



OPEN

Molecular structure, comparative and phylogenetic analysis of the complete chloroplast genome sequences of weedy rye *Secale cereale* ssp. *segetale*

Lidia Skuza^{1,2}✉, Piotr Androsiuk³, Romain Gastineau⁴, Łukasz Paukszto³, Jan Paweł Jastrzębski³ & Danuta Cembrowska-Lech^{5,6}

The complete chloroplast genome of *Secale cereale* ssp. *segetale* (Zhuk.) Roshev. (Poaceae: Triticeae) was sequenced and analyzed to better use its genetic resources to enrich rye and wheat breeding. The study was carried out using the following methods: DNA extraction, sequencing, assembly and annotation, comparison with other complete chloroplast genomes of the five *Secale* species, and multigene phylogeny. As a result of the study, it was determined that the chloroplast genome is 137,042 base pair (bp) long and contains 137 genes, including 113 unique genes and 24 genes which are duplicated in the IRs. Moreover, a total of 29 SSRs were detected in the *Secale cereale* ssp. *segetale* chloroplast genome. The phylogenetic analysis showed that *Secale cereale* ssp. *segetale* appeared to share the highest degree of similarity with *S. cereale* and *S. strictum*. Intraspecific diversity has been observed between the published chloroplast genome sequences of *S. cereale* ssp. *segetale*. The genome can be accessed on GenBank with the accession number (OL688773).

Secale cereale ssp. *segetale* is one of the many species of the genus *Secale* with a previously unknown chloroplast and mitochondrial genome. However, it can be a source of desired genes (e.g., resistance to diseases, high protein content, morphological and biochemical traits) that can enrich rye or wheat breeding^{1,2}. The lack of knowledge of phylogenetic relationships reduces the progress in rye breeding, which can be enriched with functional features derived from wild rye species³. With new biotic and abiotic stresses and climate change, there is also a need to study wild rye species, which is crucial to improving the yield and quality of this cereal⁴. Therefore, more genetic markers are needed. One of the ways to achieve this is to sequence complete chloroplast genomes. Due to their conservative and non-recombinant nature, chloroplast genomes are a solid tool in genomics and evolutionary research⁵. Certain evolutionary hotspots of the plant plastid genome, such as single nucleotide polymorphisms and insertions/deletions, may provide useful information to elucidate the phylogeny of taxonomically unresolved plant taxa^{6,7}. Thus, the availability of complete chloroplast genomes, which include new variable and informational sites, should help explain more precise phylogeny.

To participate in this effort, we have undertaken the sequencing of the complete chloroplast genomes in genus *Secale*, which are smaller and easier to analyze compared to mitochondrial genomes. So far, only the incomplete *S. cereale* cpDNA sequences (NC_021761)⁸, three sequences for *S. strictum* (KY636137, KY636138 and OL979486)⁹ and *S. sylvestre* (MW557517)¹⁰ are available. The chloroplast genome of *S. segetale* has recently been published¹¹, however a comprehensive phylogenetic analysis based on whole chloroplast genomes has not been done to date. Therefore, we presume that analysis of the complete chloroplast genome sequences of *Secale*

¹Institute of Biology, University of Szczecin, 71415 Szczecin, Poland. ²Centre for Molecular Biology and Biotechnology, Institute of Biology, University of Szczecin, 71415 Szczecin, Poland. ³Department of Plant Physiology, Genetics and Biotechnology, Faculty of Biology and Biotechnology, University of Warmia and Mazury, 10719 Olsztyn, Poland. ⁴Institute of Marine and Environmental Sciences, University of Szczecin, 70383 Szczecin, Poland. ⁵Department of Physiology and Biochemistry, Institute of Biology, University of Szczecin, Felczaka 3c St., 71412 Szczecin, Poland. ⁶Sanprobi Sp. z o. o. Sp. k., Kurza Stopka 5c St., 70535 Szczecin, Poland. ✉email: lidia.skuza@usz.edu.pl

spp., starting with *S. sylvestre*¹⁰, will be useful and cost-effective for evolutionary and phylogenetic studies, as it was suggested by our previous studies¹².

In this study, we present the complete chloroplast genome of *S. cereale* ssp. *segetale*, which will provide valuable information for genetic studies of *Secale* species.

Results

Chloroplast genome of *Secale cereale* ssp. *segetale*. Sequencing of *Secale cereale* ssp. *segetale* chloroplast genome yielded 41 653 350 raw reads, out of which 88 777 were mapped to the reference genome of *S. cereale* with 97× average coverage. The *S. cereale* ssp. *segetale* cp genome appeared as a typical circular, double-stranded molecule with the length of 137,042 bp (Fig. 1) and overall GC content of 38%. The large single copy (LSC) region is 81,060 bp long, the short single copy (SSC) region is 12,820 bp long, and each of the inverted repeat regions (IR) is 21,581 bp long. Reported cp genome contains 137 genes, including 113 unique genes and 24 genes which are duplicated in the IRs. Group of 113 unique genes features 73 protein-coding genes, 30 tRNA genes, four rRNA genes and five conserved chloroplast open reading frames (ORFs) (Table 1).

The LSC region appeared as the most abundant in genes—57 PCGs, 21 tRNA genes and two ORFs (*ycf3* and *ycf4*), whereas there are only ten PCGs and one tRNA gene in SSC. In IR there are four rRNA genes, eight tRNA genes, three ORFs (*ycf2*, *ycf15* and *ycf68*) and nine PCGs including *ndhH* located on the junction between IR and SSC region.

Repeat sequence analysis. A total of 52 repeat sequences structures with length ranging from 30 to 286 bp were revealed in the plastome of *Secale cereale* ssp. *segetale* (Table 2). The forward repeats (37) dominated over palindromic (15) repeats. Neither complementary nor reverse repeats were found. Most repeat sequences

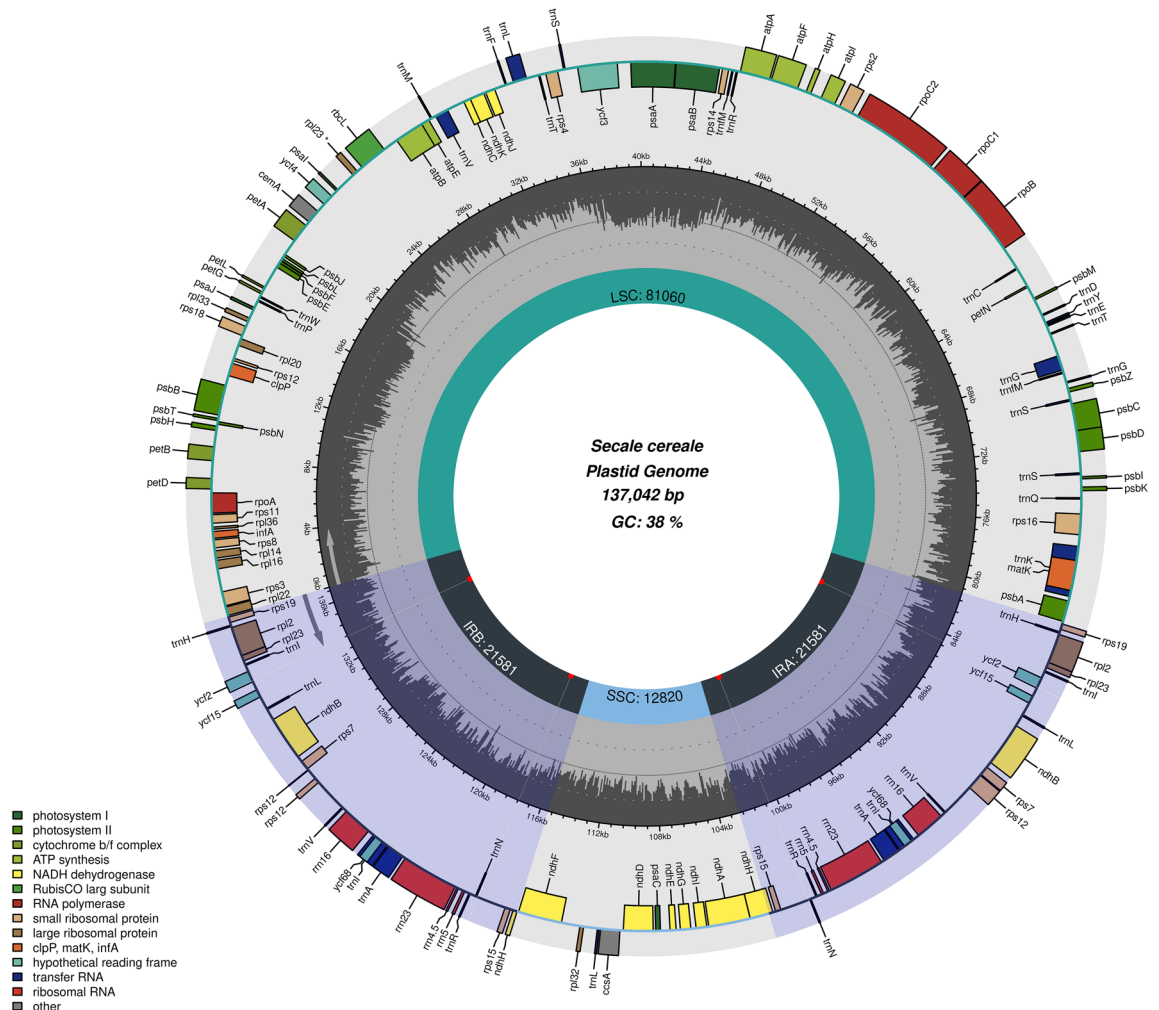


Figure 1. Map of the chloroplast genome of *Secale cereale* ssp. *segetale*. The genes inside and outside the circle are transcribed in the clockwise and counterclockwise directions, respectively. Genes belonging to different functional groups are shown in different colors. Tick lines indicate the extent of the inverted repeats (IRa and IRb) that separate the genomes into small single-copy (SSC) and large single-copy (LSC) regions. The innermost darker gray corresponds to GC-content while the lighter gray corresponds to AT content.

Category	Group of gene	Name of genes
Photosynthesis	Photosystem I	<i>psaA, psaB, psaC, psaI, psaJ</i>
	Photosystem II	<i>psbA, psbB, psbC, psbD, psbE, psbF, psbH, psbI, psbJ, psbK, psbL, psbM, psbN, psbT, psbZ</i>
	Cytochrome complex	<i>petA, petB, petD, petG, petL, petN</i>
	ATP synthase	<i>atpA, atpB, atpE, atpF, atpH, atpI</i>
	NADH dehydrogenase	<i>ndhA^c, ndhB^c (×2), ndhC, ndhD, ndhE, ndhF, ndhG, ndhH^b (2x), ndhI, ndhJ, ndhK</i>
	Large subunit of RUBISCO	<i>rbcL</i>
DNA replication and protein synthesis	Ribosomal RNA	<i>rrn4.5 (×2), rrn5 (×2), rrn16 (×2), rrn23 (×2)</i>
	Small subunit ribosomal proteins	<i>rps2, rps3, rps4, rps7 (×2), rps8, rps11, rps12^c, rps14, rps15(×2), rps16^c, rps18, rps19 (×2)</i>
	Large subunit ribosomal proteins	<i>rpl2^c (×2), rpl14, rpl16, rpl20, rpl22, rpl23^b (×3), rpl32, rpl33, rpl36</i>
	RNA polymerase subunits	<i>rpoA, rpoB, rpoC1, rpoC2</i>
	Translational initiation factor	<i>infA</i>
	Transfer RNA	<i>trnA-UGC^c (×2), trnC-GCA, trnD-GUC, trnE-UUC, trnF-GAA, trnJ^c-CAU, trnM-AUG, trnG-UCC^c (×2), trnH-GUG (×2), trnI-CAU (×2), trnI-GAU^f (×2), trnK-UUU^f, trnL-CAA (×2), trnL-CUA, trnL-UAA^c, trnM-CAU, trnN-GUU (×2), trnP-UGG, trnQ-UUG, trnR-ACG (×2), trnR-UCU, trnS-GCU, trnS-GGA, trnS-UGA, trnT-GGU, trnT-UGU, trnV-GAC (×2), trnV-UAC^c, trnW-CCA, trnY-GUA</i>
Other genes	Conserved hypothetical chloroplast ORF	<i>ycf2 (×2), ycf3^{ad}, ycf4^a, ycf15 (×2), ycf68 (×2)</i>
	Other proteins	<i>ccsA, cema, clpP, matK</i>

Table 1. Genes present in chloroplast genome of *Secale cereale* ssp. *segetale*. Genes list arranged alphabetically.

^aGenes associated with Photosystem I. ^bOne copy of the gene is a pseudogene. ^cGene containing one intron.

^dGene containing two introns. ^eTransspliced gene.

(69.3%) were detected in the LSC region, followed by IR (28.8%) and SSC regions (1.9%). 50% of these sequences were found within coding regions. The highest number of repeats were found within the sequences of the following genes: *rpoC2* (9F), *rpl23* (2F and 2P) and *rps18* (3F and 1P).

A total of 29 SSRs were detected in the *Secale cereale* ssp. *segetale* chloroplast genome (Table 3). The mononucleotide SSRs composed of A/T units were the most common, whereas hexanucleotide SSRs were not detected. 79.3% of SSRs were located within LSC region, 13.8% in IR region while only 6.9% of SSRs were found in SSC region. Most of the SSRs were identified within intergenic spacers (58.6%), while equal proportions (20.7%) were located in the introns and coding sequences.

Multigene phylogeny. Phylogeny reconstruction based on sequences of 73 protein-coding genes shared by *Secale cereale* ssp. *segetale* and 38 representatives of Pooideae subfamily appeared to be consistent with the systematic position of studied species. The BI and ML tree divided analyzed species into six major clades (Fig. 2). The first cluster contained 23 species representing Triticinae subtribe, four other clades gathered 13 species representing Hordeinae subtribe, whereas the last clad consisted of three *Littledalea* species (Littledaleeae tribe). *Secale cereale* ssp. *segetale* appeared to share the highest degree of similarity with *S. cereale* and *S. cereale* ssp. *segetale*. Mentioned above five *Secale* species form separate sub-clad within the Triticinae tribe.

Comparison with other complete chloroplast genomes of the *Secale* species. The overall sequence identity of five cp genomes of *Secale* species was plotted using mVISTA with the annotation of *S. cereale* ssp. *segetale* cp genome (obtained by new sequencing in this study) as reference (Fig. 3). The results showed that the *Secale* cp genomes exhibited a high level of sequence synteny, suggesting a conserved evolutionary pattern. The plastome sequences were fairly conserved across the four data with a few regions with a variation. The sequences of exons were nearly identical throughout the all taxa.

Discussion

The task of modern cereal breeding is to obtain new, higher-yielding varieties that have high resistance to pathogens, diseases and abiotic conditions. Unfortunately, progress in rye breeding has been limited, as the varieties used in cultivation have had limited variability due to selection. In addition, attempts to use old varieties have been unsuccessful.

A major advance in rye breeding has been the introduction of hybrid varieties, through which individual genotypes are fixed by continuing self-pollination and transferring monogenic traits into varieties¹³. However, despite the increase in yield, intermediate quality traits are subject to large annual fluctuations. Thus, despite significant increases in grain yield and decreases in protein content in the experiments, increases in grain yield did not significantly positively or negatively affect intermediate quality traits⁴.

A number of taxa in the genus *Secale* may represent a potential source of genetic variability in rye breeding³. Species such as *Secale strictum* and *Secale vavilovii* may be sources of new genetic variability, with resistance to ear fusariosis and septoria leaf blotch), while *Secale vavilovii* may also be a source of sterilizing cytoplasm (source of sterilising cytoplasm). Wild rye species and subspecies provide excellent starting material for studies

Repeat length	Strat site of repeat A	Location	Repeat A region	Strat site of repeat B	Location	Repeat B region	Repeat type
286	56561	<i>rpl23</i>	LSC	83149	<i>rpl23</i>	IR	P
286	56561	<i>rpl23</i>	LSC	134665	<i>rpl23</i>	IR	F
160	56687	<i>rpl23</i>	LSC	83149	<i>rpl23</i>	IR	P
160	56687	<i>rpl23</i>	LSC	134791	<i>rpl23</i>	IR	F
74	12628	IGS (<i>trnG-UCC-trnfM-CAU</i>)	LSC	12796	IGS (<i>trnfM-CAU-trnG-UCC</i>)	LSC	F
60	101806	IGS (<i>trnN-GUU-rps15</i>)	IR	101806	IGS (<i>trnN-GUU-rps15</i>)	IR	P
60	101806	IGS (<i>trnN-GUU-rps15</i>)	IR	116234	IGS (<i>trnN-GUU-rps15</i>)	IR	F
60	116234	IGS (<i>trnN-GUU-rps15</i>)	IR	116234	IGS (<i>trnN-GUU-rps15</i>)	IR	P
45	56364	IGS (<i>rbcL-rpl23</i>)	LSC	56364	IGS (<i>rbcL-rpl23</i>)	LSC	P
42	41556	IGS (<i>psaA-ycf3</i>)	LSC	41570	IGS (<i>psaA-ycf3</i>)	LSC	F
41	12735	<i>trnfM-CAU</i>	LSC	36344	<i>trnfM-CAU</i>	LSC	F
40	12628	IGS (<i>trnG-UCC-trnfM-CAU</i>)	LSC	36407	IGS (<i>trnfM-CAU-rps14</i>)	LSC	F
40	12796	IGS (<i>trnfM-CAU-trnG-UCC</i>)	LSC	36407	IGS (<i>trnfM-CAU-rps14</i>)	LSC	F
39	12835	IGS (<i>trnfM-CAU-trnG-UCC</i>)	LSC	36447	IGS (<i>trnfM-CAU-rps14</i>)	LSC	F
39	14437	IGS (<i>trnG-UCC-trnT-GGU</i>)	LSC	89975	IGS (<i>rps7-trnV-GAC</i>)	IR	F
39	14437	IGS (<i>trnG-UCC-trnT-GGU</i>)	LSC	128086	IGS (<i>rps7-trnV-GAC</i>)	IR	P
38	38373	<i>psaB</i>	LSC	40597	<i>psaA</i>	LSC	F
36	7542	<i>trnS-GCU</i>	LSC	44850	<i>trnS-GGA</i>	LSC	P
36	43333	I intron <i>ycf3</i>	LSC	90638	IGS (<i>rps7-trnV-GAC</i>)	IR	F
36	43333	I intron <i>ycf3</i>	LSC	127426	IGS (<i>rps7-trnV-GAC</i>)	IR	P
35	12667	IGS (<i>trnG-UCC-trnfM-CAU</i>)	LSC	36447	IGS (<i>trnfM-CAU-rps14</i>)	LSC	F
35	12719	<i>trnfM-CAU</i>	LSC	36328	<i>trnfM-CAU</i>	LSC	F
35	76754	<i>infA</i>	LSC	76772	<i>infA</i>	LSC	F
34	27191	<i>rpoC2</i>	LSC	27212	<i>rpoC2</i>	LSC	F
34	27253	<i>rpoC2</i>	LSC	27328	<i>rpoC2</i>	LSC	F
33	27113	<i>rpoC2</i>	LSC	27164	<i>rpoC2</i>	LSC	F
33	61501	IGS (<i>petA-psbJ</i>)	LSC	61501	IGS (<i>rbcL-rpl23</i>)	LSC	P
32	11271	<i>trnS-UGA</i>	LSC	44857	<i>trnS-GGA</i>	LSC	P
32	12677	IGS (<i>trnG-UCC-trnfM-CAU</i>)	LSC	36457	IGS (<i>trnfM-CAU-rps14</i>)	LSC	F
32	14794	<i>trnT-GGU</i>	LSC	46100	<i>trnT-UGU</i>	LSC	P
32	27072	<i>rpoC2</i>	LSC	27171	<i>rpoC2</i>	LSC	F
32	27219	<i>rpoC2</i>	LSC	27273	<i>rpoC2</i>	LSC	F
32	41556	IGS (<i>psaA-ycf3</i>)	LSC	41584	IGS (<i>psaA-ycf3</i>)	LSC	F
31	8407	IGS (<i>trnS-GCU-psbD</i>)	LSC	8444	IGS (<i>trnS-GCU-psbD</i>)	LSC	F
31	12843	IGS (<i>trnfM-CAU-trnG-UCC</i>)	LSC	36455	IGS (<i>trnfM-CAU-rps14</i>)	LSC	F
31	15484	IGS (<i>trnY-GUA-trnD-GUC</i>)	LSC	33828	intron <i>atpF</i>	LSC	F
31	27054	<i>rpoC2</i>	LSC	27324	<i>rpoC2</i>	LSC	F
31	27064	<i>rpoC2</i>	LSC	27259	<i>rpoC2</i>	LSC	F
31	27241	<i>rpoC2</i>	LSC	27382	<i>rpoC2</i>	LSC	F
31	27316	<i>rpoC2</i>	LSC	27382	<i>rpoC2</i>	LSC	F
31	66279	<i>rps18</i>	LSC	66300	<i>rps18</i>	LSC	F
31	80266	<i>rps3</i>	LSC	80311	<i>rps3</i>	LSC	F
31	101837	IGS (<i>trnN-GUU-rps15</i>)	IR	116265	IGS (<i>trnN-GUU-rps15</i>)	IR	F
31	105588	IGS (<i>ndhF-rpl32</i>)	SSC	105612	IGS (<i>ndhF-rpl32</i>)	SSC	F
30	12870	<i>trnG-UCC</i>	LSC	36314	<i>trnfM-CAU</i>	LSC	F
30	16756	IGS (<i>psbM-petN</i>)	LSC	16756	IGS (<i>psbM-petN</i>)	LSC	P
30	66231	<i>rps18</i>	LSC	66315	<i>rps18</i>	LSC	F
30	66298	<i>rps18</i>	LSC	66319	<i>rps18</i>	LSC	F
30	66677	<i>rps18</i>	LSC	66677	<i>rps18</i>	LSC	P
30	87638	Intron <i>ndhB</i>	IR	87638	Intron <i>ndhB</i>	IR	P
30	87638	Intron <i>ndhB</i>	IR	130432	Intron <i>ndhB</i>	IR	F
30	130432	Intron <i>ndhB</i>	IR	130432	Intron <i>ndhB</i>	IR	P

Table 2. List of repeated sequences in the chloroplast genome of *Secale cereale* ssp. *segetale*. IGS (*trnG-UCC-trnfM-CAU*) means spacer between *trnG-UCC* and *trnfM-CAU*.

Type	Repeat unit	Length	Start	End	Location	Region
Mononucleotide	A	13	7211	7223	IGS (<i>psbK-psbJ</i>)	LSC
		13	7937	7949	IGS (<i>trnS-GCU-psbD</i>)	LSC
		18	11227	11244	IGS (<i>psbC-trnS-UGA</i>)	LSC
		12	29538	29549	<i>rpoC2</i>	LSC
		12	29945	29956	IGS (<i>rpoC2-rps2</i>)	LSC
		12	33569	33580	intron <i>atpF</i>	LSC
		12	33844	33855	intron <i>atpF</i>	LSC
		13	36192	36204	IGS (<i>trnR-UCU-trnM-CAU</i>)	LSC
		12	43027	43038	II intron <i>ycf3</i>	LSC
		13	46728	46740	IGS (<i>trnT-UGU-trnL-UAA</i>)	LSC
		12	76716	76727	IGS (<i>rpl36-infA</i>)	LSC
		12	105202	105213	IGS (<i>ndhF-rpl32</i>)	SSC
Trinucleotide	AAT	15	24652	24666	<i>rpoC1</i>	LSC
	AAT	13	47511	47523	IGS (<i>trnL-UAA-trnF-GAA</i>)	LSC
	AAC	12	31552	31563	<i>atpI</i>	LSC
	AAT	12	56494	56505	IGS (<i>rbcL-rpl23</i>)	LSC
	AAG	12	65670	65681	IGS (<i>psaJ-rpl33</i>)	LSC
Tetranucleotide	AAGG	12	42927	42938	II intron <i>ycf3</i>	LSC
	AATG	12	64604	64615	IGS (<i>trnW-CCA-trnP-UGG</i>)	LSC
	AAAG	12	64889	64900	IGS (<i>trnP-UGG-psaJ</i>)	LSC
	AAAG	12	68899	68910	IGS (<i>clpP-psbB</i>)	LSC
	AACG	13	99471	99483	4.5S rRNA	IR
	AAAT	12	108269	108280	<i>ndhD</i>	SSC
	AACG	13	118618	118630	4.5S rRNA	IR
Pentanucleotide	AATAT	18	15656	15673	IGS (<i>trnY-GUA-trnD-GUC</i>)	LSC
	ACCAT	15	43805	43819	I intron <i>ycf3</i>	LSC
	AATAT	18	47154	47171	intron <i>trnL-UAA</i>	LSC
	AATAT	16	101098	101113	IGS (<i>trnN-GUU-rps15</i>)	IR
	AATAT	16	116988	117003	IGS (<i>trnN-GUU-rps15</i>)	IR

Table 3. Distribution of SSR in the *Secale cereale* ssp. *segetale* cp genome. IGS IGS (*psbK-psbJ*) means spacer between *psbK* and *psbJ*.

aimed at expanding recombination variability in cultivated rye and triticale (\times Triticosecale Wittmark). Because of their genetic distinctiveness and high trait expression, they represent a valuable source of genes in which our cultivars are deficient¹⁴. An example is the study of the efficiency of crossing the wild species *Secale vavilovii* and the rye subspecies *Secale cereale* ssp. *afghanicum*, *Secale cereale* ssp. *ancestrale*, *Secale cereale* ssp. *dighoricum*, *Secale cereale* ssp. *segetale* with the crop species *Secale cereale* ssp. *cereale*, and the resulting F1 crosses may be a potential source of variation in common rye³. Unfortunately, the lack of knowledge of phylogenetic relationships reduces the progress in rye breeding.

For understanding plant origin and evolution chloroplast genome sequences are very useful. With maternally inherited traits, a genome of relatively small size and a slow mutation rate of the genome¹⁵, analysis of the phylogenetic relationships of multiple chloroplast DNA can help understand plant phylogeny, population genetic analysis and taxonomic status at the molecular level¹⁶.

Although cp genomes of angiosperm plants are generally conservative in terms of sequence and number of genes¹⁷, levels of structural variation have been observed in the genome that vary across families and genera, such as gene duplication and large-scale rearrangements of genes, introns and IR domains (e.g. ^{18,19}).

The *S. cereale* ssp. *segetale* cp genome appeared as a typical circular, double-stranded molecule (Fig. 1) and overall GC content, which is similar to previously sequenced plastomes of *S. cereale* (137,051 bp; NCBI LC645358), *S. sylvestre* (137 116 bp)¹⁰ or within the size range of angiosperms²⁰.

The results obtained by Du et al.¹¹ are similar to ours. The size of the genome, the lengths of the LSC, SSC and IR sequences differ slightly. In contrast, larger differences are seen in the number of genes. The genome we analyzed contains 73 protein-coding genes (82 in¹¹), 30 tRNA genes (41 in¹¹) four rRNA genes (8 in¹¹) and five conserved chloroplast open reading frames (ORFs) (lack of information in¹¹).

It is difficult to say where the above-mentioned differences came from. The rich interspecific genetic diversity of *S. S. cereale* ssp. *segetale* has been previously reported (e.g.²¹). Significant differences were found between and within populations of *S. c. ssp. segetale*. A high degree of genetic variability has also been described using chromosomal markers^{22,23}. These results deserve attention and further research.

The polymorphisms found in *S. c. ssp. segetale* chloroplast genome sequences can be used e.g. to elucidate evolutionary histories such as the origin of *Secale* species or accessions at the inter- and, thanks to the research

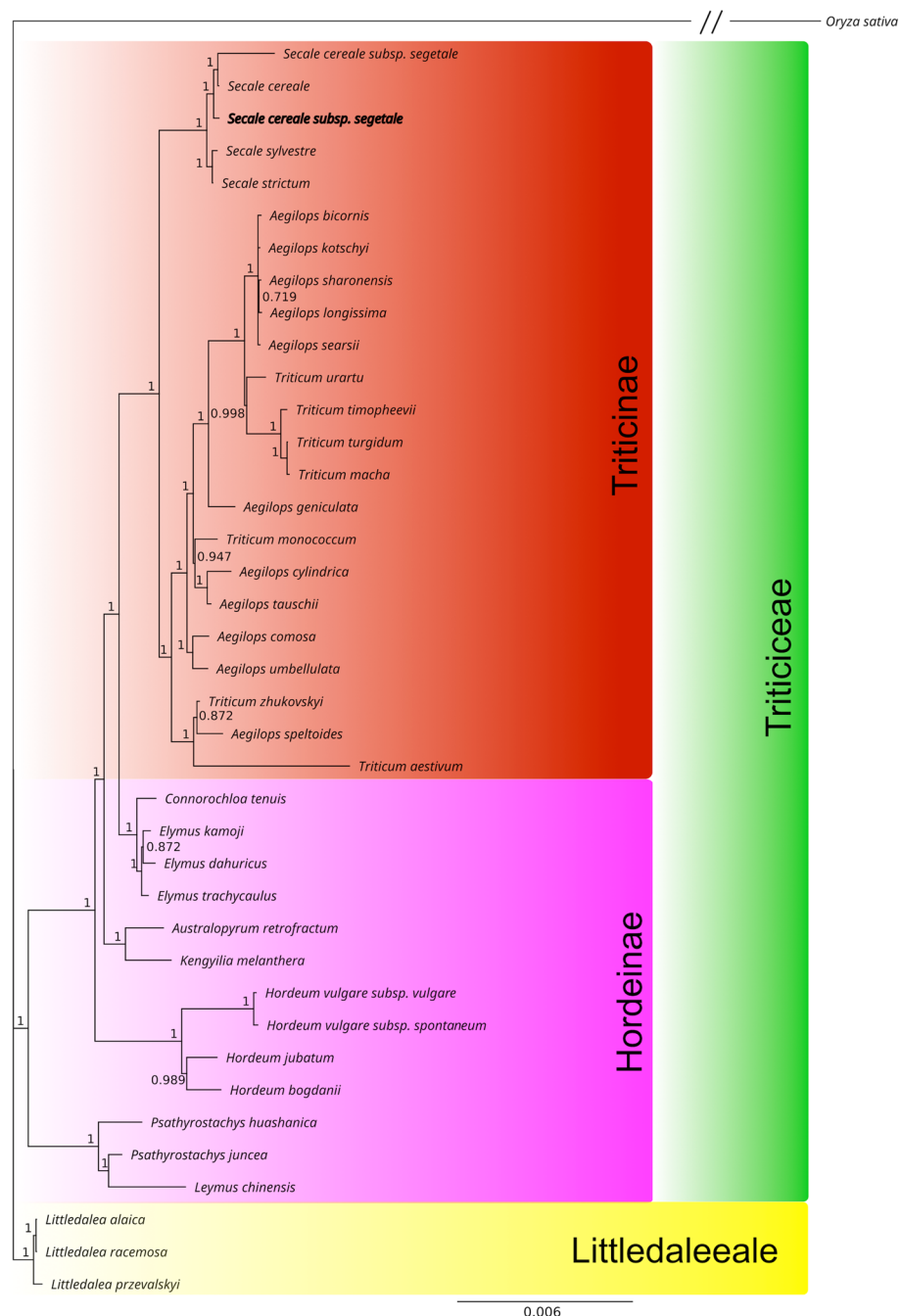


Figure 2. Cladogram illustrating the phylogenetic relationships for *Secale cereale* ssp. *segetale* based on complete cp genome sequences. Phylogenetic tree based on sequences of sheared 73 protein-coding genes from five *Secale* species and 34 other cereal lineages representing Triticoideae group within subfamily Pooideae and the cp genome of *Oryza sativa* as an outgroup, using Bayesian posterior probabilities (PP) and maximum likelihood (ML). Each node has 100% bootstrap support value. The cpDNA sequence obtained in this study is shown in bold.

described in this manuscript, intra-species level. Furthermore, the polymorphic sites promote practical applications for molecular analysis to protect *S. c. ssp. segetale* accession²⁴ and, potentially in the long term, the rye breeding industry. Unfortunately, the analyses of the genome previously published by Du et al. do not include many details, in addition to those mentioned above, which does not allow for a more detailed analysis.

Certain regions of the plastome are predisposed to indel and substitution mutations. Comparative studies of the plastome show the evolution of, among other, tandem repeats and their role in generating substitutions and indels^{25,26}. Once the composition of repeat sequences in the plastome is determined, it is possible to predict microstructural changes by analyzing the correlation between repeats, indels and substitutions. In addition to the

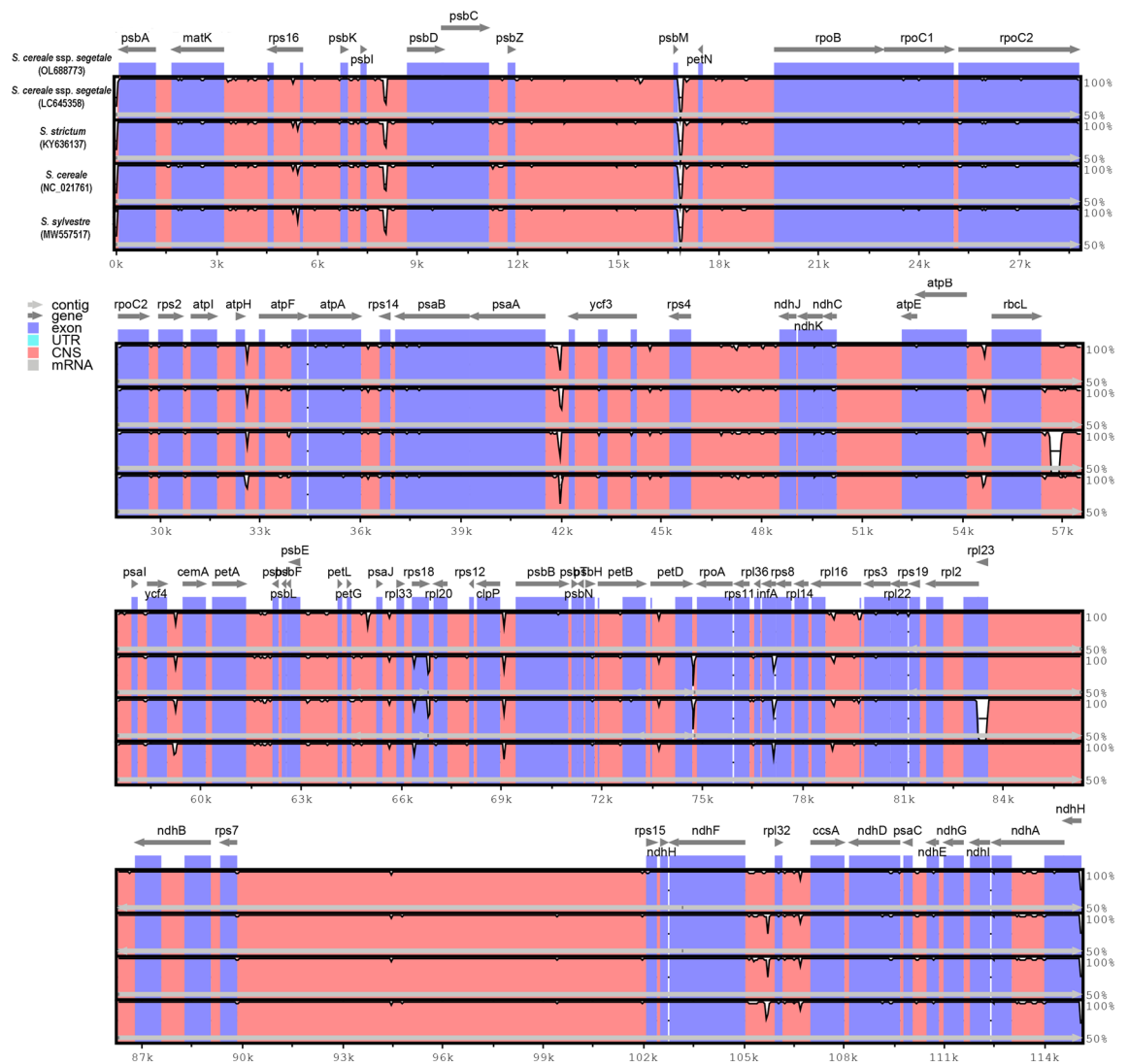


Figure 3. Percentage of sequence identity between chloroplast genomes of *Secale cereale* ssp. *segetale* and other four *Secale* species using mVISTA program. Gray arrows on the top line show transcriptional direction. The y-axis represents average percent identity between sequences of *S. cereale* ssp. *segetale* and other three *Secale* chloroplast genomes. The x-axis represents the coordinate in the chloroplast genome using *S. cereale* ssp. *segetale* as reference. Genome regions are color coded as exon, untranslated regions (UTR), and conserved non-coding sequences (CNS).

paucity of genomic resources, the phylogeny of the genus *Secale* is enigmatic (e.g.^{27,28}). Therefore, it is important to fully explore the polymorphic regions of *Secale* chloroplasts in an evolutionary context.

For the total of 52 repeat sequence structures revealed in the *Secale cereale* ssp. *segetale* plastome, the vast majority were detected in the LSC region (Table 2). The highest number of repeats was found within the sequences of the *rpoC2*, *rpl23* and *rps18* genes. Regardless of its function the *rpoC2* gene encoding the β -subunit of plastid RNA polymerase is a relatively rapidly evolving chloroplast sequence²⁹. Analogically, *rpl23* gene and its pseudogene which are observed in the grass family belong to highly polymorphic genes considered as hotspots of illegitimate recombination in cp genomes³⁰.

Chloroplast SSRs identification not only serves as a one of cp genome characteristics but also represent ideal molecular tools with various applications like investigation of domestication history, sites of origin or genetic diversity and relationships between wild and cultivated species^{31,32}. In 2016, Hagenblad et al.³³ analyzed the genetic diversity of 76 accessions of wild, feral and cultivated rye based on SNP polymorphisms. They performed an analysis of five chloroplast SSRs, derived from *Lolium* and wheat. Discriminant analysis of principal components (DAPC) of cpSSR data indicated very large genetic variation within the genus *Secale* and did not reflect taxonomic groups, except for *S. strictum* and *S. africanum*, which formed a separate cluster.

CpSSRs are mainly distributed within intergenic spacers of *Secale* plastomes; similar distribution preferences of cpSSRs have been reported in *Avena* spp., *Pseudoroegneria libanotica* and *Salvia miltiorrhiza*^{34–36}.

Phylogenetic analysis has shown that *Secale cereale* ssp. *segetale* appeared to share the highest degree of similarity with *S. cereale* and *S. cereale* ssp. *segetale*. The five *Secale* species form separate sub-clad within the Triticeae tribe, which confirms previous phylogenetic data of the genus *Secale* (e.g.³⁷).

The results showed that the *Secale* cp genomes exhibited a high level of sequence synteny, suggesting a conserved evolutionary pattern. The plastome sequences were fairly conserved across the four data with a few regions with a variation. The sequences of exons were nearly identical throughout the all taxa.

Conclusions

Here we assembled the complete, annotated chloroplast genome sequence of *Secale cereale* ssp. *segetale*. The genome is 137 042 base pair (bp) long and contains 137 genes, including 113 unique genes and 24 genes which are duplicated in the IRs. The phylogenetic analysis showed that *Secale cereale* ssp. *segetale* appeared to share the highest degree of similarity with *S. cereale* and *S. strictum*. Intraspecific diversity has been observed between the published chloroplast genome sequences of *S. cereale* ssp. *segetale*. The cp genome will provide a series of resources for evolutionary and genetic studies about species of rye. The assembled genome sequences and annotation information have been deposited in GenBank under the accession number OL688773.

Material and methods

DNA extraction, sequencing, assembly and annotation. Seeds of *Secale cereale* ssp. *segetale* introd. no. 1782/94 were obtained from the Botanical Garden of the Polish Academy of Sciences in Warsaw. Total DNA was extracted from young sprouts following Doyle and Doyle³⁸.

The chloroplast (cp) genome of *Secale cereale* ssp. *segetale* was sequenced with the use of DNBseq platform in BGI Shenzhen (China). After the quality check (FastQC tool available online at: <http://www.bioinformatics.babraham.ac.uk/projects/fastqc>) the raw reads were mapped to the reference genome of *Secale cereale* (NC_021761) in Geneious v.R7 software with default medium–low sensitivity settings³⁹. Reads aligned to the reference cpDNA genome were extracted and used for de novo assembly (K-mer—23–41, low coverage cut-off—5, minimum contig length—300). De novo contigs were extended by mapping raw reads to the generated contigs, reassembling the contigs with mapped reads, and manually scaffolding the extended contigs (minimum sequence overlap of 50 bp and 97% overlap identity). This process was iterated five times. Finally, the reduced sequences were assembled in the circular chloroplast genome. The chloroplast genome was annotated using MFannot⁴⁰ and PlasMapper⁴¹ with manual adjustments. The gene map of the annotated cp genome was developed with the OrganellarGenome DRAW tool⁴².

Repeat sequence analysis. The chloroplast simple sequence repeats (SSRs) were detected using Phobos v.3.3.12⁴³. Only perfect SSRs with a motif size of one to six nucleotide units were considered, the following thresholds for chloroplast SSRs identification were used: ≥ 12 repeat units for mononucleotide SSRs, ≥ 6 repeat units for dinucleotide SSRs, ≥ 4 repeat units for trinucleotide SSRs, and ≥ 3 repeat units for tetra-, penta- and hexanucleotide SSRs⁴⁴. Analysis of long genomic repeats, i.e. forward (F), reverse (R), palindromic (P) and complementary (C) sequences, was performed using REPuter software⁴⁵ with the following settings: (1) hamming distance of 3, (2) sequence identity $\geq 90\%$, and (3) minimum repeat size ≥ 30 bp. A single IR region was used to eliminate the influence of doubled IR regions.

Multigene phylogeny. The phylogenetic position of *Secale cereale* ssp. *segetale* within Triticeae group was also evaluated. For that purpose 73 concatenated protein-coding gene sequences shared with other 38 Poaceae species were used. The cpDNA of *Oryza sativa* was used as an outgroup (Table 4). For phylogeny reconstruction Bayesian Inference (BI) method was used. The best-fit model of sequence evolution was identified in MEGA v.7⁴⁶, and the GTR+G+I model was selected. The BI analysis was performed in MrBayes v.3.2.6⁴⁷. Parameter settings were previously described by Androsiuk et al.⁴⁸.

For multigene phylogeny maximum likelihood (ML) analyses was conducted using RAxML-NG⁴⁹ under three different strategies. (1) One of the IR regions was removed from all chloroplast genomes to reduce overrepresentation of duplicated sequences then we run RAxML-NG on the unpartitioned alignment under GTR+I+G substitution model as a single partition; (2) The same data was partitioned by gene, exon, intron and intergenic spacer regions and allowed separate base frequencies, α -shape parameters, and evolutionary rates to be estimated for each; (3) we inferred the best-fitting partitioning strategy with PartitionFinder2⁵⁰ for the alignment. The best fitting nucleotide substitution models were inferred with jModelTest2⁵¹. Phylogenetic trees were visualized and edited with FigTree 1.4.4⁵². Support for the ML tree branches was calculated by the non-parametric bootstrap method with 1000 replicates.

Comparison with other complete chloroplast genomes of the *Secale* species. The percentage of sequence identity among complete chloroplast genomes of the five *Secale*: *S. cereale* ssp. *segetale* (OL688773), *S. cereale* ssp. *segetale* (LC645358), *S. cereale* (NC_021761), *S. strictum* (KY636137), and *S. sylvestris* (MW557517) was comparatively analyzed and plotted using the program mVISTA⁵³, with alignment algorithm of LAGAN⁵⁴, a cut-off of 70% identity, and annotation of *S. cereale* ssp. *segetale* (OL688773) as reference.

Ethics approval and consent to participate. Authors confirm that the use of plants in the present study complies with international, national and/or institutional guidelines.

Species	Accession number
<i>Oryza sativa</i>	NC_008155
<i>Aegilops bicornis</i>	NC_024831
<i>Aegilops comosa</i>	NC_046697
<i>Aegilops cylindrica</i>	NC_023096
<i>Aegilops geniculata</i>	NC_023097
<i>Aegilops kotschy</i>	NC_024832
<i>Aegilops longissima</i>	NC_024830
<i>Aegilops searsii</i>	NC_024815
<i>Aegilops sharonensis</i>	NC_024816
<i>Aegilops speltoides</i>	NC_022135
<i>Aegilops tauschii</i>	NC_022133
<i>Aegilops umbellulata</i>	NC_046696
<i>Australopyrum retrofractum</i>	NC_043840
<i>Connorochloa tenuis</i>	NC_037165
<i>Elymus dahuricus</i>	NC_049159
<i>Elymus kamoji</i>	NC_051511
<i>Elymus trachycaulus</i>	NC_050404
<i>Hordeum bogdanii</i>	NC_043839
<i>Hordeum jubatum</i>	NC_027476
<i>Hordeum vulgare</i> subsp. <i>Spontaneum</i>	NC_042692
<i>Hordeum vulgare</i> subsp. <i>vulgare</i>	NC_008590
<i>Kengyilia melanthera</i>	NC_042706
<i>Leymus chinensis</i>	NC_044900
<i>Littledalea alaica</i>	NC_037519
<i>Littledalea przewalskyi</i>	NC_037497
<i>Littledalea racemosa</i>	NC_036350
<i>Psathyrostachys huashanica</i>	NC_045871
<i>Psathyrostachys juncea</i>	NC_043838
<i>Secale cereale</i>	NC_021761
<i>Secale cereale</i> subsp. <i>segetale</i>	LC645358
<i>Secale cereale</i> subsp. <i>segetale</i>	LC645358
<i>Secale strictum</i>	KY636137
<i>Secale sylvestre</i>	MW557517
<i>Triticum aestivum</i>	NC_002762
<i>Triticum macha</i>	NC_025955
<i>Triticum monococcum</i>	NC_021760
<i>Triticum timopheevii</i>	NC_024764
<i>Triticum turgidum</i>	NC_024814
<i>Triticum urartu</i>	KJ614411
<i>Triticum zhukovskyi</i>	NC_046698

Table 4. List of species used in phylogenetic studies. Species names arranged alphabetically.

Data availability

The genome can be accessed on GenBank with the accession number (OL688773).

Received: 30 November 2022; Accepted: 29 March 2023

Published online: 03 April 2023

References

1. Kubicka, H., Puchalski, J., Niedzielski, M., Łuczak, W. & Martyniszyn, A. Collection and evaluation of rye gene resources (in Polish). *Bull. Plant Breed. Accl. Inst.* **40**(241), 141–149 (2006).
2. Schittenhelm, S., Kraft, M. & Wittich, K. P. Performance of winter cereals grown on field-stored soil moisture only. *Eur. J. Agron.* **52**(B), 247–258 (2014).
3. Mikołajczyk, S., Broda, Z., Mackiewicz, D., Weigt, D. & Bocianowski, J. Biometric characteristics of interspecific hybrids in the genus *Secale*. *Biometric. Lett.* **51**(2), 153–170 (2014).
4. Laidig, F., Piepho, H. P., Rentel, D. & Huesken, A. Breeding progress, variation, and correlation of grain and quality traits in winter rye hybrid and population varieties and national on-farm progress in Germany over 26 years. *Theor. Appl. Genet.* **130**, 981–998 (2017).

5. Daniell, H., Lin, C.-S., Yu, M. & Chang, W.-J. Chloroplast genomes: Diversity, evolution, and applications in genetic engineering. *Genome Biol.* **17**, 134 (2016).
6. Eguiluz, M., Rodrigues, N. F., Guzman, F., Yuyama, P. & Margis, R. The chloroplast genome sequence from *Eugenia uniflora*, a Myrtaceae from Neotropics. *Plant Syst. Evol.* **303**, 1199–1212 (2017).
7. Ruhfel, B. R., Gitzendanner, M. A., Soltis, P. S., Soltis, D. E. & Burleigh, J. G. From algae to angiosperms—inferring the phylogeny of green plants (Viridiplantae) from 360 plastid genomes. *BMC Evol. Biol.* **14**, 23 (2014).
8. Middleton, C. P. *et al.* Sequencing of chloroplast genomes from wheat, barley, rye and their relatives provides a detailed insight into the evolution of the Triticeae tribe. *PLoS ONE* **9**(3), E85761 (2014).
9. Bernhardt, N., Brassac, J., Kilian, B. & Blattner, F. R. Dated tribe-wide whole chloroplast genome phylogeny indicates recurrent hybridizations within Triticeae. *BMC Evol. Biol.* **17**(1), 141 (2017).
10. Skuza, L., Gastineau, R. & Sielska, A. The complete chloroplast genome of *Secale sylvestre* (Poaceae: Triticeae). *J. Appl. Genet.* **63**, 115–117 (2022).
11. Du, T., Hu, Y., Sun, Y., Ye, C. & Shen, E. The complete chloroplast genome of weedy rye *Secale cereale* subsp. *segetale*. *Mitochondrial DNA B Resour.* **7**(6), 959–960 (2022).
12. Skuza, L., Szucko, I., Filip, E. & Strzała, T. Genetic diversity and relationship between cultivated, weedy and wild rye species as revealed by chloroplast and mitochondrial DNA non-coding regions analysis. *PLoS ONE* **14**(2), e0213023 (2019).
13. Miedaner, T. & Huebner, M. Quality demands for different uses of hybrid rye. in *Tagung der Vereinigung der Pflanzenzüchter und Saatgutkaufleute Österreichs 2010*. Vol. 61. 45–49 (2011).
14. Rzepka-Plevneš, D. Utility properties of hybrids *S. cereale* × *S. vavilovii* Gross. in terms of their suitability in growing rye varieties resistant to sprouting. Part I. *Bull. Plant Breed. Accl. Inst.* **37**, 69–79 (1993).
15. Palmer, J. D., Jansen, R. K., Michaels, H. J., Chase, M. W. & Manhart, J. R. Chloroplast DNA variation and plant phylogeny. *Ann. Missouri Bot. Garden* **75**, 1180–1206 (1988).
16. Alwadani, K. G., Janes, J. K. & Andrew, R. L. Chloroplast genome analysis of box-ironbark *Eucalyptus*. *Mol. Phylogenet. Evol.* **136**, 76–86 (2019).
17. Jansen, R.K., & Ruhlmann, T.A. *Plastid Genomes of Seed Plants, Genomics of Chloroplasts, and Mitochondria*. 103–126. https://doi.org/10.1007/978-94-007-2920-9_5 (Springer, 2012)
18. Guisinger, M. M., Chumley, T. W., Kuehl, J. V., Boore, J. L. & Jansen, R. K. Implications of the plastid genome sequence of *Typha* (Typhaceae, Poales) for understanding genome evolution in Poaceae. *J. Mol. Evol.* **70**, 149–166. <https://doi.org/10.1007/s00239-009-9317-3> (2010).
19. Martin, G. E. *et al.* The first complete chloroplast genome of the Genistoid legume *Lupinus luteus*: Evidence for a novel major lineage-specific rearrangement and new insights regarding plastome evolution in the legume family. *Ann. Bot.* **113**, 1197–1210. <https://doi.org/10.1093/aob/mcu050> (2014).
20. Dong, W., Xu, C., Cheng, T., Lin, K. & Zhou, S. Sequencing angiosperm plastid genomes made easy: A complete set of universal primers and a case study on the phylogeny of Saxifragales. *Genome Biol. Evol.* **5**, 989–997 (2013).
21. Che, Y. H. *et al.* Genetic diversity of *Secale cereale* subsp. *segetale* populations in Xinjiang. *J. Triticeae Crops* **28**, 755–758 (2008).
22. Yang, X. M. *et al.* Cytology and disease resistance identification of *Secale cereale* subsp. *segetale* in Xinjiang of China. *Xinjiang Agric. Sci.* **3**, 117–120 (1994).
23. Dai, M. *et al.* Karyotypes analysis of *Secale cereale* subsp. *segetale*. *J. Triticeae Crops* **33**, 440–444 (2013).
24. Che, Y. *et al.* Genetic diversity of gliadin in *Secale cereale* subsp. *segetale* from Xinjiang, China. *Genet. Resour. Crop Evol.* **63**, 1173–1179 (2016).
25. Abdullah, M. F. *et al.* Correlations among oligonucleotide repeats, nucleotide substitutions and insertion-deletion mutations in chloroplast genomes of plant family Malvaceae. *J. Syst. Evol.* **59**(2), 388–402 (2020).
26. Henriquez, C. L. *et al.* Evolutionary dynamics in chloroplast genomes of subfamily Aroideae (Araceae). *Genomics* **112**, 2349–2360 (2020).
27. Chikmawati, T., Skovmand, B. & Gustafson, J. P. Phylogenetic relationships among *Secale* species revealed by amplified fragment length polymorphisms. *Genome* **48**(5), 792–801 (2005).
28. Maraci, Ö., Özkan, H. & Bilgin, R. Phylogeny and genetic structure in the genus *Secale*. *PLoS ONE* **19**(7), e0200825 (2018).
29. Logacheva, M. D., Penin, A. A., Samigullin, T. H., Vallejo-Roman, C. M. & Antonov, A. S. Phylogeny of flowering plants by the chloroplast genome sequences: in search of a “lucky gene”. *Biochem. Mosc.* **72**, 1324–1330 (2007).
30. Lencina, F. *et al.* The *rpl23* gene and pseudogene are hotspots of illegitimate recombination in barley chloroplast mutator seedlings. *Sci. Rep.* **9**, 9960 (2019).
31. Provan, J., Powell, W. & Hollingsworth, P. M. Chloroplast microsatellites: new tools for studies in plant ecology and evolution. *Trends Ecol. Evol.* **16**, 142–147 (2001).
32. Delplancke, M. *et al.* Gene flow among wild and domesticated almond species: insights from chloroplast and nuclear markers. *Evol. Appl.* **5**, 317–329 (2012).
33. Hagenblad, J., Oliveira, H. R., Forsberg, N. E. & Leino, M. W. Geographical distribution of genetic diversity in *Secale* landrace and wild accessions. *BMC Plant Biol.* **16**, 23 (2016).
34. Liu, Q. *et al.* Comparative chloroplast genome analyses of *Avena*: Insights into evolutionary dynamics and phylogeny. *BMC Plant Biol.* **20**, 406 (2020).
35. Wu, D. D. *et al.* The complete chloroplast genome sequence of *Pseudoroegneria libanotica*, genomic features, and phylogenetic relationship with Triticeae species. *Biol. Plantarum* **62**(2), 231–240 (2018).
36. Qian, J. *et al.* The complete chloroplast genome sequence of the medicinal plant *Salvia miltiorrhiza*. *PLoS ONE* **8**(2), e57607 (2013).
37. Schreiber, M., Himmelbach, A., Borner, A. & Mascher, M. Genetic diversity and relationship between domesticated rye and its wild relatives as revealed through genotyping-by-sequencing. *Evol. Appl.* **12**(1), 66–77 (2019).
38. Doyle, J. J. & Doyle, J. L. Isolation of plant DNA from fresh tissue. *Focus (Madison)* **12**, 13–15 (1990).
39. Kearse, M. *et al.* Geneious basic: An integrated and extendable desktop software platform for the organization and analysis of sequence data. *Bioinformatics* **28**, 1647–1649 (2012).
40. MFannot. *The Robert Cedergren Centre of the Université de Montréal, France.* https://megasun.bch.umontreal.ca/cgi-bin/dev_mfa/mfannotInterface.pl.
41. Dong, X., Stothard, P., Forsythe, I. J. & Wishart, D. S. PlasMapper: A web server for drawing and auto-annotating plasmid maps. *Nucleic Acids Res.* **32**, W660–W664 (2004).
42. Lohse, M., Drechsel, O. & Bock, R. OrganellarGenomeDRAW (OGDRAW): A tool for the easy generation of high-quality custom graphical maps of plastid and mitochondrial genomes. *Curr. Genet.* **52**, 267–274 (2007).
43. Phobos v.3.3.12. *Dr. Christoph Mayer, Ruhr-Universität, Bochum.* http://www.rub.de/ecoevo/cm/cm_phobos.htm. (2010).
44. Sablok, G. *et al.* ChloroMitoSSRDDB 2.00: More genomes, more repeats, unifying SSRs search patterns and on-the-fly repeat detection. *Database* **2015**, bav084 (2015).
45. Kurtz, S. *et al.* REPuter: The manifold applications of repeat analysis on a genomic scale. *Nucleic Acids Res.* **29**, 4633–4642 (2001).
46. Kumar, S., Stecher, G. & Tamura, K. MEGA7: Molecular evolutionary genetics analysis version 7.0 for bigger datasets. *Mol. Biol. Evol.* **33**, 1870–1874 (2016).
47. Ronquist, F. & Huelsenbeck, J. P. MrBayes 3: Bayesian phylogenetic inference under mixed models. *Bioinformatics* **19**, 1572–1574 (2003).

48. Androsiuk, P. *et al.* The complete chloroplast genome of *Colobanthus apetalus* (Labill.) Druce: Genome organization and comparison with related species. *PeerJ* **6**, e4723 (2018).
49. Kozlov, A. M., Darrriba, D., Flouri, T., Morel, B. & Stamatakis, A. RAxML-NG: A fast, scalable and user-friendly tool for maximum likelihood phylogenetic inference. *Bioinformatics* **35**, 4453–4455 (2019).
50. Lanfear, R., Frandsen, P. B., Wright, A. M., Senfeld, T. & Calcott, B. PartitionFinder 2: New methods for selecting partitioned models of evolution for molecular and morphological phylogenetic analyses. *Mol. Biol. Evol.* **34**, 772–773 (2017).
51. Darrriba, D., Taboada, G. L., Doallo, R. & Posada, D. jModelTest 2: More models, new heuristics and parallel computing. *Nat. Methods* **9**, 772 (2012).
52. Rambaut, A. *FigTree v. 1.4.4*. <http://tree.bio.ed.ac.uk/software/figtree/> (2014).
53. Frazer, K. A., Pachter, L., Poliakov, A., Rubin, E. M. & Dubchak, I. VISTA: computational tools for comparative genomics. *Nucleic Acids Res.* **32**, W273–W279 (2004).
54. Brudno, M. *et al.* LAGAN and Multi-LAGAN: Efficient tools for large-scale multiple alignment of genomic DNA. *Genome Res.* **13**, 721–731 (2003).

Author contributions

L.S. study conception and design, L.S. and R.G. conducted experiments, P.A. and L.S. drafted the manuscript, bioinformatic analyses were performed by P.A., Ł.P., J.P.J. and D.C.-L.

Funding

This work was supported by a grant from the Institute of Biology, University of Szczecin, Poland.

Competing interests

The authors declare no competing interests.

Additional information

Correspondence and requests for materials should be addressed to L.S.

Reprints and permissions information is available at www.nature.com/reprints.

Publisher's note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.



Open Access This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2023