

# The extent and meaning of hybridization and introgression between Siberian spruce (*Picea obovata*) and Norway spruce (*Picea abies*): cryptic refugia as stepping stones to the west?

YOSHIAKI TSUDA,\*†<sup>1</sup> JUN CHEN,\*<sup>1</sup> MICHAEL STOCKS,\*<sup>1</sup> THOMAS KÄLLMAN,\* JØRN HENRIK SØNSTEBØ,‡ LAURA PARDUCCI,\* VLADIMIR SEMERIKOV,§ CHRISTOPH SPERISEN,¶ DMITRY POLITOV,\*\* TIINA RONKAINEN,†† MINNA VÄLIRANTA,†† GIOVANNI GIUSEPPE VENDRAMIN,† MARI METTE TOLLEFSRUD‡ and MARTIN LASCOUX\*

\*Department of Ecology and Genetics, Evolutionary Biology Centre, Uppsala University, Norbyvägen 18D, 75236, Uppsala, Sweden, †CNR, Institute of Biosciences and Bioresources, Via Madonna del Piano 10, 50019, Sesto Fiorentino, Firenze, Italy, ‡Norwegian Institute of Bioeconomy Research, Post Box 115, 1431, Ås, Norway, §Urals Division of the Russian Academy of Sciences, Institute of Plant and Animal Ecology, 8 Marta Str., 202, 620144, Ekaterinburg, Russia, ¶Swiss Federal Research Institute for Forest, Snow and Landscape Research (WSL), Zürcherstrasse 111, CH-8903, Birmensdorf, Switzerland, \*\*Vavilov Institute of General Genetics, Russian Academy of Sciences, Gubkin str. 3, 119991, Moscow, Russia, ††Environmental Change Research Unit (ECRU), Department of Environmental Sciences, University of Helsinki, PO Box 65, FI-00014, Helsinki, Finland

## Abstract

Boreal species were repeatedly exposed to ice ages and went through cycles of contraction and expansion while sister species alternated periods of contact and isolation. The resulting genetic structure is consequently complex, and demographic inferences are intrinsically challenging. The range of Norway spruce (*Picea abies*) and Siberian spruce (*Picea obovata*) covers most of northern Eurasia; yet their geographical limits and histories remain poorly understood. To delineate the hybrid zone between the two species and reconstruct their joint demographic history, we analysed variation at nuclear SSR and mitochondrial DNA in 102 and 88 populations, respectively. The dynamics of the hybrid zone was analysed with approximate Bayesian computation (ABC) followed by posterior predictive STRUCTURE plot reconstruction and the presence of barriers across the range tested with estimated effective migration surfaces. To estimate the divergence time between the two species, nuclear sequences from two well-separated populations of each species were analysed with ABC. Two main barriers divide the range of the two species: one corresponds to the hybrid zone between them, and the other separates the southern and northern domains of Norway spruce. The hybrid zone is centred on the Urals, but the genetic impact of Siberian spruce extends further west. The joint distribution of mitochondrial and nuclear variation indicates an introgression of mitochondrial DNA from Norway spruce into Siberian spruce. Overall, our data reveal a demographic history where the two species interacted frequently and where migrants originating from the Urals and the West Siberian Plain recolonized northern Russia and Scandinavia using scattered refugial populations of Norway spruce as stepping stones towards the west.

**Keywords:** divergence, Eurasia, introgression, phylogeography, spruce

Received 3 December 2015; revision received 23 February 2016; accepted 9 April 2016

Correspondence: Martin Lascoux, Fax: +46 (0) 18 471 64 57; E-mail: martin.lascoux@ebc.uu.se, Mari Mette Tollefsrud, Fax: +47 22 59 51 01; E-mail: mari.mette.tollefsrud@nibio.no and Giovanni Giuseppe Vendramin, Fax: +39 055 5225729; E-mail: giovanni.vendramin@ibbr.cnr.it

<sup>1</sup>These authors contributed equally to this work.

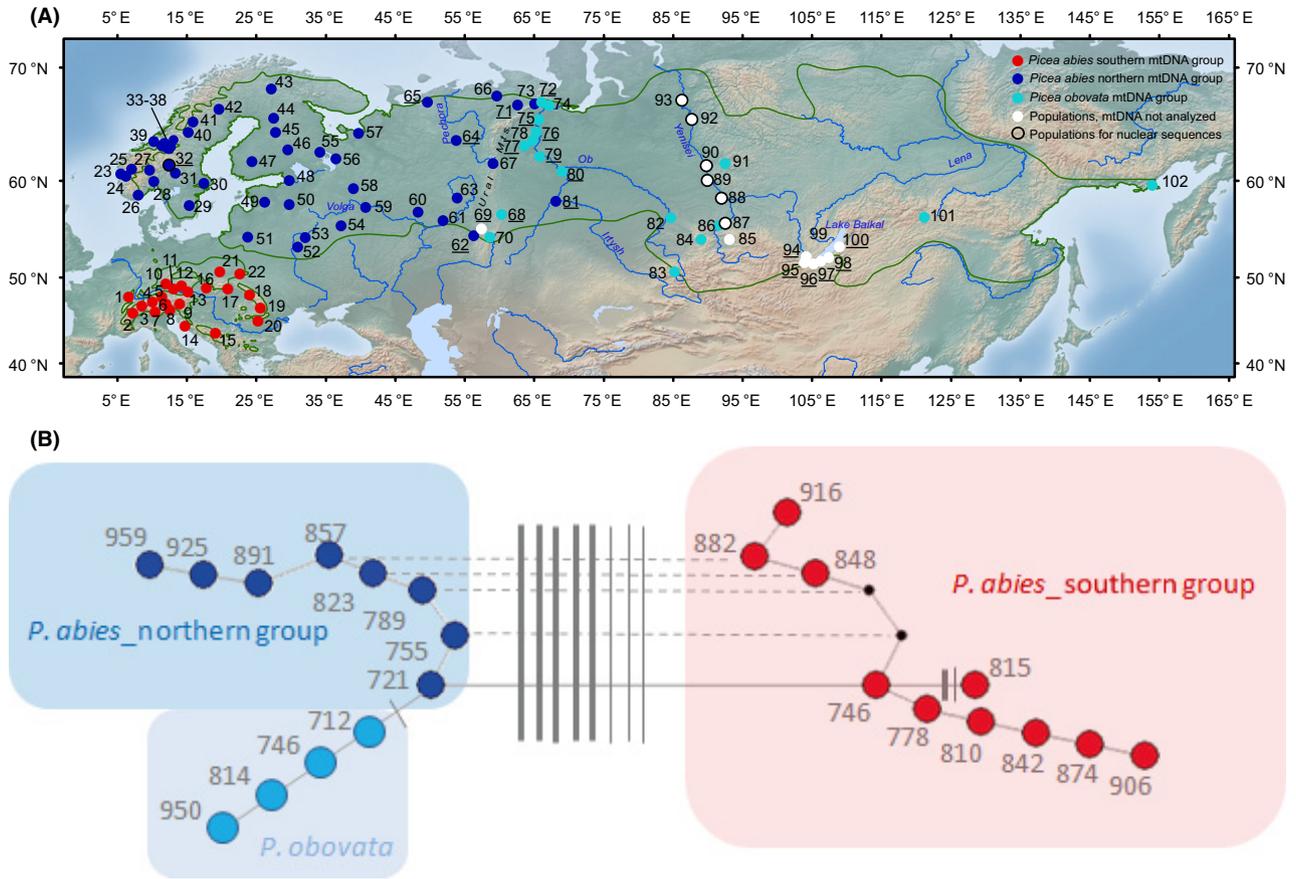
## Introduction

Hybridization and introgression are pervasive in plant species and have played a major role in their evolution (e.g. Anderson 1949; Arnold 1997; Gompert *et al.* 2012). Hybrid zones are interesting *per se*, and their genetic structure certainly contains a great deal of information about past demographics (Currat *et al.* 2008). Hybrid zones also offer a unique window on species barriers and on the role played by adaptation in speciation (Abbott *et al.* 2013; de Lafontaine *et al.* 2015). Finally, linkage disequilibrium can be higher in hybrid zones than within pure species ranges and thereby be useful for association mapping and the identification of adaptive genes (Rieseberg & Buerkle 2002). A recent surge in interest in hybridization and introgression among Asian spruce species (Du *et al.* 2009, 2011) or North American ones (Hamilton & Aitken 2013; Hamilton *et al.* 2013a,b; Haselhorst & Buerkle 2013; De La Torre *et al.* 2014a,b; De La Torre *et al.* 2015; de Lafontaine *et al.* 2015) has demonstrated that hybridization and introgression are widespread and that introgression is often asymmetrical (e.g. Du *et al.* 2011; Hamilton *et al.* 2013b). While, in some cases, this asymmetry seems to reflect primarily an initial asymmetry in the demographics of the two species, one being local and stable and the other invading and expanding (Currat *et al.* 2008; Du *et al.* 2009, 2011), it was proposed in a few studies that selection could also have played a role in the structure of the hybrid zone (Hamilton *et al.* 2013b; De La Torre *et al.* 2014a,b; Cinget *et al.* 2015). In either case, an important prerequisite to prove beyond doubt that selection played a part in shaping the hybrid zone is a good understanding of the demographic history of the two species and of their hybrid zone. More generally, information on a species history and demography is needed to reduce the number of both false positives and false negatives in genome scans for selection as recently illustrated in Swedish populations of *Arabidopsis thaliana* (Huber *et al.* 2014).

Boreal species were repeatedly exposed to ice ages and went through cycles of contraction and expansion while sister species alternated periods of contact and isolation. Recovering their demographic histories from standing genetic variation is therefore challenging, as we expect contact and hybrid zones to have moved through time. The combined range of the closely related species Norway spruce (*Picea abies*) and Siberian spruce (*Picea obovata*) constitutes one of the largest forested areas in the world, extending from the Norwegian coast in the west to the Sea of Okhotsk in the east and from the tree line in the north to the Alps and the Carpathians in the south (Fig. 1). Yet the phylogeny and history of the two species, as well as the exact limits of their respective ranges and

their hybridization zone, remain poorly understood. Morphologically, the two species are distinguished by the length of the cone and the shape of the cone scales; Norway spruce cones tend to be larger, and cone scales are more pointed and sharper than those of Siberian spruce, but in both cases, the change in the traits is more gradual across the range than discrete (Schmidt-Vogt 1974; Popov 2003; Volkova *et al.* 2014). An early study based on ten populations and 26 isozyme loci indicated the presence of a hybrid zone centred around the Urals, but the number of populations was clearly too small to delineate the extent of it (Krutovskii & Bergmann 1995). Krutovskii & Bergmann (1995) suggested that the 'two species should be considered as two closely related subspecies or two geographical races of one spruce species undergoing considerable gene exchange'. This suggestion was contradicted by the results of Tollefsrud *et al.* (2015) who showed, using a larger number of samples and chloroplastic (cp) and mitochondrial (mt) markers, that the two species constitute two clearly distinct lineages with a border between them located slightly east of the Urals. However, as cpDNA and mtDNA are easily introgressed (e.g. Du *et al.* 2009, 2011; Petit & Excoffier 2009), and as morphological traits do not permit an easy distinction of the two species, these data are not fully informative on the width of the hybrid zone and on the asymmetry of introgression. Nor did these data provide a dating of the divergence between the two species. Tollefsrud *et al.* (2015), however, combined cytoplasmic variation and palaeoecological data (both macrofossils and pollen) to draw a tentative picture of the history of the two species since the last glacial maximum (LGM). The two species appeared to have been widespread during the LGM and occupied distinct geographical areas, Norway spruce being primarily found across the East European Plain and Siberian spruce in southern Siberia and on the West Siberian Plain. Pockets of the two species were also able to survive at fairly high latitudes and close to the ice sheet that covered most of northern Europe during the last glaciations (Binney *et al.* 2009; Väli-ranta *et al.* 2011; Tollefsrud *et al.* 2015). The importance of the contribution of such remnant populations (aka 'cryptic refugia'; Stewart & Lister 2001) to the recolonization of Europe after the LGM has been and remains disputed (Tzedakis *et al.* 2013; Edwards *et al.* 2014; Huntley 2014; de Lafontaine *et al.* 2014a,b). As we shall see, the current genetic structure at nuclear loci lends support to the general picture drawn by Tollefsrud *et al.* (2015) and moreover suggests that the remnant populations played an important role in the recolonization of the west after each glaciation, though, in a rather unanticipated way.

To delineate the hybrid zone and reconstruct the demographic history of the two species, we genotyped 10 nuclear microsatellite (nSRR) loci in a large number of



**Fig. 1** Distribution map of samples and their mtDNA haplotypes. (A) Distribution map of all populations included in the study. The green outline illustrates the geographical distribution of Norway spruce and Siberian spruce in northern Eurasia after Schmidt-Vogt (1974) and Bezrukova *et al.* (2005). Populations are shown with their SAMOVA grouping based on mtDNA (the white dots correspond to populations where mtDNA was not analysed, bold circles indicate populations for which nuclear sequences were obtained). Underlined populations are additional populations compared to Tollefsrud *et al.* (2015). (B) Haplotype network of the mitochondrial haplotypes. Numbers indicate the size of the haplotypes in bp. Bars are nucleotide substitutions and lines are indels. The haplotypes are distinguished by variable copy numbers of two minisatellites, 32 bp (southern lineage) and 34 bp (southern and northern lineage) in size.

individuals from 102 populations representative of most of the natural range of the two species. We used available mtDNA data (Tollefsrud *et al.* 2008, 2009, 2015) and added new data for populations from the putative hybridization zone. Finally, in order to have an estimate of the divergence time between the two species, and more specifically between the northern domain of Norway spruce and Siberian spruce, we sequenced 22 nuclear loci in individuals from two populations located in the nonadmixed parts of each species. We then combined these data to (i) determine the location and breadth of the hybrid zone, (ii) identify the direction of introgression, (iii) estimate the divergence time between the two species and (iv) infer their past demography. The latter was achieved through two novel approaches. The first one constructs maps of effective migration rates and compares those to expectations under an isolation-by-distance model thereby allowing the identification of

barriers to gene flow. The second one estimates demographic parameters, such as effective population size and divergence time, through approximate Bayesian computation (ABC), and assess the fit of the model to the data by coalescent simulations and Bayesian clustering. Taken together, these results lay the foundations for, future work testing whether selection played a part in the establishment and maintenance of the hybrid zone or whether asymmetry in the demography of the two species is a sufficient explanation.

## Materials and methods

### Samples

For the nSSR analysis, a total of 102 populations (1299 individuals) were examined covering most of the distribution range of Norway spruce and Siberian spruce.

We divided those populations into four groups according to their geographical locations (Fig. 1A) and their general genetic structure at SSR loci (see Results section below): a southern Norway spruce domain (ALP, pop. 1–22), a northern Norway spruce domain (NOR, pop. 23–64), a hybrid zone between Norway spruce and Siberian spruce (HYZ, pop. 65–81) and, finally, a Siberian spruce group from the Ob River and eastward (OBR, pop. 82–102). We used DNA samples collected for earlier studies (Tollefsrud *et al.* 2008, 2009, 2015; Parducci *et al.* 2012; Chen *et al.* 2014; Table S1, Supporting information) and completed the sampling by adding 20 new populations from hitherto poorly characterized areas such as northeast Russia, the Ob River area and regions east of the Yenisei River (Fig. 1A; Table S1, Supporting information). To evaluate variation in the mtDNA fragment *nad1*, we used already available data from 68 populations (Tollefsrud *et al.* 2008, 2009, 2015) and added new data for 20 populations, altogether 88 populations (Fig. 1A; Table S1, Supporting information). Finally, we also analysed nuclear sequences from individuals sampled in one population located in central Sweden (12 individuals from the Fulufjället National Park) and in a set of populations located along the Yenisei in Siberia (96 individuals over six populations; Chen *et al.* 2014). Total genomic DNA of the newly collected samples was extracted from dried or fresh needles using the DNeasy plant mini kit or 96-well kit (Qiagen).

#### Genetic data

**Mitochondrial DNA.** Variation in the *nad1* gene for the newly collected populations was analysed as in Tollefsrud *et al.* (2008). As all the polymorphisms were binary [indels and single-nucleotide polymorphisms (SNPs)] or minisatellites (Tollefsrud *et al.* 2008), they were all treated as SNP-like data to estimate genetic distance among the different mitotypes.

**Nuclear microsatellites (nSSR).** A total of 10 nSSR loci were genotyped. The 10 loci were developed in previous studies: three by Scotti *et al.* (2002a) (EATC2B02, EATC1E03 and EATC2G05), one by Scotti *et al.* (2002b) (EAC2C08), one by Pfeiffer *et al.* (1997) (SPAC1F7) and five by Rungis *et al.* (2004) (WS0022.B15, WS0016.O09, WS00716.F13, WS0073.H08 and WS0092.A19). For the samples already analysed in Tollefsrud *et al.* (2009), we used the genotypes reported therein for the five loci EATC2B02, EATC1E03, EATC2G05, EAC2C08 and SPAC1F7. Except for these samples, fluorescent-dyed primer pairs of the 10 loci were mixed into three sets of multiplexes and amplified by polymerase chain reaction (PCR) in mixtures containing 1.2 µL of 1–10 ng DNA,

3.0 µL of master mix buffer (Type-it Microsatellite PCR kit; Qiagen), 1.2 µL of H<sub>2</sub>O and 0.6 µL of primer mix (with the concentration of each primer pair adjusted to 0.5 µM). Amplification was carried out using the following programme: initiation of Hot Start DNA polymerase and denaturation at 95 °C for 15 min; 32 cycles of 95 °C for 30 s, 57 °C for 30 s and 72 °C for 30 s; and a final 30-min extension step at 72 °C. The amplified PCR products were loaded onto ABI 3130 or ABI 3500 automatic sequencers (Applied Biosystems), and their sizes and genotypes were determined with the GENEMAPPER software (Applied Biosystems). To check possible allele size differences between newly obtained and previously analysed data (Tollefsrud *et al.* 2009), allele size calling was calibrated using samples shared by the two data sets as well as across the different automatic sequencers for the newly generated data.

**Nuclear DNA sequences.** Primers for the 22 loci used here were a subset of the loci designed using *Picea glauca* by Pavy *et al.* (2012). The Siberian spruce data were obtained previously through Sanger sequencing as described in Chen *et al.* (2014). The Norway spruce raw sequencing reads were obtained with the Illumina MiSeq sequencing platform (M. Stocks, J. Chen, T. Källman, J. Bousquet & M. Lascoux, unpublished data). Reads were trimmed to within 20–100 bps using *fastx\_trimmer* (A. Gordon & G. H. Hannon, unpublished data, [http://hannonlab.cshl.edu/fastx\\_toolkit](http://hannonlab.cshl.edu/fastx_toolkit)), mapped to the Siberian spruce reference sequences using *BWA* (Li & Durbin 2009) and filtered using *SAMTOOLS* (Li *et al.* 2009) such that the mapping quality must be  $\geq 29$ . Fastq files were generated from .vcf pileup files using *SAMTOOLS* and by marking sites with coverage of  $< 8$  reads as missing data. If, for any locus, a sample contained a number of missing sites covering more than around 10% of the total sequence, then this individual was excluded. The multiple sequence aligner *muscle* (Edgar 2004) was used to align Norway spruce to Siberian spruce sequences.

#### Analyses

**mtDNA.** To establish the phylogenetic relationships among the mitotypes, a minimum-spanning network was computed and drawn using *NETWORK* 4.6.1.1 (Fluxus Technology). We used spatial analysis of molecular variance (*SAMOVA* 1.0) to investigate the spatial patterns of genetic subdivision across the range of the two taxa (Dupanloup *et al.* 2002). This approach uses an iterative procedure to delineate contiguous groups of populations that are maximally differentiated. The number of inferred genetic clusters, *K*, was user-defined and set between 2 and 10 with 100 independent simulated annealing processes in each run. The optimal grouping

was considered when the differentiation among groups reached a plateau and before single populations began to be delimited.

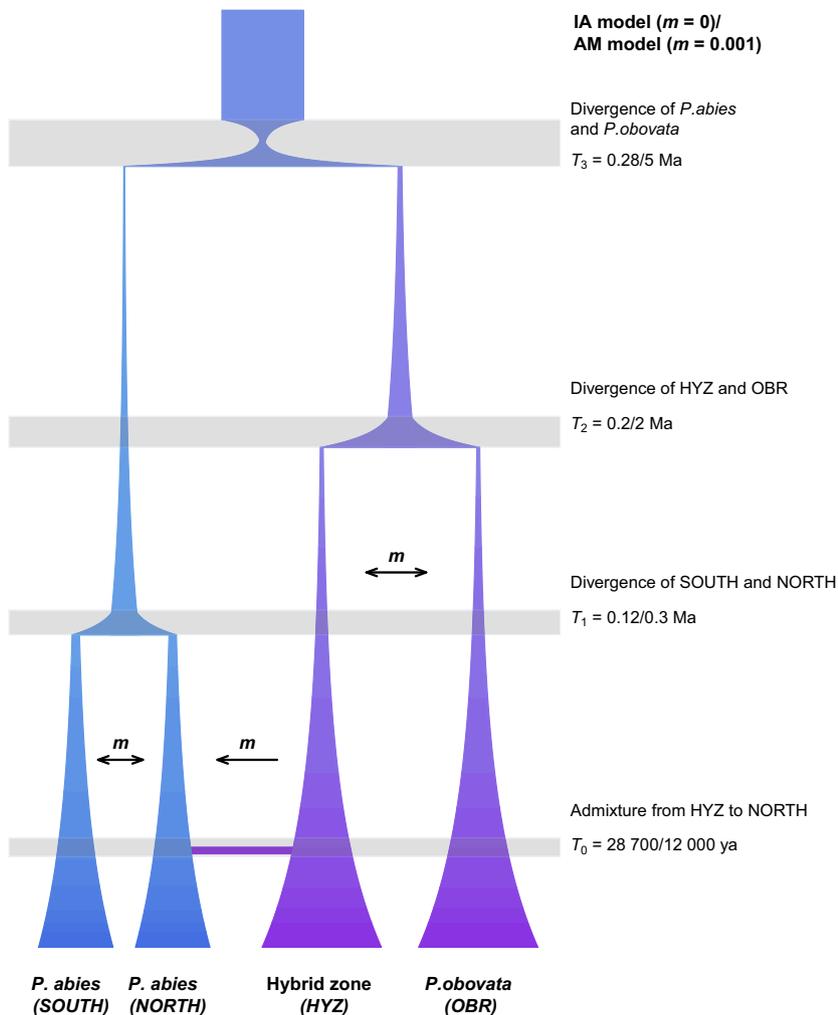
*SSR loci.  $F_{ST}$  and STRUCTURE*—Genetic divergence among populations, evaluated by Wright's fixation index,  $F_{ST}$  (Weir & Cockerham 1984), was estimated using  $F_{STAT}$  (Goudet 1995). Confidence intervals (95% and 99%) were obtained by bootstrapping, using 1000 replicates. The genetic structure was investigated with the model-based clustering algorithm implemented in the software STRUCTURE v. 2.3.4 (Pritchard *et al.* 2000; Hubisz *et al.* 2009). The number of clusters ( $K$ ) varied from 1 to 10, and we used a correlated allele frequencies model with 50 000 burn-in Markov chain Monte Carlo (MCMC) repetitions and 250 000 subsequent repetitions. The probabilities of the data for each  $K$  were averaged over 10 runs. STRUCTURE HARVESTER (Earl & vonHoldt 2012) was employed to calculate the probability of the data for each  $K$  (LnPD; Pritchard *et al.* 2000) and to calculate  $\Delta K$  (Evanno *et al.* 2005). Bar charts representing the proportion of cluster membership in each individual were edited and modified using CLUMP (Jakobsson & Rosenberg 2007) and DISTRICT (Rosenberg 2004).

*Estimated effective migration surfaces*—The estimated effective migration surfaces (EEMS) method is a new approach for analysing population structure from georeferenced genetic samples (Petkova *et al.* 2016). Briefly, the method assumes that the data follows roughly a model of isolation by distance (IBD). However, because landscapes are nonhomogeneous, in some areas, the data will conform well to the IBD model whereas in others, the genetic similarity will decay faster or slower than expected. The method identifies those areas and thereby allows the identification of barriers to gene flow. The EEMS are calculated by adjusting the migration rates so that the genetic differences obtained under a stepping-stone model closely match the observed genetic differences. The expected distances are calculated by integrating overall possible migration paths between populations, migration being more likely between closely located ones. The method is model based and less sensitive to sampling than principal component analysis (Petkova *et al.* 2016), and the resulting maps are easier to interpret and more informative about barriers to gene flow. We ran the programme with 30 000 burn-in MCMC steps and 50 000 subsequent iterations. To reduce the potential influence of the grid size, we averaged the results over runs of different number of demes (800, 1000, and 1500) as suggested by Petkova *et al.* (2016).

*ABC analysis with nSSR*—To infer the joint demographic history of Norway spruce and Siberian spruce, and in

particular the creation of an hybrid and introgression zone between them, we built an 'isolation with admixture' (IA) and an 'admixture with migration' (AM) model (Fig. 2) and estimated their parameters using ABC. The two models, isolation with admixture (IA) and admixture with migration (AM), are similar but the former is, except for the admixture event, a pure population split model whereas the latter populations are connected through migration. Importantly, we should stress that our prime goal was not to retrieve the detailed demographic history of either species, a task beyond our reach with the number of loci we had. Instead our aim, here was to test whether simple models of admixture and migration could retrieve the pattern of introgression revealed by the joint use of mtDNA and nSSR loci. The demographic models were suggested by palaeoecological data (Giesecke & Bennett 2004; Binney *et al.* 2009; Tollefsrud *et al.* 2015) and by the current distribution of mtDNA and nSSR genotypes inferred with SAMOVA and STRUCTURE, respectively (see Results, Fig. 3). To check the fit of the models to the data, and as we were primarily interested in explaining the current distribution of genetic variation, we tested whether data sets generated by coalescent simulations based on the model parameters and analysed with STRUCTURE retrieved the STRUCTURE plot obtained from the nSSR data (details are given in the Supporting information).

The coalescent simulations were performed using FASTSIMCOAL2 version 2.5.1 (Excoffier *et al.* 2013). We applied a generalized stepwise-mutation (GSM) model on 10 nSSR loci, with a mutation rate of  $2 \times 10^{-4}$  per locus and 60% of the mutations changing the allele size by more than one step. The choice of a mutation model and a mutation rate is always a bit arbitrary. In the present case, we choose this specific mutation rate and model based on previous studies (Marriage *et al.* 2009; Kuchma *et al.* 2011), the fact that a GSM seemed to offer a better fit to the observed allele distribution pattern than a simple stepwise-mutation model (SMM) model and also, as it was shown that model misspecification, in particular the adoption of a simple SMM when a GSM is the true model, can have a significant impact on the estimation of demographic parameters (Gonser *et al.* 2000; Leblois *et al.* 2014). Four million sets of multivariate parameters were drawn from their prior distributions and used in simulations for both the IA and AM models. Summary statistics were calculated using ARLSUMSTAT ver. 3.5.1.3 (Excoffier & Lischer 2010). We estimated the posterior distributions of the parameters using the neutral network regression algorithm implemented in the R package 'abc' with a tolerance ratio of  $5 \times 10^{-5}$  (Csilléry *et al.* 2012). The goodness of fit of each model was evaluated by posterior prediction on



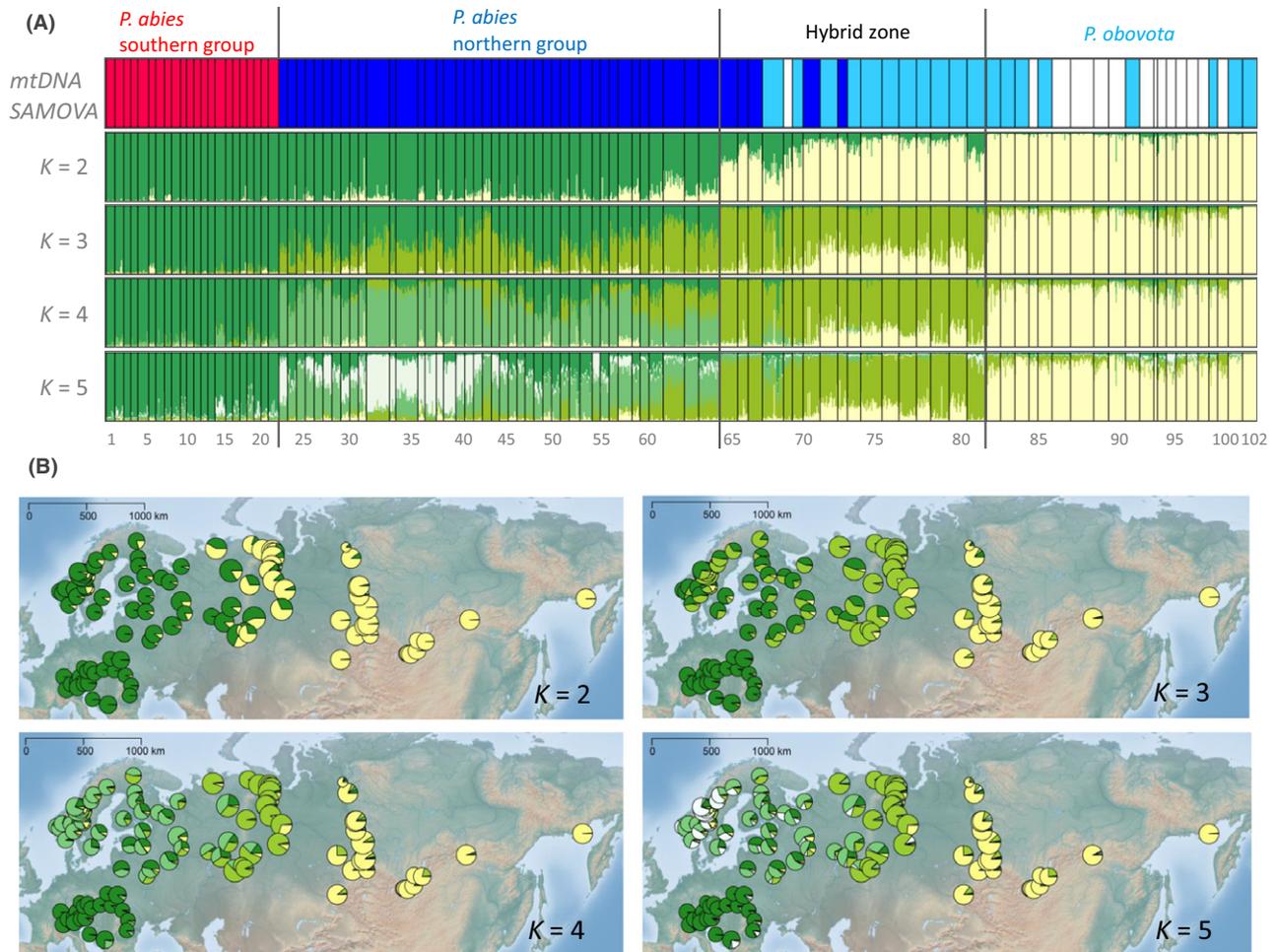
**Fig. 2** Approximate Bayesian computation (ABC) models used to retrieve the clustering results based on the SSR loci. The two models, isolation with admixture (IA) and admixture with migration (AM), are similar but the former is, except for the admixture event, a pure population split model whereas the latter populations are linked by migration.

summary statistics. A total of 1000 sets of parameters were randomly drawn from the posterior distributions and used in coalescent simulations for both models. Posterior predictions of each summary statistics were calculated from those 1000 simulated data sets and compared to the observed values.

We further computed the posterior model probabilities by comparing observed and simulated summary statistics and used Bayes factors to compare the two models. However, Bayes factor alone can be misleading if comparing against a naive null model. Thus, as we are mostly interested in knowing whether our model can produce the observed population genetic structure, we also created 1000 data sets using parameters from the posterior distributions and analyse those data sets with STRUCTURE. We ran STRUCTURE analysis on 1000 simulated data sets of both ABC models with 10 000 steps of burn-in and 50 000 MCMC iterations. We calculated the probabilities of every individual being assigned to four ( $K = 4$ ) or five clusters ( $K = 5$ ) separately. We chose to consider  $K = 4$  as  $K = 3-4$  was the most likely number

of clusters in the STRUCTURE analysis when using Pritchard *et al.* (2000)'s initial ad hoc procedure, and jointly, it fits well with the mtDNA data and with what one would expect from the palaeoecological data. Although it was slightly less likely than  $K = 4$ , we also considered  $K = 5$  for which a cluster centred on Scandinavia occurred. To measure the goodness of fit of the simulated STRUCTURE clustering generated by either model, we calculated the proportions of individuals in each of inferred clusters for each population. Two-tail  $P$ -values were computed by comparing the observed proportions to the ranked simulated values and were adjusted using false discovery rate correction.

*ABC analysis of nuclear DNA sequences*—All summary statistics for the observed data were calculated using the python library Evolib (<https://github.com/mspopgen/Evolib>). ABC was used to fit population split models to two populations, Fulufjället, which belongs to the northern Norway spruce domain, and Yenisei, a group of populations located along the Yenisei River in Siberia, which



**Fig. 3** Genetic clustering analysis conducted in SAMOVA and STRUCTURE (A) Bar plots of clustering results based on mtDNA (white colour represents the lack of mtDNA haplotypes) and nSSR ( $K = 2$  to  $K = 5$ ) (B) Pie charts showing the geographical distribution map of the STRUCTURE clusters across the different populations.

show no genetic differentiation (Chen *et al.* 2014), and are part of the Siberian spruce range. Two models were considered: (i) a simple isolation model (IS) where two populations split at some time point in the past and (ii) an isolation with migration model (IM) where symmetrical migration is allowed between the two diverging populations. In total, we selected 13 summary statistics for parameter estimation (see details of the ABC analysis in the Supporting information).

## Results

### Mitochondrial DNA variation

Twenty-two haplotypes were detected in the 88 populations surveyed for mtDNA variation. All of them have been previously described in Tollefsrud *et al.*'s study (2008, 2015). Five nucleotide substitutions and three indels defined a northern and a southern lineage

(Fig. 1B). The northern lineage consisted of eight mitotypes restricted to the northern range of Norway spruce and four mitotypes restricted to the range of Siberian spruce. The 10 mitotypes of the southern lineage were restricted to the southern range of Norway spruce. Hence, for mitochondrial DNA, Norway spruce is not monophyletic (Fig. 1B). SAMOVA delineated three highly differentiated groups ( $F_{CT\ mtDNA} = 0.744$ ; Fig. 3A), corresponding to the southern and the northern domains of Norway spruce and to the proposed range of Siberian spruce in Siberia (c.f. Tollefsrud *et al.* 2015). The northern domain of Norway spruce and Siberian spruce is closely related as reflected also in the phylogenetic network.

### Nuclear SSR

STRUCTURE and  $F_{ST}$ . The STRUCTURE analysis based on nSSR loci gave a dramatically different genetic structuring compared to that obtained from mitochondrial

DNA (Fig. 3). The most likely number of clusters in the STRUCTURE analysis was  $K = 2$  when using Evanno *et al.* (2005)'s  $\Delta K$  criterion and  $K = 3-4$  using Pritchard *et al.* (2000)'s initial ad hoc procedure (Fig. S2, Supporting information). At  $K = 2$ , the range of the two species was divided into two large clusters. A first group comprises Norway spruce populations from both the southern and the northern domains (pop. 1–64), and a second one includes Siberian spruce populations starting in the hybridization zone west of the Urals and stretching out east all the way to the Sea of Okhotsk (Fig. 3B). Populations in the hybrid zone along the Ural Mountains (pop. 61–81) hold an increasing share of ancestry from the Siberian spruce group as one moves east (pop. 61–82). At  $K = 3$ , two noteworthy changes occur. First, the southern domain emerges as a homogeneous group (pop. 1–22). Second, a new cluster centred along the Urals (pop. 65–81) appears with extension westwards across all the northern domain of Norway spruce. Notably, almost all the populations in the northern domain of Norway spruce (pop. 23–64) show a shared ancestry between these two clusters. At  $K = 4$ , the northern domain of the Norway spruce range (pop. 23–64) emerges as a separate cluster, though, with admixed ancestries. In particular, the populations from the Baltics, southern Finland and southwestern Russia (pop. 44–54) exhibit a significant level of admixture from the southern group and the ancestry of the populations to the west of the Urals (pop. 60–64) traces back to the hybrid group and the northern Norway spruce group. Finally, at  $K = 5$ , a new cluster emerges, which is centred on central Scandinavia. Examination of the STRUCTURE plot indicates that all new clusters roughly stem from admixed populations as one moves from  $K = 2$  to  $K = 5$ . It is also important to reiterate that all individuals from the northern part of the Norway spruce range (pop. 23–64) and a few from the hybrid zone (pop. 66, 67, 71, 73 and 81) belong to the northern Norway spruce group in the SAMOVA analysis of the mitochondrial polymorphism.

$F_{ST}$  values estimated from 10 nSSR loci showed limited genetic divergence between populations from the four main regions identified by STRUCTURE (the southern and northern domains of Norway spruce, the hybrid zone and the Siberian spruce region; Table 1):  $F_{ST}$  increases slightly with the geographical distance between groups and the most divergent pair is the northern domain of Norway spruce and Siberian spruce ( $F_{ST} = 0.143$ ). The expected heterozygosity varies from 0.68 (Siberian spruce) to 0.74 (southern domain) among the four main regions (Table S2, Supporting information).

*Estimated effective migration surfaces.* Figure 4 gives the EEMS for our data. The EEMS map is dominated by the

**Table 1**  $F_{ST}$  at nSSR loci among the four main clusters defined by STRUCTURE: the southern domain of Norway spruce (pop. 1–22), the northern domain of Norway spruce (pop. 23–64), the hybrid zone between Norway spruce and Siberian spruce (pop. 65–81) and Siberian spruce (pop. 82–102)

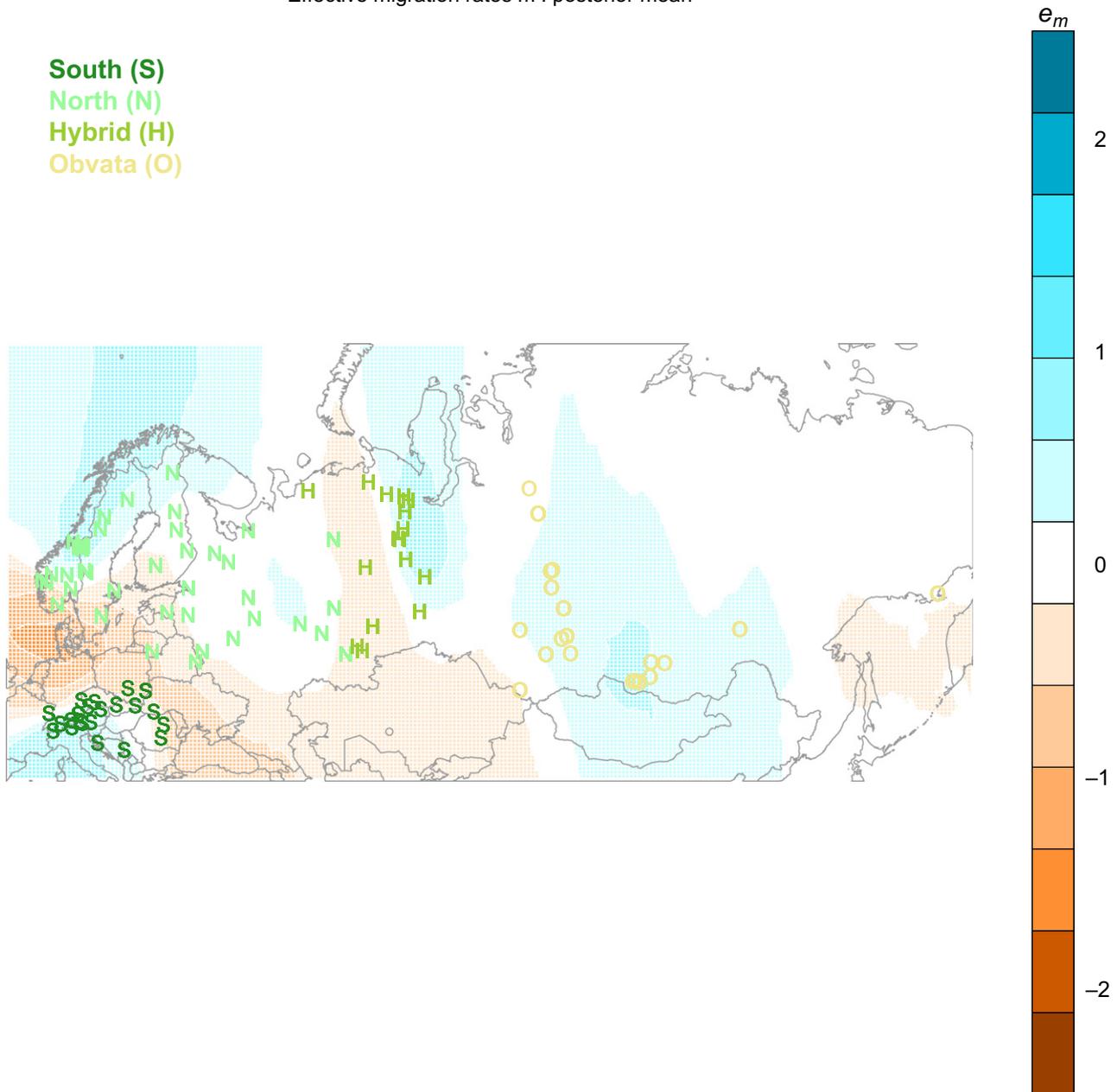
|                    | Southern<br>(pop. 1–22) | Northern<br>(pop. 23–64) | Hybrid<br>(pop. 65–81) | Siberian<br>spruce<br>(pop.<br>82–102) |
|--------------------|-------------------------|--------------------------|------------------------|--|
| Southern           | 0.000                   |                          |                        |  |
| Northern           | 0.024                   | 0.000                    |                        |  |
| Hybrid             | 0.068                   | 0.034                    | 0.000                  |  |
| Siberian<br>spruce | 0.103                   | 0.143                    | 0.122                  | 0.000                                  |

presence of two large barriers to gene flow: the first one separates the northern and southern domains of Norway spruce while the second one corresponds to the Urals Mountains and the hybrid zone between Norway spruce and Siberian spruce. These barriers separate areas where the decay of genetic differences across geographical distances is in line with the expectations of an IBD model. There are also areas (in blue in Fig. 4), such as Central Russia, Scandinavia or along the River Ob, where the effective migration rates tend to be slightly higher than expected under an IBD model.

*ABC analysis of the hybrid zone using nSSR data.* We used nSSR data to test whether admixture alone (from HYZ to NOR, the IA model) can generate the introgressive clustering patterns observed in NOR and HYZ populations, or whether migrations between groups is necessary (the AM model). The model comparisons based on Bayes factor ( $BF = 34$ ) and the posterior prediction of summary statistics (Figs S3 and S4, Supporting information) both support the presence of migration (the AM model, detailed descriptions of all parameters and model comparison are given in Table S3, Supporting information). More importantly, STRUCTURE analyses based on 1000 posterior runs showed that simulated data generated with the AM model parameters led to a much closer fit to the observed than simulated data generated with the IA model parameters at both  $K = 4$  and  $K = 5$  (Table 2 and Figs S5–S8, Supporting information). In general, the IA model failed to generate the correct genetic clusters in HYZ populations (Clusters 3 and 4) while the AM model had a very satisfying match in this case.

#### Nuclear sequences

*Within and between species summary statistics.* Measures of the population mutation rate are similar between the

Effective migration rates  $m$  : posterior mean

**Fig. 4** The estimated effective migration surfaces plot shows the effective migration rates between all populations. Populations in the blue areas are connected by higher migration rates than expected under isolation by distance (IBD) while the ones in the orange areas have lower migration rates than expected and are interpreted as migration barriers. In white areas, the effective migration surface is close to the one expected under IBD.

two species (Table S4, Supporting information), with Siberian spruce showing a slightly higher value according to Watterson's estimate (Norway spruce:  $\theta_W = 0.00388$ ; Siberian spruce:  $\theta_W = 0.00403$ ) but with Norway spruce showing a higher value according to the average number of pairwise nucleotide differences (Norway spruce:  $\pi = 0.00392$ ; Siberian spruce:  $\pi = 0.00333$ ). This is reflected in Siberian spruce with a

more negative Tajima's  $D$  ( $-0.544$ ) than in Norway spruce ( $-0.132$ ). In total, we observed 171 segregating sites in Norway spruce and 255 in Siberian spruce. This difference could be due to the difference in sample size between Siberian spruce (96) and Norway spruce (12). In total, we observed 332 segregating sites in the two species, 76 private to Norway spruce, 160 to Siberian spruce and 95 occur in both Norway spruce and

**Table 2** Proportions of individuals in each STRUCTURE cluster ( $K = 4$  and  $K = 5$ ) for all four groups used in the approximate Bayesian computation (ABC) analysis. Comparison of observed values (first row) and predicted values (second row) based on STRUCTURE analyses carried out using data sets generated from the parameters estimated in the ABC analysis with the IA and AM models (first and second figure within parentheses, respectively)

|         | Cluster 1 <sup>†</sup>                      | Cluster 2                       | Cluster 3                        | Cluster 4                        | Cluster 5                      |
|---------|---|---------------------------------|----------------------------------|----------------------------------|--------------------------------|
| $K = 4$ |   |                                 |                                  |                                  |                                |
| ALP     | 0.905 (0.963, 0.965)                        | 0.062 (0.032, 0.021)            | 0.018 ( <b>0.002**</b> , 0.008)  | 0.015 ( <b>0.003***</b> , 0.006) |                                |
| NOR     | 0.187 ( <b>0.012***</b> , 0.068)            | 0.655 (0.460, 0.663)            | 0.119 (0.523, 0.230)             | 0.039 ( <b>0.005*</b> , 0.039)   |                                |
| HYZ     | 0.023 ( <b>0.003*</b> , <b>0.002***</b> )   | 0.048 (0.137, <b>0.004***</b> ) | 0.774 (0.845, 0.851)             | 0.156 ( <b>0.015***</b> , 0.143) |                                |
| OBR     | 0.026 ( <b>0.003**</b> , <b>0.001***</b> )  | 0.038 (0.027, <b>0.003***</b> ) | 0.087 ( <b>0.009***</b> , 0.037) | 0.849 (0.962, 0.959)             |                                |
| $K = 5$ |   |                                 |                                  |                                  |                                |
| ALP     | 0.884 (0.908, 0.909)                        | 0.051 (0.011, 0.029)            | 0.017 ( <b>0.002***</b> , 0.009) | 0.015 ( <b>0.002**</b> , 0.006)  | 0.034 (0.077, 0.047)           |
| NOR     | 0.183 ( <b>0.008***</b> , 0.058)            | 0.479 (0.451, 0.487)            | 0.090 (0.521, 0.197)             | 0.037 ( <b>0.004*</b> , 0.037)   | 0.212 ( <b>0.016*</b> , 0.220) |
| HYZ     | 0.020 ( <b>0.002***</b> , <b>0.002***</b> ) | 0.067 (0.055, 0.018)            | 0.754 (0.729, 0.795)             | 0.138 ( <b>0.007***</b> , 0.145) | 0.020 (0.206, 0.040)           |
| OBR     | 0.02 ( <b>0.002**</b> , <b>0.001***</b> )   | 0.03 (0.012, 0.005)             | 0.093 ( <b>0.006***</b> , 0.040) | 0.782 (0.913, 0.942)             | 0.070 (0.067, 0.011)           |

Numbers in bold font show the simulated values that significantly deviated from the observed ones.

<sup>†</sup>The observed value is presented followed by the mean values summarized from the IA and AM models in parentheses.

\* $P$ -value < 0.05, \*\* $P$ -value < 0.01, \*\*\* $P$ -value < 0.001.  $P$ -values were adjusted for multiple tests using false discovery rate.

Siberian spruce. There was only one single fixed site between the two species.

*ABC analysis.* Parameters for both the isolation (IS) and the isolation-with-migration (IM) models (Fig. S9, Supporting information) were well estimated by the neural network ABC method. Estimates of the population mutation rate ( $\theta$ ) were similar for both the IS and IM models (Table S5, Supporting information: median  $\theta$  IS = 0.00395; median  $\theta$  IM = 0.00365). Estimates of the mutation rate per site and per year are  $0.7\text{--}1.3 \times 10^{-9}$  (Willyard *et al.* 2007) in *Pinus* and  $0.6\text{--}1.1 \times 10^{-9}$  in *Picea abies* (Chen *et al.* 2012). Adjusting to an assumed generation time comprised between 25 and 50 years, the per site per generation mutation rate was set at approximately  $2.5 \times 10^{-8}$  (Bodare *et al.* 2013). This gives estimates of the effective population size for the two species ( $N_e$ ) of 42 586 and 36 163 for the IS and IM models, respectively (see Table 3 for confidence intervals). The time of divergence (split, as measured in coalescent time units) estimated between the two species for the two models did, however, show an approximately twofold difference between the two models (median IS model split = 0.07668; median IM model split = 0.15274). Assuming the above mutation rate, this translates divergence times into generations of 6563 and 11 891 for the IS and IM models, respectively. Assuming a 50-year generation time (e.g. Chen *et al.* 2010), we estimate a divergence time of 328 161 and 594 551 years for the IS and IM model, respectively. Assuming instead a generation time of 25 years per generation gives estimates of 164 081 and 297 276 years for the IS and IM models, respectively. For the IM model, we estimate the population migration rate at 1.237,

corresponding to a symmetrical migration rate of  $9.33 \times 10^{-6}$ . We find no statistical difference in the fit of the models to the data using both the rejection and mnlogistic methods and a slightly higher model probability for the IS model compared to the IM model (rejection: IS = 0.512, IM = 0.488; mnlogistic: IS = 0.5419, IM = 0.4581).

## Discussion

In the present study, we combined mitochondrial DNA, nuclear microsatellites and SNP data to infer the joint demographic history of Norway spruce and Siberian spruce. Combining information from genomes with different effective population sizes and mutation rates, as well as contrasted modes of inheritance and dispersal, allowed us to capture demographic events on different time scales and to infer the major features of population movements and gene flow across the combined ranges of the two species. The emerging picture is that Norway spruce and Siberian spruce are two distinct species with a joint and dynamic demographic history, involving repeated recolonizations from the east that were accompanied by extensive hybridization and introgression. This is today reflected by a complex population genetic structure, in particular within Norway spruce that was more strongly affected than Siberian spruce by glaciations (Chen *et al.* 2012, 2014; Tollefsrud *et al.* 2015). The new insights emerging from our analyses are the following: (i) the range of the two species is divided by two main barriers: one corresponding to the hybrid zone between the two species and another one separating the southern and northern domains of Norway spruce, (ii) the hybrid zone between the two species is

**Table 3** Median estimates and confidence intervals (2.5%, 97.5%) for the neural network approximate Bayesian computation analysis estimates of the effective population size ( $N$ ), the time of divergence in generations ( $T_g$ ), the time of divergence in years (Ty25 and Ty50, assuming 25- and 50-year generation times, respectively) and the migration rate ( $m$ ). Calculations assume a per site, per generation mutation rate of  $2.5 \times 10^{-8}$

|                                | $N$                        | $T_g$                    | Ty25                          | Ty50                            | $m$                              |
|--------------------------------|----------------------------|--------------------------|-------------------------------|---------------------------------|----------------------------------|
| Isolation model                | 42 586<br>(36 292, 51 685) | 6563<br>(4276, 9807)     | 164 081<br>(106 906, 245 175) | 328 161<br>(213 813, 490 351)   | —                                |
| Isolation with migration model | 36 163 (29 714, 45 047)    | 11 891<br>(5274, 24 791) | 297,276<br>(131 847, 619 783) | 594,551<br>(263 695, 1 239 566) | 9.33e-06<br>(2.47e-06, 2.24e-05) |

rather large, extending roughly from the Ob River east of the Ural Mountains to the Russian plain in the west with even some impact of Siberian spruce towards northern Fennoscandia, (iii) comparison of the distribution maps of mitochondrial and nuclear variation (Figs 1 and 3), clearly indicates an introgression of mtDNA from Norway spruce into Siberian spruce in the hybrid zone and (iv) coalescent simulations based on the parameters of a model with admixture and introgression allow the reconstruction of the main genetic clusters. The last two results, taken together, support a scenario under which Norway spruce populations were invaded by a migration front originating from the Urals or even the West Siberian Plain after the LGM and most likely after preceding glacial periods. These new results, their limitations and their implications for the history of the two species are discussed below.

#### *Barriers and differentiation at nuclear and mitochondrial genes*

The two main barriers retrieved by EEMS, the one separating the southern and northern domains of Norway spruce and the one corresponding to the hybrid zone between the two species, are consistent with previous population genetic and phylogeographical studies in Norway spruce (Lagercrantz & Ryman 1990; Heuertz *et al.* 2006; Tollefsrud *et al.* 2008, 2009) or in both species (Krutovskii & Bergmann 1995; Tollefsrud *et al.* 2015). Krutovskii & Bergmann (1995) isozyme study established that the contact zone between the two species was centred on the Urals and Tollefsrud *et al.* (2015) confirmed it using cytoplasmic markers and a larger number of populations. Furthermore, Tollefsrud *et al.* (2015) contrasted the population genetic structure of mitochondrial and chloroplastic variation. For both markers, the main border between the major haplotypes was centred on the Urals. However, whereas the area around the Urals was a rather sharp border for mitochondrial haplotypes, the main Siberian spruce chloroplastic haplotype was still found at high frequency in Scandinavian populations suggesting an impact of

Siberian spruce much beyond the Urals. Haplotype sharing can be a result of introgression and/or retention of ancestral polymorphisms (e.g. Wang *et al.* 2011). Also, like nuclear DNA, cpDNA is spread through both seeds and pollen and if we believe that the pattern observed at nuclear markers reflects past population movements, these population movements should, at least in part, be also responsible for the current distribution of chloroplast DNA variation. A hybrid zone centred around the Urals is also in agreement with Popov (2003) morphological study that indicates that Norway spruce-like characteristics increase progressively as one moves west of the Urals/Ob River area while the Siberian spruce-like characteristics increase progressively as one moves east. Finally, it is worth noting that most phylogeographical studies indicate that the Urals, or the area around it, constituted a major contact zone for many plant and animal species (e.g. Semerikov & Lascoux 2003; Fedorov *et al.* 2008; Eidesen *et al.* 2013; Semerikov *et al.* 2013).

Recently, Lockwood *et al.* (2013), using results from a phylogenetic study based on a limited number of cytoplasmic and nuclear markers, have argued that 'northern Norway spruce may be part of an Asian species complex that diverged from the southern Norway spruce lineage in the Upper Miocene, some 6 Ma, which can explain the deep genetic gap noted in phylogeographical studies of Norway spruce'. This interpretation is consistent with the pattern that we observe with mtDNA, where the number of mutations separating the southern and northern domains is larger than the number of mutations separating the latter and Siberian spruce (Fig. 1B), but seems at variance with the  $F_{ST}$  values obtained with nSSR loci.  $F_{ST}$  values are generally low (0.028–0.143) and are lowest between the southern and northern domains of Norway spruce and highest between Siberian spruce and the two Norway spruce domains. Also, in the STRUCTURE analysis, the separation between the northern and the southern domains occurs after the separation between the two species (at  $K = 4$  and  $K = 2$ , respectively) and also indicates the presence of a contact zone with some amount of admixture

around the Baltics, southern Finland and southern Russia. In part, this apparent discrepancy could be due to the fact that the Lockwood *et al.* (2013)'s data primarily reflect mitochondrial polymorphism, the only nuclear sequence used having limited variation. How can one then reconcile the difference in divergence for the two genomes, or more precisely between the two genomes and the different types of markers? The key probably lies in the different mode of inheritance and level of introgression of the two genomes and in the history of the two domains. A scenario compatible with the data could be the following. The northern and the southern domains diverged some 6 Ma. The southern domain became very restricted and was more recently recolonized by the same source as the northern domain. During the recolonization process, the southern mitochondrial DNA lineage introgressed into the incoming trees as predicted by the model proposed by Currat *et al.* (2008), something that has been generally observed, in particular in spruce species (e.g. Du *et al.* 2009). The two domains then diverged and have since experienced limited gene flow and probably most of it via pollen. Of course, this model is certainly a strongly simplified picture of what actually happened, especially as it does not consider the multiple cycles of glaciations that have affected Eurasia during the Quaternary, but it depicts a scenario that could generate patterns consistent with both our and Lockwood *et al.* (2013) results. It should also be pointed out that both our estimates and that of Lockwood *et al.* (2013) are based on a rather limited number of loci and that a larger number of loci would undoubtedly improve those estimates.

Finally, regarding the breadth of the hybrid zone, as inferred from the STRUCTURE plot or the estimated effective migration surfaces analysis, it is worth pointing out that our results depend rather strongly on sampling. In particular, we have a large gap in sampling between the Ob and Yenisei rivers. This gap creates well-separated clusters between the hybrid zone (pop. 65–81) and the Siberian spruce cluster. It may well be that, had samples from the vast area separating the two river areas been available, the change in admixture would have been more progressive. As a matter of fact, an early analysis without data from the Ob River led to a very different interpretation, with an even sharper border. Overall, however, the results fit very well with those obtained from chloroplast DNA (Tollefsrud *et al.* 2015) and with the introgression zone suggested by morphological data (Popov 2003, 2010).

#### *Introgression and demographic inferences*

The general picture of the demographic history emerging from the present study complements previous

genetic studies that focused on a more limited geographical area (Lagercrantz & Ryman 1990; Krutovskii & Bergmann 1995; Potokina *et al.* 2013) or were only based on cytoplasmic markers (Mudrik *et al.* 2015; Potokina *et al.* 2015; Tollefsrud *et al.* 2015). It also corroborates Binney *et al.* (2009)'s survey of macrofossils across Eurasia: macrofossil remains of *Picea* dating to prior to the LGM were found along the southern edge of the West Siberian Plain and on the East European Plain (see also Tollefsrud *et al.* 2015). *Picea* was also present in sites close to the edge of the Scandinavian ice sheet when the latter reached its maximum extent, that is west of the Urals in the Pechora region (around pop. 64–66 in the present study, Fig. 1A) (Väliranta *et al.* 2011). The pattern of mtDNA introgression from Norway spruce into incoming Siberian spruce in the Pechora region observed in the present study suggests that these sites could have contributed to population expansion, though, both indirectly and directly: as conditions improved and ice retreated the remnants populations of Norway spruce scattered across the East European Plain were swamped by pollen and seeds originating from refugial populations located east of the Urals or in the West Siberian Plain where Siberian spruce today dominates. The importance of pollen flow from Siberian spruce into the Norway spruce domain is supported by the presence of Siberian spruce chloroplast haplotypes in Scandinavia, in individuals all having, otherwise, Norway spruce mitochondrial haplotypes (Tollefsrud *et al.* 2015). Both fossil record and genetic data are limited in western Siberia, yet all evidence point to the survival of trees in large pockets or along river banks (Semerikov *et al.* 2013 and references therein). Hence, the cryptic refugia scattered at high latitudes west of the Urals may have had a limited direct contribution to the migration front but they certainly did contribute through introgression. In other words, they probably served as 'genetic stepping stones' as the migration front moved westwards. The presence of these 'stepping stones' contributing both directly, through seeds, and indirectly, through introgression, to the recolonization process, implies that simple demographic models such as classical stepping stone, isolation-by-distance or diffusion models would be difficult to fit to the data. The observed introgression of mtDNA from Norway spruce into Siberian spruce is indeed predicted by Currat *et al.* (2008) model. Under this model, introgression occurs almost exclusively from the local to the invading species. Introgression is more important for markers that have low levels of gene flow within species, for instance mtDNA (maternally inherited) in conifers, than for markers with high levels of gene flow, that is nuclear markers (biparentally inherited). This model has been shown to provide a simple explanation

for the rather general observation of asymmetric introgression between markers with different modes of inheritance and/or dispersal. As in the present study, most analyses of hybrid zones in conifers have found a higher level of introgression for mitochondrial DNA than for nuclear or chloroplast DNA (e.g. Du *et al.* 2009, 2011).

Currat *et al.* (2008) model implies that the current population genetic structure is the result of a primarily westwards migration starting within the current distribution of Siberian spruce or, at least, from the hybrid zone. Do demographic inferences from the variation at microsatellite loci support this model? To test the plausibility of the proposed scenario we used, to our knowledge for the first time, a combination of ABC and posterior predictive simulations to see whether demographic models with and without migration could recreate the observed population genetic structure. A demographic model including admixture and migration, along the lines suggested by the comparison of mtDNA and nSSR variation, was much more likely than a model without migration and did recreate the observed STRUCTURE pattern while a model without migration failed to do so. These results therefore demonstrate that a rather simple model of westwards migration and admixture could explain the observed STRUCTURE pattern, despite its apparent complexity. The timing of these events, however, is more difficult to estimate. The choice of the model had a strong impact on the estimates of divergence times and effective population sizes (Table S3, Supporting information), and those were quite different from the estimates obtained with nuclear sequence data. This may feel somewhat problematic but there are some caveats. First, the models were not strictly equivalent. Second the data were in both cases limited and estimates have fairly large confidence intervals. Third, and more generally, one has to realize that effective population size and divergence time, apart from being dependent on strong assumptions on mutation rates and generation times, are intrinsically extremely model dependent and are rather meaningless taken out of context. This goes at the heart of the concept of effective population, which was initially defined as the size of a Wright–Fisher population that would have the same rate of inbreeding (inbreeding population size) as the observed population. So, by definition, the effective population size is truly a scaling factor of the amount of random genetic drift for a given demographic model and has little relation with census size, a point that is often missed (see e.g. Ewens 2004, p. 128). It is therefore not surprising that different models, applied to different data (nSSR vs. DNA sequences), lead to different estimates of demographic parameters. In any case, with all these caveats in mind, it is

noteworthy that all analyses point at a split between the two species, and the creation of a hybrid zone, predating the LGM. Considering that many glacials and interglacials took place over the last million years and that each cycle was probably accompanied by similar recolonization processes, there is no reason to expect to detect solely a signature of the population movement induced by the LGM. Indeed, under the admixture model (AM), we have a fairly ancient split between the hybrid zone (HYZ) and Siberian spruce (OBR) *c.* 2 Ma, which was accompanied by much more recent evidence of admixture (*c.* 12 000 years ago) (Table S3, Supporting information).

## Conclusion

The present study demonstrates that the joint use of nuclear and cytoplasmic markers, combined with the use of new approaches such as EEMS and STRUCTURE-based posterior predictives, allows a much finer reconstruction of the joint demographic history of species. In particular, maternally inherited markers alone, such as mitochondrial DNA, failed to reveal the breadth of the hybrid zone between the two species and the direction of introgression could not be inferred from them. Chloroplast DNA did suggest it (Tollefsrud *et al.* 2015) but have the limitation inherent to nonrecombining DNA. What is to be done next? We see at least four main avenues to enrich our understanding of the evolution of spruce species across Eurasia. First, inferences would undoubtedly benefit from a more extensive and representative sampling in Northern Russia and between the Ob and the Yenisei Rivers, both of which are vast and difficult to access. Second, the power of demographic inferences is highly dependent on the number of independent gene genealogies that can be estimated and genomic data should allow much finer demographic reconstructions (Wakeley 2008). Third, those demographic inferences could benefit from the inclusion of ancient DNA (aDNA) samples. Previous studies based on pollen fossils have suggested that spruce populations survived the LGM in Scandinavia (Parducci *et al.* 2012) and new methods are developed to include aDNA in demographic inferences (e.g. Skoglund *et al.* 2014). Another promising avenue, especially in cold areas such as northern Russia and Siberia, is aDNA obtained from sediments, aka environmental DNA, eDNA (Pedersen *et al.* 2015). eDNA could for instance help to characterize the fluctuations of the tree limit through time. Finally, we have focused here on neutral processes and it would be interesting to investigate whether selection also plays a part in the establishment and maintenance of the hybrid zone. This would entail estimating fitness components of individuals with

different levels of hybridization a serious challenge in any organisms, especially long-lived ones such as spruce.

## Acknowledgements

ML and GGV thank FORMAS and the BioDiversa projects Linktree and Tiptree for funding. GGV was also supported by a grant of the European Commission through the FP7-project FORGER (KBBE-289119). Yoshiaki Tsuda was supported by the Japan Society for the Promotion of Science (JSPS). MV acknowledges funding from the Academy of Finland. ML and VS thank people who helped them during sampling along the Ob and Yenisei rivers. We thank Yoshihisa Suyama from Tohoku University, Japan, for providing samples from the Baikal region.

## References

- Abbott R, Albach D, Ansell S *et al.* (2013) Hybridization and speciation. *Journal of Evolutionary Biology*, **26**, 229–246.
- Anderson E (1949) *Introgressive Hybridization*. Wiley, New York, New York.
- Arnold ML (1997) *Natural Hybridization and Evolution*. Oxford University Press, Oxford.
- Bezrukova EV, Abzaeva AA, Letunova PP *et al.* (2005) Post-glacial history of Siberian spruce (*Picea obovata*) in the Lake Baikal area and the significance of this species as a paleo-environmental indicator. *Quaternary International*, **136**, 47–57.
- Binney HA, Willis KJ, Edwards ME *et al.* (2009) The distribution of late-Quaternary woody taxa in northern Eurasia: evidence from a new macrofossil database. *Quaternary Science Reviews*, **28**, 2445–2464.
- Bodare S, Stocks M, Yang J-C, Lascoux M (2013) Origin and demographic history of the endemic Taiwan spruce (*Picea morrisonicola*). *Ecology and Evolution*, **3**, 3320–3333.
- Chen J, Källman T, Gyllenstrand N, Lascoux M (2010) New insights on the speciation history and nucleotide diversity of three boreal spruce species and a Tertiary relict. *Heredity*, **104**, 3–14.
- Chen J, Uebbing S, Gyllenstrand N, Lagercrantz U, Lascoux M, Källman T (2012) Sequencing of the needle transcriptome from Norway spruce (*Picea abies* Karst L.) reveals lower substitution rates, but similar selective constraints in gymnosperms and angiosperms. *BMC Genomics*, **13**, 589.
- Chen J, Tsuda Y, Stocks M *et al.* (2014) Clinal variation at phenology-related genes in spruce: parallel evolution in FTL2 and Gigantea? *Genetics*, **197**, 1025–1038.
- Cinget B, de Lafontaine G, Gérardi S, Bousquet J (2015) Integrating phylogeography and paleoecology to investigate the origin and dynamics of hybrid zones: insights from two widespread North American firs. *Molecular Ecology*, **24**, 2856–2870.
- Csilléry K, François O, Blum MGB (2012) abc: an R package for approximate Bayesian computation (ABC). *Methods in Ecology and Evolution*, **3**, 475–479.
- Currat M, Ruedi M, Petit RJ, Excoffier L (2008) The hidden side of invasions: massive introgression by local genes. *Evolution*, **62**, 1908–1920.
- De La Torre AR, Roberts DR, Aitken SN (2014a) Genome-wide admixture and ecological niche modelling reveal the maintenance of species boundaries despite long history of inter-specific gene flow. *Molecular Ecology*, **23**, 2046–2059.
- De La Torre AR, Wang T, Jaquish B, Aitken SN (2014b) Adaptation and exogenous selection in a *Picea glauca* × *Picea engelmannii* hybrid zone: implications for forest management under climate change. *New Phytologist*, **201**, 687–699.
- De La Torre A, Ingvarsson PK, Aitken SN (2015) Genetic architecture and genomic patterns of gene flow between hybridizing species of *Picea*. *Heredity*, **115**, 153–164.
- Du FK, Petit RJ, Liu JQ (2009) More introgression with less gene flow: chloroplast vs. mitochondrial DNA in the *Picea asperata* complex in China, and comparison with other Conifers. *Molecular Ecology*, **18**, 1396–1407.
- Du FK, Peng XL, Liu JQ, Lascoux M, Hu FS, Petit RJ (2011) Direction and extent of organelle DNA introgression between two spruce species in the Qinghai-Tibetan Plateau. *New Phytologist*, **192**, 1024–1033.
- Dupanloup I, Schneider S, Excoffier L (2002) A simulated annealing approach to define the genetic structure of populations. *Molecular Ecology*, **11**, 2571–2581.
- Earl DA, vonHoldt BM (2012) STRUCTURE HARVESTER: a website and program for visualizing STRUCTURE output and implementing the Evanno method. *Conservation Genetics Resources*, **4**, 359–361.
- Edgar RC (2004) MUSCLE: multiple sequence alignment with high accuracy and high throughput. *Nucleic Acids Research*, **32**, 1792–1797.
- Edwards ME, Armbruster WS, Elias SC (2014) Constraints on post-glacial boreal tree expansion out of far-northern refugia. *Global Ecology and Biogeography*, **23**, 1198–1208.
- Eidesen PB, Ehrich D, Bakkestuen V *et al.* (2013) Genetic roadmap of the Arctic: plant dispersal highways, traffic barriers and capitals of diversity. *New Phytologist*, **200**, 898–910.
- Evanno G, Regnaut S, Goudet J (2005) Detecting the number of clusters of individuals using the software STRUCTURE: a simulation study. *Molecular Ecology*, **14**, 2611–2620.
- Ewens WJ (2004) *Mathematical Population Genetics. I. Theoretical Introduction*. Springer, New York, New York, 417 p.
- Excoffier L, Lischer HEL (2010) Arlequin suite ver 3.5: a new series of programs to perform population genetics analyses under Linux and Windows. *Molecular Ecology Resources*, **10**, 564–567.
- Excoffier L, Dupanloup I, Huerta-Sanchez E, Sousa VC, Foll M (2013) Robust demographic inference from genomic and SNP data. *Plos Genetics*, **10**, e1003905.
- Fedorov VB, Goropashnaya AV, Boeskorov GG, Cook JA (2008) Comparative phylogeography and demographic history of the wood lemming (*Myopus schisticolor*): implications for late Quaternary history of the taiga species in Eurasia. *Molecular Ecology*, **17**, 598–610.
- Giesecke T, Bennett K (2004) The Holocene spread of *Picea abies* (L.) Karst. in Fennoscandia and adjacent areas. *Journal of Biogeography*, **31**, 1523–1548.
- Gompert Z, Parchman TL, Buerkle CA (2012) Genomics of isolation in hybrids. *Philosophical Transactions of the Royal Society of London B: Biological Sciences*, **367**, 439–450.
- Gonser R, Donnelly P, Nicholson G, Di Rienzo A (2000) Microsatellite mutations and inferences about human demography. *Genetics*, **154**, 1793–1807.

- Goudet J (1995) FSTAT (version 1.2): a computer program to calculate F-statistics. *Journal of Heredity*, **86**, 485–486.
- Hamilton JA, Aitken SN (2013) Genetic and morphological structure of a spruce hybrid (*Picea sitchensis* × *P. glauca*) zone along a climatic gradient. *American Journal of Botany*, **100**, 1651–1662.
- Hamilton JA, Lexer C, Aitken SN (2013a) Differential introgression reveals candidate genes for selection across a spruce (*Picea sitchensis* × *P. glauca*) hybrid zone. *New Phytologist*, **197**, 927–938.
- Hamilton JA, Lexer C, Aitken SN (2013b) Genomic and phenotypic architecture of a spruce hybrid zone (*Picea sitchensis* × *P. glauca*). *Molecular Ecology*, **22**, 827–841.
- Haselhorst MSH, Buerkle C (2013) Population genetic structure of *Picea engelmannii*, *P. glauca* and their previously unrecognized hybrids in the central Rocky Mountains. *Tree Genetics & Genomes*, **9**, 669–681.
- Heuertz M, De Paoli E, Källman T *et al.* (2006) Multilocus patterns of nucleotide diversity, linkage disequilibrium and demographic history of Norway spruce [*Picea abies* (L.) Karst]. *Genetics*, **174**, 2095–2105.
- Huber CD, Nordborg M, Hermisson J, Hellmann I (2014) Keeping it local: evidence for positive selection in Swedish *Arabidopsis thaliana*. *Molecular Biology and Evolution*, **31**, 3026–3039.
- Hubisz MJ, Falush D, Stephens M, Pritchard JK (2009) Inferring weak population structure with the assistance of sample group information. *Molecular Ecology Resources*, **9**, 1322–1332.
- Huntley B (2014) Extreme temporal interpolation of sparse data is not a sufficient basis to substantiate a claim to have uncovered Pleistocene forest microrefugia. *New Phytologist*, **204**, 447–449.
- Jakobsson M, Rosenberg NA (2007) CLUMPP: a cluster matching and permutation program for dealing with label switching and multimodality in analysis of population structure. *Bioinformatics*, **23**, 1801–1806.
- Krutovskii KV, Bergmann F (1995) Introgressive hybridization and phylogenetic relationships between Norway, *Picea abies* (L.) Karst., and Siberian, *P. obovata* Ledeb., spruce species studied by isozyme loci. *Heredity*, **74**, 464–480.
- Kuchma O, Vornam B, Finkeldey R (2011) Mutation rates in Scots pine (*Pinus sylvestris* L.) from the Chernobyl exclusion zone evaluated with amplified fragment-length polymorphisms (AFLPs) and microsatellite markers. *Mutation Research*, **725**, 29–35.
- de Lafontaine G, Amasifuen Guerra CA, Ducouso A, Petit RJ (2014a) Cryptic no more: soil macrofossils uncover Pleistocene forest microrefugia within a periglacial desert. *New Phytologist*, **204**, 715–729.
- de Lafontaine G, Amasifuen Guerra CA, Ducouso A, Sanchez-Goni MF, Petit RJ (2014b) Beyond skepticism: uncovering cryptic refugia using multiple lines of evidence. *New Phytologist*, **204**, 450–454.
- de Lafontaine G, Prunier J, Gérardi S, Bousquet J (2015) Tracking the progression of speciation: variable patterns of introgression across the genome provide insights on the species delimitation between progenitor-derivative spruces (*Picea mariana* × *P. rubens*). *Molecular Ecology*, **24**, 5229–5247.
- Lagercrantz U, Ryman N (1990) Genetic structure of Norway spruce (*Picea abies*)—concordance of morphological and allozymic variation. *Evolution*, **44**, 38–53.
- Leblois R, Pudlo P, Néron J *et al.* (2014) Maximum-likelihood inference of population size contractions from microsatellite data. *Molecular Biology and Evolution*, **31**, 2805–2823.
- Li H, Durbin R (2009) Fast and accurate short read alignment with Burrows-Wheeler transform. *Bioinformatics*, **25**, 1754–1760.
- Li H, Handsaker B, Wysoker A *et al.* (2009) The sequence alignment/map format and SAMtools. *Bioinformatics*, **25**, 2078–2079.
- Lockwood JD, Aleksić JM, Zou J, Wang J, Liu J, Renner SS (2013) A new phylogeny for the genus *Picea* from plastid, mitochondrial, and nuclear sequences. *Molecular Phylogenetics and Evolution*, **69**, 717–727.
- Marriage TN, Hudman S, Mort ME, Orive ME, Shaw RG, Kelly JK (2009) Direct estimation of the mutation rate at dinucleotide microsatellite loci in *Arabidopsis thaliana* (Brassicaceae). *Heredity*, **103**, 310–317.
- Mudrik EA, Polyakova TA, Shatokhina AV, Bondarenko GN, Politov DV (2015) Spatial distribution of intron 2 of nad1 gene haplotypes in populations of Norway and Siberian spruce (*Picea abies*–*P. obovata*) species complex. *Russian Journal of Genetics*, **51**, 957–965.
- Parducci L, Jørgensen TJ, Tollefsrud MM *et al.* (2012) Glacial survival of boreal trees in Northern Scandinavia. *Science*, **335**, 1083–1086.
- Pavy N, Namroud M-C, Gagnon F, Isabel N, Bousquet J (2012) The heterogeneous levels of linkage disequilibrium in white spruce genes and comparative analysis with other conifers. *Heredity*, **108**, 273–284.
- Pedersen MW, Overballe-Petersen S, Ermini L *et al.* (2015) Ancient and modern environmental DNA. *Philosophical Transactions of the Royal Society of London. Series B, Biological Sciences*, **370**, 1660.
- Petkova D, Novembre J, Stephens M (2016) Visualizing spatial population structure with estimated effective migration surfaces. *Nature Genetics*, **48**, 94–100.
- Petit RJ, Excoffier L (2009) Gene flow and species delimitation. *Trends in Ecology & Evolution*, **24**, 386–393.
- Pfeiffer A, Olivieri AM, Morgante M (1997) Identification and characterization of microsatellites in Norway spruce (*Picea abies* K.). *Genome*, **40**, 411–419.
- Popov PP (2003) Structure and differentiation of spruce populations in Eastern Europe and Western Siberia. *Russian Journal of Ecology*, **34**, 27–33.
- Popov PP (2010) Form structure and geographic differentiation of spruce populations in northwestern Russia. *Russian Journal of Ecology*, **41**, 336–343.
- Potokina EK, Orlova LV, Vishnevskaya MS, Alekseeva EA, Potokin AF, Egorov AA (2013) Genetic differentiation of spruce populations in Northwest Russia according to the results of microsatellite loci analysis. *Russian Journal of Genetics: Applied Research*, **3**, 352–360.
- Potokina EK, Kiseleva AA, Nikolaeva MA, Ivanov SA, Ulianich PS, Potokin AF (2015) Analysis of the polymorphism of organelle DNA to elucidate the phylogeography of Norway spruce in the East European Plain. *Russian Journal of Genetics: Applied Research*, **4**, 430–439.
- Pritchard J, Stephens M, Donnelly P (2000) Inference of population structure using multilocus genotype data. *Genetics*, **155**, 945–959.

- Rieseberg LH, Buerkle CA (2002) Genetic mapping in hybrid zones. *The American Naturalist*, **159**, S36–S50.
- Rosenberg NA (2004) DISTRUCT: a program for the graphical display of population structure. *Molecular Ecology Notes*, **4**, 137–138.
- Rungis D, Bérubé Y, Zhang J, et al. (2004) Robust simple sequence repeat markers for spruce (*Picea spp.*) from expressed sequence tags. *Theoretical and Applied Genetics*, **109**, 1283–1294.
- Schmidt-Vogt H (1974) Das natürliche Verbreitungsgebiet der Fichte (*Picea abies* [L.] Karst.) in Eurasien. *Allgemeine Forst- und Jagdzeitung*, **145**, 185–197.
- Scotti I, Magni F, Paglia GP, Morgante M (2002a) Trinucleotide microsatellites in Norway spruce (*Picea abies*): their features and the development of molecular markers. *Theoretical and Applied Genetics*, **106**, 40–50.
- Scotti I, Paglia P, Magni G, Morgante M (2002b) Efficient development of dinucleotide microsatellite markers in Norway spruce (*Picea abies* Karst.) through dot-blot selection. *Theoretical and Applied Genetics*, **104**, 1035–1041.
- Semerikov V, Lascoux M (2003) Nuclear and cytoplasmic variation within and between Eurasian *Larix* (Pinaceae) species. *American Journal of Botany*, **90**, 1113–1123.
- Semerikov VL, Semerikova SA, Polezhaeva MA, Kosintsev PA, Lascoux M (2013) Southern montane populations did not contribute to the recolonization of western Siberian plains by Siberian larch (*Larix sibirica*): evidence from the first range-wide analysis of cytoplasmic markers. *Molecular Ecology*, **22**, 4958–4971.
- Skoglund P, Sjödin P, Skoglund T, Lascoux M, Jakobsson M (2014) Investigating population history using temporal genetic differentiation. *Molecular Biology and Evolution*, **31**, 2516–2527.
- Stewart JR, Lister AM (2001) Cryptic northern refugia and the origins of the modern biota. *Trends in Ecology & Evolution*, **16**, 608–613.
- Tollefsrud MM, Kissling R, Gugerli F et al. (2008) Genetic consequences of glacial survival and postglacial colonization in Norway spruce: combined analysis of mitochondrial DNA and fossil pollen. *Molecular Ecology*, **17**, 4134–4150.
- Tollefsrud MM, Sønstebo JH, Brochmann C, Johnsen Ø, Skråppa T, Vendramin GG (2009) Combined analysis of nuclear and mitochondrial markers provide new insight into the genetic structure of North European *Picea abies*. *Heredity*, **102**, 549–562.
- Tollefsrud MM, Latalowa M, van der Knaap WO, Brochmann C, Sperisen C (2015) Late Quaternary history of North Eurasian Norway spruce (*Picea abies*) and Siberian spruce (*Picea obovata*) inferred from macrofossils, pollen and cytoplasmic DNA variation. *Journal of Biogeography*, **42**, 1431–1442.
- Tzedakis PC, Emerson BC, Hewitt GM (2013) Cryptic or mystic? Glacial tree refugia in northern Europe. *Trends in Ecology & Evolution*, **28**, 696–704.
- Väliranta M, Kaakinen A, Kuhry P, Kultti S, Salonen JS, Seppä H (2011) Scattered late-glacial and early Holocene tree populations as dispersal nuclei for forest development in north-eastern European Russia. *Journal of Biogeography*, **38**, 922–932.
- Volkova P, Shipunov A, Borisova P, Moseng R, Ivens I (2014) In search of hybridity: the case of Karelian spruces. *Silva Fennica*, **48**, 2.
- Wakeley J (2008) *Coalescent Theory. An Introduction*. Roberts and Co, Greenwood Village, Colorado.
- Wang J, Abbott RJ, Peng YL, Du FK, Liu JQ (2011) Species delimitation and biogeography of two fir species (*Abies*) in central China: cytoplasmic DNA variation. *Heredity*, **107**, 362–370.
- Weir BS, Cockerham CC (1984) Estimating F-statistics for the analysis of population structure. *Evolution*, **38**, 1358–1370.
- Willyard AW, Syring J, Gernandt DS, Liston A, Cronn R (2007) Fossil calibration of molecular divergence infers a moderate mutation rate and recent radiations for pinus. *Molecular Biology and Evolution*, **24**, 90–101.

---

M.L., M.M.T., Y.T. and G.G.V. conceived and designed the study; Y.T. performed the SSR laboratory work and part of the SSR analyses. M.M.T., C.S., L.P. and Y.T. generated and analysed the mtDNA data. J.C. performed the SSR ABC analysis. J.H.S. contributed to the analysis of the SSR data. M.S. generated the Norway spruce nuclear sequences and carried out sequence data analysis; D.P., V.S., M.V. and T.R. contributed materials. All co-authors discussed the results. M.L. and M.M.T. wrote the manuscript with critical input from other co-authors.

---

## Data accessibility

Data have been deposited to the Dryad Digital Repository Dryad with 'doi: 10.5061/dryad.6bf38'.

## Supporting information

Additional supporting information may be found in the online version of this article.

### Appendix S1 Methods and results.

**Fig. S1** The correlation of summary statistics and sample sizes using the rarefaction test.

**Fig. S2** The tests of the most likely number of clusters in the STRUCTURE analysis using SSR loci based on  $\Delta K$  (Evanno et al. 2005) or  $\ln P(D)$  (Pritchard et al. 2000).

**Fig. S3** The posterior prediction of summary statistics based on the IA model.

**Fig. S4** The posterior prediction of summary statistics based on the AM model.

**Fig. S5** The proportion of each inferred genetic cluster for four subpopulation groups based on STRUCTURE analysis ( $K = 4$ ) using simulated dataset under the IA model.

**Fig. S6** The proportion of each inferred genetic cluster for four subpopulation groups based on STRUCTURE analysis ( $K = 5$ ) using simulated dataset under the IA model.

**Fig. S7** The proportion of each inferred genetic cluster for four subpopulation groups based on STRUCTURE analysis ( $K = 4$ ) using simulated dataset under the AM model.

**Fig. S8** The proportion of each inferred genetic cluster for four subpopulation groups based on STRUCTURE analysis ( $K = 5$ ) using simulated dataset under the AM model.

**Fig. S9** Posterior distributions for parameters estimated under the Isolation model (A) and the Isolation-with-migration model (B).

**Table S2** SSR loci polymorphism for all four subpopulations: ALP (pop. 1–22), NOR (pop. 23–64), HYZ (pop. 65–81), and OBR (pop. 82–102).

**Table S3** The mode and 95% confident interval of marginal posterior distribution for parameters estimated under the IA and AM models using 10 SSR loci.

**Table S4** Observed summary statistics for the total number of segregating sites ( $S$ ) and confidence intervals (2.5%, 97.5%) and means for Watterson's theta ( $\theta_W$ ), the average number of pairwise nucleotide differences ( $\theta_{\pi}$ ) and Tajima's  $D$  calculated both within *Picea abies* and *Picea obovata* and across both species.

**Table S5** Median parameter estimates and confidence intervals (2.5%, 97.5%) for the neural network ABC analysis on sequence data.

**Table S1** Summary for references of datasets of nuclear and mitochondrial (mt) DNA variation of 102 populations examined in this study.