

# Complete Chloroplast Genome Sequence of a Major Allogamous Forage Species, Perennial Ryegrass (*Lolium perenne* L.)

KERSTIN Diekmann<sup>1,2,\*</sup>, TREVOR R. Hodkinson<sup>2</sup>, KENNETH H. Wolfe<sup>3</sup>, ROB VAN DEN Bekerom<sup>1,4</sup>, PHILIP J. Dix<sup>4</sup>, and SUSANNE Barth<sup>1</sup>

*Teagasc Crops Research Centre, Oak Park, Carlow, Ireland*<sup>1</sup>; *School of Natural Sciences, Trinity College Dublin, Dublin 2, Ireland*<sup>2</sup>; *Smurfit Institute of Genetics, Trinity College Dublin, Dublin 2, Ireland*<sup>3</sup> and *Department of Biology, National University of Ireland Maynooth, Maynooth, Ireland*<sup>4</sup>

(Received 10 February 2009; accepted 12 April 2009; published online 4 May 2009)

## Abstract

*Lolium perenne* L. (perennial ryegrass) is globally one of the most important forage and grassland crops. We sequenced the chloroplast (cp) genome of *Lolium perenne* cultivar Cashel. The *L. perenne* cp genome is 135 282 bp with a typical quadripartite structure. It contains genes for 76 unique proteins, 30 tRNAs and four rRNAs. As in other grasses, the genes *accD*, *ycf1* and *ycf2* are absent. The genome is of average size within its subfamily Pooideae and of medium size within the Poaceae. Genome size differences are mainly due to length variations in non-coding regions. However, considerable length differences of 1–27 codons in comparison of *L. perenne* to other Poaceae and 1–68 codons among all Poaceae were also detected. Within the cp genome of this outcrossing cultivar, 10 insertion/deletion polymorphisms and 40 single nucleotide polymorphisms were detected. Two of the polymorphisms involve tiny inversions within hairpin structures. By comparing the genome sequence with RT-PCR products of transcripts for 33 genes, 31 mRNA editing sites were identified, five of them unique to *Lolium*. The cp genome sequence of *L. perenne* is available under Accession number AM777385 at the European Molecular Biology Laboratory, National Center for Biotechnology Information and DNA DataBank of Japan.

**Key words:** chloroplast genome; *Lolium perenne*; Poaceae; chloroplast DNA variation; RNA editing

## 1. Introduction

Chloroplasts (cps), plant cell organelles derived from independent living cyanobacteria,<sup>1–3</sup> contain their own small genome averaging 150 kb in flowering plants. The cp genome molecules can be circular or linear, mono- or multimeric,<sup>4</sup> but the genome can be represented by a monomeric circular map containing two copies of an inverted repeat (IR) region (~23 kb) which separate a small single copy (SSC)

region (~18 kb) from a large single copy (LSC) region (~84 kb). In most angiosperm species, the cp genome contains ~113 different genes<sup>5</sup> that primarily encode for proteins and RNAs for the photosynthetic system and that are generally highly conserved in terms of content and order among plant families.<sup>6</sup> Cp genomes are usually inherited maternally,<sup>7</sup> and this property is useful for several applications such as for defining cytoplasmic breeding pools in plant breeding, and tracking parentage in interspecific hybrids (e.g. *Arabidopsis suecica*<sup>8</sup>). Cp genetic engineering is also an ideal approach for minimizing the risk of spreading transgenes into wild plants via pollen.<sup>9</sup> In comparison with nuclear genetic engineering, much higher expression of the transgenic insertion can also be obtained because of

Edited by Satoshi Tabata

\* To whom correspondence should be addressed. Tel. +353 59-9170-243. Fax. +353 59-9142-423. E-mail: kerstin.diekmann@teagasc.ie

© The Author 2009. Kazusa DNA Research Institute.

The online version of this article has been published under an open access model. Users are entitled to use, reproduce, disseminate, or display the open access version of this article for non-commercial purposes provided that: the original authorship is properly and fully attributed; the Journal and Oxford University Press are attributed as the original place of publication with the correct citation details given; if an article is subsequently reproduced or disseminated not in its entirety but only in part or as a derivative work this must be clearly indicated. For commercial re-use, please contact journals.permissions@oxfordjournals.org

the high copy number of cp genomes within a single plant cell. Cp genome sequences are also highly suitable for phylogenetic studies.<sup>10</sup>

To date (February 2009), entire cp genome sequences of 117 streptophytic species are publicly available ([http://www.ncbi.nlm.nih.gov/genomes/ORGANELLES/plastids\\_tax.html](http://www.ncbi.nlm.nih.gov/genomes/ORGANELLES/plastids_tax.html)). Only 18 of these genome sequences belong to the monocot group of angiosperms, and of these 13 are from the grass family Poaceae. Poaceae include the most important agricultural plant species from a socio-economic perspective as they contain the cereals and forage species.<sup>11</sup> *Lolium perenne* (perennial ryegrass) is globally one of the most important grassland species especially for the northern hemisphere (<http://www.worldseed.org>). In 2006–2007, more than one-third of world grass seed production was from *L. perenne*. Thus *L. perenne* has the highest economic impact as a forage and grassland crop. It is a cross-pollinating species and cultivar populations consist of a heterozygous nuclear genome background.

Several methods exist for obtaining complete cpDNA sequences. The *Arabidopsis thaliana* cp genome, for example, was sequenced using cpDNA clones found as ‘contaminations’ in genomic libraries.<sup>12</sup> The cp genome of *Nicotiana sylvestris*, the maternal genome donor of *Nicotiana tabacum*, was obtained by sequencing extracted high-purity cpDNA that was cloned into sequencing vectors.<sup>13</sup> A commonly used method involves amplifying the cp genome by rolling circle amplification and then cloning this product into sequencing vectors.<sup>14</sup> Recently, consensus cpDNA sequencing primers have become available for sequencing cp genomes using a primer walking strategy.<sup>15</sup> For sequencing the cp genome of *L. perenne*, we extracted high-purity cpDNA which we amplified with a whole genome amplification kit and used a shotgun sequencing approach. Thus each region of the genome was sequenced several fold from independent clones, which allowed us to detect SNPs and indels.

Few studies have examined variation of the cp genome within a population of a species. However, McGrath et al.<sup>16</sup> discovered more than 500 haplotypes within 1575 individual plants of *Lolium*, *Festulolium* and *Festuca* populations. We hope to add to this information by assessing cp genome variation within a *Lolium* cultivar by detecting SNPs and indels. This assessment should reveal highly variable regions in the *Lolium* cp genome, from which markers can be designed for assessing cytoplasmic breeding pools and to add to population genetic and phylogenetic studies.

In this study, we also analysed RNA editing sites in *L. perenne* cp transcripts. RNA editing is a repair mechanism that alters the genetic information of

land plant organelles at the transcript level. It is a post-transcriptional modification (mostly C to U conversion) of the nucleotide sequence of pre-mRNAs by inserting, deleting or substituting nucleotides in order to yield functional RNA species.<sup>17,18</sup> Editing in cps was first discovered by Hoch et al.<sup>19</sup> for the cp *rpl2* gene in maize, where it creates a start codon and hence restores the functionality of the *rpl2* gene. Knowledge about RNA editing sites is essential for describing the functional capability of cp genes, characterizing different species and obtaining a better understanding of how these sites have evolved.

## 2. Materials and methods

### 2.1. Sequencing, assembling and annotating the cp genome

cpDNA was isolated from the *L. perenne* cultivar Cashel following a protocol from Diekmann et al.<sup>20</sup> Approximately 400 g of 3-week-old leaf material derived from ca. 200 g of a heterozygous Cashel seed population was used. Sequencing of the cpDNA was sourced to a commercial company (GATC Biotech/Germany). A shotgun sequencing approach was used resulting in 2179 trace files. A pre-assembly was carried out with the program PHRAP (<http://www.phrap.org/index.html>). The final assembly was based on the contiguous sequences obtained from PHRAP and done in comparison with the cpDNA sequence of *Agrostis stolonifera* (bentgrass).<sup>21</sup> For three genome regions with low coverage, primers were designed to re-sequence these regions (*trnL-trnF*: forward primer (FP): AGTTGTGAGGGTTC AAGTCC and reverse primer (RP): GAACTGGTGACAC GAGGATT; *atpB*: FP: GTTCGTTGCCAACAACTCTA and RP: AGGTAGCTCTAGTCTATGGC; *atpB-rbcL*: FP: TGTGG AAGATCTGTGCCTAC and RP: GCTGAGGAGTTACTCG GAAT).

The annotation of the cp genome was based on two online available programs: DOGMA (<http://dogma.cccb.utexas.edu/>) and tRNA-Scan SE (<http://lowelab.ucsc.edu/tRNAscan-SE/>) using the default settings. Intron positions were determined following Sugita and Sugiura.<sup>22</sup> The circular cp genome map was drawn using the GenomeVx program.<sup>23</sup> Differences between the available cp genomes were analysed based on gene, intergenic spacer (IGS) and intron lengths which were extracted from the published cp genome sequences (*A. stolonifera*: EF115543; *Brachypodium distachyon*: EU325680; *Hordeum vulgare*: EF115541; *Oryza nivara*: AP006728; *Oryza sativa indica*: AY522329; *Oryza sativa japonica*: X15901; *Sorghum bicolor*: EF115542; *Saccharum officinarum*: AP006714; *Triticum aestivum*: AB042240; *Zea mays*: X86563).

## 2.2. SNP and indel analysis

Because the cpDNA had been extracted from a population of plants belonging to the cultivar Cashel, several SNPs and indels could be detected. A thorough SNP and indel analysis was carried out by manually checking the alignment of the read and trace files from which the genome assembly was undertaken using the programme Lasergene (DNASTAR, Inc., Madison, Wisconsin). Only SNPs and indels supported by trace files with low background and clear, distinguishable peaks were recorded. Indels were only taken into account if they were supported by at least two trace files and not located in coding regions where they would cause a frame shift. This way the possibility of cloning and sequencing artefacts was considered.

## 2.3. RNA editing analysis

Thirty-three genes (*atpA*, *atpB*, *atpF*, *clpP*, *matK*, *ndhA*, *ndhB*, *ndhD*, *ndhF*, *ndhG*, *ndhI*, *ndhK*, *petA*, *petB*, *psaA*, *psaB*, *psaJ*, *psbC*, *psbD*, *psbE*, *psbJ*, *psbL*, *psbZ*, *rpl2*, *rpl20*, *rpoA*, *rpoB*, *rpoC1*, *rpoC2*, *rps14*, *rps2*, *rps8*, *ycf3*) were analysed for RNA editing sites. Of these, 22 were chosen for study because they had been previously reported to be edited in other monocot plants,<sup>24–28</sup> and 11 were included because of observed differences from existing expressed sequence tags (EST) in Poaceae,<sup>21</sup> but no information was previously available for *Lolium*. Primers (Supplementary Table S1) for these genes were designed using Primer Express (version 2.0, Applied Biosystems, Foster City, CA, USA) and Primer3 software (<http://frodo.wi.mit.edu/>). For genes > 700 bp, several primer pairs were designed to cover the complete gene region. Primers were designed in the untranslated regions (UTR) to ensure complete coverage of genes. Since the length of the UTR of genes was not known, the primers were designed in the 30 bp region before and after each gene.

cDNA was used as template for the RT–PCRs. Total RNA was extracted using TRI Reagent® Solution (Ambion Inc., Austin, TX, USA) following the supplier's protocol ([http://www.ambion.com/techlib/prot/bp\\_9738.pdf](http://www.ambion.com/techlib/prot/bp_9738.pdf)) with the following modifications: the incubation of the homogenate was extended to 10 min; instead of 100 µl bromochloropropane, 200 µl of ice cold chloroform was used; the steps including the addition of ice cold chloroform, followed by incubation at room temperature and centrifugation at 12 000g were repeated once; in addition to the 500 µl isopropanol, 0.5 µl Glycogen (Sigma-Aldrich, St Louis, Missouri, USA) was added to enhance the RNA yield; the centrifugation following the addition of isopropanol was extended to

10 min. The RNA was finally dissolved in nuclease free water and treated with DNA-free™ (Ambion Inc., Austin, TX, USA) following the manufacturer's instructions to remove possible DNA contamination. First strand cDNA was synthesized using SuperScript™ III Reverse Transcriptase (Invitrogen™ Corporation, Carlsbad, CA, USA) following the manufacturer's instructions.

For each gene region, two independent RT–PCR reactions were set up using the following components per 30 µl PCR reaction: 3 µl cDNA, 3 µl 10 x Thermo Buffer (New England Biolabs, Inc., Ipswich, MA, USA), 0.6 µl FP, 0.6 µl RP, 0.6 µl dNTPs (metabion international AG, Martinsried, Germany) (10 mM), 21.9 µl ddH<sub>2</sub>O, 0.3 µl Taq-Polymerase (New England Biolabs, Inc.). The PCR programme settings were 95°C 5 min, (95°C 1 min, 55°C 1 min, 72°C 1 min) 35 cycles, 72°C 10 min. The annealing temperature was adjusted according to the optimal primer requirements. The resulting RT–PCR products were sequenced twice using both forward and reverse primers. The analysis of the editing sites was carried out in MEGA 3.1<sup>29</sup> by aligning the cDNA sequence results with the corresponding DNA sequences and checking visually for SNPs.

## 3. Results and discussion

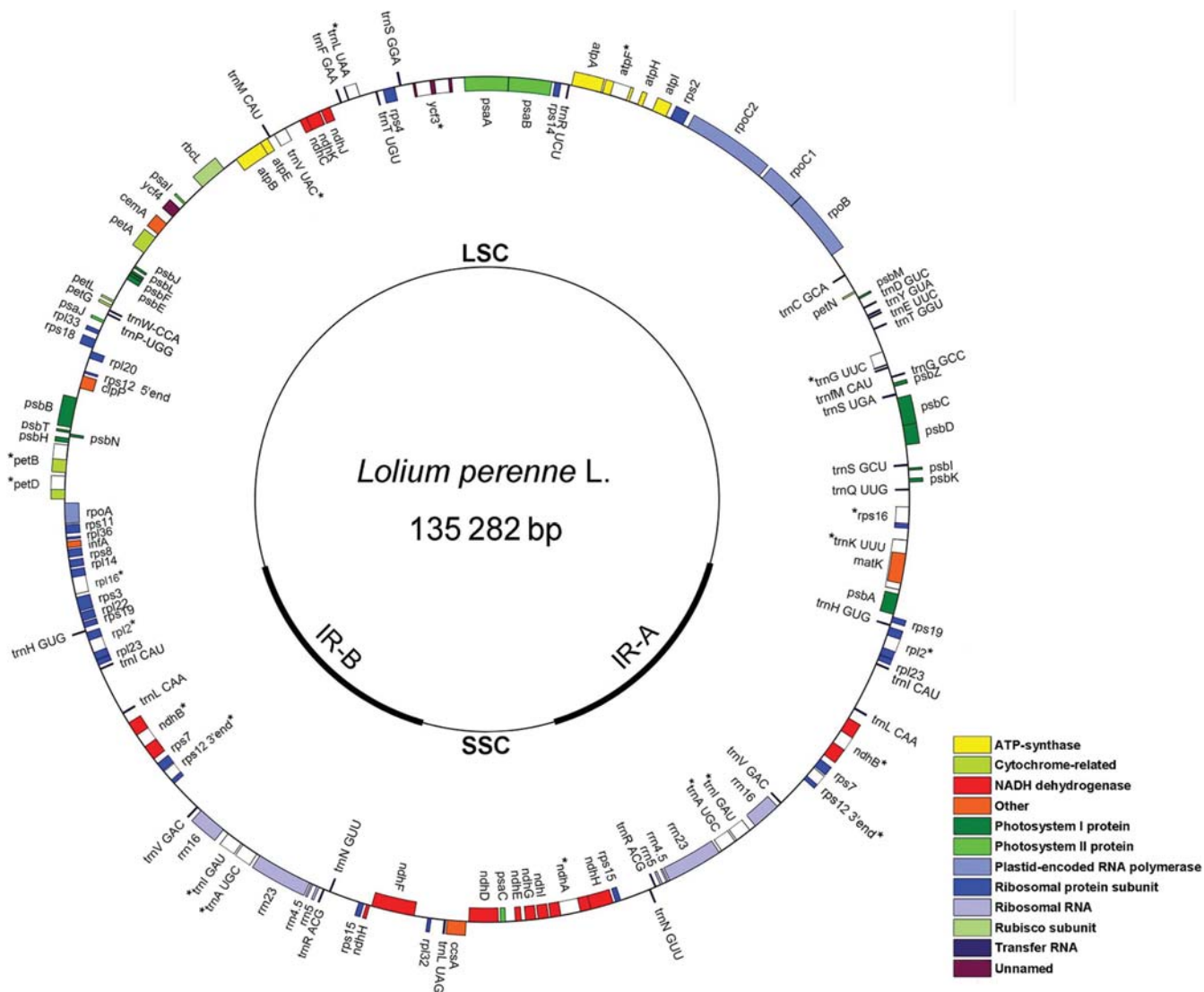
Using the shotgun sequencing approach an average eightfold genome coverage was achieved. The cp genome of *L. perenne* has a total length of 135 282 bp with a quadripartite structure typical of angiosperms. The LSC region consists of 79 972 bp, the SSC of 12 428 bp and the IRs of 21 441 bp each (Fig. 1). The genome has a GC content of 38% and codes for 128 genes of which 18 are duplicated in the IR region. The genome contains 264 simple sequence repeats (SSRs) with mononucleotide repeats of 7–16 bp in length. The cp genome sequence of *L. perenne* is deposited at the European Molecular Biology Laboratory under Accession number AM777385.

### 3.1. Comparison to other species

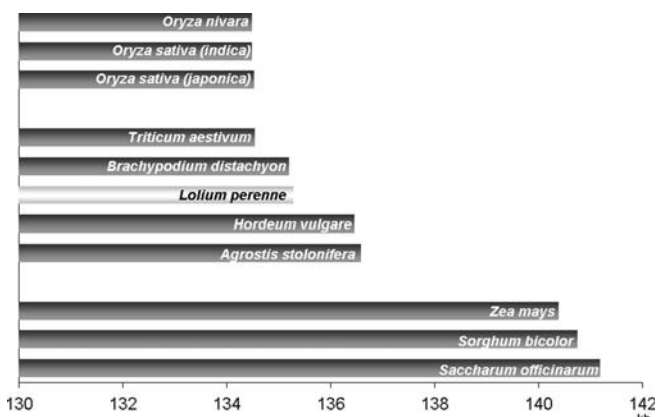
The average size of publicly available Poaceae cp genomes is 137 091 bp. The subfamily Ehrhartoideae has the smallest genome with an average size of 134 505 bp; subfamily Panicoideae has the largest genome with an average size of 140 876 bp. The subfamily Pooideae, to which *L. perenne* belongs, has an average size of 135 614 bp. Thus *L. perenne* is of average size within Pooideae and of medium size within Poaceae (Fig. 2).

The gene content and intron content of *L. perenne* cpDNA are the same as that of other





**Figure 1.** Circular structure of the chloroplast genome of *Lolium perenne*. Genes written on the outside are transcribed clockwise, genes on the inside counter-clockwise, annotated genes are colour coded according to their function, genes containing introns are highlighted with an asterisk; LSC, large single copy region; SSC, small single copy region; IR, inverted repeat.



**Figure 2.** Chloroplast genome sizes of 11 different Poaceae species grouped by taxonomic sub-families.

grasses,<sup>21,24,27,30,31,32</sup> with 76 protein-coding genes, 30 tRNA genes and four rRNA genes. Eighteen genes are completely duplicated within the IR, as are the 3' exons of the *trans*-spliced gene *rps12* and the 5' part of *ndhH* which overlaps the IR/SSC junctions. When compared with the standard set of genes in angiosperm cp genomes, the genes *accD*, *ycf1* and *ycf2* are absent. After our analysis was completed, the cp genome sequence of the very closely related species *Festuca arundinacea* became available in GenBank (Cahoon et al., unpublished data; accession number FJ466687). Rather surprisingly, in addition to the expected absences of *accD*, *ycf1* and *ycf2*, the *Festuca* sequence also lacks intact copies of the genes *psbF*, *rps14*, *rps18* and *ycf4*. All four of these

**Table 1.** Length variation of > 100 bp for intergenic spacer and intron regions in Poaceae chloroplast genomes

Genome region	Ehrhartoideae			Pooideae			Panicoideae			Variation (bp) <sup>a</sup>
	<i>Oryza nivara</i>	<i>Triticum aestivum</i>	<i>Brachipodium distachyon</i>	<i>Lolium perenne</i>	<i>Hordeum vulgare</i>	<i>Agrostis stolonifera</i>	<i>Zea mays</i>	<i>Sorghum bicolor</i>	<i>Saccharum officinarum</i>	
<b>Large single copy region</b>										
<i>matK</i> – <i>trnK</i> -UUU	692	599	<b>444</b>	680	690	691	<b>695</b>	693	688	251
<i>rps16</i> – <i>trnQ</i> -UUG	1061	1062	784	815	<b>772</b>	787	1169	1521	<b>1530</b>	758
<i>trnS</i> -UGA– <i>psbZ</i>	347	358	359	356	357	<b>362</b>	<b>261</b>	344	344	101
<i>trnG</i> -GCC– <i>trnJ</i> M-CAU	434	449	438	451	449	438	<b>494</b>	<b>378</b>	492	116
<i>trnG</i> -UCC– <i>trnT</i> -GGU	1306	1200	<b>782</b>	1175	1194	1284	1874	<b>2013</b>	1956	1231
<i>trnT</i> -GGU– <i>trnE</i> -UUC	519	<b>306</b>	507	<b>537</b>	453	470	529	462	535	231
<i>trnD</i> -GUC– <i>psbM</i>	<b>381</b>	774	799	474	778	698	1052	<b>1059</b>	1052	678
<i>psbM</i> – <i>petN</i>	761	628	797	<b>279</b>	717	287	799	808	<b>811</b>	532
<i>petN</i> – <i>trnC</i> -GCA	<b>414</b>	949	935	917	725	921	955	933	<b>962</b>	548
<i>trnC</i> -GCA– <i>rpoB</i>	<b>1084</b>	1165	1196	1185	1166	1173	<b>1273</b>	1204	1213	189
<i>atpI</i> – <i>atpH</i>	795	572	<b>387</b>	570	568	531	818	756	<b>820</b>	433
<i>trnT</i> -UGU– <i>trnL</i> -UAA	764	<b>613</b>	<b>829</b>	825	624	819	813	797	801	216
<sup>b</sup> <i>trnL</i> -UAA– <i>trnL</i> -UAA	542	<b>589</b>	543	550	569	<b>424</b>	459	450	453	165
<i>trnL</i> -UAA– <i>trnF</i> -GAA	<b>245</b>	355	357	349	321	341	364	365	<b>366</b>	121
<i>trnF</i> -GAA– <i>ndhJ</i>	495	<b>448</b>	575	586	587	584	<b>591</b>	591	570	143
<i>ndhC</i> – <i>trnV</i> -UAC	<b>706</b>	911	816	844	727	926	<b>941</b>	924	929	235
<i>rbcL</i> – <i>psaI</i>	<b>1683</b>	880	<b>462</b>	1184	1604	1561	889	862	945	1221
<i>ycf4</i> – <i>cemA</i>	420	<b>475</b>	426	460	470	455	<b>330</b>	373	372	145
<i>petA</i> – <i>psbJ</i>	<b>1006</b>	821	835	<b>796</b>	821	806	900	900	900	210
<i>psbE</i> – <i>petL</i>	1197	1169	1277	1281	<b>1162</b>	<b>1286</b>	1237	1265	1214	124
<sup>b</sup> <i>petB</i> – <i>petB</i>	<b>814</b>	749	809	747	<b>697</b>	759	699	702	758	117
<sup>b</sup> <i>rpl16</i> – <i>rpl16</i>	1056	1044	1050	<b>868</b>	1064	1045	1043	1072	<b>1080</b>	212
<b>Inverted repeat</b>										
<i>trnI</i> -CAU– <i>trnL</i> -CAA	<b>1498</b>	1498	2497	2452	2380	2487	3630	3632	<b>3633</b>	2135
<i>rps12_3end</i> – <i>trnV</i> -GAC	1724	1726	<b>1642</b>	1646	1726	1756	1758	1767	<b>1767</b>	125
<sup>b</sup> <i>trnI</i> -GAU– <i>trnI</i> -GAU	948	807	806	808	802	<b>801</b>	950	948	<b>952</b>	151
<i>rps15</i> – <i>ndhF</i>	343	421	399	404	<b>422</b>	413	<b>107</b>	<b>107</b>	125	315
<b>Small single copy region</b>										
<i>ndhF</i> – <i>rpl32</i>	715	<b>919</b>	846	<b>703</b>	848	897	839	814	856	216
<i>rpl32</i> – <i>trnL</i> -UAG	547	690	697	663	<b>725</b>	659	531	<b>522</b>	523	203
<i>ndhG</i> – <i>ndhI</i>	243	<b>264</b>	251	<b>116</b>	250	252	184	184	184	148

Bold numbers show the shortest length for that intergenic spacer/intron. Bold and underlined numbers show the largest length for that intergenic spacer/intron.

<sup>a</sup>Difference between smallest and largest values.

<sup>b</sup>Highlights introns.

**Table 2.** Variation in length of different chloroplast genes

Gene name	Length in		Ehrhartoideae	Pooideae					Panicoideae			Codon variation
	bp	codons		<i>Oryza nivara</i>	<i>Triticum aestivum</i>	<i>Brachipodium distachyon</i>	<i>Lolium perenne</i>	<i>Hordeum vulgare</i>	<i>Agrostis stolonifera</i>	<i>Zea mays</i>	<i>Sorghum bicolor</i>	
<b>Large single copy region</b>												
<i>atpA</i>	1515	505	3	—	—	3	—	1	3	3	3	3
<i>atpF</i>	552	184	—	—	—	—	—	3	—	—	—	3
<i>infA</i>	324	108	—	6	—	—	6	—	—	—	—	6
<i>matK</i>	1536	512	—	—	31	—	—	—	2	4	2	31
<i>ndhK</i>	678	226	—	20	20	21	20	20	2	2	2	21
<i>petB</i>	648	216	—	—	—	—	19	—	19	19	—	19
<i>psaB</i>	2205	735	—	—	—	—	—	—	1	—	—	1
<i>psaJ</i>	129	43	2	—	—	—	—	—	—	—	—	2
<i>psbK</i>	183	61	1	1	—	1	1	2	1	1	1	2
<i>psbT</i>	102	34	2	5	2	5	5	5	—	—	—	5
<i>rbcL</i>	1431	477	1	1	—	1	3	1	—	—	—	3
<i>rpl16</i>	450	150	—	—	—	—	1	—	—	—	—	1
<i>rpl22</i>	444	148	2	1	2	—	2	2	1	1	1	2
<i>rpoA</i>	1014	338	—	2	—	4	2	2	2	2	2	4
<i>rpoB</i>	3228	1076	—	1	1	1	1	1	—	—	—	1
<i>rpoC1</i>	2031	677	6	7	6	—	6	6	7	7	7	7
<i>rpoC2</i>	4401	1467	47	13	37	—	36	—	61	54	68	68
<i>rps16</i>	189	63	—	23	23	27	23	23	23	23	23	27
<i>rps18</i>	471	157	7	14	7	—	14	13	14	7	7	14
<i>rps3</i>	675	225	15	15	15	15	15	15	—	—	—	15
<b>Inverted repeat</b>												
<i>rps12</i>	363	121	4	1	—	4	4	4	4	4	4	4
<i>rps15</i>	237	79	12	12	12	12	12	12	—	—	—	12
<b>Small single copy region</b>												
<i>ccsA</i>	960	320	2	3	3	—	3	—	2	2	2	3
<i>ndhD</i>	1503	501	—	—	—	2	—	—	—	—	—	2
<i>ndhF</i>	2205	735	—	5	7	7	5	5	4	4	4	7
<i>rpl32</i>	180	60	—	4	—	—	—	2	—	—	—	4

If not otherwise stated numbers shown refer to amount of additional codons; —, no variation to the smallest length observed.

genes are intact and apparently functional in *L. perenne*.

Differences in the cp genome size of *L. perenne* compared with other Poaceae species are mainly due to length variations of IGS regions and introns (Table 1) and this finding was consistent with previous observations.<sup>27,32</sup> The length of IGS regions and introns varies widely from only a few base pairs up to several hundred. Twenty-five IGS regions and four introns were found to vary in length by more than 100 bp (Table 1). The highest variation in size (given in brackets) was found in the *trnI-CAU-trnL-CAA* IGS (2135 bp), the *trnG-UCC-trnT-GGU* IGS (1231 bp) and the *rbcL-psaI* IGS (1221 bp). The *trnI-CAU-trnL-CAA* IGS and *rbcL-psaI* IGS are sites

that contain pseudogenes for *ycf2* and *accD*, respectively, in Poaceae.<sup>27</sup> Both these pseudogenes and a *ycf1* pseudogene were detected in *L. perenne*. The *trnG-UCC - trnT-GGU* IGS is part of a 'divergence hotspot' described by Maier et al.<sup>27</sup> whose variability is caused by a large number of deletion/insertion events.

A comparison between *L. perenne* and the other Poaceae species showed differences in gene length for 26 genes (Table 2). The majority of these genes is in the LSC region. Length variations of more than ten codons were observed in eight genes (codon variation): *matK* (31), *ndhK* (21), *petB* (19), *rpoC2* (68), *rps3* (15), *rps15* (12), *rps16* (27) and *rps18* (14). The variation in gene length for the *rpoC2* gene was

more than twice that found in any other gene. *L. perenne* and *A. stolonifera* have the shortest *rpoC2* genes (each 4 401 bp). The *rps18* gene in *L. perenne* is up to 14 codons shorter than in the other species. The *ndhK* and *rps16* genes are 21 and 27 codons, respectively, longer in *L. perenne* than in *O. nivara*.

The length variations observed in *rps18* and *rpoC2* are noteworthy. In both cases, *L. perenne* showed the shortest of all sequences. Sequence variation between monocots and dicots for *rps18* has been described by Weglöhner et al.,<sup>33</sup> based on the occurrence of different numbers of the heptapeptide repeat SKQPFK near the N terminus of the protein. Our study revealed that length differences among Poaceae *rps18* sequences are mainly based on the same heptapeptide repeat (S/F)K(Q/K)(P/T)F(R/L/H/S/N)(K/R) as described by Weglöhner et al.<sup>33</sup> (Supplementary Fig. S1). This motif is present six times in *rps18* of *L. perenne*, *B. distachyon*, *O. sativa*, *O. nivara*, *S. bicolor* and *S. officinarum* and seven times in *rps18* of *A. stolonifera*, *H. vulgare*, *T. aestivum* and *Z. mays*. The *L. perenne rps18* gene is the shortest detected so far, because it has undergone an additional deletion of seven amino acids near its C terminus. The deletions do not result in the creation of stop codons and we expect the *L. perenne rps18* gene to be functional.

The largest length variation in Poaceae genes was found in *rpoC2*, of up to 68 codons difference between *L. perenne* and *S. officinarum*, and is due to several insertion/deletion events (data not shown). Comparisons of the *rpoC2* gene from dicots and monocots revealed that Poaceae have a unique insertion of ~400 bp in the middle of this gene.<sup>34–36</sup> Cummings et al.<sup>36</sup> demonstrated that this region is highly variable compared with its flanking regions and is rich with tandem repeats. Nearly, all the variations found between *L. perenne* and the panicoids are located in this specific insertion region. Analysing cytoplasmic male sterile (CMS) lines of *Sorghum*, Chen et al.<sup>37</sup> discovered a 165 bp deletion in this insertion region that suggests a possible relation between this deletion in *rpoC2* and the CMS-system.<sup>38</sup> So far this deletion was only observed in *Sorghum* but sequence comparisons (data not shown) revealed that one deletion that results in the shorter *L. perenne rpoC2* gene is located in the same region where the deletion occurs in *Sorghum*. Hence a higher susceptibility to variation in this gene region could be indicated and an investigation of *L. perenne* CMS lines in regard to variation to fertile lines may prove valuable for improving future *Lolium* breeding schemes.

### 3.2. Indel/SNP analysis

A total of 10 indels (Table 3) and 40 SNPs (Table 4) were found to be polymorphic among our sequencing

reads. All indels are located in intergenic regions. Indels occurred in microsatellite regions, resulting in both shortening (one occurrence) and lengthening (nine occurrences) of the sequenced region compared with the length that was observed in the majority of the trace files. Knowledge gained about the sequence variability of these regions can be used to design primers around those microsatellites for population genetic and phylogenetic studies and can be also used to support breeding schemes via defining cytoplasmic breeding pools. This will be of especially high value for breeding schemes based on inter-specific crosses between *Lolium* and *Festuca*.

Nineteen SNPs were found within IGS regions and introns and 21 within coding regions (Table 4). Most of the SNPs are due to transition mutations (20 A↔G and 8 C↔T), with 12 transversions. Closer analysis of the SNPs found at position 100 655 and 100 656 (*trnN-rps15* IGS) revealed that these SNPs are caused by a tiny inversion of two nucleotides which are flanked by an IR of 29 bp length forming a stable hairpin secondary structure (Fig. 3). The small inversion of TG within the *trnN-rps15* region in the IR is found in 13 of the 29 trace files covering the region. We also noticed another small inversion that was supported by only one sequence read and caused SNPs at position 18, 20, 21 and 23 (*rps19-psbA* IGS). This inversion spans six nucleotides (TTCTAG) that are flanked by an IR of 25 bp length (Fig. 3).

Small inversions like the ones revealed by our study have been found between species,<sup>39,40</sup> between genera<sup>39,41</sup> and also within populations of one other species, the conifer *Abies*.<sup>42</sup> The two inversions found in the current study lead, together with the

**Table 3.** Indels observed in the cp genome of *Lolium perenne*

Position	Nucleotide		Region	Trace files	
	Major	Minor		Absolute	%
8258	—	T	<i>trnS-psbD</i>	2	50.00
18191	—	T	<i>trnC-rpoB</i>	3	50.00
31190	T	—	<i>atpI-atpH</i>	3	27.27
62835	—	A	<i>psbE-petL</i>	3	42.86
62836	—	A	<i>psbE-petL</i>	3	42.86
63107	—	T	<i>psbE-petL</i>	4	50.00
66367	—	A	<i>rpl20-rps12</i>	3	30.00
66368	—	A	<i>rpl20-rps12</i>	3	30.00
80295	—	T	<i>rps19-trnH</i>	2	28.57
93161	—	G	<i>rrn16-trnI</i>	3	16.67

Major, most commonly found nucleotide; minor, least commonly found nucleotide; absolute and % columns refer to the amount of trace files containing the under-represented nucleotide.

**Table 4.** Single nucleotide polymorphisms in the chloroplast genome of *Lolium perenne*

Position	Nucleotide		IUPAC	Amino acid change	Region	Trace files	
	Major	Minor				Absolute	%
1618	A	T	W		<i>trnK</i> intron	1	20.00
19560	T	C	Y	I→T	<i>rpoB</i>	1	25.00
27177	G	A	R	S→N	<i>rpoC2</i>	1	25.00
28829	G	A	R	A→T	<i>rpoC2</i>	1	14.29
34720	G	A	R	—	<i>atpA</i>	4	36.36
37506	T	C	B	—	<i>psaB</i>	2	25.00
38978	C	A	M	—	<i>psaA</i>	1	7.14
40609	A	G	R	G→V	<i>psaA</i>	5	33.33
42894	A	G	R		<i>ycf3</i> intron	2	25.00
43270	G	A	R		<i>ycf3</i> intron	1	10.00
54360	C	A	M	Q→K	<i>rbcL</i>	1	50.00
61647	A	C	M	—	<i>psbE</i>	2	28.57
65631	C	T	Y	P→L	<i>rps18</i>	1	7.69
69066	G	A	R	A→T	<i>psbB</i>	1	16.67
86203	G	A	R	A→V	<i>ndhB</i> exon	1	5.56
94732	G	A	R		<i>trnA</i> intron	1	4.00
95307	G	A	R		<i>rrn23</i>	1	4.17
96920	C	A	M		<i>rrn23</i>	2	4.76
96968	C	G	S		<i>rrn23</i>	2	5.13
10390	A	G	R	—	<i>psaC</i>	6	42.86
109007	G	A	R	A→V	<i>ndhE</i>	1	10.00
18 <sup>a</sup>	C	T	Y		<i>rps19-psbA</i>	1	16.67
20 <sup>a</sup>	A	C	M		<i>rps19-psbA</i>	1	16.67
21 <sup>a</sup>	G	T	K		<i>rps19-psbA</i>	1	16.67
23 <sup>a</sup>	A	G	R		<i>rps19-psbA</i>	1	16.67
45874	A	C	M		<i>trnT-trnL</i>	2	50.00
47636	T	C	Y		<i>trnF-ndhJ</i>	1	20.00
51379	A	G	R		<i>trnV-trnM</i>	1	16.67
62341	T	G	K		<i>psbE-petL</i>	3	42.86
62521	A	C	M		<i>psbE-petL</i>	4	50.00
63360	G	A	R		<i>petL-petG</i>	3	50.00
73849	C	T	Y		<i>petD-rpoA</i>	1	8.33
82491	A	G	R		<i>rpl23-trnI</i>	1	5.88
83207	G	T	K		<i>trnI-trnL</i>	1	7.14
85260	G	A	R		<i>trnL-ndhB</i>	1	5.88
100655 <sup>a</sup>	C	T	Y		<i>trnN-rps15</i>	13	43.33
100656 <sup>a</sup>	A	G	R		<i>trnN-rps15</i>	13	43.33
103870	A	G	R		<i>ndhF-rpl32</i>	2	18.18
105222	C	T	Y		<i>rpl32-trnL</i>	1	14.29
108689	G	A	R		<i>psaC-ndhE</i>	1	7.69

Major, most commonly found nucleotide; minor, least commonly found nucleotide; absolute and % columns refer to the amount of trace files containing the under-represented nucleotide.

<sup>a</sup>Inversions.

level of observed SNPs, to the conclusion that the cp genome of *L. perenne* cv Cashel consists of at least two haplotypes but potentially scores more.

McGrath et al.<sup>16</sup> detected five haplotypes in 16 individuals of Cashel using a set of ten primers<sup>43</sup> amplifying eight different regions in the cp genome and sizing





**Table 5.** RNA editing sites found in the chloroplast genome of *Lolium perenne* in comparison with the editing sites found in other monocots

Gene	Site	Codon position	Editing sites	Edited codon	Amino acid change	<i>Lolium perenne</i>	<i>Hordeum vulgare</i> <sup>26</sup>	<i>Oryza sativa</i> <sup>25</sup>	<i>Saccharum officinarum</i> <sup>24</sup>	<i>Zea mays</i> <sup>27</sup>
<i>atpA</i>	1	383	35112	tCa	S→L	+		+	+	+
<i>matK</i>	1	420	1993	Cat	H→Y	+ <sup>a</sup>	+ <sup>46</sup>	+ <sup>50</sup>		
<i>ndhA</i>	1	17	111 250	tCa	S→L	+	+ <sup>47</sup>	(—)	+	+
	2	158	112 355	tCa	S→L	+	+ <sup>47</sup>	+	+	+
	3	188	112 777	tCa	S→L	+	+ <sup>47</sup>	+ <sup>50</sup>	+	+
	4	357		tCc	S→F	(—)	+ <sup>47</sup>	+	+	+
<i>ndhB</i>	1	50	87743	tCa	S→L	+	+	(—)	(—)	(—)
	2	156	87425	cCa	P→L	+	+	+	+	+
	3	196	87306	Cat	H→Y	+	+	+	+	+
	4	204	87281	tCa	S→L	+	+	+ <sup>a</sup>	+	+
	5	235	87188	tCc	S→F	+	+	+	(—)	(—)
	6	246	87155	cCa	P→L	+	+	+	+	+
	7	277	86347	tCa	S→L	+	+	+	+	+
	8	279	86341	tCa	S→L	+	+	+	(—)	(—)
	9	494	85696	cCa	P→L	+	+	+ <sup>a</sup>	+	+
<i>ndhD</i>	1	295 (293)	107 165	tCa	S→L	+ <sup>a</sup>	+ <sup>48</sup>	+	+	+
<i>ndhF</i>	1	21	103 675	tCa	S→L	+		+	+ <sup>a</sup>	+
<i>ndhG</i>	1	116	109 624	cCa	P→L	+		+ <sup>50</sup>		(—)
5'UTR		−10				(—)	(—) <sup>49</sup>	+	+	+
<i>ndhK</i>	1	2	49367	<i>gtC</i>	V→V	+				
	2	43	49245	<i>cCa</i>	P→L	+				
<i>petB</i>	1	204	72259	cCa	P→L	+		(—)	+	+
<i>psbJ</i>	1	20	61111	<i>cCt</i>	P→L	+				
<i>psbL</i>	1	37	61339	<i>ttC</i>	F→F	+ <sup>a</sup>		(—)		
<i>rpl2</i>	1	1	82030	aCg	T→M	+ <sup>a</sup>	+	+ <sup>a</sup>	+ <sup>a</sup>	+
<i>rpl20</i>	1	103	66009	tCa	S→L	+ <sup>a</sup>		(—)	+ <sup>a</sup>	+
	1	156	19737	tCa	S→L	+ <sup>a</sup>	+ <sup>28</sup>	+ <sup>a</sup>	+	+
	2	182	19815	tCa	S→L	+ <sup>a</sup>	+ <sup>28</sup>	+ <sup>a</sup>	+	+
	3	187	19830	tCg	S→L	+ <sup>a</sup>	+ <sup>28</sup>	+ <sup>a</sup> <sub>uCa</sub>	+	+
4	206		cCg	P→L	—	− <sup>28</sup>	—	—	+	
<i>rpoC2</i>	1	925		tCg	S→L	(—)		(—)	+	+
	2	1320	28731	<i>tCa</i>	S→L	+				—
<i>rps8</i>	1	61	76422	tCa	S→L	+		+	+	+
<i>rps14</i>	1	27		tCa	S→L	(—)		+	+ <sup>a</sup>	+
<i>ycf3</i>	1	15	43599	tCc	S→F	+		(—)	—	+ <sup>51a</sup>
	2	62	42 700	aCg	T→M	+		+ <sup>a</sup>	+	+

—, editing although C encoded in DNA; (—), no editing, U encoded in DNA; blank space, editing not yet determined/no information available; italic text: unique for *Lolium perenne*.

<sup>a</sup>Partially edited.

correctly as editing sites and not accidentally mistaken as SNPs or vice versa.

**Acknowledgements:** We thank Dr Jiri Ködding and Dr Gavin Conant for all their support while sequencing, assembling and annotating the chloroplast genome. K.H.W. is supported by Science Foundation Ireland. K.D. and R.vdB. were financed under the Teagasc Walsh Fellowship Scheme.

**Supplementary data:** Supplementary data are available at [www.dnaresearch.oxfordjournals.org](http://www.dnaresearch.oxfordjournals.org).

## Funding

The project funding was obtained from the Teagasc 'Vision' programme.

## References

1. Gray, M. W. 2004, The evolutionary origins of plant organelles, In: Daniell and Chase, H. (eds.), *Molecular Biology and Biotechnology of Plant Organelles*, Springer: Dordrecht, pp. 87–108.

2. Keeling, P. J. 2004, Diversity and evolutionary history of plastids and their hosts, *Am. J. Bot.*, **91**, 1481–1493.
3. Margulis, L. 1970, Origin of eukaryotic cells, Yale University Press: New Haven.
4. Lilly, J. W., Havey, M. J., Jackson, S. A. and Jiang, J. 2001, Cytogenomic analyses reveal the structural plasticity of the chloroplast genome in higher plants, *Plant Cell*, **13**, 245–254.
5. Sugiura, M. 1992, The chloroplast genome, *Plant Mol. Biol.*, **19**, 149–168.
6. Palmer, J. D. 1987, Chloroplast DNA evolution and bio-systematic uses of chloroplast DNA variation, *Am. Nat.*, **130**, S6.
7. Corriveau, J. L. and Coleman, A. W. 1988, Rapid screening method to detect potential biparental inheritance of plastid DNA and results for over 200 angiosperm species, *Am. J. Bot.*, **75**, 1443–1458.
8. Säll, T., Jakobsson, M., Lind-Halldén, C. and Halldén, C. 2003, Chloroplast DNA indicates a single origin of the allotetraploid *Arabidopsis suecica*, *J. Evol. Biol.*, **16**, 1019–1029.
9. Daniell, H., Datta, R., Varma, S., Gray, S. and Lee, S. B. 1998, Containment of herbicide resistance through genetic engineering of the chloroplast genome, *Nat. Biotechnol.*, **16**, 345–348.
10. Wu, F. -H., Kan, D. -P., Lee, S. -B., et al. 2009, Complete nucleotide sequence of *Dendrocalamus latiflorus* and *Bambusa oldhamii* chloroplast genomes, *Tree Physiol.*, doi:10.1093/treephys/tpp015.
11. Hodkinson, T. R., Waldren, S., Parnell, J. A. N., Kelleher, C. T., Salamin, K. and Salamin, N. 2007, DNA banking for plant breeding, biotechnology and biodiversity evaluation, *J. Plant Res.*, **120**, 17–29.
12. Sato, S., Nakamura, Y., Kaneko, T., Asamizu, E. and Tabata, S. 1999, Complete structure of the chloroplast genome of *Arabidopsis thaliana*, *DNA Res.*, **6**, 283–90.
13. Yukawa, M., Tsudzuki, T. and Sugiura, M. 2006, The chloroplast genome of *Nicotiana sylvestris* and *Nicotiana tomentosiformis*: complete sequencing confirms that the *Nicotiana sylvestris* progenitor is the maternal genome donor of *Nicotiana tabacum*, *Mol. Genet. Genomics*, **275**, 367–373.
14. Jansen, R. K., Raubeson, L. A., Boore, J. L., et al. 2005, Methods for obtaining and analyzing whole chloroplast genome sequences, *Methods Enzymol.*, **395**, 348–384.
15. Chung, S. M., Gordon, V. S. and Staub, J. E. 2007, Sequencing cucumber (*Cucumis sativus* L.) chloroplast genomes identifies differences between chilling-tolerant and-susceptible cucumber lines, *Genome*, **50**, 215–225.
16. McGrath, S., Hodkinson, T. R. and Barth, S. 2007, Extremely high cytoplasmic diversity in natural and breeding populations of *Lolium* (Poaceae), *Heredity*, **99**, 531–544.
17. Bock, R. 2000, Sense from nonsense: how the genetic information of chloroplasts is altered by RNA editing, *Biochimie*, **82**, 549–557.
18. Tsudzuki, T., Wakasugi, T. and Sugiura, M. 2001, Comparative analysis of RNA editing sites in higher plant chloroplasts, *J. Mol. Evol.*, **53**, 327–332.
19. Hoch, B., Maier, R. M., Appel, K., Igloi, G. L. and Kössel, H. 1991, Editing of a chloroplast mRNA by creation of an initiation codon, *Nature*, **353**, 178–180.
20. Diekmann, K., Hodkinson, T. R., Fricke, E. and Barth, S. 2008, An optimized chloroplast DNA extraction protocol for grasses (Poaceae) proves suitable for whole plastid genome sequencing and SNP detection, *PLoS ONE*, **3**, e2813.
21. Sasaki, C., Lee, S. B., Fjellheim, S., et al. 2007, Complete chloroplast genome sequences of *Hordeum vulgare*, *Sorghum bicolor* and *Agrostis stolonifera*, and comparative analyses with other grass genomes, *Theor. Appl. Genet.*, **115**, 591.
22. Sugita, M. and Sugiura, M. 1996, Regulation of gene expression in chloroplasts of higher plants, *Plant Mol. Biol.*, **32**, 315–326.
23. Conant, G. C. and Wolfe, K. H. 2008, GenomeVx: simple web-based creation of editable circular chromosome maps, *Bioinformatics*, **24**, 861–862.
24. Calsa Júnior, T., Carraro, D. M., Benatti, M. R., Barbosa, A. C., Kitajima, J. P. and Carrer, H. 2004, Structural features and transcript-editing analysis of sugarcane (*Saccharum officinarum* L.) chloroplast genome, *Curr. Genet.*, **46**, 366–373.
25. Corneille, S., Lutz, K. and Maliga, P. 2000, Conservation of RNA editing between rice and maize plastids: are most editing events dispensable, *Mol. Gen. Genet.*, **264**, 419–424.
26. Freyer, R., Hoch, B., Neckermann, K., Maier, R. M. and Kössel, H. 1993, RNA editing in maize chloroplasts is a processing step independent of splicing and cleavage to monocistronic mRNAs, *Plant J.*, **4**, 621–629.
27. Maier, R. M., Neckermann, K., Igloi, G. L. and Kössel, H. 1995, Complete sequence of the maize chloroplast genome: gene content, hotspots of divergence and fine tuning of genetic information by transcript editing, *J. Mol. Biol.*, **251**, 614–628.
28. Zeltz, P., Hess, W. R., Neckermann, K., Börner, T. and Kössel, H. 1993, Editing of the chloroplast *rpoB* transcript is independent of chloroplast translation and shows different patterns in barley and maize, *EMBO J.*, **12**, 4291–4296.
29. Kumar, S., Tamura, K. and Nei, M. 2004, MEGA3: integrated software for molecular evolutionary genetics analysis and sequence alignment, *Brief. Bioinform.*, **5**, 150–163.
30. Ogihara, Y., Isono, K., Kojima, T., et al. 2002, Structural features of a wheat plastome as revealed by complete sequencing of chloroplast DNA, *Mol. Genet. Genomics*, **266**, 740–746.
31. Bortiri, E., Coleman-Derr, D., Lazo, G. R., Anderson, O. D. and Gu, Y. Q. 2008, The complete chloroplast genome sequence of *Brachypodium distachyon*: sequence comparison and phylogenetic analysis of eight grass plastomes, *BMC Res. Notes*, **31**, 61.
32. Palmer, J. D. 1991, Plastid chromosomes: structure and evolution, In: Vasil, I. K. and Bogorad, L. (eds.), *Cell Culture and Somatic Cell Genetics of Plants, vol 7A, The Molecular Biology of Plastids*, Academic Press: San Diego, pp. 5–53.

33. Weglöhner, W. and Subramanian, A. R. 1991, A heptapeptide repeat contributes to the unusual length of chloroplast ribosomal protein S18. Nucleotide sequence and map position of the *rpl33-rps18* gene cluster in maize, *FEBS Lett.*, **279**, 193–197.
34. Igloi, G. L., Meinke, A., Döry, I. and Kössel, H. 1990, Nucleotide sequence of the maize chloroplast *rpo B/C1/C2* operon: comparison between the derived protein primary structures from various organisms with respect to functional domains, *Mol. Gen. Genet.*, **221**, 379–394.
35. Shimada, H., Fukuta, M., Ishikawa, M. and Sugiura, M. 1990, Rice chloroplast RNA polymerase genes: the absence of an intron in *rpoC1* and the presence of an extra sequence in *rpoC*, *Mol. Gen. Genet.*, **221**, 395–402.
36. Cummings, M. P., King, L. M. and Kellogg, E. A. 1994, Slipped-strand mispairing in a plastid gene: *rpoC2* in grasses (Poaceae), *Mol. Biol. Evol.*, **11**, 1–8.
37. Chen, Z., Muthukrishnan, S., Liang, G. H., Schertz, K. F. and Hart, G. E. 1993, A chloroplast DNA deletion located in RNA polymerase gene *rpoC2* in CMS lines of sorghum, *Mol. Gen. Genet.*, **236**, 251–259.
38. Chen, Z., Schertz, K. F., Mullet, J. E., DuBell, A. and Hart, G. E. 1995, Characterization and expression of *rpoC2* in CMS and fertile lines of sorghum, *Plant Mol. Biol.*, **28**, 799–809.
39. Kim, K. J. and Lee, H. L. 2004, Complete chloroplast genome sequences from *Korean Ginseng (Panax schin-seng* Nees) and comparative analysis of sequence evolution among 17 vascular plants, *DNA Res.*, **11**, 247–261.
40. Kim, K. J. and Lee, H. L. 2005, Widespread occurrence of small inversions in the chloroplast genomes of land plants, *Mol. Cells*, **19**, 104–113.
41. Kelchner, S. A. and Wendel, J. F. 1996, Hairpins create minute inversions in non-coding regions of chloroplast DNA, *Curr. Genet*, **30**, 259–262.
42. Tsumura, Y., Suyama, Y. and Yoshimura, K. 2000, Chloroplast DNA inversion polymorphism in populations of *Abies* and *Tsuga*, *Mol. Biol. Evol.*, **17**, 1302–1312.
43. McGrath, S., Hodkinson, T. R., Salamin, N. and Barth, S. 2006, Development and testing of novel chloroplast microsatellite markers for *Lolium perenne* and other grasses (Poaceae) from de novo sequencing and *in silico* sequences, *Mol. Ecol. Notes*, **6**, 449–452(4).
44. Daniell, H., Lee, S. B., Grevich, J., et al. 2006, Complete chloroplast genome sequences of *Solanum bulbocastanum*, *Solanum lycopersicum* and comparative analyses with other Solanaceae genomes, *Theor. Appl. Genet.*, **112**, 1503–1518.
45. Timme, R. E., Kuehl, J. V., Boore, J. L. and Jansen, R. K. 2007, A comparative analysis of the *Lactuca* and *Helianthus* (Asteraceae) plastid genomes: identification of divergent regions and categorization of shared repeats, *Am. J. Bot.*, **94**, 302–312.
46. Vogel, J., Hübschmann, T., Börner, T. and Hess, W. R. 1997, Splicing and intron-internal RNA editing of *trnK-matK* transcripts in barley plastids: support for *matK* as an essential splice factor, *J. Mol. Biol.*, **270**, 179–187.
47. Lopez, C., Freyer, R., Guera, A., et al. 1997, Sequence of *ndhA* gene of barley (*Hordeum vulgare* L.) plastid (accession nos. Y13729 & Y13730). Transcript editing in Graminean organs, *Plant Physiol.*, **115**, 313.
48. del Campo, E. M., Sabater, B. and Martin, M. 1997, Plastid *ndhD* gene of barley (*Hordeum vulgare* L.), sequence and transcript editing (Accession no. Y12258), *Plant Physiol.*, **114**, 748.
49. Drescher, A., Hupfer, H., Nickel, C., et al. 2002, C-to-U conversion in the intergenic *ndhI/ndhG* RNA of plastids from monocot plants: conventional editing in an unconventional small reading frame, *Mol. Genet. Genomics*, **267**, 262–269.
50. Inada, M., Sasaki, T., Yukawa, M., Tsudzuki, T. and Sugiura, M. 2004, A systematic search for RNA editing sites in pea chloroplasts: an editing event causes diversification from the evolutionarily conserved amino acid sequence, *Plant Cell Physiol.*, **45**, 1615–1622.
51. Ruf, S. and Kössel, H. 1997, Tissue-specific and differential editing of the two *ycf3* editing sites in maize plastids, *Curr. Genet.*, **32**, 19–23.