

1 **Haplotype-resolved *de novo* genome assemblies of four coniferous tree species**

2

3 Kenta Shirasawa^{1#}, Kentaro Mishima^{2#}, Hideki Hirakawa¹, Tomonori Hirao³, Miyoko Tsubomura³,
4 Soichiro Nagano³, Taiichi Iki², Sachiko Isobe¹, and Makoto Takahashi³

5

6 ¹Kazusa DNA Research Institute, Kisarazu, Chiba 292-0818, Japan

7 ²Tohoku Regional Breeding Office, Forest Tree Breeding Center, Forestry and Forest Products Research
8 Institute, Forest Research and Management Organization, Takizawa, Iwate 020-0621, Japan

9 ³Forest Tree Breeding Center, Forestry and Forest Products Research Institute, Forest Research and
10 Management Organization, Hitachi, Ibaraki 319-1301, Japan

11

12 #These two authors contributed equally to this work.

13 *Correspondence: shirasaw@kazusa.or.jp

14

15 **Abstract**

16 Coniferous trees in gymnosperm are an important source of wood production. Because of their long
17 lifecycle, the breeding programs of coniferous tree are time- and labor-consuming. Genomics could
18 accelerate the selection of superior trees or clones in the breeding programs; however, the genomes of
19 coniferous trees are generally giant in size and exhibit high heterozygosity. Therefore, the generation of
20 long contiguous genome assemblies of coniferous species has been difficult. In this study, we optimized
21 the DNA library preparation protocols and employed high-fidelity (HiFi) long-read sequencing technology
22 to sequence and assemble the genomes of four coniferous tree species, *Larix kaempferi*, *Chamaecyparis*
23 *obtusa*, *Cryptomeria japonica*, and *Cunninghamia lanceolata*. Genome assemblies of the four species
24 totaled 13.5 Gb (*L. kaempferi*), 8.5 Gb (*C. obtusa*), 9.2 Gb (*C. japonica*), and 11.7 Gb (*C. lanceolata*),
25 which covered 99.6% of the estimated genome sizes on average. The contig N50 value, which indicates
26 assembly contiguity, ranged from 1.2 Mb in *C. obtusa* to 16.0 Mb in *L. kaempferi*, and the assembled
27 sequences contained, on average, 89.2% of the single-copy orthologs conserved in embryophytes.
28 Assembled sequences representing alternative haplotypes covered 70.3–95.1% of the genomes, suggesting
29 that the four coniferous tree genomes exhibit high heterozygosity levels. The genome sequence information
30 obtained in this study represents a milestone in tree genetics and genomics, and will facilitate gene discovery,
31 allele mining, phylogenetics, and evolutionary studies in coniferous trees, and accelerate forest tree

1 breeding programs.

2

3 **Keywords:** Breeding, coniferous trees, genome, gymnosperm, long-read sequencing

4

5 **Introduction**

6 Forests cover approximately 31% of the global land area (FAO 2020). In Japan, forest area is as high as
7 more than 60%, of which approximately 40% is accounted for by artificial forests, mainly coniferous trees,
8 which have been planted for applications in the forestry industry (Forestry Agency 2022). *Cryptomeria*
9 *japonica*, a Cupressaceae family member endemic to Japan, is the most important forestry species in the
10 country, occupying >40% of the artificial forests (Forestry Agency 2022). *Chamaecyparis obtusa*, another
11 member of the Cupressaceae family, is also an important conifer species with superior-quality wood that
12 has been used for the construction of buildings in Japan since the ancient times (Tsumura et al. 2007).
13 *Cunninghamia lanceolata* (Cupressaceae) was introduced from China and Taiwan, and has been recognized
14 as a new woody resource in Japan in recent years because of its fast growth and desirable wood properties
15 (Fujisawa 2017). Additionally, *Larix kaempferi* (Pinaceae) is an important breeding tree species in the
16 northern part of Japan (Kurinobu 2005). Timber tree germplasm and breeding programs have been used to
17 improve wood production and quality. However, the breeding programs of coniferous trees are generally
18 time-consuming because their lifecycle can be as long as or longer than 50 years. Furthermore, the highly
19 heterozygous genomes and outbreeding mating system of coniferous tree species complicate the breeding
20 systems (Burdon and Wilcox 2011).

21 Since the genomes of coniferous trees are often giant in size (>10 Gb), the analysis of the genome
22 sequence data of these species remains challenging (Neale and Wheeler 2019). Owing to the recent
23 advancements in long-read sequencing technologies, the genomes of gymnosperm species belonging to six
24 families, including Cupressaceae, Cycadaceae, Ginkgoaceae, Pinaceae, Taxaceae, and Welwitschiaceae,
25 have become available (Wan et al. 2022). The giant genomes of gymnosperm species are rich in repetitive
26 sequences such as transposable elements (Ohri 2021). The high-fidelity (HiFi) long-read sequencing
27 technology generates reads spanning the repetitive sequences, resulting in an assembly that covers most of
28 the gene spaces in the genome (Hon et al. 2020).

29 Coniferous tree genomics using the greatly advanced sequencing technologies promises the
30 acceleration of not only the breeding programs but also the phylogenetic study of coniferous species. In
31 this study, we determined the genome sequences of four coniferous tree species, including *L. kaempferi*, *C.*
32 *obtusa*, *C. japonica*, and *C. lanceolata*, using the HiFi long-read sequencing technology. The genome

1 sequence information obtained in this study could serve as a useful resource for breeding coniferous trees,
2 and for understanding the physiology of forest trees and the population genetics and phylogenetics of
3 gymnosperms.

4 **Materials and methods**

6 *Plant materials and DNA extraction*

7 Four coniferous tree species were used in this study: *Larix kaempferi* (GFE32203) was planted at the
8 Tohoku Regional Breeding Office, Forest Tree Breeding Center, Forestry and Forest Products Research
9 Institute, Forest Research and Management Organization in Iwate, Japan; *Chamaecyparis obtusa*
10 (GFB00119) and *Cryptomeria japonica* (GFA01029) were planted at the Forest Tree Breeding Center,
11 Forestry and Forest Products Research Institute, Forest Research and Management Organization in Ibaraki,
12 Japan; and *Cunninghamia lanceolata* (GFHN00090) was also planted at Forest Tree Breeding Center,
13 Forestry and Forest Products Research Institute, Forest Research and Management Organization in Ibaraki,
14 Japan. The original *C. lanceolata* tree was planted at the Kiyosumi Work Station in the University of Tokyo
15 Chiba Forest in Chiba, Japan.

16 Genomic DNA was extracted from the young leaves of each tree species using Genome-tips
17 (Qiagen, Hilden, Germany). DNA concentration was measured using the Qubit dsDNA BR assay kit
18 (Thermo Fisher Scientific, Waltham, MA, USA), and DNA fragment length was evaluated by agarose gel
19 electrophoresis with Pippin Pulse (Sage Science, Beverly, MA, USA).

21 *DNA library preparation and sequencing*

22 DNA libraries were prepared using six protocols (Table 1). Genomic DNA was sheared either by six
23 centrifugations at $1,600 \times g$ in the g-Tube (Covaris, Woburn, MA, USA) or in Megaruptor 2 (Deagenode,
24 Liege, Belgium) with the Large Fragment Hydropore mode and mean fragment sizes of 20, 30, or 40 kb.
25 Then, DNA library construction was performed with the SMRTbell Express Template Prep Kit 2.0 (PacBio),
26 according to the manufacturer's instructions. The obtained DNA libraries were fractionated with BluePippin
27 (Sage Science) to eliminate fragments shorter than 15, 20, or 25 kb in length. The fractionated DNA libraries
28 were sequenced on SMRT cells on the Sequel II and Sequel IIE system (PacBio). HiFi reads were
29 constructed with the CCS pipeline (<https://ccs.how>).

31 *Genome sequence assembly*

1 Genome sizes of the four tree species were estimated with Genomic Character Estimator (GCE) (Liu et al.
2 2013), based on the k -mer frequency ($k = 21$) calculated with Jellyfish version 2.3.0 (Marçais and Kingsford
3 2011). The reads were assembled using Hifiasm version 0.16.1 (Cheng et al. 2021) with default parameters.
4 Assembly completeness was evaluated using Benchmarking Universal Single-Copy Orthologs (BUSCO)
5 version 5.2.2, with default parameters (Simão et al. 2015), using the lineage dataset embryophyta_odb10
6 (eukaryota, 2020-09-10).

7

8 **Results**

9 *DNA sequencing*

10 We constructed the DNA libraries of *L. kaempferi*, *C. obtusa*, *C. japonica*, and *C. lanceolata* using six
11 library preparation protocols (Protocol1–Protocol6), with slight modification to the sheared DNA fragment
12 sizes and cutoff values to eliminate short DNA in libraries (Table 1). The resultant libraries were sequenced
13 on a total of 38 SMRT cells: 2 cells, Protocol1; 23 cells, Protocol2; 1 cell, Protocol3; 3 cells, Protocol4; 2
14 cells, Protocol5; and 7 cells, Protocol6.

15 The number, N50 value, and total length of HiFi reads per cell varied among the protocols (Figure
16 1). The total length of HiFi reads per cell ranged, on average, from 23.2 Gb in Protocol2 to 35.8 Gb in
17 Protocol6, while the read N50 value of these reads ranged, on average, from 16.3 kb in Protocol4 to 27.8
18 kb in Protocol6. The read N50 was expectedly long in long-insert libraries prepared using Protocol3 and
19 Protocol6.

20

21 *Genome assembly of L. kaempferi*

22 A total of 251.4 Gb HiFi reads (median N50 = 20.9 kb) were obtained from 11 SMRT cells (Figure 1). In
23 the k -mer distribution analysis, two peaks were detected at k -mer multiplicities of 10 and 19 (Figure 2); the
24 former and latter multiplicities corresponded to heterozygous regions of the diploid genome and
25 homozygous regions of haploid genome, respectively. The haploid genome size of *L. kaempferi* was
26 estimated to be 13.2 Gb (Figure 2). Accordingly, the genome coverage of the HiFi reads was calculated as
27 $19\times$ (= 251.4 Gb/13.2 Gb), which was sufficient for *de novo* genome assembly. The reads were assembled
28 with a haplotype-resolved *de novo* assembly method to obtain primary contigs (long continuous stretches
29 of contiguous sequences from one haplotype) and alternate contigs (contiguous sequences including
30 sequence and structural variants from another haplotype). The primary contigs included 4,655 sequences
31 spanning 13.5 Gb, with an N50 value of 16.0 Mb, which corresponded to the estimated size (Table 2). The

1 complete BUSCO score was 89.1% (single-copy BUSCO score = 77.5%, duplicated BUSCO score =
2 11.6%) (Table 2). On the other hand, the total length of alternate contigs was 9.5 Gb (N50 = 1.7 Mb) (Table
3 2). The ratio of the length of alternative contigs to that of primary contigs was 70.3% (= 9.5 Gb/13.5 Gb).

4 5 *Genome assembly of C. obtusa*

6 A total of 254.5 Gb HiFi reads (median N50 = 23.0 kb) were obtained from ten SMRT cells (Figure 1). Two
7 peaks were detected at *k*-mer multiplicities of 16 (diploid) and 32 (haploid), and the genome size of *C.*
8 *obtusa* was estimated to be 7.9 Gb (Figure 2). The genome coverage of HiFi reads was 30× (= 254.5 Gb/7.9
9 Gb). The HiFi reads were assembled into 15,466 primary contigs spanning 8.5 Gb, with an N50 value of
10 1.2 Mb (Table 2). The total contig length corresponded to the estimated genome size. The complete BUSCO
11 score was 89.4% (single-copy BUSCO = 82.3%, and duplicated BUSCO = 7.1%) (Table 2). Alternate contig
12 length was 7.8 Gb, with an N50 value of 0.5 Mb (Table 2). The ratio of alternate contig length to primary
13 contig length was 91.3% (= 7.8 Gb/8.5 Gb).

14 15 *Genome assembly of C. japonica*

16 A total of 237.6 Gb HiFi reads (median N50 = 17.3 kb) were obtained from nine SMRT cells (Figure 1).
17 Two peaks were detected at *k*-mer multiplicities of 12 (diploid) and 24 (haploid), and the genome size of *C.*
18 *japonica* was estimated to be 9.9 Gb (Figure 2). The genome coverage of HiFi reads was 24× (= 237.6
19 Gb/9.9 Gb). The HiFi reads were assembled into 2,741 primary contigs spanning 9.2 Gb, with an N50 value
20 of 8.3 Mb (Table 2). The total contig length covered 93.4% of the estimated genome size. The complete
21 BUSCO score was 89.0% (single-copy BUSCO = 82.3%, duplicated BUSCO = 6.7%) (Table 2). The
22 alternate contig length was 8.7 Gb, with an N50 value of 2.1 Mb (Table 2). The ratio of alternate contig
23 length to primary contig length was 94.3% (= 8.7 Gb/9.2 Gb).

24 25 *Genome assembly of C. lanceolata*

26 A total of 289.1 Gb HiFi reads (median N50 = 27.8 kb) were obtained from eight SMRT cells (Figure 1).
27 Two peaks were detected at *k*-mer multiplicities of 12 (diploid) and 24 (haploid), and the genome size of *C.*
28 *lanceolata* was estimated to be 12.0 Gb (Figure 2). The genome coverage of HiFi reads was 24× (= 289.1
29 Gb/12.0 Gb). The HiFi reads were assembled into 2,472 primary contigs spanning 11.5 Gb, with an N50
30 value of 11.7 Mb (Table 2). The total contig length covered 95.9% of the estimated genome size. The
31 complete BUSCO score was 89.1% (single-copy BUSCO = 80.8%, duplicated BUSCO = 8.3%) (Table 2).

1 The alternate contig length was 11.0 Gb, with an N50 value of 2.9 Mb (Table 2). The ratio of alternate
2 contig length to primary contig length was 95.1% (= 8.7 Gb/9.2 Gb).

3

4 **Discussion**

5 Here, we report the genome assemblies of four gymnosperm coniferous tree species, *L. kaempferi*, *C. obtusa*,
6 *C. japonica*, and *C. lanceolata* (Table 2). Since the genomes of these species were giant in size (~10 Gb),
7 similar to those of other gymnosperm species, we modified the DNA library preparation protocols to
8 maximize the data production efficiency. Among the six protocols (Table 1), Protocol6 was the most
9 effective with respect to the data yield (median 35.8 Gb) and read N50 length (median 27.8 kb), although
10 the number of reads obtained using Protocol6 was less than that obtained using Protocol4 and Protocol5
11 (Figure 1). In addition, the genome coverage of HiFi reads differed among the four species. In *C. lanceolata*,
12 HiFi read N50 length was the highest (27.8 kb; Figure 1), and genome coverage was the second highest
13 (24×) among the four species. On the other hand, the HiFi read N50 length of *L. kaempferi* (20.9 kb) was
14 lower than that of *C. lanceolata* (Figure 1), and genome coverage (19×) in *L. kaempferi* was the lowest
15 among the four species. However, unexpectedly, *L. kaempferi* showed the longest sequence contiguity,
16 followed by *C. lanceolata* (Table 2). The heterozygosity level of genomes might affect the assembled
17 sequence contiguity. The ratio of alternate contig length to primary contig length (70.3% in *L. kaempferi*
18 and 95.1% in *C. lanceolata*) might support this assumption. However, the *C. obtusa* assembly was
19 remarkably fragmented, even though its HiFi read N50 (23.0 kb) (Figure 1) and genome coverage (30×)
20 were comparable with those of *C. lanceolata*. This suggests that not only the heterozygosity level but also
21 other genome features, such as the length and/or distribution patterns of repetitive elements, affect the
22 contiguity of the assembled sequence.

23 In the phylogenetic tree of seed plants, gymnosperms occupy a basal position (Chase et al. 1993)
24 that branches out into the different clades of angiosperms. Therefore, genetic and genomic studies on
25 gymnosperms could shed light on the evolutionary history of plants. However, the genomic investigation
26 of gymnosperms has been lagging behind that of angiosperms because gymnosperms possess large-sized
27 genomes (≥ 10 Gb) and are relatively less commercially important than angiosperms (Wan et al. 2022). The
28 advent of sequencing technologies has steered this situation in favor of gymnosperm genomics. Long-read
29 sequencing technologies have enabled the sequencing and assembly of giant gymnosperm genomes rich in
30 repetitive sequences, which was not possible with short-read sequencing technologies (Wan et al. 2022).
31 The genome assemblies of four gymnosperm species generated in this study, in addition to those of other

1 species sequenced previously (Wan et al. 2022), will provide new insights into the genome evolution of
2 plants.

3 The four coniferous species sequenced in this study are important for the forestry. Genome
4 assemblies of these four coniferous tree species could be used as references for the identification of
5 sequence and structural variants in the genomes of divergent cultivars and breeding materials belonging to
6 the same species. Furthermore, in previous genetic and genomic studies, the utilization of a genome
7 sequence as a reference also enabled the identification of genes of interest (Neale and Kremer 2011). The
8 variant and gene information obtained in this study could be used for the development of DNA markers to
9 facilitate genetic studies and breeding programs (Muranty et al. 2014). Genome prediction might be another
10 powerful tool for the selection of elite tree lines from breeding programs (Lebedev et al. 2020; Grattapaglia
11 2022), which usually require a long time. In addition, gene editing technology could also be used as an
12 effective breeding strategy for shortening the duration of tree breeding programs (Bewg et al. 2018;
13 Goralogia et al. 2021).

14 The genome sequence information obtained in this study could contribute to breeding programs,
15 gene discovery, and allele mining in coniferous tree species with giant genomes. Since the coverage of
16 genome sequences was as high as ~90%, according to BUSCO evaluation, the assemblies could be used as
17 reference sequences in transcriptome analysis (Mishima et al. 2022). The protocols presented in this study
18 would contribute to and accelerate the genome sequence analysis of coniferous species with giant genomes.
19 Moreover, chromosome-level assemblies, which would enable phylogenetics and evolutionary studies
20 based on comparative genomics, could also be established through further genomic and genetic analyses in
21 the near future.

22

23 **Data availability**

24 Raw sequence reads were deposited in the Sequence Read Archive (SRA) database of the DNA Data Bank
25 of Japan (DDBJ) under the accession numbers DRA014993 (*L. kaempferi*), DRA014992 (*C. obtusa*),
26 DRA014994 (*C. japonica*), and DRA014995 (*C. lanceolata*). The assembled sequences are available at
27 DDBJ (accession numbers: BSBM01000001-BSBM01004655 [*L. kaempferi*]; BSBK01000001-
28 BSBK01015466 [*C. obtusa*]; BSBL01000001-BSBL01002741 [*C. japonica*]; and BSBN01000001-
29 BSBN01002472 [*C. lanceolata*]), BreedingTrees-by-Genes (<http://btg.kazusa.or.jp>), and Plant GARDEN
30 (<https://plantgarden.jp>).

31

1 **Acknowledgments**

2 We thank the University of Tokyo Chiba Forest for their support and cooperation during the collection of
3 *C. lanceolata* leaf samples. We also thank Y. Kishida, M. Kohara, C. Minami, K. Ozawa, H. Tsuruoka, and
4 A. Watanabe (Kazusa DNA Research Institute) for technical assistance. This study was supported in part
5 by the MAFF commissioned project study on “Development of efficient breeding technique aiming at
6 forestry trees with superior carbon storage capacity” (Grant Number JPJ009841), JSPS KAKENHI
7 (22H05172 and 22H05181), and the Kazusa DNA Research Institute Foundation.

8
9 **Competing interests**

10 The authors have no competing interests to declare that are relevant to the content of this article.

11
12 **References**

- 13 Bewg WP, Ci D, Tsai C-J (2018) Genome Editing in Trees: From Multiple Repair Pathways to Long-Term
14 Stability. *Front Plant Sci* 9:1732
- 15 Burdon RD, Wilcox PL (2011) Integration of Molecular Markers in Breeding. In: Plomion C, Bousquet J,
16 Kole C (eds) *Genetics, Genomics and Breeding of Conifers*. CRC Press, p 47
- 17 Chase MW, Soltis DE, Olmstead RG, et al (1993) Phylogenetics of Seed Plants: An Analysis of Nucleotide
18 Sequences from the Plastid Gene *rbcL*. *Ann Mo Bot Gard* 80:528–580
- 19 Cheng H, Concepcion GT, Feng X, et al (2021) Haplotype-resolved de novo assembly using phased
20 assembly graphs with hifiasm. *Nat Methods* 18:170–175
- 21 FAO (2020) *Global forest resources assessment 2020*. FAO
- 22 Forestry Agency (2022) *Annual Report on Forest and Forestry in Japan*. Ministry of Agriculture, Forestry
23 and Fisheries, Japan
- 24 Fujisawa Y (2017) The future of forestry and Chinese fir. *Forest Genetics and Tree Breeding* 6:132–136
- 25 Goraloglia GS, Redick TP, Strauss SH (2021) Gene editing in tree and clonal crops: progress and challenges.
26 *In Vitro Cellular & Developmental Biology - Plant* 57:683–699
- 27 Grattapaglia D (2022) Twelve Years into Genomic Selection in Forest Trees: Climbing the Slope of
28 Enlightenment of Marker Assisted Tree Breeding. *For Trees Livelihoods* 13:1554
- 29 Hon T, Mars K, Young G, et al (2020) Highly accurate long-read HiFi sequencing data for five complex

- 1 genomes. *Scientific Data* 7:1–11
- 2 Kurinobu S (2005) Forest Tree Breeding for Japanese larch. *Eurasian J For Res* 8:127–137
- 3 Lebedev VG, Lebedeva TN, Chernodubov AI, Shestibratov KA (2020) Genomic Selection for Forest Tree
4 Improvement: Methods, Achievements and Perspectives. *For Trees Livelihoods* 11:1190
- 5 Liu B, Shi Y, Yuan J, et al (2013) Estimation of genomic characteristics by analyzing k-mer frequency in
6 de novo genome projects. *arXiv:1308.2012*
- 7 Marçais G, Kingsford C (2011) A fast, lock-free approach for efficient parallel counting of occurrences of
8 k-mers. *Bioinformatics* 27:764–770
- 9 Mishima K, Hirakawa H, Iki T, et al (2022) Comprehensive collection of genes and comparative analysis
10 of full-length transcriptome sequences from Japanese larch (*Larix kaempferi*) and Kuril larch
11 (*Larix gmelinii* var. *japonica*). *BMC Plant Biol* 22:470
- 12 Muranty H, Jorge V, Bastien C, et al (2014) Potential for marker-assisted selection for forest tree breeding:
13 lessons from 20 years of MAS in crops. *Tree Genet Genomes* 10:1491–1510
- 14 Neale DB, Kremer A (2011) Forest tree genomics: growing resources and applications. *Nat Rev Genet*
15 12:111–122
- 16 Neale DB, Wheeler NC (2019) Gene and Genome Sequencing in Conifers: Modern Era. In: Neale DB,
17 Wheeler NC (eds) *The Conifers: Genomes, Variation and Evolution*. Springer International
18 Publishing, Cham, pp 43–60
- 19 Ohri D (2021) Variation and Evolution of Genome Size in Gymnosperms. *Silvae Genet* 70:156–169
- 20 Simão FA, Waterhouse RM, Ioannidis P, et al (2015) BUSCO: assessing genome assembly and annotation
21 completeness with single-copy orthologs. *Bioinformatics* 31:3210–3212
- 22 Tsumura Y, Matsumoto A, Tani N, et al (2007) Genetic diversity and the genetic structure of natural
23 populations of *Chamaecyparis obtusa*: implications for management and conservation. *Heredity*
24 99:161–172
- 25 Wan T, Gong Y, Liu Z, et al (2022) Evolution of complex genome architecture in gymnosperms.
26 *Gigascience* 11:giac078
- 27

1 **Table 1** HiFi library preparation protocols used in this study

Protocol	Shearing method	Shearing condition	Library cutoff size
1	g-Tube	Six centrifugations, each at $6,000 \times g$	15 kb
2	g-Tube	Six centrifugations, each at $6,000 \times g$	20 kb
3	g-Tube	Six centrifugations, each at $6,000 \times g$	25 kb
4	Megaruptor 2	20 kb	15 kb
5	Megaruptor 2	30 kb	15 kb
6	Megaruptor 2	40 kb	20 kb

2

1 **Table 2** Genome assembly statistics of four timber tree species

	<i>Larix kaempferi</i>	<i>Chamaecyparis obtusa</i>	<i>Cryptomeria japonica</i>	<i>Cunninghamia lanceolata</i>
Estimated genome size (bp)	13,216,140,982	7,946,916,792	9,890,457,748	12,037,454,484
Genome coverage of HiFi reads (×)	19.0	32.0	24.0	24.0
Primary contigs	LKA_r1.0	COB_r1.0	CJA_r1.0	CLA_r1.0
Total contig size (bp) (a)	13,492,429,495	8,513,066,914	9,239,540,489	11,548,046,079
No. of contigs	4,655	15,466	2,741	2,472
Contig N50	15,952,621	1,153,237	8,324,041	11,743,988
Gaps (bp)	0	0	0	0
GC content (%)	38.1	35.4	36.3	36.8
Complete BUSCOs, single-copy (%)	77.5	82.3	82.3	80.8
Complete BUSCOs, duplicated (%)	11.6	7.1	6.7	8.3
Fragmented BUSCOs (%)	6.8	5.5	6.1	6.3
Missing BUSCOs (%)	4.1	5.1	4.9	4.6
Alternate contigs	LKA_r1.0a	COB_r1.0a	CJA_r1.0a	CLA_r1.0a
Total contig size (bp) (b)	9,480,875,637	7,771,030,914	8,709,242,091	10,979,409,264
No. of contigs	24,731	35,816	13,386	11,396
Contig N50	1,663,422	496,648	2,129,829	2,899,092
Gaps (bp)	0	0	0	0
GC content (%)	38.0	35.4	36.3	36.8
Ratio of alternate contig size to primary contig size (%) (b/a)	70.3	91.3	94.3	95.1

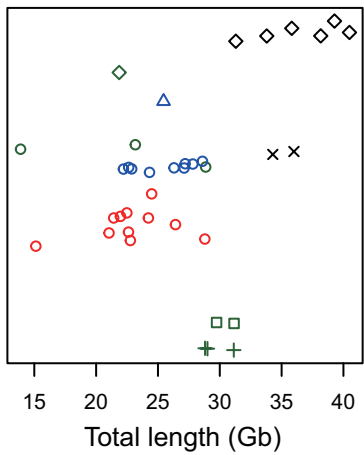
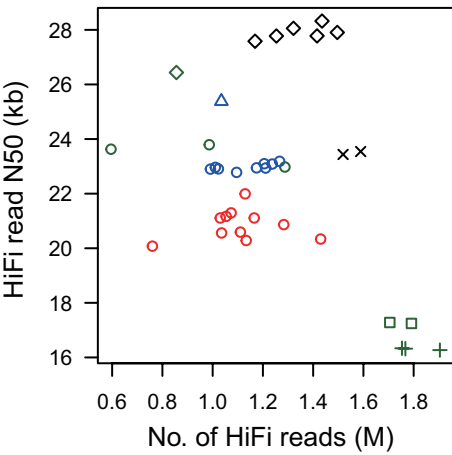
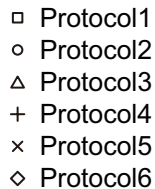
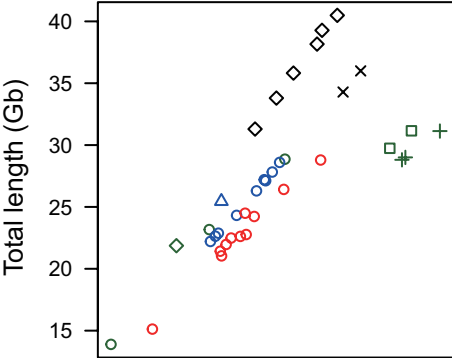
2

1 **Figure legends**

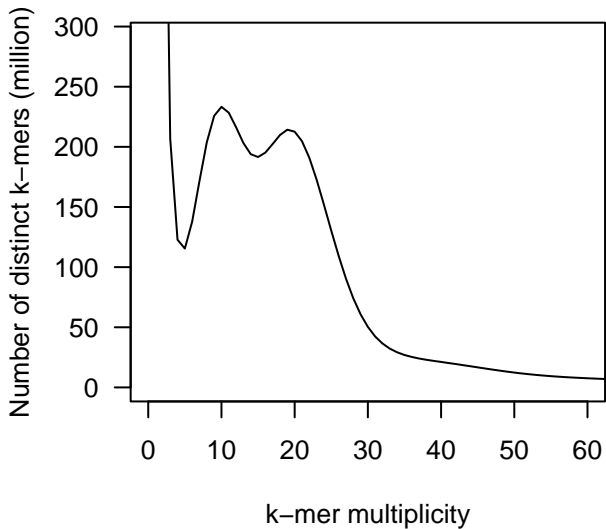
2 **Figure 1** Total length, N50 value, and number of HiFi reads obtained from DNA libraries prepared using
3 six different protocols. The different symbols (square, circle, triangle, plus, cross, and diamond) indicate
4 the different protocols (1, 2, 3, 4, 5, and 6, respectively; detailed in Table 1), and different colors (red, blue,
5 green, and black) indicate the different tree species (*Larix kaempferi*, *Chamaecyparis obtusa*, *Cryptomeria*
6 *japonica*, and *Cunninghamia lanceolata*, respectively).

7 **Figure 2** Estimation of the genome sizes of *L. kaempferi*, *C. obtusa*, *C. japonica*, and *C. lanceolata* by *k*-
8 mer distribution analysis. The *k*-mer size of 21 was used for the four timber tree species.

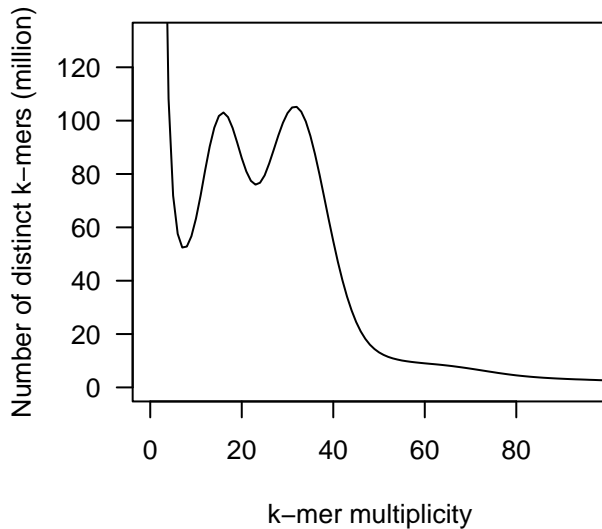
9



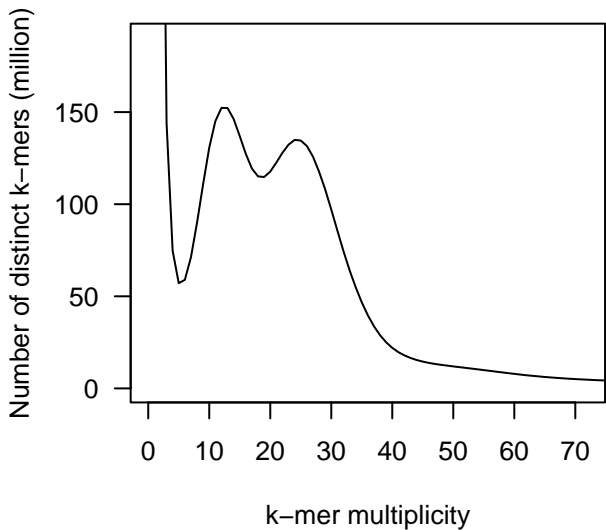
L. kaempferi



C. obtusa



C. japonica



C. lanceolata

