary history than the
7 chromosomal genes.

8 We applied additional genetic differentiation tests (K_s^* , K_{st} , Z , Z^* , S_{nn}) to the
9 pSym and chromosomal sequences to test whether the two IS and *nodC* have a lower
10 level of genetic differentiation than the chromosomal sequences which would further
11 support a recent origin within the three populations. The results of the genetic
12 differentiation tests for the two housekeeping genes were highly significant ($p < 0.001$),
13 supporting the idea of genetic differentiation among the three populations (Table 4 and
14 Supp. Table 3). In contrast, the three pSym sequences had much lower levels of genetic
15 differentiation across the three populations. For example, comparison of the pSym
16 sequences from the two Mexican populations (those from Puebla and Guanajuato-
17 Michoacan) yielded K_{st} values that were either non-significant or were just barely
18 significant at $p < 0.05$ (Table 4). Although the chromosomal housekeeping genes showed
19 evidence of genetic differentiation, we did not find evidence of fixed polymorphisms;
20 this may indicate that the three populations also diverged recently. Given that the pSym
21 sequences had low π_s values, their Dxy and π values were not well differentiated, had
22 several shared alleles, and the genetic differentiation tests were not significant, it is not
23 surprising that these sequences also lacked fixed polymorphisms.

24 Next, we analyzed whether there was some degree of gene flow across the three
25 populations. If a F_{st} value above 0.25 and a N_m value > 1 were taken to represent a

1 significant between-population gene flow (37), most of the analyzed genes could be
2 considered to have evidence of gene flow. The pSym sequences often had higher Fst
3 values than the chromosomal genes (Table 4). A comparison between the Mexican and
4 Spanish populations also yielded values indicative of gene flow. Given the geographical
5 distance between these populations, it seems conceivable that these values could reflect
6 a recent divergence of the three populations and a recent origin for the pSym genes.

7 There is published evidence indicating that genetic exchange occurred within *R.*
8 *etli* populations from a single field; this process involved both plasmid and
9 chromosomal loci (36). Accordingly, we hypothesized that the number of recombination
10 events would be different between the chromosomal and pSym genes if their divergence
11 times were clearly distinct. To assess this, we used the RDP3 software package to
12 implement eight different recombination tests to search for recombination events across
13 the three populations (Supp. Table 4). The chromosomal genes, *dnaB* and *glyA*, showed
14 evidence of four and two intragenic recombination events, involving ten and five
15 strains, respectively. In contrast, only one pSym sequence (ISRel4) showed a
16 recombination event, even though the low sequence divergence make it harder to detect
17 recombination. Interestingly, some of the same recombination events were found in
18 strains of both *R. etli* and *R. gallicum*; moreover, four recombination events found
19 between populations could explain some intermingled strains in the phylogeny of the
20 three populations (Fig. 1 and Supp. Table 4).

21 Next, we used a phylogenetic network approach, consisting of split
22 decomposition and neighbor net analyses to test whether the inconsistencies in the
23 phylogenetic reconstructions could be due to recombination events. The split
24 decomposition analysis of the two chromosomal genes showed a partition in the data
25 between sequences from the *R. etli* and *R. gallicum* species. In the case of *glyA*, we

1 found some splits within the *R. etli* strains; these could come from the detected
2 recombination events. ISRel4 was the only pSym sequence for which we obtained a
3 split that clearly represents the unique recombination event we detected in the *R. etli*
4 and *R. gallicum* strains. The neighbor net analysis, which is more sensitive to the
5 conflicting sequence data that may represent recombination events, showed that the two
6 chromosomal genes had several splits, some of which were related to the detected
7 recombination events (Supp. Fig. 1). ISRel4 yielded the same split we had found in the
8 split decomposition analysis. The other two pSym sequences failed to reveal splits with
9 either analysis; this was consistent with the lack of any recombination events detected
10 by our RDP3 analysis of these sequences. The relative lack of splits and recombination
11 events in the pSym sequences supports our hypothesis that these genes (and probably
12 the symbiotic plasmid itself) originated relatively recently in the three populations.

13 Another possibility that could explain the high DNA conservation of the pSym
14 sequences and the conserved genomic contexts of ISRel2 and ISRel4 is that selection
15 pressures may inhibit variation at these sequences. However, Tajima's D neutrality tests
16 showed that most of the observed nucleotide substitution patterns in the chromosomal
17 and pSym sequences from the three populations were under a neutral equilibrium model
18 (i.e., all statistical results were non-significant; Table 2 and Supp. Table 2). The neutral
19 equilibrium model was rejected only for *nodC*. However, this could be due to the effect
20 of a decrease in population size and/or balancing selection.

21

22 **IS elements in *Rhizobium etli* symbiotic plasmids**

23 We recently sequenced another *Rhizobium etli* strain, the Costa Rican strain
24 CIAT652, which has one chromosome, three plasmids, and 22 IS elements distributed
25 mainly in the chromosome and the symbiotic plasmid (15). Comparative analysis of the

1 CFN42 and CIAT652 symbiotic plasmids showed that their sequences were highly
2 conserved, with a nucleotide identity value of 99% (Fig. 3). In contrast, the identities of
3 the chromosome and other plasmid sequences ranged from 90 to 95%. The two pSyms
4 differed mainly in the presence/absence of certain genomic regions, the diversity of IS
5 families, and presence/absence the IS elements in a given genomic context. This
6 indicates that the major evolutionary events modifying these two symbiotic plasmids
7 were the gains/losses of genomic regions and the losses/translocations of IS elements. In
8 contrast, the chromosomes and other plasmids did not contain IS in the same genomic
9 contexts. Comparison of the IS elements in the chromosomes and plasmids indicated
10 that the symbiotic plasmids may have originated externally, and harbored IS that
11 thereafter transposed to the chromosomes. Both CFN42 and CIAT652 harbored IS
12 elements that were 100% identical in sequence between the chromosome and pSym. If
13 these plasmids have a recent origin (as suggested by our population genetic analyses),
14 this would seem to indicate that the IS had recently moved to the chromosome in both
15 strains, in the manner of an infection-expansion-extinction cycle (43).

16

17 **DISCUSSION**

18

19 In order to examine the evolutionary dynamics of IS elements in relation to the
20 evolutionary history of the chromosomal and symbiotic plasmid genes, we herein
21 compared the distribution and genetic diversity of these sequences in three *R. etli*
22 populations. First, we used the concatenated alignment of the chromosomal *glyA* and
23 *dnaB* genes to perform a phylogenetic reconstruction, and found that the three tested *R.*
24 *etli* populations were almost completely differentiated according to their geographic
25 origins. Moreover, they were clearly differentiated from the *R. gallicum* strains, which

1 formed a fourth clade. The close phylogenetic relationship between the Michoacan-
2 Guanajuato and Spanish populations was highlighted by the individual phylogenies of
3 *dnaB* and *glyA* especially the one based on *glyA*, where several strains appeared to be
4 intermingled in both populations. This pattern might be explained by recombination
5 events, which were assessed using the RDP3 software (Supp. Table 4). However,
6 another possibility is that the intermingled strains from Guanajuato could be closely
7 related to migrant strains isolated from the Spanish population. The phylogeny made
8 with the concatenated alignment of *dnaB* and *glyA* was used as the reference for
9 mapping the IS distributions.

10 Comparison of the IS presence/absence profiles of the *R. etli* populations to that
11 of the CFN42 reference strain showed that the most conserved IS were on the symbiotic
12 plasmid. Because the majority of the IS elements did not maintain their genomic
13 context, the presence/absence profiles were inconsistent with the housekeeping genes
14 phylogeny; some strains that were closely related in the phylogeny had very different IS
15 profiles. For example, most of the Puebla strains had more similar IS profiles to that of
16 the model strain, CFN42, which belonged to the Guanajuato-Michoacan clade. These
17 data demonstrate that the distribution of IS elements between and within populations is
18 strain-specific in concordance with the HGT and transpositional dynamics of IS (4).

19 We analyzed two of the IS elements, ISRe12 and ISRel4 (both from pSym) in
20 detail, as they were the most conserved, in terms of their genomic context, across the
21 three *R. etli* populations. This pattern may indicate that some selective advantage is
22 conferred by the presence of these IS elements, or it could be a consequence of the
23 apparently recent origin of the pSym plasmid. Other studies on the evolutionary
24 dynamics of IS elements have offered different explanations for the maintenance of
25 such elements in populations or strains of the same species. For example, a report on IS

1 in strains of *Helicobacter pylori* from different geographical origins suggested that IS
2 are ancient components of the *H. pylori* gene pool, and they evolve at approximately the
3 same rate as normal chromosomal genes (18). A paper on IS within a single population
4 of the hyperthermophilic archaeon, *Pyrococcus*, suggested that the high frequencies of
5 some IS were due to genetic drift (8). Neither of these explanations seems to be
6 applicable to the IS in *R. etli* examined herein. Instead, we think that the IS have
7 undergone rapid turnover in the population, and that the contextual and sequence
8 conservation of ISRel2 and ISRel4 should be viewed within the evolutionary history of
9 pSym.

10 Due to the transpositional nature of IS elements, they are not expected to be
11 found at the same genomic site in populations of different geographical origin.
12 However, we herein observed such conservation for ISRel2 and ISRel4. One
13 explanation for our observation is that that the IS might have been frequently found in
14 the same genomic context as consequence of the small population size of the *R. etli*
15 collections analyzed in this work. The probability of finding several strains with the
16 same IS in the same location is lower for large populations than for small populations
17 (38). In our case, however, this may not be the best explanation given the characteristics
18 of the three populations (see Methods) and the geographical distances separating the
19 collection sites. Another possibility is that these IS elements were recently acquired by
20 the studied *R. etli* strains. Comparisons among these IS elements, *nodC* and the
21 chromosomal genes suggested that the pSym sequences recently entered the *R. etli* gene
22 pool of the three populations. Population genetic analyses of the pSym sequences
23 showed that they have the following features: i) low π values; ii) Dxy and π values that
24 imply no clear differentiation; iii) several different shared alleles; and iv) non-
25 significant results from their genetic differentiation tests. In contrast, opposing results

1 were obtained for the chromosomal genes, *glyA* and *dnaB*. Thus, the differences noted
2 in our population genetic analyses suggest that the chromosome and pSym have
3 different evolutionary histories in the studied *R. etli* populations. These findings, along
4 with the high sequence conservation found in the pSym sequences of CFN42 and
5 CIAT652, support the proposition that pSym could have a recent origin in *Rhizobium*
6 *etli* (14). However, the pSym sequences from CFN42 and CIAT652 have clear
7 differences in terms of the presence/absence of certain genomic regions (some as large
8 as 50 kb) and IS elements; notably, the gains and losses of IS elements appear to be the
9 main events differentiating these two pSyms.

10 A prior analysis of the complete genomes of *R. etli* CFN42 and *R. etli* CIAT652
11 demonstrated that the IS do not disrupt genes (with one exception), and these strains
12 appear to harbor the same pSym plasmid (15). Horizontal transfer of this plasmid to *R.*
13 *etli* CFN42 and *R. etli* CIAT652 could explain the asymmetric distribution of IS, which
14 were only found in these plasmids and in the chromosome. Some IS elements probably
15 moved from the plasmids to the chromosome, but they do not appear to have moved to
16 other plasmids. The asymmetrical distribution of IS in the chromosome and plasmids
17 could be due the size of the different replicons in the *R. etli* CFN42 genome (41). In the
18 chromosome of CFN42 there are six probable recent transposition events of IS elements
19 (100% identity between the chromosomal IS and another IS in the pSym or conjugative
20 plasmid). This implies a transposition rate of one IS per 730 kb, which is higher than the
21 size of any other replicon of CFN42 (14), so the asymmetrical distribution could be the
22 product of the chromosome size.

23 The complete genomes of the two strains of *R. etli*, CFN42 and CIAT652,
24 revealed the presence of different IS family members in different copy numbers. It has
25 been proposed that when new IS elements enter the genome, they actively transpose to

1 different genomic positions (43). This hypothesis may help explain the behavior of
2 some of the IS in *R. etli* CFN42. For instance, members of the IS66 family, which is a
3 very common IS family in rhizobial species (27), were found in several copies in the
4 pSym plasmids of CFN42 and CIAT652, and IS66 copies with 100% nucleotide
5 sequence identity can also be found in the chromosome. The most plausible explanation
6 for this finding is that the IS originally entered the genome via plasmids and were then
7 transposed to the chromosome.

8 The present work provides useful new insights into differences of the
9 distribution and genomic context maintenance of IS elements in natural populations of
10 *Rhizobium etli*. Although the pSym in these populations seems to be of relatively recent
11 origin, the harbored IS appear to be active elements that may participate in the genomic
12 plasticity of these organisms. Our comparisons of the two genomes of *R. etli* and across
13 the natural populations provides further evidence that different IS can rapidly expand
14 when they arrive in a new host genome, as recently proposed by Wagner (2006).

15

16 **ACKNOWLEDGMENTS**

17

18 We thank Santiago Castillo for critical reading of the manuscript. We thank Susana Brom
19 (Universidad Nacional Autónoma de México) for providing the Spanish strains; Valeria
20 Souza for providing those from Puebla and Michoacan-Guanajuato; and J. Espíritu for
21 technical and computational assistance. This work was supported by the grants from
22 CONACyT (grant U4633), PAPIIT-UNAM (grant IN223005) and the Natural
23 Environment Research Council.

24 L.L. was responsible for data analysis and manuscript preparation. L.L. and V.G. were
25 responsible for the experimental design. I.H-G was responsible for bioinformatic

1 analysis. P.B. and R.I.S. were responsible for DNA sequencing. L.L., V.S., J.P.W.Y.,
2 G.D. and V.G. were responsible for discussion of the data.

3

4 **REFERENCES**

5

- 6 1. **Altschul, S., Madden, T., Schaffer, A., Zhang, J., Zhang, Z., Millar, W. and**
7 **D. Lipman.** 1997. Gapped BLAST and PSI-BLAST: a new generation of
8 protein database search programs. *Nucleic Acids Res.* 25(17):3389-402.
- 9 2. **Amadou, C., Pascal, G., Mangenot, S., Glew, M., Bontemps, C., Capela, D.,**
10 **Carrere, S., Cruveiller, S., Dossat, C., Lajus, A., Marchetti, M., Poinot, V.,**
11 **Rouy, B. Servin, Z., Saad, M., Schenowitz, C., Barbe, V., Batut, J., Medigue**
12 **C. and C. Masson-Boivin.** 2008. Genome sequence of the beta-rhizobium
13 *Cupriavidus taiwanensis* and comparative genomics of rhizobia. *Genome Res.*
14 18(9): 1472-83.
- 15 3. **Bailly, X., Olivieri, I., De Mita, S., Cleyet-Marel J. C. and G. Bena.** 2006.
16 Recombination and selection shape the molecular diversity pattern of nitrogen-
17 fixing *Sinorhizobium* sp. associated to *Medicago*. *Mol. Ecol.* 15(10): 2719-34.
- 18 4. **Chandler, M. and J. Mahillon.** 2002. Insertion Sequences revisited. Pp.305-
19 366. In Craig, L. et al. (Eds). *Mobile DNA II*, ASM Press, Washington D. C.
- 20 5. **Crossman L. C., Castillo-Ramírez S., McAnnula C., Lozano L., Vernikos G.**
21 **S., Acosta J. L., Ghazoui Z. F., Hernández-González I., Meakin G., Walker**
22 **A. W., Hynes M. F., Young J. P., Downie J. A., Romero D., Johnston A. W.,**
23 **Dávila G., Parkhill J. And V. González.** 2008. A common genomic framework
24 for a diverse assembly of plasmids in the symbiotic nitrogen fixing bacteria.
25 *PloS ONE* 3(7): e2567.

- 1 6. **Doolittle, W. F. and C. Sapienza.** 1980. Selfish genes, the phenotype paradigm
2 and genome evolution. *Nature*. 284:601-603.
- 3 7. **Edgar, R.** 2004. MUSCLE: a multiple sequence alignment method with reduced
4 time and space complexity. *BMC Bioinformatics*. 5:113.
- 5 8. **Escobar-Páramo, P., Ghosh, S. and J. DiRuggiero.** 2005. Evidence for
6 genetic drift in the diversification of a geographically isolated population of the
7 hyperthermophilic archaeon *Pyrococcus*. *Mol. Biol. Evol.* 22 (11):2297-2303.
- 8 9. **Filée, J., Siguier, P. and M. Chandler.** 2007. Insertion sequence diversity in
9 Archea. *Microbiol. Mol. Biol. Rev.* 71 (1): 121-157.
- 10 10. **Flores, M., Morales, L., Avila, M., Bustos, P., González, V., Garcia, D.,**
11 **Mora, Y., Guo, X., Collado-Vides, J., Piñero, D., Davila, G Mora, J. and R.**
12 **Palacios.** 2005. Diversification of DNA sequences in the symbiotic genome of
13 *Rhizobium etli*. *J. Bacteriol.* 187(21):7185-92.
- 14 11. **Freiberg, C., Fellay, R., Bairoch, A., Broughton, W. J., Rosenthal A. and X.**
15 **Perret.** 1997. Molecular basis of symbiosis between *Rhizobium* and legumes.
16 *Nature* 387(6631): 394-401.
- 17 12. **Galibert, F., Finan, T. M., Long, S. R., Puhler, A., Abola, P., Ampe, F.,**
18 **Barloy-Hubler, F., Barnett, M. J., Becker, A., Boistard, P., Bothe, G.,**
19 **Boutry, M., Bowser, L., Buhrmester, J., Cadieu, E., Capela, D., Chain, P.,**
20 **Cowie, A., Davis, R. W., Dreano, S., Federspiel, N. A., Fisher, R. F., Gloux,**
21 **S., Godrie, T., Goffeau, A., Golding, B., Gouzy, J., Gurjal, M., Hernandez-**
22 **Lucas, I., Hong, A., Huizar, L., Hyman, R. W., Jones, T., Kahn, D., Kahn,**
23 **M. L., Kalman, S., Keating, D. H., Kiss, E., Komp, C., Lelaure, V., Masuy,**
24 **D., Palm, C., Peck, M. C., Pohl, T. M., Portetelle, D., Purnelle, B.,**
25 **Ramsperger, U., Surzycki, R., Thebault, P., Vandenbol, M., Vorholter, F. J.,**

- 1 **Weidner, S., Wells, D. H., Wong, K., Yeh K. C. and J. Batut.** 2001. The
2 composite genome of the legume symbiont *Sinorhizobium meliloti*. *Science*
3 293(5530): 668-72.
- 4 13. **Gasca, J.** Genética de poblaciones de *Rhizobium etli* asociado a *Phaseolus*
5 *vulgaris* en dos localidades del bajío mexicano. Tesis de Licenciatura. Instituto de
6 Ecología. UNAM. 102 pag.
- 7 14. **González, V., Santamaría, R., Bustos, P., Hernandez-González, I.,**
8 **Medrano-Soto, A., Moreno-Hagelsieb, G., Janga, S., Ramirez, M., Jiménez-**
9 **Jacinto, V., Collado-Vides, J. and G. Davila.** 2006. The partitioned *Rhizobium*
10 *etli* genome: genetic and metabolic redundancy in seven interacting replicons.
11 PNAS 103(10):3834-3839.
- 12 15. **González, V., Acosta, J.L., Santamaría, R., Bustos, P., Fernandez, J.L.,**
13 **Hernandez-González, I., Diza, R., Flores, M., Palacios, R., Mora, J. and G.**
14 **Davila.** 2010. Conserved symbiotic plasmid DNA sequences in the
15 Multireplicon pangenomic structure of *Rhizobium etli*. *Appl. Environ.*
16 *Microbiol.* 76(5):1604-1614.
- 17 16. **Guindon, S. and O. Gascuel.** 2002. Efficient biased estimation of evolutionary
18 distances when substitution rates vary across sites. *Mol Biol Evol.* (4):534-43.
- 19 17. **Huson, D. H. and D. Bryant.** 2006. Application of phylogenetic networks in
20 evolutionary studies. *Mol. Biol. Evol.* 23(2): 254-67.
- 21 18. **Kalia, A., A. K. Mukhopadhyay, A. K., Dailide, G., Ito, Y., Azuma, T.,**
22 **Wong B. C. and D. E. Berg.** 2004. Evolutionary dynamics of insertion
23 sequences in *Helicobacter pylori*. *J. Bacteriol.* 186(22): 7508-20.
- 24 19. **Kaneko, T., Nakamura, Y., Sato, S., Asamizu, E., Kato, T., Sasamoto, S.,**
25 **Watanabe, A., Idesawa, K., Ishikawa, A., Kawashima, K., Kimura, T.,**

- 1 **Kishida, Y., Kiyokawa, C., Kohara, M., Matsumoto, M., Matsuno, A.,**
2 **Mochizuki, Y., Nakayama, S., Nakazaki, N., Shimpo, S., Sugimoto, M.,**
3 **Takeuchi, C., Yamada M. and S. Tabata.** 2000. Complete genome structure of
4 the nitrogen-fixing symbiotic bacterium *Mesorhizobium loti*. *DNA. Res.* 7(6):
5 331-8.
- 6 20. **Kaneko, T., Nakamura, Y., Sato, S., Minamisawa, K., Uchiumi, T.,**
7 **Sasamoto, S., Watanabe, A., Idesawa, K., Iriguchi, M., Kawashima, K.,**
8 **Kohara, M., Matsumoto, M., Shimpo, S., Tsuruoka, H., Wada, T., Yamada**
9 **M. and S. Tabata.** 2002. Complete genomic sequence of nitrogen-fixing
10 symbiotic bacterium *Bradyrhizobium japonicum* USDA110. *DNA. Res.* 9(6):
11 189-97.
- 12 21. **Kidwell, M.G., and D. R. Lisch.** (2002) Transposable elements as sources of
13 genomic variation. Pp.305-366. In Craig, L. et al. (Eds). *Mobile DNA II*, ASM
14 Press, Washington D. C.
- 15 22. **Korber B.** 2000. HIV Signature and Sequence Variation Analysis.
16 Computational Analysis of HIV Molecular Sequences, Chapter 4, pages 55-72.
17 Allen G. Rodrigo and Gerald H. Learn, eds. Dordrecht, Netherlands: Kluwer
18 Academic Publishers.
- 19 23. **Kurtz, S., Phillippy, A., Delcher, A., Smoot, M., Shumway, M., Antonescu,**
20 **C. and S. Salzberg.** 2004. Versatile and open software for comparing large
21 genomes. *Genome Biol.* 5:R12.
- 22 24. **Kwak, M. and P. Gepts.** 2009. Structure of genetic diversity in the two major
23 gene pools of common bean (*Phaseolus vulgaris* L., Fabaceae). *Theor. Appl.*
24 *Genet.* 118(5): 979-92.

- 1 25. **Lenski, R., Winkwirth, C., and M. Riley.** 2003. Rates of DNA sequence
2 evolution in experimental populations of *Escherichia coli* during 20,000
3 generations, *J. Mol. Evol.* 56:498-508.
- 4 26. **Librado, P. and J. Rozas.** 2009. DnaSP v5: a software for comprehensive
5 analysis of DNA polymorphism data. *Bioinformatics* 25(11): 1451-2.
- 6 27. **Mahillon, J., Léonard, C., and M. Chandler.** 1999. IS elements as constituents
7 of bacterial genomes. *Res Microbiol.* 150:675-687.
- 8 28. **Martin, D. P., Williamson, C. and D. Posada.** 2005. RDP2: recombination
9 detection and analysis from sequence alignments. *Bioinformatics* 21(2): 260-2.
- 10 29. **Papadopoulos, D., Schneider, D., Meier-Eis, J., Arber, W., Lenski, R. E.**
11 **and Blot, M.** et al. 1999. Genomic evolution during a 10,000-generation
12 experiment with bacteria. *Proc. Natl. Acad. Sci.* 96:3807-3812.
- 13 30. **Perez-Mendoza, D., Dominguez-Ferreras, A., Munoz, S., Soto, M. J.,**
14 **Olivares, J., Brom, S., Girard, L., Herrera-Cervera J. A. and J. Sanjuan.**
15 2004. Identification of functional mob regions in *Rhizobium etli*: evidence for
16 self-transmissibility of the symbiotic plasmid pRetCFN42d. *J. Bacteriol.*
17 186(17): 5753-61.
- 18 31. **Posada, D. and K. A. Crandall.** 1998. MODELTEST: testing the model of
19 DNA substitution. *Bioinformatics* 14(9): 817-8.
- 20 32. **Ramsay, J. P., Sullivan, J. T., Stuart, G. S., Lamont I. L. and C. W. Ronson.**
21 2006. Excision and transfer of the *Mesorhizobium loti* R7A symbiosis island
22 requires an integrase IntS, a novel recombination directionality factor RdfS, and
23 a putative relaxase RlxS. *Mol. Microbiol.* 62(3): 723-34.

- 1 **33. Rodríguez-Navarro, D.,** Buendía, A. Camacho, M., Lucas M. and C.
2 Santamaria. 2000. Characterization of *Rhizobium* spp. Bean isolates from south-
3 west Spain. *Soil Biol. Biochem.* 32:1601-1613.
- 4 **34. Schneider, D., and R. Lenski.** 2004. Dynamics of insertion sequence elements
5 during experimental evolution of bacteria. *Res. Microbiol.* 155:319-327.
- 6 **35. Siguier, P., Filée, J., and M. Chandler.** 2006. Insertion sequences in
7 prokaryotic genomes. *Curr. Opin. Microbiol.* 9:1-6.
- 8 **36. Silva, C., Vinuesa, P., Eguiarte, L. Martínez-Romero, E. and V. Souza.**
9 2003. *Rhizobium etli* and *Rhizobium gallicum* nodulate common bean
10 (*Phaseolus vulgaris*) in a traditionally managed milpa plot in Mexico:
11 Population genetics and Biogeographic implications. *Appl. Env. Microbiol.*
12 69:884-893.
- 13 **37. Silva, C., Vinuesa, P., Eguiarte, L., Souza, V. and Martínez-Romero, E.**
14 2005. Evolutionary genetics and biogeographic structure of *Rhizobium gallicum*
15 sensu lato, a widely distributed bacterial symbiont of diverse legumes. *Mol.*
16 *Ecol.* 14:4033-4050.
- 17 **38. Slatkin, M.** 1985. Genetic differentiation of transposable elements under
18 mutation and unbiased gene conversion. *Genetics* 110(1): 145-58.
- 19 **39. Sullivan, J. T. and C. W. Ronson.** 1998. Evolution of rhizobia by acquisition
20 of a 500-kb symbiosis island that integrates into a phe-tRNA gene. *Proc. Natl.*
21 *Acad. Sci. U S A* 95(9): 5145-9.
- 22 **40. Sullivan, J. T., Trzebiatowski, J. R., Cruickshank, R. W., Gouzy, J., Brown,**
23 **S. D., Elliot, R. M., Fleetwood, D. J., McCallum, N. G., Rossbach, U., Stuart,**
24 **G. S., Weaver, J. E., Webby, R. J., De Bruijn F. J. and C. W. Ronson.** 2002.

- 1 Comparative sequence analysis of the symbiosis island of *Mesorhizobium loti*
2 strain R7A. *J. Bacteriol.* 184(11): 3086-95.
- 3 41. **Touchon, M., and Rocha, E.** 2007. Causes of insertion sequences abundance in
4 prokaryotic genomes. *Mol. Biol. Evol.* 4:969-981.
- 5 42. **Young, J. P., Crossman, L. C., Johnston, A. W., Thomson, N. R., Ghazoui,**
6 **Z. F., Hull, K. H., Wexler, M., Curson, A. R., Todd, J. D., Poole, P. S.,**
7 **Mauchline, T. H., K. East, M. A., Quail, A., Churcher, C., Arrowsmith, C.,**
8 **Cherevach, I., Chillingworth, T., Clarke, K., Cronin, A., Davis, P., Fraser,**
9 **A., Hance, Z., Hauser, H., Jagels, K., Moule, S., Mungall, K., Norbertczak,**
10 **H., Rabinowitsch, E., Sanders, M., Simmonds, M., Whitehead S. and J.**
11 **Parkhill.** 2006. The genome of *Rhizobium leguminosarum* has recognizable
12 core and accessory components. *Genome Biol* 7(4): R34.
- 13 43. **Wagner, A.** 2006. Periodic extinctions of transposable elements in bacterial
14 lineales: Evidence from intragenomic variation in multiple genomes. *Mol. Biol.*
15 *Evol.* 23 (4):723-733.
- 16 44. **Zeigler, D. R.** 2003. Gene sequences useful for predicting relatedness of whole
17 genomes in bacteria. *Int. J. Syst. Evol. Microbiol.* 53(Pt 6): 1893-900.

18

19 Figure 1. Phylogenetic reconstruction representing the evolutionary relationships across
20 the three collections. The circles inside the branches represent bootstrap support >70%.
21 The external circles represent *R. etli* strains intermingled from one of the other
22 collections. The arrows represent other species intermingled (4872, GR42 and GR60 are
23 *R. gallicum*; GR93 is *R. giardinii*; and GR84 is *R. leguminosarum*). The underlined
24 strains represent recombinant strains for *dnaB* or *glyA* between the three populations
25 and the *R. gallicum* strain.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25

Figure 2. Presence/absence profile of the IS elements present in the chromosome, conjugative and symbiotic plasmids across the three populations. The bottom line represents CFN42 control strain. Red bars represent positive confirmation of the IS element (a DNA fragment of similar size to that in CFN42), blue bars represent a negative result for the presence of the IS element (a DNA fragment equivalent to the distance between the two sequences near the insertion site but lacking the IS), no colour represents no PCR product and yellow bars represent DNA fragments bigger in size than expected.

Figure 3. Dot plot graphic of pSym from *R. etli* CFN42 and CIAT652. Diagonal red lines represent DNA regions that were aligned between the two replicons. The blue and red points represent small DNA regions, in many cases repeated sequences, that were aligned in different regions of each replicon.

Table 1. Rhizobium isolates and reference strains used in this study

Table 2. DNA divergence in *R. etli* IS elements, plasmid genes and chromosomal genes

Table 3. Comparison of mean genetic divergence between populations (D_{xy}) and mean genetic divergence (π) within each population

Table 4. Genetic differentiation and gene flow tests in *R. etli* IS elements, plasmid genes and chromosomal genes

1 Supp. Table 1. Complete set of sequences obtained from PCR reactions
2
3 Supp. Table 2. DNA divergence and neutrality test in *R. etli* IS elements, plasmid genes
4 and chromosomal genes within populations
5
6 Supp. Table 3. Genetic differentiation analyses of *R. etli* IS elements, plasmid genes and
7 chromosomal genes
8
9 Supp. Table 4. Recombination events in plasmid and chromosomal genes between
10 populations. Each row represents a different recombination event for each gene. The
11 underlined strains represent recombinant strains for *dnaB* or *glyA* between the three
12 populations and the *R. gallicum* strains that could explain the intermingled strains in the
13 phylogeny reconstruction of the three populations (Fig. 1).
14
15 Supp. Figure 1. Phylogenetic network analyses with Neighbor Net for a) *dnaB* and b)
16 *glyA*. Underlined labels represent strains participating in recombination events
17 determined with RDP3 (Supp. Table 4): bold labels represent recombinant strains and
18 italic labels represents parental strains.
19

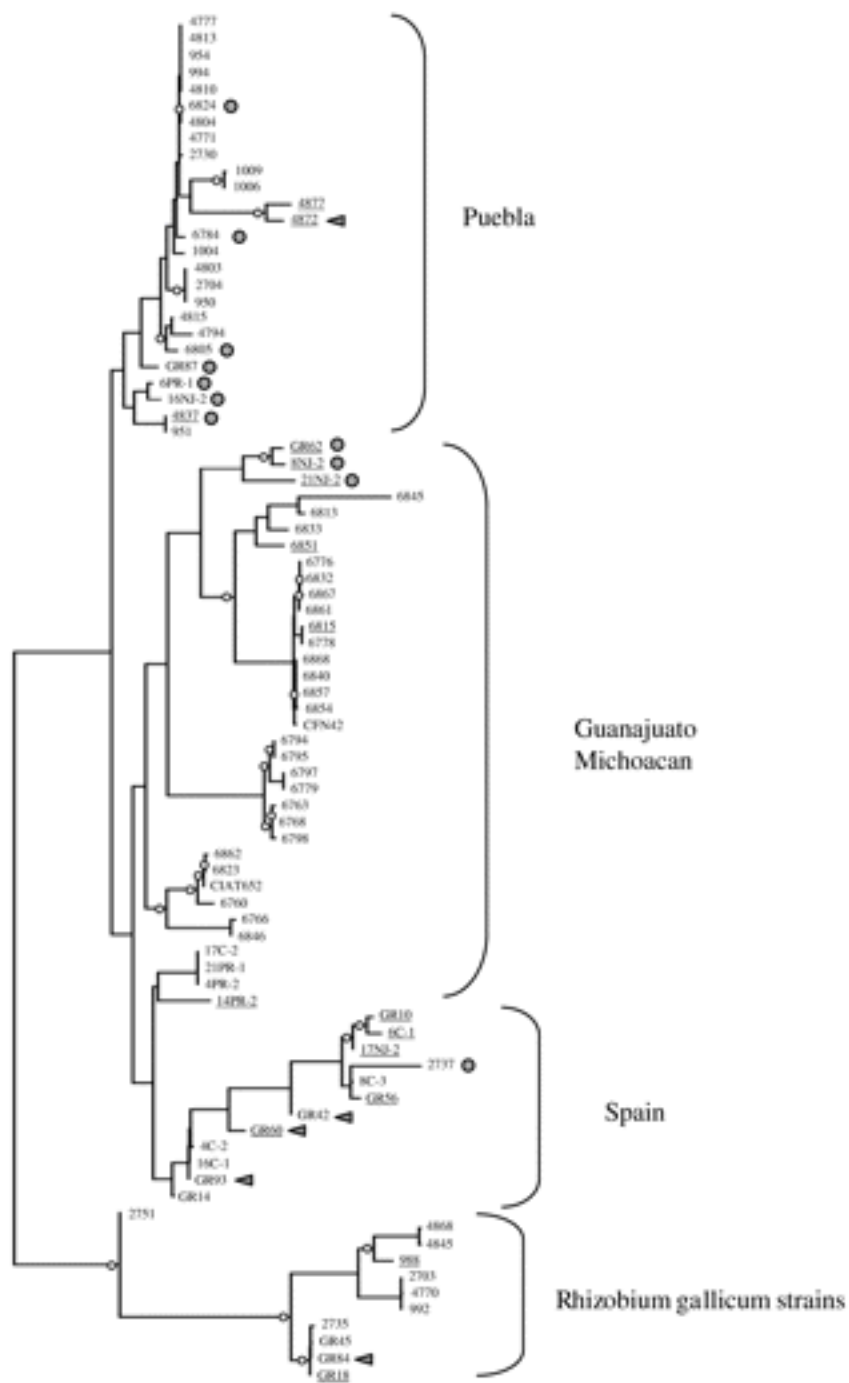


Figure 1. Phylogenetic reconstruction representing the evolutionary relationships across the three collections. The circles inside the branches represent bootstrap support >70%. The external circles represent *R. etli* strains intermingled from one of the other collections. The arrows represent other species intermingled (4872, GR42 and GR60 are *R. gallicum*; GR93 is *R. giardinii*; and GR84 is *R. leguminosarum*). The underlined strains represent recombinant strains for *dnaB* or *glyA* between the three populations and the *R. gallicum* strain.

CHROMOSOME

CONJUGATIVE PLASMID

SYMBIOTIC PLASMID

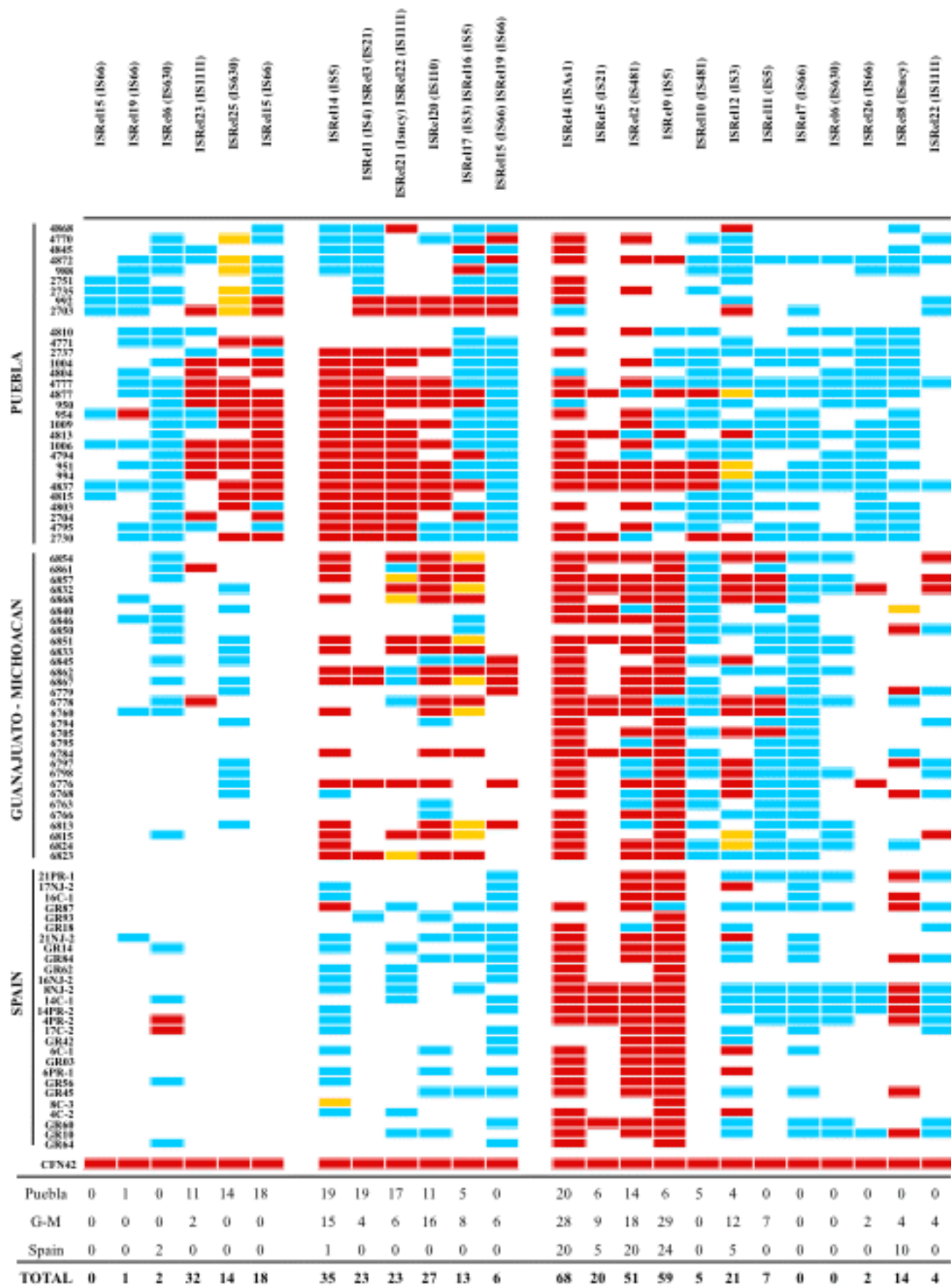


Figure 2. Presence/absence profile of the IS elements present in the chromosome, conjugative and symbiotic plasmids across the three populations. The bottom line represents CFN42 control strain. Red bars represent positive confirmation of the IS element (a DNA fragment of similar size to that in CFN42), blue bars represent a negative result for the presence of the IS element (a DNA fragment equivalent to the distance between the two sequences near the insertion site but lacking the IS), no colour represents no PCR product and yellow bars represent DNA fragments bigger in size than expected.

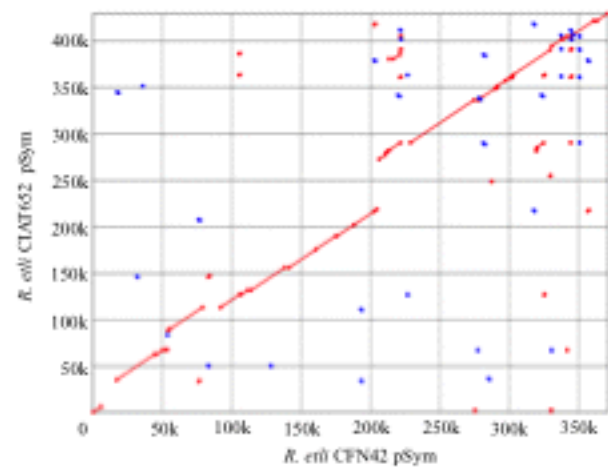


Figure 3. Dot plot graphic of pSym from *R. etli* CFN42 and CIAT652. Diagonal red lines represent DNA regions that were aligned between the two replicons. The blue and red points represent small DNA regions, in many cases repeated sequences, that were aligned in different regions of each replicon.

Table 1. Rhizobium isolates and reference strains used in this study

| Origin | No. Strains | Species | Identifier | Reference |
|-----------------------------|-------------|-------------------------|---|-----------|
| Rhizobium Isolates | | | | |
| Puebla, Mexico | 30 | Rhizobium etli | 1009, 1006, 4771, 1004, 2737, 4815, 4803, 4777, 950, 2704, 4813, 4794, 951, 994, 4837, 4804, 954, 4877, 4795, 2730, 4810 | (34) |
| | | Rhizobium gallicum | 4868, 4872, 4770, 4845, 988, 2751, 2735, 992, 2703 | |
| Guanajuato, Mexico | 30 | Rhizobium etli | Domesticated: 6854, 6861, 6857, 6832, 6868, 6867, 6840, 6846, 6850, 6851, 6833, 6845, 6862 Wild: 6779, 6778, 6760, 6794, 6805, 6795, 6766, 6784, 6797, 6798, 6776, 6768, 6763 Weedy: 6815, 6824, 6823, 6813 | (12) |
| Valle Guadalquivir, Spain | 27 | Rhizobium etli | 21PR-1, 16C-1, 21NJ-2, 8C-3, 8NJ-2, 14PR-2, 17NJ-2, 16NJ-2, 14C-1, 4PR-2, 17C-2, 6PR-1, 6C-1, 4C-2, GR10, GR62, GR14, GR87, GR56 | (31) |
| | | Rhizobium gallicum | GR18, GR45, GR60, GR42 | |
| | | Rhizobium giardinii | GR93, GR03 | |
| | | Rhizobium fredii | GR64 | |
| | | Rhizobium leguminosarum | GR84 | |
| Rhizobium Reference Strains | | | | |
| Mexico | 1 | Rhizobium etli | CFN42 | (13) |
| Costa Rica | 1 | Rhizobium etli | CIAT652 | (9) |

| Gene | Polymorphic sites ^a | dn/ds | Θ | Nucleotide diversity π | Nucleotide diversity at synonymous sites π_S | Nucleotide diversity at nonsynonymous sites π_{NS} | Tajima's D |
|--------|--------------------------------|-------|----------|-------------------------------|---|---|------------|
| dnaB | 263 (23.6) | 0.019 | 0.047 | 0.066 | 0.237 | 0.009 | 1.39 |
| glyA | 278 (24.7) | 0.029 | 0.070 | 0.088 | 0.122 | 0.075 | 0.84 |
| nodC | 19 (3.3) | 0.355 | 0.007 | 0.014 | 0.028 | 0.009 | 0.72 |
| ISReI2 | 6 (0.6) | 0.267 | 0.001 | 0.002 | 0.001 | 0.002 | 2.72 |
| ISReI4 | 13 (1.3) | 0.325 | 0.002 | 0.002 | 0.003 | 0.007 | 1.50 |

^a The percentages of polymorphic sites per sequence length are shown in parentheses.

Table 2. DNA divergence in *R. etli* IS elements, plasmid genes and chromosomal genes

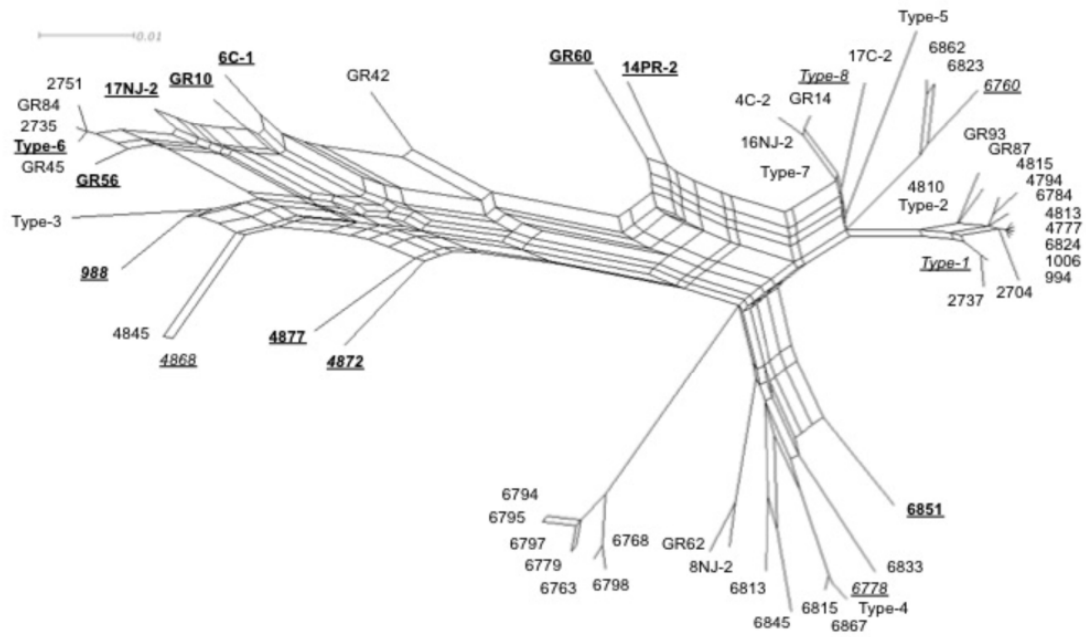
| | Dxy | π |
|---------------|------------|-------------------------|
| dnaB | | |
| Pueb - Guanj | 0.072 | 0.056 - 0.052 |
| Pueb - Spain | 0.068 | 0.056 - 0.064 |
| Guanj - Spain | 0.073 | 0.052 - 0.064 |
| glyA | | |
| Pueb - Guanj | 0.097 | 0.080 - 0.066 |
| Pueb - Spain | 0.100 | 0.080 - 0.077 |
| Guanj - Spain | 0.089 | 0.066 - 0.077 |
| ISRel2 | | |
| Pueb - Guanj | 0.003 | 0.002 - 0.002 |
| Pueb - Spain | 0.002 | 0.002 - 0.001 |
| Guanj - Spain | 0.002 | 0.002 - 0.001 |
| ISRel4 | | |
| Pueb - Guanj | 0.002 | 0.001 - 0.003 |
| Pueb - Spain | 0.002 | 0.001 - 0.002 |
| Guanj - Spain | 0.002 | 0.003 - 0.002 |
| nodC | | |
| Pueb - Guanj | 0.003 | 0.001 - 0.005 |
| Pueb - Spain | 0.026 | 0.001 - 0.008 |
| Guanj - Spain | 0.025 | 0.005 - 0.008 |

Table 3. Comparison of mean genetic divergence between populations (Dxy) and mean genetic divergence (π) within each population

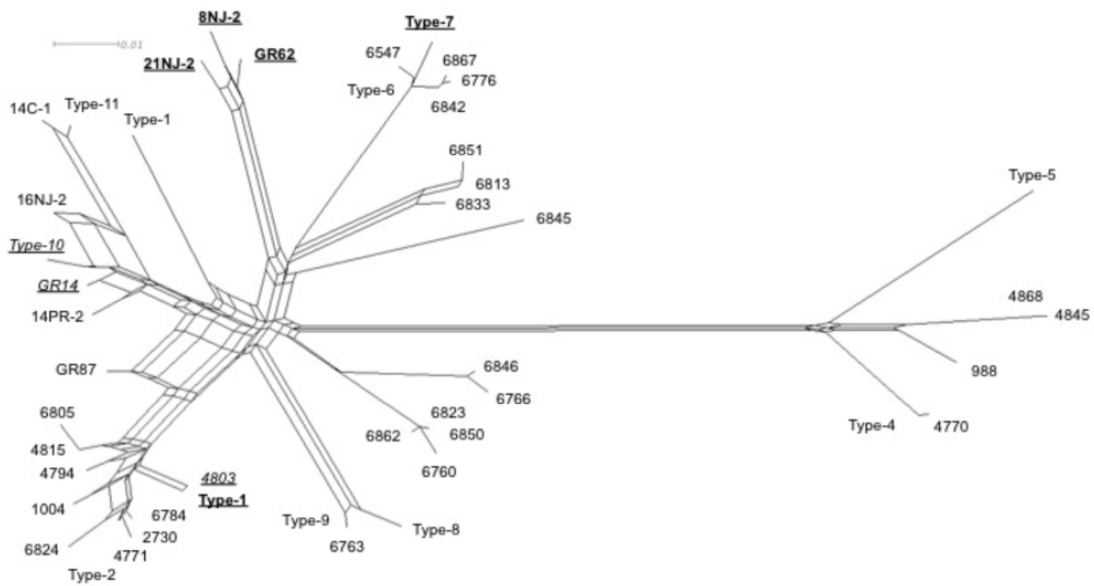
| Gene | K_{ST} | F_{ST} | Nm |
|---------------|-----------------------|-----------------------|-----------|
| dnaB | | | |
| Pueb - Guanj | 0.146*** | 0.252 | 1.48 |
| Pueb - Spain | 0.064*** | 0.118 | 3.72 |
| Guanj - Spain | 0.118*** | 0.208 | 1.90 |
| glyA | | | |
| Pueb - Guanj | 0.144*** | 0.249 | 1.51 |
| Pueb - Spain | 0.142*** | 0.246 | 1.54 |
| Guanj - Spain | 0.129*** | 0.226 | 1.71 |
| ISRel2 | | | |
| Pueb - Guanj | 0.007 | 0.025 | 19.75 |
| Pueb - Spain | 0.652** | 0.845 | 0.09 |
| Guanj - Spain | 0.607** | 0.750 | 0.17 |
| ISRel4 | | | |
| Pueb - Guanj | 0.056 | 0.217 | 1.80 |
| Pueb - Spain | 0.310** | 0.634 | 0.29 |
| Guanj - Spain | 0.156** | 0.270 | 1.35 |
| nodC | | | |
| Pueb - Guanj | 0.184* | 0.374 | 0.84 |
| Pueb - Spain | 0.118* | 0.258 | 1.44 |
| Guanj - Spain | 0.225** | 0.357 | 0.90 |

Table 4. Genetic differentiation and gene flow tests in *R. etli* IS elements, plasmid genes and chromosomal genes. Asterisks represent the probability obtained by the permutation test with 1000 replicates (ns, not significant; *, 0.01<P<0.05; **, 0.001<P<0.01; ***, P<0.001).

A)



B)



Supp. Figure 1. Phylogenetic network analyses with Neighbor Net for a) *dnaB* and b) *glyA*. Underlined labels represent strains participating in recombination events determined with RDP3 (Supp. Table 4): bold labels represent recombinant strains and italic labels represents parental strains.

| Gene | No. of sequences | Sequence length (bp) | Puebla | Guanajuato | Spain |
|--------|------------------|----------------------|--------|------------|-------|
| dnaB | 87 | 1113 | 30 | 30 | 25 |
| glyA | 88 | 1125 | 29 | 31 | 28 |
| ISRel2 | 45 | 981 | 5 | 16 | 20 |
| ISRel4 | 53 | 1035 | 6 | 23 | 21 |
| nodC | 44 | 576 | 7 | 22 | 15 |

Supp. Table 1. Complete set of sequences obtained from PCR reactions

| Gene | θ | Nucleotide diversity π | Nucleotide diversity at synonymous sites π_S | Nucleotide diversity at nonsynonymous sites π_{NS} | Tajima's D |
|---------------|---------------|----------------------------|--|--|------------|
| dnaB | 0.047 | 0.066 | 0.237 | 0.009 | 1.39 |
| Puebla | 0.040 | 0.056 | 0.191 | 0.011 | 1.48 |
| Guanajuato | 0.040 | 0.052 | 0.194 | 0.005 | 1.08 |
| Spain | 0.042 | 0.064 | 0.227 | 0.010 | 2.15* |
| glyA | 0.070 | 0.088 | 0.122 | 0.075 | 0.84 |
| Puebla | 0.065 | 0.080 | 0.118 | 0.066 | 0.88 |
| Guanajuato | 0.057 | 0.066 | 0.091 | 0.057 | 0.57 |
| Spain | 0.070 | 0.077 | 0.098 | 0.061 | 0.07 |
| nodC | 0.007 | 0.014 | 0.028 | 0.009 | 2.72** |
| Puebla | 0.001 | 0.001 | 0.002 | 0 | -1.01 |
| Guanajuato | 0.007 | 0.005 | 0.012 | 0.002 | -1.16 |
| Spain | 0.010 | 0.008 | 0.015 | 0.005 | -0.87 |
| ISRel2 | 0.001 | 0.002 | 0.001 | 0.002 | 0.72 |
| Puebla | 0.001 | 0.002 | 0.002 | 0.001 | 0.70 |
| Guanajuato | 0.001 | 0.002 | 0.001 | 0.002 | 0.29 |
| Spain | 0.001 | 0.001 | 0 | 0.001 | 0.52 |
| ISRel4 | 0.002* | 0.002 | 0.003 | 0.002 | -0.62 |
| Puebla | 0.001 | 0.001 | 0 | 0.001 | -1.13 |
| Guanajuato | 0.003 | 0.003 | 0.005 | 0.002 | -0.14 |
| Spain | 0.002 | 0.002 | 0.003 | 0.001 | -0.65 |

Supp. Table 2. DNA divergence and neutrality test in *R. etli* IS elements, plasmid genes and chromosomal genes within populations.

| Gene | K_S | K_S* | Z | Z* | S_{nn} | PM test P-value |
|---------------|----------------------|-----------------------|----------|-----------|-----------------------|----------------------------|
| dnaB | | | | | | |
| Pueb - Guanj | 59.84 | 3.53 | 786.71 | 6.22 | 0.95 | <0.001 |
| Pueb - Spain | 66.33 | 3.65 | 699.60 | 6.14 | 0.88 | <0.01 |
| Guanj - Spain | 63.80 | 3.79 | 672.31 | 6.10 | 0.93 | <0.001 |
| glyA | | | | | | |
| Pueb - Guanj | 56.90 | 3.52 | 752.71 | 6.18 | 0.93 | <0.001 |
| Pueb - Spain | 59.18 | 3.37 | 682.83 | 6.08 | 0.92 | <0.001 |
| Guanj - Spain | 53.62 | 3.52 | 729.46 | 6.15 | 0.86 | <0.001 |
| ISRel2 | | | | | | |
| Pueb - Guanj | 2.19 | 0.58 | 205.22 | 5.20 | 0.62 | <0.05 |
| Pueb - Spain | 3.05 | 0.58 | 79.22 | 4.22 | 0.82 | <0.05 |
| Guanj - Spain | 3.42 | 0.71 | 207.62 | 5.14 | 0.89 | <0.01 |
| ISRel4 | | | | | | |
| Pueb - Guanj | 2.73 | 1.14 | 199.90 | 5.04 | 0.76 | <0.05 |
| Pueb - Spain | 2.13 | 0.90 | 120.35 | 4.52 | 0.81 | <0.01 |
| Guanj - Spain | 2.82 | 1.12 | 404.66 | 5.64 | 0.78 | <0.01 |
| nodC | | | | | | |
| Pueb - Guanj | 1.63 | 0.78 | 90.11 | 4.27 | 0.71 | ns |
| Pueb - Spain | 1.13 | 0.52 | 138.43 | 4.76 | 0.88 | <0.001 |
| Guanj - Spain | 1.29 | 0.60 | 273.51 | 5.40 | 0.69 | <0.001 |

Supp. Table 3. Genetic differentiation analyses of *R. etli* IS elements, plasmid genes and chromosomal genes

| Gene | Recombinant strains | Parental strains | Methods | P-value |
|-------------|--|-------------------------|----------------|----------------|
| dnaB | 4877, <u>4872</u> , 988 | 4868 - 4795 | 3 | 0.564 |
| | 6851 | 6760 - 6778 | 4 | 0.509 |
| | 14PR-2, <u>GR60</u> | <u>6PR-1</u> - 4872 | 3 | 0.634 |
| | GR18, 17NJ-2, GR10, 6C-1, GR56 | 988 – <u>6PR-1</u> | 4 | 0.621 |
| glyA | <u>4837</u> | 4803 – <u>GR42</u> | 2 | 0.479 |
| | 6815, <u>8NJ-2</u> , <u>21NJ-2</u> , <u>GR62</u> | GR14 - 4803 | 2 | 0.510 |
| ISRel4 | 954 | GR87 - 4795 | 4 | 0.549 |

Supp. Table 4.

Recombination events in plasmid and chromosomal genes between populations. Each row represents a different recombination event for each gene. The underlined strains represent recombinant strains for dnaB or glyA between the three populations and the *R. gallicum* strains that could explain the intermingled strains in the phylogeny reconstruction of the three populations (Fig. 1).