



**HAL**  
open science

## Structural peculiarities of linear megaplasmid, pLMA1, from interfere with pyrosequencing reads assembly

Martin Wagenknecht, Julián R. Dib, Andrea Thürmer, Rolf Daniel, María E. Farías, Friedhelm Meinhardt

### ► To cite this version:

Martin Wagenknecht, Julián R. Dib, Andrea Thürmer, Rolf Daniel, María E. Farías, et al.. Structural peculiarities of linear megaplasmid, pLMA1, from interfere with pyrosequencing reads assembly. *Biotechnology Letters*, 2010, 32 (12), pp.1853-1862. 10.1007/s10529-010-0357-y . hal-00610056

**HAL Id: hal-00610056**

**<https://hal.science/hal-00610056>**

Submitted on 21 Jul 2011

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

1 **Section Microbial and Enzyme Technology**

2  
3 **Structural peculiarities of the linear megaplasmid, pLMA1, from *Micrococcus luteus***  
4 **interfere with pyrosequencing reads assembly**

5  
6 Martin Wagenknecht, Julián R. Dib, Andrea Thürmer, Rolf Daniel, María E. Farías, Friedhelm Meinhardt\*

7  
8  
9 M. Wagenknecht, F. Meinhardt:

10 Institut für Molekulare Mikrobiologie und Biotechnologie, Westfälische Wilhelms -Universität Münster,  
11 Corrensstr. 3, D-48149 Münster, Germany

12  
13 J.R. Dib, M.E. Farías:

14 Laboratorio de Investigaciones Microbiológicas de Lagunas Andinas (LIMLA), Planta Piloto de Procesos  
15 Industriales Microbiológicos-CONICET, Av. Belgrano y Pje. Caseros, (4000) Tucumán, Argentina

16  
17 A. Thürmer, R. Daniel:

18 Göttingen Genomics Laboratory, Institut für Mikrobiologie und Genetik, Georg-August-Universität Göttingen,  
19 Grisebachstr. 8, D-37077 Göttingen, Germany

20  
21  
22 \*Corresponding author:

23 Phone: +49 251 8339825. Fax: +49 251 8338388. E-mail address: meinhar@uni-muenster.de

24  
25  
26  
27  
28 **Keywords:** 454 sequencing; Inverted repeats; Iteron; Linear plasmid; *Micrococcus* Repetitive sequences;  
29 Transposase

31

32 **Abstract**

33

34 Different strains of *Micrococcus luteus*, isolated from high-altitude Argentinean wetlands, were recently reported  
35 to harbour the linear plasmids pLMA1, pLMH5 and pLMV7, all of which with 5'-covalently attached terminal  
36 proteins. The link between pLMA1 and the host's erythromycin resistance as well as further presumptive  
37 qualities prompted us to perform a detailed characterization. When the 454 technology was applied for direct  
38 sequencing of gel-purified pLMA1, assembly of the reads was impossible. However, combined Sanger/454  
39 sequencing of cloned pLMA1 fragments, covering altogether 23 kb of the 110-kb spanning plasmid, allowed  
40 numerous sequence repeats of varying in lengths to be identified thus rendering an explanation for the above 454  
41 assembly failure. A large number of putative transposase genes were identified as well. Furthermore, a region  
42 with five putative iteron sequences is possibly involved in pLMA1 replication.

43

44

45 **Introduction**

46

47 Microbial linear plasmids are rather widespread, particularly among various Gram-positive bacteria. They have  
48 been found in many *Streptomyces* spp., several rhodococci and mycobacteria, but also in *Arthrobacter*  
49 *nitroguajacolicus* Rüd61a, *Planobispora rosea*, and the plant pathogen *Clavibacter michiganensis*. Such linear  
50 replicons belong to a class of genetic elements, which are characterized by terminal inverted repeats (TIRs) and  
51 terminal proteins (TPs) attached to each 5'-end. Most of the linear megaplasmids are conjugative and display  
52 rather low-copy numbers (see Wagenknecht and Meinhardt, 2010, and references therein). Linear plasmids often  
53 encode nonessential functions but – under certain conditions – they may provide advantages attributes, such as  
54 heavy metal resistance or specific catabolic traits. Antibiotic resistance is rather seldomly found to be linear  
55 plasmid-encoded, as for methylenomycin of *Streptomyces* linear elements. [For a monograph on microbial linear  
56 plasmids see Meinhardt and Klassen (2007).]

57 We recently reported the isolation and characterization of the first linear plasmids in different strains of  
58 *Micrococcus luteus* (Dib et al. 2010a). All of them (pLMA1, pLMH5, and pLMV7) possess 5'-attached TPs.  
59 Host strains were isolated from high-altitude Argentinean wetlands, which are considered extreme and pristine  
60 environments characterized by high UV radiation, arsenic concentration, and salinity. Such bacteria displayed

61 high UV tolerance, heavy metal resistance, and, unexpectedly, resistance against a number of antibiotics (Dib et  
62 al. 2008). As there is already a proven link between pLMA1 occurrence and erythromycin resistance (Dib et al.  
63 2010a) we decided to focus on such plasmid.

64 One of the most powerful new sequencing technologies is the 454 pyrosequencing as such parallel  
65 noncloning pyrosequencing-based system is capable of delivering sequence information much faster than current  
66 Sanger sequencing platforms. 454 sequencers provide shorter reads (ideally ~350 bp on average versus ~800 bp  
67 for Sanger reads) but at greatly reduced per-base costs. Moreover, the 454 technology was shown to be capable  
68 of resolving hard stops for which Sanger sequencing is ineffective (Goldberg et al. 2006).

69 When we sequenced pLMA1 with the 454 technology, attempts to assemble the reads encountered  
70 insurmountable obstacles. However, applying both, the Sanger sequencing and the 454 technology for cloned  
71 pLMA1 fragments, long sequence stretches of the linear plasmid were obtained. Peculiar attributes, such as a  
72 large number of repetitive sequences and transposase encoding genes became obvious.

73

74

## 75 **Materials and methods**

76

77 Bacterial strains, plasmids, and growth conditions

78

79 *Micrococcus luteus* strains were cultivated as described in Dib et al. (2010a). Erythromycin  $100 \mu\text{g ml}^{-1}$  was  
80 added to the medium to counteract loss of pLMA1. *Escherichia coli* NEB 5-alpha (New England Biolabs) and  
81 pUC18 were used for cloning *Pst*I restriction fragments.

82

83 DNA isolation

84

85 Extraction of *Pst*I-digested DNA fragments from agarose gels was achieved using the QIAquick Gel Extraction  
86 Kit (Qiagen). Isolation of bulk and pLMA1 DNA was done as described previously (Dib et al. 2010a, 2010b).

87

88 Labeling of DNA probes and Southern hybridization

89

90 Probe fragments were PCR-amplified using Phusion Hot Start High-Fidelity DNA Polymerase (Finnzymes),  
91 primer pairs op15-revw5/op15-uniw3-1 (2315 bp), op34-revw3-1/op34-uniw4 (1597 bp), op54-revw2/op54-  
92 uniw2 (1607 bp), op62-revw2/op62-uniw2 (999 bp), and op74-revw2/op74-uniw1 (1456 bp), and pP86-15,  
93 pP62-34, pP45-54, pP37-62, and pP37-74, respectively, as template (for primer sequences see Supplementary  
94 Table S1). Probe labeling, capillary blotting, and hybridization/detection were done as described in  
95 Wagenknecht and Meinhardt (2010).

96

97 454 sequencing of pLMA1, assembly, and mapping

98

99 Sequencing of gel-purified plasmid pLMA1 was done using the Genome Sequencer FLX system (Roche Applied  
100 Science). A single-stranded DNA shotgun library (ssDNA library) was generated from approximately 5 µg of  
101 isolated pLMA1 DNA. The DNA was fragmented by nebulization for 30 s and 1 bar. Further steps were done  
102 according to the Roche protocol. The size selection of the ssDNA library resulted in an average length of 469 bp.  
103 A total of 84302 reads were achieved. The GS *De Novo* Assembler (Roche Applied Science) software package  
104 was used for sequence assembly. The GS Reference Mapper was used to align the reads from the 454  
105 sequencing run to the cloned *Pst*I restriction fragments.

106

107 Sequencing of cloned *Pst*I restriction fragments by primer walking

108

109 Sequencing of the pUC18-cloned *Pst*I restriction fragments was done using BigDye Terminator 3.1 chemistry  
110 and an ABI Prism 3700 DNA Analyzer (Applied Biosystems). For sequencing primers see Table S1 in the  
111 Supplementary material. Hard-stop events that occurred during primer walking were bypassed by an additional  
112 10 min denaturation step combined with 1 M betaine prior to the sequencing reaction.

113

114 Prediction and annotation of ORFs and identification of repetitive sequences

115

116 An initial set of predicted protein-coding regions was identified using Artemis V11.22 (Sanger Institute). ORFs  
117 shorter than 50 amino acids and with overlaps of higher scoring regions were eliminated. Annotation was done  
118 considering a Blast (Altschul et al. 1990) search against Swissprot (Pearson 1994) and GenBank nr (Benson et

119 al. 2004). Putative protein functions were assigned for hits with a full length alignment and appropriate similarity  
120 and confirmed by protein domain hits.

121 Repetitions within the nucleotide sequence of the cloned *Pst*I fragments of pLMA1 and the reference  
122 sequences were identified and analysed using Spectral Repeat Finder (Sharma et al. 2004) and Clone Manager  
123 (Sci-Ed Software).

124

125

## 126 **Results and discussion**

127

128 454 pyrosequencing of pLMA1 and read assembly

129

130 Approx. 5 µg pLMA1 DNA, isolated by electroelution from preparative pulsed-field (PF) gels, were used to  
131 generate the single-stranded DNA shotgun library that was subsequently subjected to the 454 sequencing  
132 procedure. The sequencing runs yielded 84302 single reads generating total sequence information of 8703704  
133 bases. Considering 110 kb as the electrophoretically determined size of pLMA1 (Dib et al. 2010a), the obtained  
134 sequence data theoretically correspond to a 79-fold depth coverage. However, by *de novo* assembling of the  
135 single reads, it turned out that only 29% of all the sequences could be assembled. Moreover, from the 955  
136 contigs obtained, the majority (~80%) displayed remarkable short sizes, ranging from 100–200 bp. About 16%  
137 of all contigs ranged from 201–400 bp; the maximum length was 684 bp only. All bioinformatic attempts to  
138 increase the contig lengths which, after repeated assembly, would allow for creating longer stretches of  
139 continuous pLMA1 sequences, failed. Thus, subsequent sequence analyses, such as ORF prediction and  
140 annotation, were virtually impossible.

141 Chromosomal DNA contamination of the plasmid DNA sample could have been the reason for above  
142 findings. However, test digests of the pLMA1 DNA sample with different restriction endonucleases revealed a  
143 distinct band pattern with a negligible smear in the gel (Supplementary Fig. S1), indicative of only faint amounts  
144 of contaminating chromosomal DNA.

145 Repetitive sequences may severely hamper *de novo* assembly of sequence data, as short read lengths,  
146 characteristic for 454 pyrosequencing, may make it impossible to span such repetitive elements. This situation is  
147 known from genome sequencing projects of bacteria (Goldberg et al. 2006) and in particular those of eukaryotes  
148 with highly repetitive genomes, such as barley (Wicker et al. 2006). The more and the longer sequence repeats

149 are present in a given DNA molecule, such issue gains in importance. Though the employed Genome Sequencer  
150 FLX System (Standard Series, Roche) theoretically allows read lengths of 200–300 bases, the average read  
151 length obtained for pLMA1 was 103 bases, making the assembly of nucleotide sequences rich in repetitive  
152 elements even more difficult. It is a fairly common experience that a high GC content may lead to short read  
153 lengths, which agrees in the case of pLMA1 with the calculated average GC content of 74.2% of all reads  
154 obtained.

155

156 Cloning and Sanger sequencing of selected restriction fragments of pLMA1

157

158 Since the above 454 pyrosequencing did not yield long continuous stretches of pLMA1 nucleotide sequences  
159 and for checking that repeats caused the failure of the assembly, we cloned restriction fragments of pLMA1 and  
160 subjected them to Sanger sequencing.

161 Five of the fragments produced by *Pst*I cleavage (8.6, 6.2, 5.8, 4.5, and 3.7 kb; Supplementary Fig. S1), were  
162 excised from a preparative gel and cloned in *Pst*I-linearized pUC18. Plasmid DNA isolated from *E. coli*  
163 transformants was cut with *Pst*I, which released the insert. Each of the transformation reactions – with one  
164 exception – yielded transformants with plasmids exhibiting the expected insert sizes. Plasmids designated pP86-  
165 15, pP62-34, pP45-54, and pP37-62, harboring the 8.6, 6.2, 4.5, and the 3.7 kb *Pst*I restriction fragment,  
166 respectively, were used for insert sequencing by primer walking starting with the standard primers ‘uni’ and  
167 ‘rev’. Regrettably, though repeatedly attempted, no hybrid plasmid containing the 5.8 kb insert was obtained.

168 As linear plasmids of actinomycetes possess TPs covalently attached to each 5′-end, proven also for pLMA1  
169 (Dib et al., 2010a), and treatment with proteinase K (done prior to PFGE to allow the DNA to enter the gel) does  
170 not remove the linking amino acid residue of the TP (Yang et al. 2002), the 5.8 kb *Pst*I restriction fragment  
171 possibly carries one of the pLMA1 termini.

172 Plasmids isolated from PF gels routinely contain trace amounts of chromosomal DNA. Thus, before  
173 continuing insert sequencing by primer walking, sequence data obtained from the initial sequencing reactions  
174 using ‘uni’ and ‘rev’ were subjected to a BlastN analysis to verify the origin of the cloned *Pst*I fragments. The  
175 chromosomal sequence of *M. luteus* NCTC 2665 (Accession no. CP001628) served as the reference. No, or only  
176 partial, similarity was seen for the inserts of pP86-15, pP62-34, and pP45-54; however, the insert sequence of  
177 plasmid pP37-62 revealed a chromosomal origin of the cloned 3.7-kb *Pst*I fragment. A PCR analysis (not  
178 shown), using insert-specific primers and total DNA of the wild-type strain *M. luteus* A1 and of the pLMA1-

179 deficient strain *M. luteus* A1-M1 as templates, confirmed above findings. In an additional Southern analysis  
180 (Fig. 1), we PCR-amplified probes deduced from the respective insert sequences (Fig. 2) and hybridized them  
181 with *Pst*I-digested total DNA of *M. luteus* strains A1 and A1-M1, and with likewise-cut pLMA1 DNA.  
182 Observed hybridization signals affirmed inserts of pP86-15, pP62-34, and pP45-54 being fragments of pLMA1  
183 (Fig. 1a, b, and c), whereas the insert of pP37-62 indeed is a chromosomal fragment (Fig. 1e). We therefore  
184 checked remaining transformants obtained from the corresponding transformation experiment by isolating  
185 plasmids and sequencing the terminal regions of their inserts with the backbone primers. Analysing insert  
186 sequences in a BlastN search, a couple of plasmids exhibited identical inserts that did not show similarity with  
187 the reference chromosome. One of them, pP37-74, was selected for entire insert sequencing and subjected to  
188 further PCR and Southern analysis (Fig. 1d), which confirmed its origin from pLMA1.

Approx.  
position of  
Fig. 1

189 During the sequencing reactions of pP86-15, pP62-34, and pP37-74 hard stops arose that were resolved as  
190 described in Material and methods. Regions with a potential to form such secondary structures are known to  
191 cause difficulties in Sanger sequencing, predominantly in DNA molecules with a high GC content (Goldberg et  
192 al. 2006).

193 Unexpectedly, hybridization of the probe deduced from the insert of pP45-54 revealed an additional signal of  
194 chromosomal origin with a size of approximately 2.4 kb (Fig. 1c). Comparison of nucleotide sequences of the  
195 probe and the *M. luteus* NCTC 2665 reference chromosome disclosed a homology stretch of 625 bp with 99.7%  
196 identity, present also on the chromosome of the pLMA1 host strain *M. luteus* A1.

197 Action of the restriction endonuclease *Pst*I is influenced by the nucleotide sequences neighbouring its  
198 recognition site in a way that adjacent GC runs significantly impede cleavage (Armstrong and Bauer 1982),  
199 rendering an explanation for the large but less intense signals observed in lanes containing pLMA1 DNA in  
200 addition to the expected ones (Fig. 1a–d, lane 6). Since the lanes in which total DNA of *M. luteus* A1 was  
201 separated contained lesser amounts of pLMA1, such signals are hardly observed (Fig. 1a–d, lane 4).

202

203 454 reads map on Sanger-sequenced pLMA1 fragments

204

205 After finishing primer walking on the cloned *Pst*I fragments, we were now able to align numerous reads of the  
206 454 sequencing to the Sanger-sequenced pLMA1 fragments. Since each of the pLMA1 fragments was  
207 abundantly covered by 454 reads, that also overlapped, the depth coverage obtained reached average values of  
208 56–87-fold. The mapping results are summarized in Table 1. Thus, we could affirm the accuracy of the

Approx.  
position of  
Table 1

209 nucleotide sequences initially obtained by primer walking. Being consistent with a chromosomal origin, only a  
210 total of 4 reads mapped on the 3721-bp fragment of pP37-62 (Table 1). Also, we aligned all 454 reads (84302) to  
211 the *M. luteus* NCTC 2665 reference chromosome; only 7432 (8.8%) of them were found to map. Thus,  
212 chromosomal DNA contamination can most likely be ruled out as the reason for the failure of the 454 reads  
213 assembly.

214 Analysis of the base composition revealed a GC content of 64.7%, 71.3%, 67.6%, and 67.4% for the 8621-  
215 bp, 6195-bp, 4498-bp, and 3681-bp insert, respectively. The GC content of the 3721-bp chromosomal *M. luteus*  
216 A1 fragment was calculated to be 74%, which is identical to the value reported for *M. luteus* (Ohama et al. 1989)  
217 and matches nicely with the GC content of 73% of the chromosome of reference strain *M. luteus* NCTC 2665.

218

219 pLMA1 is rich in ORFs involved in transposition

220

221 The obtained sequence data of the cloned pLMA1 fragments were subjected to a gene prediction analysis. On the  
222 8621-bp, 6195-bp, 4498-bp, and 3681-bp insert a total of 12, 7, 5, and 5 ORFs were identified, respectively (Fig.  
223 2), covering 76.3% of the sequence of all four pLMA1 fragments. Possible functions of the putative proteins,  
224 based on amino acid sequence similarities, are summarized in Supplementary Table S2. For 12 of the predicted  
225 ORFs, no function could be attributed.

Approx.  
position of  
Fig. 2

226 One remarkable and unexpected feature of the four pLMA1 fragments is the high proportion of genes  
227 involved in transposition; particularly, the 3681-bp pLMA1 fragment is rich in such genes (Fig. 2). The function  
228 of the closest relative suggests ORF 3 of pP45-54 to be a DNA replication protein; other closely related proteins  
229 are annotated as transposase helpers. A copy of this ORF (99% nucleotide sequence identity) was found on the  
230 chromosome of *M. luteus* NCTC 2665. Such ORF is probably also present on the chromosome of *M. luteus* A1,  
231 as suggested by above Southern analysis (Fig. 1c). ORF 1 of pP86-15, ORFs 2 and 7 of pP62-34, ORF 4 of  
232 pP45-54, and ORFs 1, 3, 4, and 5 of pP37-74 (gray shaded in Fig. 2) share significant similarities with known  
233 transposases (Supplementary Table S2). Interestingly, nucleotide sequence analysis revealed similarities to  
234 putative transposases MC9, MC13, MC19, MC34, and MC35 of pSD10, a 50-kb circular cryptic plasmid from a  
235 marine *Micrococcus* strain (Zhong et al. 2002). A region starting 43 bp upstream of ORF 7 (pP62-34) and  
236 extending through the ORF is similar to a 1007 bp block consisting of the first 894 nucleotides of MC13 and 113  
237 nucleotides of its upstream sequence; both share 84% nucleotide sequence identity. It is not surprising that such  
238 a similar region is again found for MC19, as MC13 and MC19 are identical except a 1-bp difference (Zhong et

239 al. 2002). The first 1067 nucleotides of ORF 1 of the pP37-74 insert display 88% sequence identity to  
240 nucleotides 84–1150 of MC9. ORFs 4 and 5 of the same insert located adjacent to each other and overlapping by  
241 four bp, exhibit a genetic organization identical to MC34 and MC35 of pSD10. Moreover, a 1056-bp region,  
242 consisting of ORFs 4 and 5 (pP37-74) and including 44 bp of upstream sequence of ORF 4, shares 89%  
243 nucleotide sequence identity with a region starting 44 bp upstream of MC34 and covering MC34 as well as the  
244 first 713 nucleotides of MC35. It is to be expected that the high degree of nucleotide sequence identity of the  
245 three truncated putative transposase genes (ORFs 1 and 5 of pP37-74, ORF 7 of pP62-34) continues throughout  
246 the missing coding sequences of these genes.

247 High similarities of above transposase genes and adjacent non-coding regions and their presence in different  
248 *Micrococcus* strains points to horizontal gene transfer as well as exchange of genetic information between the  
249 plasmid and the host chromosome. Furthermore, the high number of ORFs on pLMA1 involved in transposition,  
250 which probably increases upon further sequencing (the total sequence of 22995 bp obtained so far represents  
251 only one fifth of the entire length of the plasmid), indicates that pLMA1 is presumably very flexible with respect  
252 to its genetic composition.

253 A number of transposases are associated with insertion sequence (IS) elements, which differ in size and  
254 sequence composition and are flanked by short inverted repeats (IR), usually between 10 and 40 bp in size  
255 (Mahillon and Chandler 1998). Based on IR sequences reported for pSD10 (Zhong et al. 2002), we checked the  
256 pLMA1 fragments and identified a 24-bp region upstream of ORF 4 of pP37-74 (5'-  
257 GGACTGACGCACGTGTAGGTGACA-3') differing in six (underlined) positions from the IR (5'-  
258 GGACTGGTGTACACATAGGT-GACA-3') found upstream of MC35. Despite the lack of the corresponding  
259 IR, which is probably located downstream of ORF 5 (pP37-74) within yet unknown pLMA1 sequence, this  
260 finding suggests ORFs 4 and 5 (pP37-74) being part of an IS element. Although the cloned pLMA1 fragments  
261 are rich in genes encoding putative transposases, further IRs were not identified. However, for some bacterial  
262 genomes and plasmids IS elements lacking IRs have been reported (Burland et al. 1998).

263 IS elements and transposons are frequently associated with antibiotic resistance genes (Varaldo et al. 2009).  
264 *M. luteus* A1, hosting pLMA1, is resistant to a number of antibiotics, particularly to macrolides, such as  
265 erythromycin, for which plasmid curing experiments suggested a link to pLMA1 (Dib et al. 2010a). Resistance  
266 to erythromycin is either conferred by a 23S rRNA methylase-mediated target site modification or by an efflux  
267 system; the latter resulting in low-level resistance (Varaldo et al. 2009). Since the minimal inhibitory  
268 concentration of erythromycin for *M. luteus* A1 wild-type reached a level  $>256 \mu\text{g ml}^{-1}$  (Dib et al. 2010a), a

269 methylase-mediated resistance mechanism is likely to be encoded. As we were not able to identify an ORF  
270 coding for a 23S rRNA methylase, such function is presumably located on the yet unknown pLMA1 sequence.

271

272 pLMA1 is laced with short repetitive sequences

273

274 To verify that repetitive sequences are responsible for the failure of the 454 read assembly, we looked for repeats  
275 within the nucleotide sequences of the cloned pLMA1 fragments. Indeed, a huge number of repetitive sequences  
276 (mostly ranging from 6–12 nucleotides) were found. Such repetitions occur at high frequency and are spread  
277 quite evenly over the pLMA1 fragments; partially they overlap each other. Conspicuously arranged repeats were  
278 seen on the 6195-bp *Pst*I fragment; we identified five direct repeats (putative iterons) starting at nucleotide  
279 positions 2510, 2571, 2632, 2693, and 2754, respectively. Iterons 1–4 are perfect direct repetitions sharing a 29-  
280 bp consensus sequence (5'-GGAAGCCCCGCGCAGCAGGGATGAGCCC-3'); only iteron 5 (5'-  
281 GGAAGCTCCGCCCGCAGGGATGAGCCA-3') differs at four (underlined) positions. All five iterons are  
282 regularly spaced by 32 bp.

283 Such iteron sequences are characteristic for replication origins of linear plasmids such as pCLP of  
284 *Mycobacterium celatum* (Le Dantec et al. 2001) and many *Streptomyces* linear plasmids such as pSLA2 (Chang  
285 et al. 1996) and pRL2 (Zhang et al. 2008). Iterons, as the initiation sites of plasmid DNA replication, usually are  
286 located adjacent to *rep* genes: *rep1* encoding a DNA-binding protein and *rep2* a DNA helicase (Chang et al.  
287 1996, Zhang et al. 2008); though, other recently described *rep* genes are clearly discriminable from the above  
288 (Zhang et al. 2006). For pLMA1, ORFs adjacent to the iterative sequences do not show any similarity to known  
289 *rep* loci. However, iterons and genes involved in plasmid replication may be separated by nonessential genes, as  
290 for the *Streptomyces* linear plasmid pRL2 (Zhang et al. 2008), or by noncoding sequences as for pSHK1 (Zhang  
291 et al. 2008). Indeed, studies on the structure of replication loci of *Streptomyces* linear plasmids in general  
292 revealed an unexpected variety of components and their positions (Zhang et al. 2008). However, the role of the  
293 putative iteron sequences of pLMA1 as well as other functions located outside of the sequenced fragments,  
294 possibly involved in plasmid replication, needs to be elucidated.

295 The arrangement of repetitive sequences is exemplarily shown for the 3681-bp pLMA1 fragment (pP37-74)  
296 in comparison to the 3721-bp chromosomal fragment of *M. luteus* A1 (pP37-62) (Fig. 3). Though the  
297 chromosome also has sequence repeats, their number is lower and they are irregularly arranged (Fig. 3). To  
298 appraise the observed frequency of the single identified repetitive sequences, we calculated their expected

Approx.  
position of  
Fig. 3

299 frequency and compared both values (Table 2). This calculation, based on length and nucleotide composition of  
300 the DNA fragment, clearly shows an over-representation of the sequence repeats; deviation from the expectation  
301 becomes higher the longer the repeat is (Table 2).

302 Such dispersed repetitions along with the short read lengths (see above) can be considered the major reason  
303 for the failure of the 454 read assembly. Additives for high-GC DNA as well as an optimized sequencing system  
304 (Titanium Series, Roche), becoming only quite recently available, may allow for longer reads of up to 400–500  
305 bases. However, despite such improvements, assembly of reads of pLMA1 still remains uncertain. As  
306 Sanger/pyrosequencing hybrid approaches were already successfully used for the sequencing of bacterial  
307 genomes, we decided to apply conventional Sanger sequencing for pLMA1 which currently is in progress.

308

### 309 **Acknowledgments**

310

311 We gratefully acknowledge the support by PICT-Agencia Nacional de Promoción Científica y Tecnológica.  
312 Julián R. Dib was a recipient of a CONICET and DAAD fellowship.

313

### 314 **References**

315

316 Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ (1990) Basic local alignment search tool. *J Mol Biol*  
317 215(3):403–410

318

319 Armstrong K, Bauer WR (1982) Preferential site-dependent cleavage by restriction endonuclease *Pst*I. *Nucleic*  
320 *Acids Res* 10(3):993–1007

321

322 Benson DA, Karsch-Mizrachi I, Lipman DJ, Ostell J, Wheeler DL (2004) Genbank: Update. *Nucleic Acids Res*  
323 32(Database issue):D23–26

324

325 Burland V, Shao Y, Perna NT, Plunkett G, Sofia HJ, Blattner FR (1998) The complete DNA sequence and  
326 analysis of the large virulence plasmid of *Escherichia coli* O157:H7. *Nucleic Acids Res* 26(18):4196–  
327 4204

328

329 Chang PC, Kim ES, Cohen SN (1996) *Streptomyces* linear plasmids that contain a phage-like, centrally located,  
330 replication origin. *Mol Microbiol* 22(5):789–800  
331

332 Dib J, Motok J, Zenoff VF, Ordonez O, Farias ME (2008) Occurrence of resistance to antibiotics, UV-B, and  
333 arsenic in bacteria isolated from extreme environments in high-altitude (above 4400 m) Andean wetlands.  
334 *Curr Microbiol* 56(5):510–517  
335

336 Dib JR, Wagenknecht M, Hill RT, Farias ME, Meinhardt F (2010a) First report of linear megaplas mids in the  
337 genus *Micrococcus*. *Plas mid* 63(1):40–45  
338

339 Dib JR, Wagenknecht M, Hill RT, Farias ME, Meinhardt F (2010b) Novel linear megaplas mid from  
340 *Brevibacterium* sp. isolated from extreme environment. *J Basic Microbiol* 50(3):280–284  
341

342 Goldberg SM, Johnson J, Busam D, Feldblyum T, Ferriera S, Friedman R et al. (2006) A sanger/pyrosequencing  
343 hybrid approach for the generation of high-quality draft assemblies of marine microbial genomes. *Proc*  
344 *Natl Acad Sci U S A* 103(30):11240–11245  
345

346 Le Dantec C, Winter N, Gicquel B, Vincent V, Picardeau M (2001) Genomic sequence and transcriptional  
347 analysis of a 23-kilobase mycobacterial linear plasmid: Evidence for horizontal transfer and identification  
348 of plasmid maintenance systems. *J Bacteriol* 183(7):2157–2164  
349

350 Mahillon J, Chandler M (1998) Insertion sequences. *Microbiol Mol Biol Rev* 62(3):725–774  
351

352 Margulies M, Egholm M, Altman WE, Attiya S, Bader JS, Bemben LA et al. (2005) Genome sequencing in  
353 microfabricated high-density picolitre reactors. *Nature* 437(7057):376–380  
354

355 Meinhardt F, Klassen R (2007) *Microbial linear plasmids*. Springer-Verlag, Berlin Heidelberg  
356

357 Ohama T, Muto A, Osawa S (1989) Spectinomycin operon of *Micrococcus luteus*: Evolutionary implications of  
358 organization and novel codon usage. *J Mol Evol* 29(5):381–395

359  
360 Pearson WR (1994) Using the FASTA program to search protein and DNA sequence databases. *Methods Mol*  
361 *Biol* 25:365–389  
362  
363 Sharma D, Issac B, Raghava GP, Ramaswamy R (2004) Spectral repeat finder (SRF): Identification of repetitive  
364 sequences using fourier transformation. *Bioinformatics* 20(9):1405–1412  
365  
366 Varaldo PE, Montanari MP, Giovanetti E (2009) Genetic elements responsible for erythromycin resistance in  
367 streptococci. *Antimicrob Agents Chemother* 53(2):343–353  
368  
369 Wagenknecht M, Meinhardt F (2010) Copy number determination, expression analysis of genes potentially  
370 involved in replication, and stability assays of pAL1 – the linear megaplasmid of *Arthrobacter*  
371 *nitroguajacolicus* Rü61a. *Microbiol Res*. doi:10.1016/j.micres.2009.12.005  
372  
373 Wicker T, Schlagenhauf E, Graner A, Close TJ, Keller B, Stein N (2006) 454 sequencing put to the test using the  
374 complex genome of barley. *BMC Genomics* 7:275  
375  
376 Yang CC, Huang CH, Li CY, Tsay YG, Lee SC, Chen CW (2002) The terminal proteins of linear streptomyces  
377 chromosomes and plasmids: A novel class of replication priming proteins. *Mol Microbiol* 43(2):297–305  
378  
379 Zhang R, Xia H, Guo P, Qin Z (2009) Variation in the replication loci of *Streptomyces* linear plasmids. *FEMS*  
380 *Microbiol Lett* 290(2):209–216  
381  
382 Zhang R, Yang Y, Fang P, Jiang C, Xu L, Zhu Y et al. (2006) Diversity of telomere palindromic sequences and  
383 replication genes among *Streptomyces* linear plasmids. *Appl Environ Microbiol* 72(9):5728–5733  
384  
385 Zhong Z, Caspi R, Mincer T, Helinski D, Knauf V, Boardman K et al. (2002) A 50-kb plasmid rich in mobile  
386 gene sequences isolated from a marine *Micrococcus*. *Plasmid* 47(1):1–9  
387

388 **Legends to figures**

389

390 **Fig. 1** Origin of the cloned *Pst*I restriction fragments by Southern analysis. Total DNA of the wild-type strain *M.*  
391 *luteus* A1, of the plasmid-deficient strain *M. luteus* A1-M1, and isolated pLMA1 DNA was digested using  
392 restriction endonuclease *Pst*I, separated on 1.0% agarose gels (left panels), and transferred onto nylon  
393 membranes. The DNA was hybridized with probes deduced and PCR-amplified from the cloned *Pst*I restriction  
394 fragments (see also Fig. 2). Right panels display corresponding sections of exposed X-ray films. Only relevant  
395 sections of gels/films are shown. Arrows point to pLMA1 restriction fragments with sizes of 8.6 kb (a), 6.2 kb  
396 (b), 4.5 kb (c), and 3.7 kb (d, e), respectively, that were selected for cloning and supposed to give a signal upon  
397 hybridization. M, DNA size standard; 1 and 4, *M. luteus* A1; 2 and 5, *M. luteus* A1-M1; 3 and 6, pLMA1

398

399 **Fig. 2** Schematic representation of annotated ORFs on the sequenced *Pst*I restriction fragments of pLMA1.  
400 Designations of the cloning plasmids harboring the respective pLMA1 fragment are shown on the left, the sizes  
401 of the pLMA1 fragments are given in brackets on the right. Predicted ORFs are shown as arrows. Arrows  
402 missing the beginning or the arrowhead (serrated lines) indicate interrupted ORFs. The direction of the arrows  
403 corresponds to the transcriptional directions. ORFs encoding for proteins involved in transposition are gray  
404 shaded. Probes used in Southern analysis (see Fig. 1) correspond to thick, black bars with their respective  
405 localisation given below the pLMA1 fragments. Recognition sites of restriction endonuclease *Pst*I are indicated  
406 as such and by short vertical lines. The sequences of the cloned 8621-bp, 6195-bp, 4498-bp, and 3681-bp  
407 pLMA1 fragments, and the 3721-bp chromosomal fragment have been deposited in the EMBL nucleotide  
408 sequence database under accession numbers FN692038, FN692039, FN692040, FN692041, and FN692042,  
409 respectively

410

411 **Fig. 3** Schematic map indicating identified repetitive sequences. The chromosomal fragment of *M. luteus* A1  
412 (pP37-62) and a fragment of pLMA1 (pP37-74) are shown as horizontal, black lines. Annotated ORFs of the  
413 pLMA1 fragment are denoted as open arrows. Arrows missing the beginning or the arrowhead (serrated lines)  
414 show interrupted ORFs. Vertical lines in different greyscales indicate repetitive sequences. In the enlarged  
415 section shown below, such sequences are exhibited in more detail. Length and orientation of the repetitive  
416 sequences are marked by gray arrows. Repeat lengths are earmarked by bold numbers. Identical repeats are

417 designated by the same subscript. The 0.5-kb size bar is for pP37-62 and pP37-74, the 50-bp bar is for the  
418 enlarged section only