

ORIGINAL RESEARCH

## Exploring the evolution of the proteins of the plant nuclear envelope

Axel Poulet<sup>a,b</sup>, Aline V. Probst<sup>b</sup>, Katja Graumann<sup>a</sup>, Christophe Tatout <sup>b</sup>, and David Evans <sup>a</sup>

<sup>a</sup>Department of Biological and Medical Sciences, Oxford Brookes University, Oxford, UK; <sup>b</sup>UMR CNRS 6293 INSERM U1103 Clermont Université, GREM, Aubière, France

### ABSTRACT

In this study, we explore the plasticity during evolution of proteins of the higher plant nuclear envelope (NE) from the most ancestral plant species to advanced angiosperms. The higher plant NE contains a functional Linker of Nucleoskeleton and Cytoskeleton (LINC) complex based on conserved Sad1-Unc84 (SUN) domain proteins and plant specific Klarsicht/Anc1/Syne homology (KASH) domain proteins. Recent evidence suggests the presence of a plant lamina underneath the inner membrane and various coiled-coil proteins have been hypothesized to be associated with it including Crowded Nuclei (CRWN; also termed LINC and NMCP), Nuclear Envelope Associated Protein (NEAP) protein families as well as the CRWN binding protein KAKU4. SUN domain proteins appear throughout with a key role for mid-SUN proteins suggested. Evolution of KASH domain proteins has resulted in increasing complexity, with some appearing in all species considered, while other KASH proteins are progressively gained during evolution. Failure to identify CRWN homologs in unicellular organisms included in the study and their presence in plants leads us to speculate that convergent evolution may have occurred in the formation of the lamina with each kingdom having new proteins such as the Lamin B receptor (LBR) and Lamin-Emerin-Man1 (LEM) domain proteins (animals) or NEAPs and KAKU4 (plants). Our data support a model in which increasing complexity at the nuclear envelope occurred through the plant lineage and suggest a key role for mid-SUN proteins as an early and essential component of the nuclear envelope.

### ARTICLE HISTORY

Received 29 June 2016  
Revised 22 August 2016  
Accepted 7 September 2016

### KEYWORDS



Chromatin; higher plant; LINC complex; KASH domain; nucleoskeleton; nucleus; SUN domain

### Introduction


The nuclear envelope is a key component of eukaryotic cells and may be considered to be composed of 3 elements, the nuclear membrane, nuclear pore complexes and the nuclear lamina.<sup>13,19</sup> These structural components are essential for many processes including nuclear morphology, nuclear migration, chromatin organization and regulation of gene expression.<sup>15</sup> Significant progress has been made in describing novel plant nuclear envelope proteins.<sup>29,35,43</sup> In *Arabidopsis thaliana* (*A. thaliana*), these include components of the Linker of Nucleoskeleton and Cytoskeleton (LINC) complex for which functional data is slowly being revealed. *Arabidopsis* contains proteins of the inner nuclear envelope of the SUN domain family including Cter-Sad1-Unc84 (Cter-SUN)<sup>16,17,28</sup> as well as mid-SUN domain proteins in which a SUN-domain homologous to that of the C-ter SUNs is located

centrally within the protein.<sup>18</sup> It also contains proteins of the outer nuclear envelope, of the KASH domain protein family including WPP Domain Interacting Proteins [WIPs], SUN interacting Nuclear Envelope Proteins [SINs] and *Arabidopsis thaliana* Toll Interleukin Receptor domain KASH protein [TIK].<sup>41,18,40</sup> In addition, plant proteins proposed to form the nuclear lamina - Crowded Nuclei (CRWNs;<sup>9,37</sup>) and CRWN-interacting proteins such as KAKU4<sup>14</sup>) as well as Nuclear Envelope Associated Proteins (NEAPs), which may be associated with the lamina,<sup>30</sup> have been described in *A. thaliana* and shown to localize to the nuclear periphery.

Sequence data now available permits comparison of components of the nuclear envelope between algae, mosses, gymnosperms and angiosperms with the components of *A. thaliana*. Functional analysis of these genes is challenging because they belong to small gene

**CONTACT** David Evans  [deevans@brookes.ac.uk](mailto:deevans@brookes.ac.uk)  Faculty of Health and Life Sciences, Oxford Brookes University, Headington Campus, Oxford OX3 0BP UK.

Color versions of one or more of the figures in this article can be found online at [www.tandfonline.com/kncl](http://www.tandfonline.com/kncl).

 Supplemental data for this article can be accessed on the [publisher's website](#).

families as a consequence of gene and whole-genome duplication (WGD) creating duplicate genes and thus gene redundancy.<sup>12,33</sup> Whole-genome duplication (WGD) is recognized as an important event for genome evolution in animals, plants and fungi and to drive key new features, with resulting increased complexity and speciation.<sup>33</sup> Following WGD, massive gene loss can occur restoring the diploid state for most duplicated loci while few duplicated genes remain and may provide new evolutionary innovation including structures (e.g. floral organs) and adaptations.<sup>21</sup> Previous analyses of plant genomes have shown that seed plants share an ancient WGD event, zeta.<sup>20</sup> A second WGD, epsilon, has been detected shortly before the diversification of angiosperms. These two WGDs are suggested to play a role in the origin and rapid diversification of the angiosperms.<sup>20</sup> Finally, the gamma WGD occurred after eudicot/monocot diversification, followed by several partial or complete duplication events.

In this study, we have selected 20 representative species based on the revised classification of eukaryotes,<sup>1</sup> their available genome sequences and gene expression description, in order to explore the evolution of components of the plant nuclear envelope. Availability of genome sequence data for the core eudicot *Amborella trichopoda* provides an opportunity to explore the nuclear envelope of a primitive angiosperm. *Amborella*, a New Caledonian shrub, has been suggested as the sole surviving sister species of all other angiosperms and is unique in sequenced plant genomes in showing no evidence of recent, lineage-specific genome duplications.<sup>32</sup> The *Amborella* genome therefore offers an opportunity to explore the composition of an ancestral plant nuclear envelope and the effect of genomic changes after polyploidy in other angiosperms.<sup>32</sup> Finally, RNAseq data for each species were used to describe expression levels within species to establish that the genes are active and not pseudogenes and to demonstrate gene activity and when possible tissue specific expression patterns.

The aim of the work presented in this paper was to explore the evolution of nuclear envelope proteins in unicellular algae and multicellular plants and to provide evidence for the composition of the simplest functional plant LINC complex. This study provides valuable information for mutant and other functional studies by identifying potential redundancy and

specialization in nuclear envelope and lamina-like components.

## Material and methods

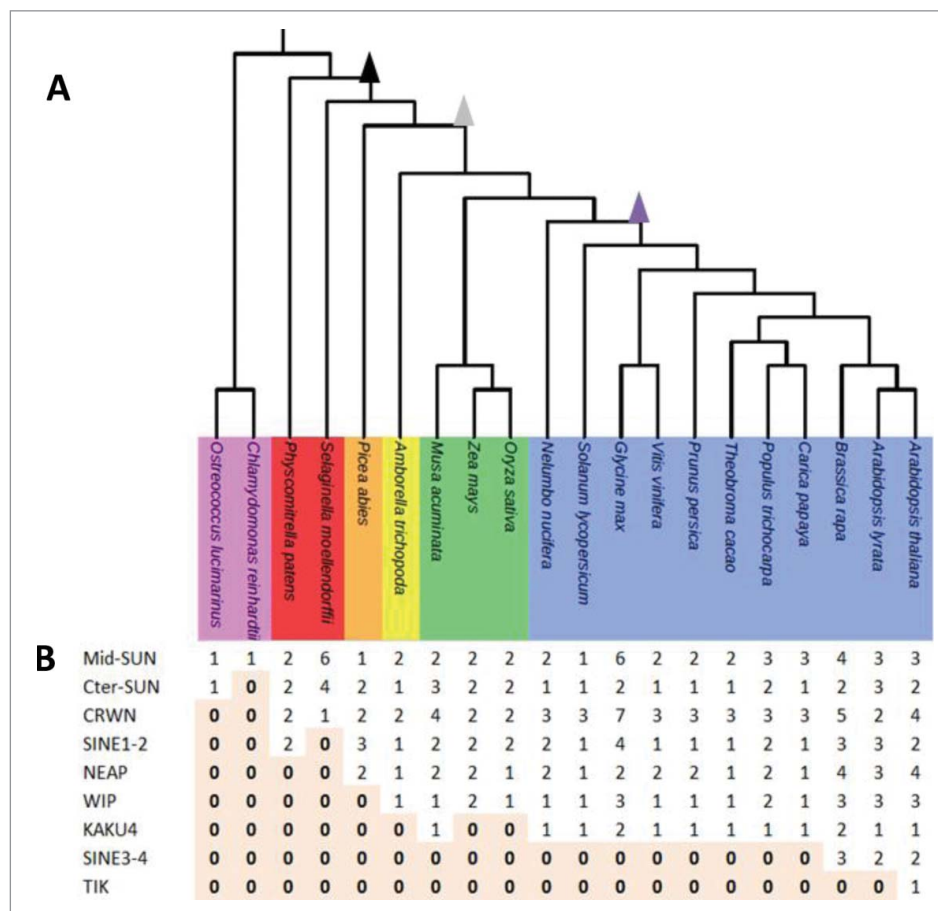
### Homologous LINC complex and lamin-like protein detection

A Perl script was developed and applied to proteomic data to identify KASH domain proteins. The program tests the presence of the trans-membrane (TM) domain and 4 specific amino acids at the C-terminus, which are characteristic for KASH proteins. The position of the TM domain is variable and the script searches this TM domain up to 40 amino acids away from the KASH-specific C-terminal motifs detected in *A. thaliana* (either VIPT, VVPT, AVPT, PLPT, TVPT, LVPT or PPPS.<sup>41,18,42</sup>). The identification of the TM domain is based on the Kyte-Doolittle method.<sup>25</sup> Only proteins, which possess a TM domain and the 4 KASH specific amino acids in the C-terminus of the protein, were selected.

For all proteins of interest a Basic Local Alignment Search Tool protein (BLASTp) was used with default parameters as well as HMMER (Hidden Markov Model-based sequence alignment tool; <http://hmmer.org>). The best hits were retained and used for phylogenetic analysis.<sup>1990</sup> The proteome of each species was used as reference for the BLASTp (Fig. 1'), and the protein sequences of the LINC complex as well as the putative lamina of *A. thaliana* were used as queries (Supplementary Table 1). BLASTp results are given as supplementary Table 2 (mid-SUN), 3 (Cter-SUN), 4 (WIP), 5 (SINE), CRWN (6), NEAP (7) and KAKU4 (8). Reciprocal BLASTp was used to verify the relevance of all identified orthologs.

### Phylogenetic reconstruction

Selected sequences were first aligned with MUSCLE, a multiple sequence alignment tool,<sup>11</sup> using default parameters. The alignment was then refined using Gblocks<sup>34</sup> Fast-Tree was then applied with default parameters, for the construction of the phylogenetic tree.<sup>31</sup> Fast-Tree infers approximately-maximum-likelihood phylogenetic trees from alignments. Finally, phylogenetic trees were drawn using the Interactive Tree Of Life ITOL.<sup>26</sup>



**Figure 1.** Distribution of components of plant nuclear envelope in the plant kingdom. (A) Selected plant lineages used in this study from left to right: Unicells Algae (pink), Moss and Club Moss (red), Gymnosperm (orange), Basal Angiosperms (yellow), Monocots (green) and Eudicots (blue). Zeta epsilon and gamma WGDs are indicated as arrow heads respectively in black, gray and purple. (B) Distribution of the 9 protein families (rows) in the 20 species (columns). Absence (0) of a given protein is highlighted in light orange.

### RNA sequencing data and analysis

Data used for the RNA-seq analysis was obtained from the NCBI (<http://www.ncbi.nlm.nih.gov/geo/browse/>) or from the Amborella Genome Database, respectively (<http://amborella.huck.psu.edu/>). Five different tissues (leaves, roots, flowers, flower buds, and seeds/siliques) as well as total seedling were chosen for the analysis of the expression patterns of the genes of interest (Supplementary Table 1). The expression was analyzed for 10 species (Supplementary Table 9). Reads from RNA-Seq libraries were mapped onto the candidate gene sequences allowing no mismatches using TOPHAT v 2.0.14<sup>22</sup> with standard settings and maximum of multihits set at 1, minimum intron length set at 15 bp, and maximum intron length set as 6,000 bp. Reads were added together for each gene using HTseq-count with the overlap resolution mode set as intersection-non empty and with no strand-specific protocol.<sup>3</sup> Transcription levels in Reads Per Kilobase

of transcript per Million mapped reads (RPKM) were normalized to *AtSAND* (At2g28390;<sup>7</sup> and Supplementary Table 10). *SAND* was chosen due to its constant gene expression levels across different tissues at developmental stages in *Arabidopsis thaliana*.<sup>7</sup> For each species, the *SAND* homolog with the closed sequence identity to *AtSAND* was chosen. Furthermore, absolute *SAND* expression levels (in RPKM) in different species were comparable to expression of *Arabidopsis AtSAND*.

### Results and discussion

In order to gain an insight into the evolutionary development of known plant nuclear envelope proteins, we reconstructed the phylogenetic distribution of the LINC complex (SUN and KASH) and plant lamina (CRWN, NEAP and KAKU4) components by exploring 20 representative species including unicellular photosynthetic algae, lycophytes, mosses,

gymnosperms and angiosperms for which genome sequences and gene expression data are available (Fig. 1; Supplementary Table 11).

### Phylogenetic analysis of inner nuclear membrane proteins

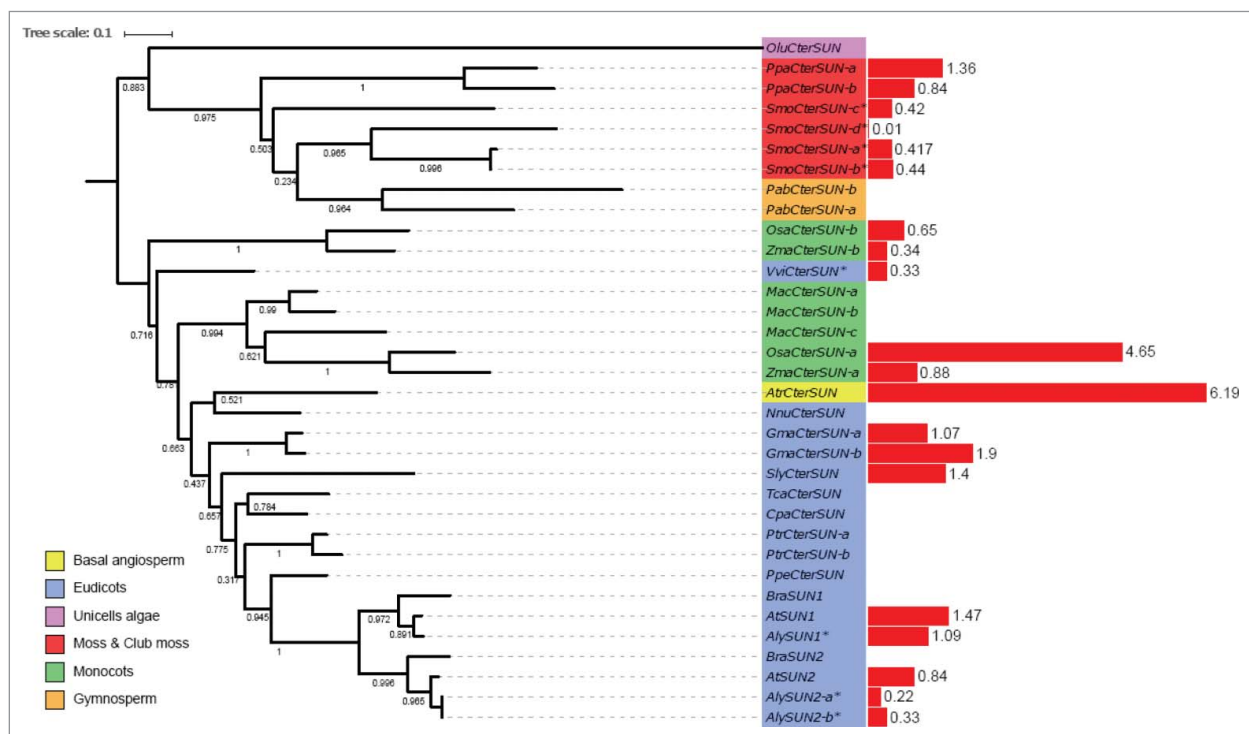
#### Cter-SUN proteins

The SUNs are divided into 2 subfamilies according to the position of the SUN domain: the mid-SUN with a central SUN domain and the Cter-SUNs having a SUN domain at the C-terminus. The potential origin of the 2 classes of SUN domain proteins remains obscure and is discussed in Graumann et al., 2015. The SUNs are key members of the LINC complex and expressed in all tissues.<sup>27,18</sup> Blastp and HMMER analysis revealed 3three Cter-SUN proteins across all the 19 species studied, other than *Chlamydomonas reinhardtii*, where no Cter-SUN protein was detected (Fig. 1). Monocots and eudicots form 2 paraphyletic groups, with the *Vitis vinifera* homolog showing greater similarity to the monocot Cter-SUN sequence. The *Brassicaceae* form a monophyletic group and the

duplication of the Cter-SUN gene seems to have occurred late in evolution because duplicated Cter-SUNs remain grouped within a given species (Fig. 2).

Expression data for the Cter-SUNs shows a similar transcript level for all the tissue analyzed in different species (Supplementary Fig. 1). In some cases, one of the 2 Cter-SUNs is more strongly expressed in the seedling (e.g.,: AtSUN1 more strongly expressed than AtSUN2; OsaSUN-a more strongly expressed than OsaSUN-b) ((WIPs)). *A. trichopoda* encodes only one Cter-SUN that is highly expressed in all tissues. The simplest functional LINC complex may therefore be based on a single Cter-SUN, and strengthens the suggestion that duplication of the Cter-SUN gene occurred after speciation.

One or 2 Cter-SUN proteins were identified in most plants and the moss, although 4 close homologues were identified for the club moss *Selaginella molen-dorfii*. In *A. thaliana*, SUN1 and SUN2 share almost the same activity and localization.<sup>17</sup> This is in contrast to mammals, where 5 Cter-SUN orthologues have clearly differentiated functions. It appears that the gene duplication resulting in these orthologues



**Figure 2.** Phylogenetic tree of Cter-SUN proteins and gene expression levels. Left: maximum likelihood tree of Cter-SUN protein homologues constructed from an alignment. Bootstrap values are presented. The color of the label shows the lineage of the plant. The gene label is constructed with the 3 letters from the species name (supplementary Table 4) and the gene name of the *A. thaliana* homologues. Right: red bar represents the value of the transcription level in seedlings expressed in RPKM, except for species indicated by \*, the RNA-seq data was obtained from leaf tissue (Supplementary Table 2).

occurred earlier in the evolution of mammals. One likely consequence is the lack of specificity of function of plant Cter-SUN homologues; for example, a disruption of a single SUN gene results in an infertility phenotype in animals,<sup>8</sup> but in *A. thaliana*, a single Cter-SUN deletion does not affect meiosis or fertility whereas the double mutant *atsun1 atsun2* impacts fertility and cell division.<sup>36</sup> This suggests a significant redundancy in Cter-SUN function in plants and that double knock-out or knock-down mutants are required for recognizable phenotypes to be obtained.

### mid-SUN proteins

All the species considered contain at least a mid-SUN protein and overall 50 mid-SUN homologues were identified during our bioinformatic screen. The mid-SUN angiosperm homologues are clustered in 2

groups, SUN3/SUN4 and SUN5. In each mid-SUN homologous group, the basal angiosperm, monocots and eudicots form monophyletic groups. This suggests that mid-SUN (3, 4) gene duplication occurred after speciation between angiosperms and gymnosperms (Fig. 3). In all tissue analyzed the *SUN3/SUN4* group tends to be more ubiquitously and highly expressed than the *SUN5* group, this is also true for the *A. trichopoda* homologues (Supplementary Fig. 1). It has been suggested that *AtSUN5* has a meiotic function<sup>18</sup> and this is also true for maize with *ZmaSUN5*,<sup>27</sup> although while the double mutants of SUN3, SUN4 and SUN5 are viable, a *sun3 sun4 sun5* triple mutant is lethal.<sup>18</sup> *A. trichopoda* has 2 mid-SUN proteins, one SUN3/SUN4 homolog and a SUN5 homolog. This suggests that the simplest LINC complex has 2 mid-SUNs each with a specific or partially overlapping function.

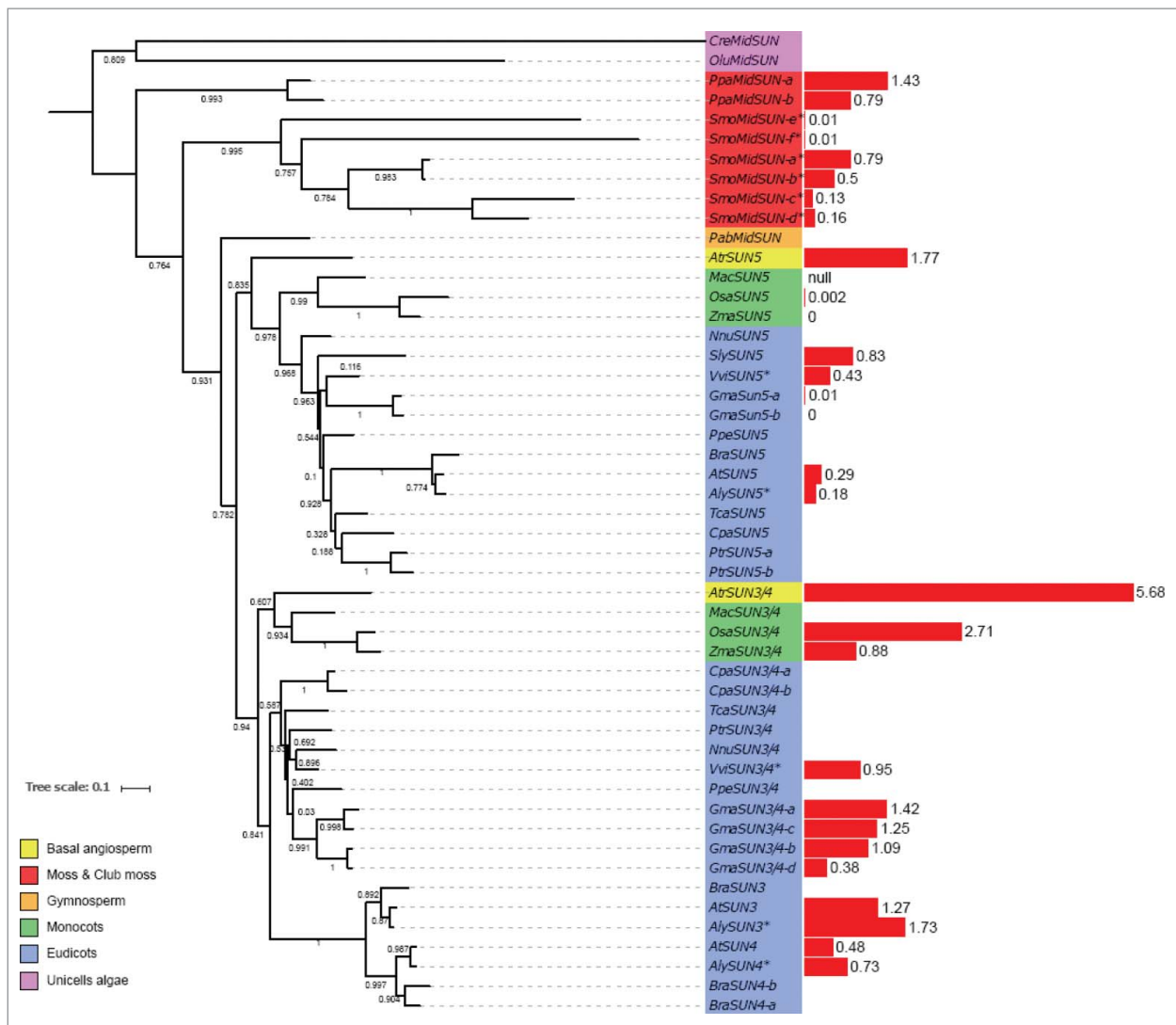


Figure 3. Phylogenetic tree of mid-SUN proteins.

In summary, the majority of the 20 species possess at least one mid-SUN and one Cter-SUN protein except for *Chlamydomonas reinhardtii*, which has only one mid-SUN protein. Interestingly, in common with Cter SUNs, the club moss *Selaginella* has the highest number of mid-Sun protein homologues (6). These results are in good agreement with previous studies that have highlighted the conservation of both Cter- and mid-SUN proteins in most eukaryotes<sup>27,18</sup> and suggest that the LINC complex was present in the Last Evolutionary Common Ancestor (LECA;<sup>23</sup>). Mid-SUN homologues and Cter-SUN proteins were detected in the unicellular algae examined, suggesting that SUN emergence pre-dates the evolution of multicellularity. The evolutionary relationship between Cter-SUN and mid-SUN proteins has yet to be described. This study suggests that SUN domain proteins may also be among the earliest evolving components of the plant nuclear envelope and that mid-SUN proteins may have significance nuclear function in the absence of Cter-SUN.

### **KASH protein homologues**

KASH proteins are diverse in sequence and structure<sup>40</sup> but possess a conserved C-terminal region with a TM domain and a conserved motif of 4 amino acids at the extreme C-terminus. For the detection of the KASH domain protein homologues, 2 strategies were used. The first was a BLASTp analysis based on known *A. thaliana* KASH domain proteins (Supplementary Table 1). This analysis permitted detection of KASH protein homologues in all the organisms studied, except for the unicellular algae where no KASH protein was detected (Supplementary Table 12). Using this method, 32 SINE homologues were found, whereas WIP and potential TIK proteins [or TIK-like proteins] (Supplementary Table 12) were much less common and were found mainly in eudicots. An exception is *Brassicaceae*, where several potential WIP (3 in *A. lyrata*, 4 in *Brassica rapa*) and TIK (2 in *A. lyrata* and 1 in *B. rapa*) homologues were identified, these were also detected in *Glycine max* (2 WIP, 1 TIK), *Prunus persica* (1 TIK), *Carica papaya* (1 WIP), *Musa acuminata* (1 WIP), *A. trichopoda* (1 TIK) and the gymnosperm *Picea abies* (1 TIK) (Supplementary Table 12). To expand the data collected by Blastp, a script was developed to detect proteins with the TM domain and C-terminal motif.

All the identified plant KASH domain proteins have been divided into 3 groups: SINEs, WIPs and TIK.<sup>40</sup> Six KASH protein clusters were revealed (Supplementary Fig. 2). One includes WIP proteins detected in the monocotyledons and the basal angiosperms (Supplementary Table 12), as well as 7 new putative WIP proteins to those detected previously by BLASTp. For SINE proteins, 3 clusters were detected, for SINE1/2, SINE3 and SINE4 adding respectively 2, 6 and 12 SINE proteins to those already identified. The high number of proteins in the SINE3 and SINE4 cluster found only by the script was due to weak conservation of these proteins. One much smaller cluster includes the TIK-like proteins. Only four putative homologues were added but these were shown subsequently to lack either the TIR domain or the C-terminal TM domain, therefore suggesting that the TIK protein<sup>18</sup> may be unique to *A. thaliana*. An additional cluster (other) had low sequence similarity and was not included subsequently. The three WIP proteins in *A. thaliana* show previously described properties<sup>38,39</sup> of a cytoplasmic domain at the N-terminus, with AtWIP1 and AtWIP2 having 3 coiled-coil domains but AtWIP3 only one. The C-terminal region is well conserved and the coiled-coil domains align, with all proteins detected as homologues having a C-terminal predicted TM domain and KASH motif, except AlyWIP1, which lacks homology at the C-terminal region but is well conserved at the N-terminus.

All SINEs have a typical KASH TM domain and C-terminal amino acid motif.<sup>42</sup> The SINE gene family comprises 4 genes in *A. thaliana*, with similarity between AtSINE1 and AtSINE2, characterized by an Armadillo repeat domain near the N-terminus; and between AtSINE3 and AtSINE4. The Perl script added only 2 sequences in the SINE1/SINE2 group while it added 14 proteins to the SINE3/SINE4 cluster. In this case the Blastp approach was less efficient than the Perl script because of the absence of well-conserved domains in the N-terminus. After removal of sequences with the lowest similarity or without the conserved domain, SINE1/SINE2 proteins are present in all species except the unicellular algae and club moss, while SINE3/SINE4 were absent (Fig. 5). In summary, the SINE1/SINE2 cluster and WIP proteins are detected in basal angiosperms whereas the TIK protein is detected only in *A. thaliana*.

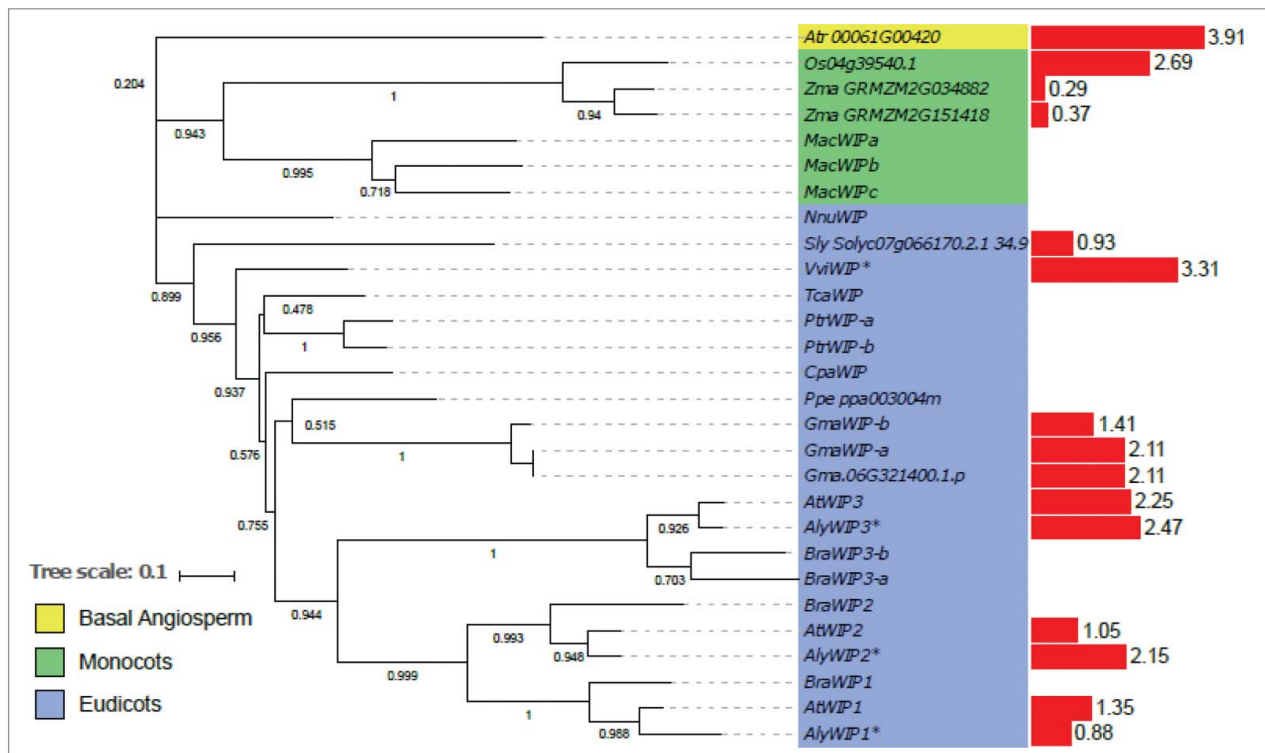


Figure 4. Phylogenetic tree of WIP proteins.

### Phylogenetic analysis of the outer nuclear membrane proteins

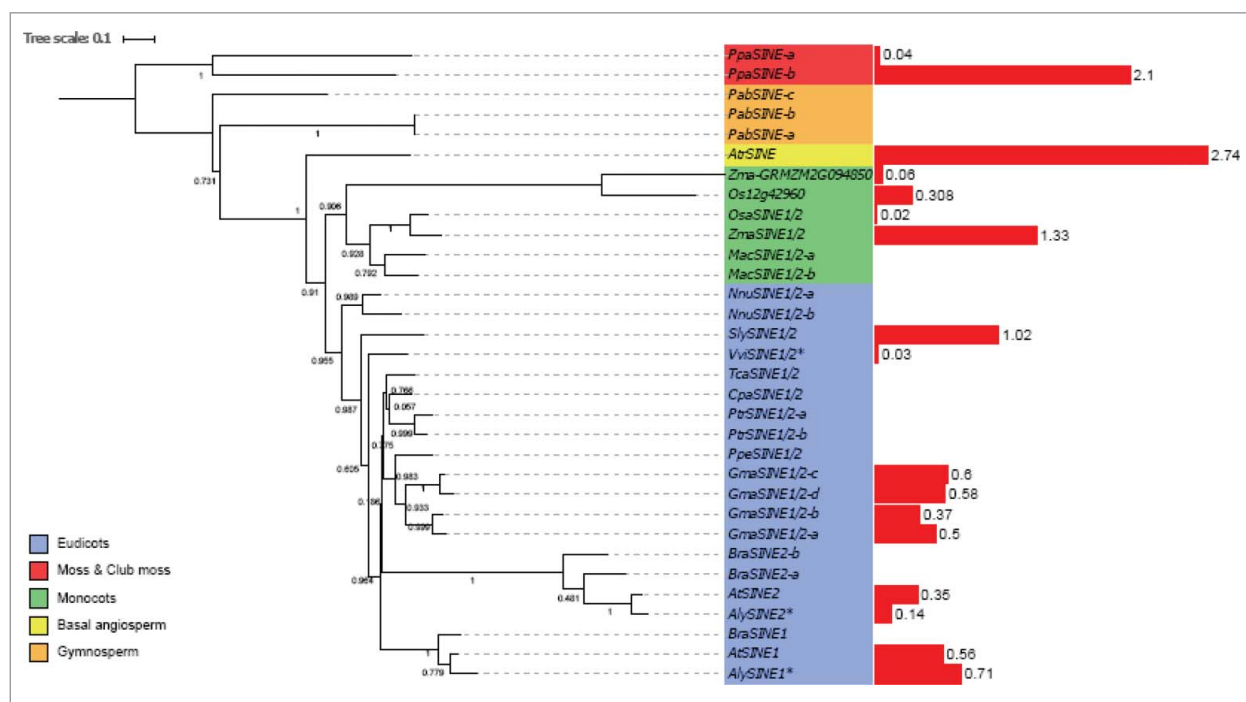
#### WIP proteins

The WIP protein family was the first KASH family detected in *A. thaliana*<sup>41</sup>). WIP proteins were not detected in unicellular algae, moss, club moss or gymnosperms; suggesting that they are angiosperm specific proteins. One WIP homolog was detected for *A. trichopoda*. The monocots form a monophyletic group, with one protein for rice, 2 and 3 for maize and *Musa acuminata* suggesting gene duplication (Fig. 4). The eudicots form a paraphyletic group because the WIP homolog of *Nelumbo nucifera* differs from, and is positioned outside, the WIPs of eudicots. The *Brassicaceae* on the other hand, form a monophyletic group (Fig. 4). This suggests that an ancestral duplication in the *Brassicaceae* ancestor gave rise to WIP1/WIP2 and WIP3, and then WIP1 and WIP2 resulted from a more recent gene duplication. All three genes are expressed in all the tissues analyzed. In *A. thaliana* *AtWIP3* transcripts are more abundant than *AtWIP1* and *AtWIP2* in all tissues. This may be due to redundancy in *AtWIP1* and *AtWIP2* function, and in *A. trichopoda*, the WIP homolog is highly expressed (Supplementary Fig. 1).

#### SINE proteins

SINEs in *A. thaliana*<sup>42</sup> comprise 2 groups, SINE1/SINE2 and SINE3/SINE4. *AtSINE1* is more expressed in guard cells, and its armadillo domain forms F-actin-associated fibers involved in nuclear positioning while *AtSINE2* is suggested to be involved in the immunity response of leaves.<sup>42</sup> No expression and activity data was available for *AtSINE3* and *AtSINE4*.

SINE1/SINE2 proteins were not found in unicellular algae and in club moss, but in contrast to WIPs, 2 and 3 SINE homologues were found in moss and gymnosperms, respectively (Fig. 5). The angiosperms form a monophyletic group and one SINE1/SINE2 homolog was detected for *A. trichopoda* and positioned at the base of the angiosperm group (Fig. 5). The phylogenetic analysis of SINE3 and SINE4 is not possible due to the low similarity between sequences and a lack of conserved domains. Although SINE3 and SINE4 are detected in the *Brassicaceae* group, the other sequences are divergent. In the monocots, 2 protein homologues were detected for *Musa acuminata*, *Oryza sativa* and *Zea mays*. However, the phylogeny suggests the presence of recent gene duplication in *Musa*



**Figure 5.** Phylogenetic tree of SINE1, SINE2 homologues proteins.

*acuminata* (Fig. 5). In contrast, the gene duplication between the 2 other monocots seems to have occurred before their speciation. All the eudicots possess at least one SINE1/SINE2 homolog. Four homologues that group together were found in *Glycine max*, suggesting a recent gene duplication. As for WIPs, *Brassicaceae* proteins cluster together, and one group of homologues is detected for each of SINE1 and SINE2. The organization between the 2 groups suggests a gene duplication to form SINE1 and SINE2.

In *A. thaliana*, *AtSINE1* and *AtSINE2* are expressed at the same level in all tissues, but at a higher level than *AtSINE3* and *AtSINE4*. However, SINE1/SINE2 homologues in most other species show the lowest level of expression of all KASH proteins for all tissues analyzed expect for maize, rice and *A. trichopoda* (Supplementary Fig. 1). In these species WIP and SINE expression is at the same level for all tissues. In *A. thaliana*, *AtWIPs* are more highly expressed than *AtSINE*.

#### Phylogenetic analysis of the putative nuclear lamina and nuclear-envelope associated proteins

Proteins of the lamin family are restricted to animals.<sup>5</sup> However, to date, 3 protein families have been

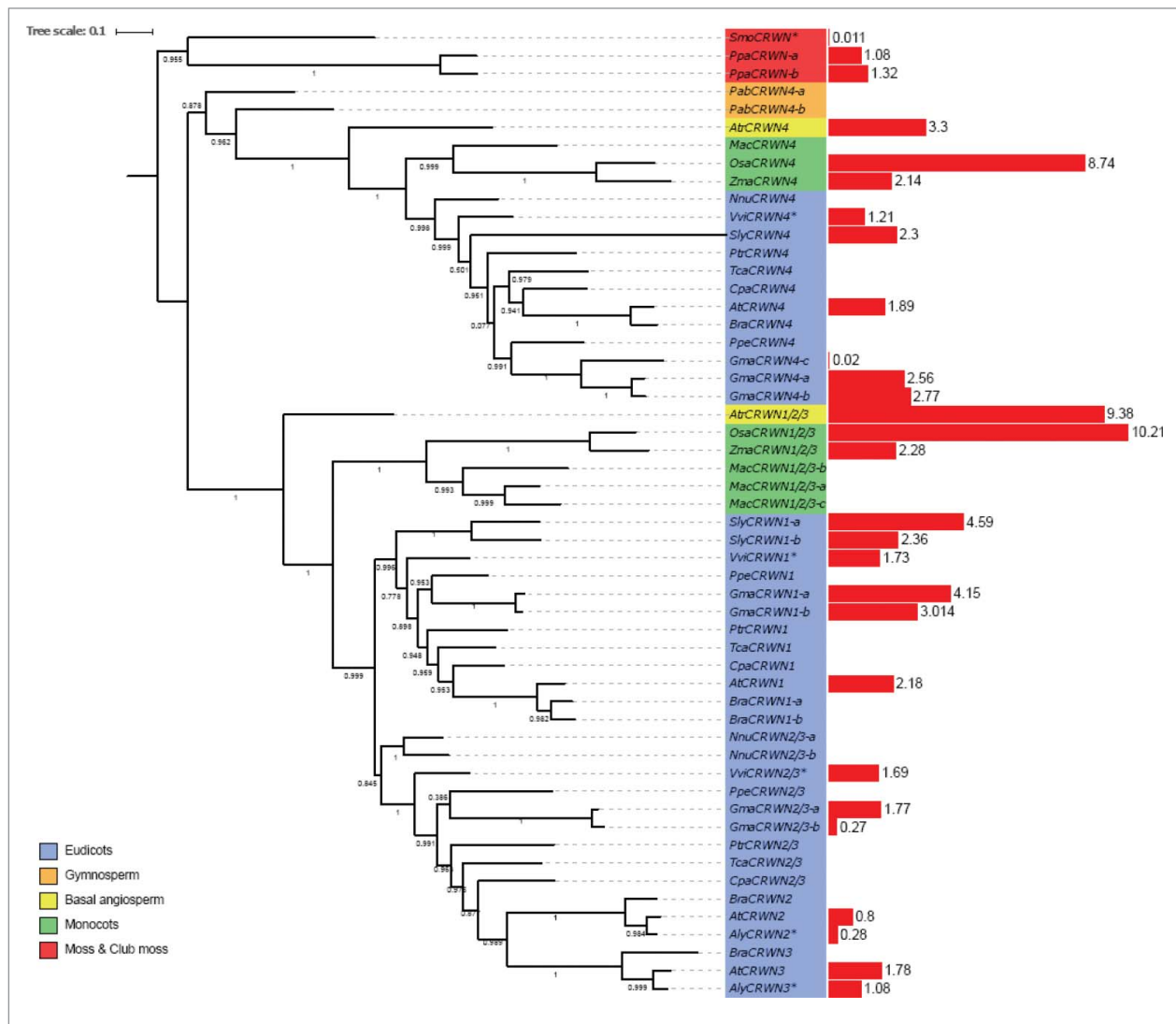
suggested to be components of the putative lamina in *A. thaliana*, CRWN,<sup>9,37</sup> KAKU4<sup>14</sup> and a novel nuclear envelope associated protein family, NEAP<sup>30</sup>.

#### CRWN proteins

The CRWN gene family is imported into the nucleus through an NLS and extensive coiled-coil domains reminiscent of the animal lamins are hypothesized to allow polymerisation of the protein to form the plant lamina. Fifty CRWN proteins were detected by BLASTp and pHMMER in all multicellular plants but are absent from unicellular algae. In most species 2 homologues were detected for each species. Two clusters of CRWN proteins were defined in a previous publication<sup>6</sup> and were also identified here as 2 main phylogenetic groups: CRWN1/CRWN2/CRWN3 and CRWN4. The clusters of CRWN4 homologues constitute monophyletic groups and only one protein was found for all species except for *Glycine max* (Fig. 6). Gymnosperm homologues seem to have only the CRWN4 lineage while *A. lyrata* has lost the CRWN4 lineage meaning that some functional redundancy exists between the 2 monophyletic groups.

For the second group made up of the homologues of the 3 other CRWN proteins, the same organization was found and only one homolog in *A. trichopoda* was detected (Fig. 6). In the monocot group, only *Musa*





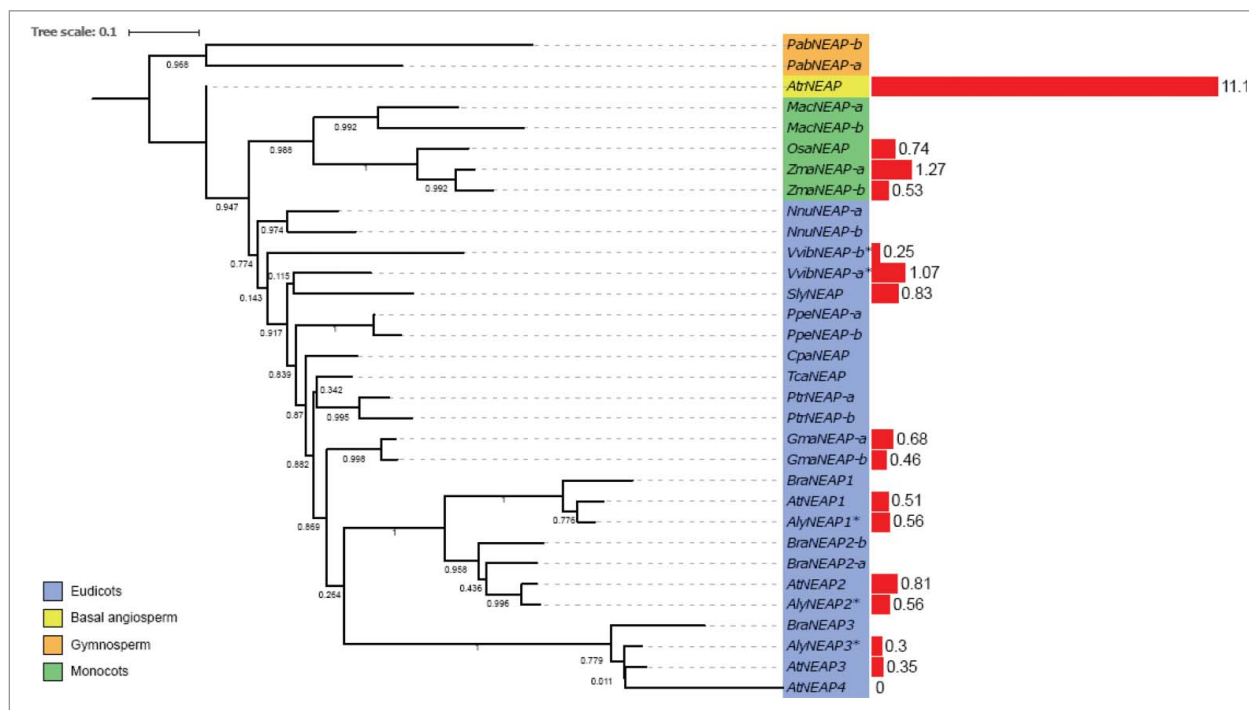
**Figure 6.** Phylogenetic tree of CRWN proteins.

*acuminata* possesses 3 homologues, the other monocots possessing only one (Fig. 6). In the eudicot group, 2 clusters can be distinguished: one for the homologues of AtCRWN1 and the other for AtCRWN2/AtCRWN3. This reveals a gene duplication, which occurred after the speciation creating monocots and eudicots. The other duplication, which gave rise to CRWN2 and CRWN3, occurred after *Brassicaceae* speciation and formed a monophyletic group. The genes belonging to the cluster CRWN1/CRWN2/CRWN3 show higher expression in comparison to CRWN4. Other than in the *Brassicaceae*, CRWN2 is less expressed than CRWN1 and CRWN3 for all the tissues analyzed. Surprisingly, no lamin-like proteins were detected in the chlorophyte unicellular algae. Previous studies have shown the presence of other lamin-like proteins in unicells like NE81 and

NUP1.<sup>10,24</sup> It is likely that several proteins have evolved in different systems to fulfil a similar role and this would reward further study.

### NEAP proteins

The NEAP proteins are characterized by a TM region at the C-terminus, a functional NLS and extensive coiled-coil domains.<sup>30</sup> NEAP1, NEAP2, and NEAP3 were identified in gymnosperms and angiosperms and 28 proteins were detected while NEAPs are absent from the more ancestral species moss, club moss and unicellular algae (Supplementary Table 11). The monocots form a monophyletic group with 2 potential specific gene duplications for *Musa acuminata* and *Zea mays* (Fig. 7). As for monocots, the eudicots form a monophyletic group (Fig. 7), and the gene duplication seems specific to species. So the 3 NEAP genes in



**Figure 7.** Phylogenetic tree of NEAP proteins.

*Brassicaceae* appear to result from a duplication event during the speciation of *Brassicaceae*. The single NEAP gene in *A. trichopoda* is expressed at very high level. Lack of expression of *AtNEAP4* and absence of protein homologues in other species imply that it is a pseudogene. The other NEAP genes are expressed in seedlings and in other tissues but at a low level (Supplementary Fig. 1).

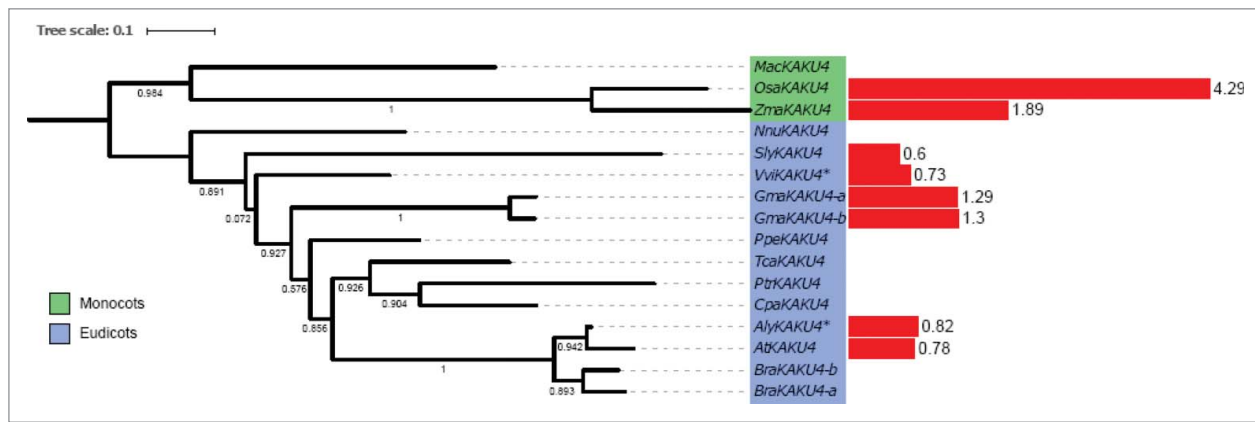
### KAKU4 proteins

KAKU4 homologues are only detected in angiosperms. Only one KAKU4 homologue is detected in each species except for *Glycine max* and *Brassica rapa*. KAKU4 is therefore a recent addition specific to angiosperms and as it interacts with CRWN1 and CRWN4, it could link CRWN proteins to other components at the nuclear periphery. Analysis of KAKU4 phylogeny reveals 2 monophyletic groups, for the monocot and eudicot homologues (Fig. 8). The protein was not detected in basal angiosperms, gymnosperms, moss, club moss and unicellular algae. Either KAKU4 is a protein with specific function in angiosperms, or was not detected due to a high variability between species. KAKU4 homologues are expressed to comparable levels in all tissues (Supplementary Fig. 1).

### Discussion

The results presented reveal functional conservation of the proteins of the plant nuclear envelope with those of other kingdoms, but surprising diversity in protein sequence. Table 1 summarizes the occurrence of each of the components in the study, together with their function. SUN domain proteins, lamina component CRWN and the KASH domain proteins SINE1-2 involved in actin binding are present before the Zeta WGD (though C-ter SUNs and CRWN are absent from the Chlamydomonas); putative NE- anchored lamina components NEAPs are first found after the Zeta WGD. Binding of RanGAP to the NE by the KASH proteins designated WIP originates with the gamma WGD of the angiosperms; the mechanism of anchorage of RanGAP in gymnosperms and mosses therefore warrants further study. CRWN interacting KAKU4 and KASH protein TIK appear to be of later origin and were only detected in the *Brassicaceae*, suggesting specialist functions.

Data from *A. trichopoda* suggests a minimal angiosperm LINC complex, with 2 KASH domain proteins (one WIP and one SINE), 3 SUN domain proteins (one C-ter-SUN and 2 mid-SUNs) and putative lamina constituents (2 CRWNs) together with one NEAP. The moss *P. patens* has 4 SUNs, 2 KASH and



**Figure 8.** Phylogenetic tree of KAKU4 proteins.

2 putative lamina constituents. The gymnosperm *P. abies* has 3 SUNs, 3 KASH (SINES) and putative lamina components (2 CRWN) together with 2 NEAPs (Fig. 1). Evolution of SUN, KASH and lamina constituents appears to have accompanied WGD and partial genome duplication events and to have resulted in a range of homologues during angiosperm speciation. The results presented also indicate a plant nuclear envelope which has developed significant complexity and redundancy through gene duplication explaining the need for multiple knockout mutants; for instance

the double mutant *atsun1 atsun2*<sup>41</sup> or the quintuple mutant *wifi (atwip1 atwip2 atwip3 atwit1 atwit2)*<sup>44</sup> before strong phenotypes are observed.

The data also suggests evolution from an ancestral LINC complex, in which SUN domain proteins are multifunctional, to a more multifaceted LINC complex containing an increasing number of KASH and lamin-like proteins. A key role for mid-SUN proteins is suggested. This is commensurate with the demonstration that SUNs play a fundamental role in chromatin interaction with the nuclear envelope during its reformation in plant

**Table 1.** protein classes and their origins and function as derived in this study.

	First appearance	WGD	Function	Location	Reference
C-ter SUN	Moss		LINC complex component; binds KASH; Nuclear shape and size; meiosis	INM	Graumann et al. <sup>17</sup> ; Oda and Fukuda, 2011
Mid-SUN	Alga		LINC complex component; binds KASH; Nuclear shape and size; fertility	INM and ER	Graumann et al. <sup>18</sup>
CRWN	Moss		nucleoskeleton; nuclear size and shape; heterochromatin organization	nuclear periphery and nucleoplasm	Dittmer et al. <sup>09</sup> ; Wang et al., 2013
SINE1-2	Moss		LINC complex component; KASH; interacts with actin cytoskeleton; nuclear positioning in guard cells (SINE1); innate immunity response (SINE2)	ONM	Zhou et al. <sup>41, 42</sup>
NEAP	Basal angiosperm, Gymnosperm	Zeta	NE anchor, SUN binding, chromatin interactor; root growth; nuclear morphology	INM	Pawar, 2015
WIP	Basal angiosperm	Epsilon	SUN binding; anchors RanGAP to NE; nuclear morphology; pollen tube termination; nuclear movement	ONM; RanGAP anchorage	Xu et al. <sup>38</sup> ; Zhao et al., 2008
KAKU 4	Monocot	Epsilon	CRWN binding; nuclear size and shape	Nucleoskeleton	Goto et al. <sup>14</sup>
SINE 3-4	Eudicot ( <i>Brassicaceae</i> )	Gamma	LINC complex component; KASH; SUN binding	ONM	Zhou et al. <sup>42</sup>
TIK	Eudicot ( <i>Arabidopsis</i> )	Gamma	SUN binding; nuclear size and shape; root growth	ONM	Graumann et al. <sup>18</sup>

mitosis<sup>15</sup> and in telomere attachment in meiosis.<sup>36</sup> KASH domain proteins appear to be evolving, with SINEs preceding WIPs, with TIK only identified in *Arabidopsis*. It is suggested that increasing specialization accompanies the acquisition of additional KASH homologues and that specific functions of the later evolving proteins (for instance RanGTP anchorage and nuclear movement in the pollen tube) are undertaken by other nuclear envelope components in their absence. A similar pattern of evolution of KASH proteins is suggested in ophisthokonts;<sup>42</sup> commenting on novel plant KASH proteins noted that while some are highly conserved (e.g. Nesprin 1 and 2, ANC-1 and MSP-300) others are restricted in distribution (e.g., Klarsicht homologues to insects and KDP-1 to nematodes) suggesting origins after SUN domain proteins and rapid evolution linked to diversifying function. Finally, higher plants have evolved a lamina-like structure based, like animal cells, on coiled-coil proteins. This appears to have arisen with the CRWN proteins present in mosses and clubmosses (lycophytes) and with KAKU4 arising later. These data are consistent with previous reports suggesting that SUN domain proteins were the first nuclear envelope proteins linking chromatin to the nuclear envelope,<sup>5</sup> predating the evolution of lamins. Indeed lamins are prone to rapid evolution as they interact with fewer partners than components of the LINC complex.<sup>23</sup> One of the striking results of our analysis is the absence of CRWN in unicellular species despite the fact that the lamina is involved in basic function such as nuclear morphology and chromatin organization which are both important for the regulation of gene expression. Although we cannot exclude if lamins and CRWN are subjected to fast evolution leading to unsuccessful recovery of homologs by Blast and HMMER analyses, it is tempting to speculate that convergent evolution occurred in animals and plants, with increasing functionality and complexity through the introduction of LBR and LEM proteins in animal and KAKU4 in plants. Similar observations were recently presented for the PRC1 polycomb group complex which is a conserved function but with poor sequence homology despite the presence of the conserved RING-domain, again suggesting convergent evolution between plant and animals. Interestingly, the PRC1 complex is involved in the regulation of gene expression through the binding of trimethylated Histone H3 at lysine 27 (H3K27me3) a well-known repressive epigenetic mark also enriched in Lamina-Associated Domains (LADs).<sup>4</sup> Exploring the possible connection between the nuclear envelope components and the PRC1

repressive complex will lead to a better understanding of the functions of the nuclear envelope in the regulation of chromatin organization and gene expression.

## Abbreviations

<i>A. thaliana</i>	<i>Arabidopsis thaliana</i>
BLASTp	Basic Local Alignment Search Tool protein
CRWN	Crowded Nuclei (also termed LINC for Little Nuclei and NMCP for Nuclear Matrix Constituent Protein)
HMMER	Hidden Markov Model-based sequence alignment tool
KASH	Klarsicht/Anc1/Syne homology
LBR	Lamin B receptor
LEM	Lamin-Emerin-Man1
LINC	Linker of Nucleoskeleton and Cytoskeleton
NEAP	Nuclear Envelope Associated Protein
RPKM	Reads Per Kilobase of transcript per Million mapped reads
SUN	Sad1-Unc84
SINEs	SUN interacting Nuclear Envelope Proteins
TIK	Toll Interleukin Receptor domain KASH protein
TM	trans-membrane
WGD	whole-genome duplication
WIPs	WPP Domain Interacting Proteins

## Disclosure of potential conflicts of interest

No potential conflicts of interest were disclosed.

## Funding

The work was supported by the CNRS, INSERM, Blaise Pascal and Oxford Brookes Universities. KG is a Leverhulme Trust Early Career Fellow.

## ORCID

Christophe Tatout  <http://orcid.org/0000-0001-5215-2338>

David Evans  <http://orcid.org/0000-0001-6248-1899>

## References

- [1] Adl SM, Simpson AGB, Lane CE, Lukeš J, Bass D, Bowser SS, Brown MW, Burki F, Dunthorn M, Hampl V, et al. The revised classification of eukaryotes. *J Eukaryot Microbiol* 2012; 59:429-514; PMID:23020233; <http://dx.doi.org/10.1111/j.1550-7408.2012.00644.x>

- [2] Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ. Basic local alignment search tool. *J Mol Biol* 1990; 215:403-410; PMID:2231712; [http://dx.doi.org/10.1016/S0022-2836\(05\)80360-2](http://dx.doi.org/10.1016/S0022-2836(05)80360-2)
- [3] Anders S, Pyl PT, Huber W. HTSeq—a Python framework to work with high-throughput sequencing data. *Bioinformatics* 2015; 31:166-9; PMID:25260700; <http://dx.doi.org/10.1093/bioinformatics/btu638>
- [4] Bickmore WA, van Steensel B. Genome architecture: domain organization of interphase chromosomes. *Cell* 2013; 152:1270-84; PMID:23498936
- [5] Cavalier-Smith T. Kingdoms protozoa and chromista and the eozoan root of the eukaryotic tree. *Biol Lett* 2010; 6:342-5; PMID:20031978; <http://dx.doi.org/10.1098/rsbl.2009.0948>
- [6] Ciska M, Moreno Diaz de la Espina S. NMCP/LINC proteins: putative lamin analogs in plants? *Plant Signal Behav* 2013; 8:e26669; PMID:24128696; <http://dx.doi.org/10.4161/psb.26669>
- [7] Czechowski T, Stitt M, Altmann T, Udvardi MK, Scheible WR. Genome-Wide Identification and Testing of Superior Reference Genes for Transcript Normalization in Arabidopsis. *Plant Physiol* 2005; 139:5-17; PMID:16166256; <http://dx.doi.org/10.1104/pp.105.063743>
- [8] Ding X, Xu R, Yu J, Xu T, Zhuang Y, Han M. SUN1 Is Required for Telomere Attachment to Nuclear Envelope and Gametogenesis in Mice. *Dev Cell* 2007; 12:863-72
- [9] Dittmer TA, Stacey NJ, Sugimoto-Shirasu K, Richards EJ. LITTLE NUCLEI Genes Affecting Nuclear Morphology in Arabidopsis thaliana. *Plant Cell* 2007; 19:2793-803; PMID:17873096; <http://dx.doi.org/10.1105/tpc.107.053231>
- [10] DuBois KN, Alsford S, Holden JM, Buisson J, Swiderski M, Bart JM, Ratushny AV, Wan Y, Bastin P, Barry JD, et al. NUP-1 Is a Large Coiled-Coil Nucleoskeletal Protein in Trypanosomes with Lamin-Like Functions. *PLoS Biol* 2012; 10:e1001287; PMID:22479148; <http://dx.doi.org/10.1371/journal.pbio.1001287>
- [11] Edgar RC. MUSCLE: a multiple sequence alignment method with reduced time and space complexity. *BMC Bioinformatics* 2004; 5:113; PMID:15318951; <http://dx.doi.org/10.1186/1471-2105-5-113>
- [12] Gaut BS, Ross-Ibarra J. Selection on Major Components of Angiosperm Genomes. *Science* 2008; 320:484-6; PMID:18436777; <http://dx.doi.org/10.1126/science.1153586>
- [13] Gerace L, Burke B. Functional organization of the nuclear envelope. *Annu Rev Cell Biol* 1988; 4:335-74; PMID:2461721; <http://dx.doi.org/10.1146/annurev.cb.04.110188.002003>
- [14] Goto C, Tamura K, Fukao Y, Shimada T, Hara-Nishimura I. The novel nuclear envelope protein KAKU4 modulates nuclear morphology in Arabidopsis. *Plant Cell* 2014; 26:2143-55; PMID:24824484; <http://dx.doi.org/10.1105/tpc.113.122168>
- [15] Graumann K, Evans DE. The plant nuclear envelope in focus. *Biochem Soc Trans* 2010a; 38:307-11; <http://dx.doi.org/10.1042/BST0380307>
- [16] Graumann K, Evans DE. Plant SUN domain proteins: Components of putative plant LINC complexes? *Plant Signal. Behav* 2010b; 5:154-6; <http://dx.doi.org/10.4161/psb.5.2.10458>
- [17] Graumann K, Runions J, Evans DE. Characterization of SUN-domain proteins at the higher plant nuclear envelope. *Plant J* 2010; 61:134-44; PMID:19807882; <http://dx.doi.org/10.1111/j.1365-313X.2009.04038.x>
- [18] Graumann K, Vanrobays E, Tutois S, Probst AV, Evans DE, Tatout C. Characterization of two distinct subfamilies of SUN-domain proteins in Arabidopsis and their interactions with the novel KASH-domain protein AtTIK. *J Exp Bot* 2014; 65:6499-512; PMID:25217773
- [19] Hetzer MW. The nuclear envelope. *Cold Spring Harb Perspect Biol* 2010; 2:a000539; PMID:20300205
- [20] Jiao Y, Wickett NJ, Ayyampalayam S, Chanderbali AS, Landherr L, Ralph PE, Tomsho LP, Hu Y, Liang H, Soltis PS, et al. Ancestral polyploidy in seed plants and angiosperms. *Nature* 2011; 473:97-100; PMID:21478875; <http://dx.doi.org/10.1038/nature09916>
- [21] Kellis M, Birren BW, Lander ES. Proof and evolutionary analysis of ancient genome duplication in the yeast *Saccharomyces cerevisiae*. *Nature* 2004; 428:617-24; PMID:15004568; <http://dx.doi.org/10.1038/nature02424>
- [22] Kim D, Pertea G, Trapnell C, Pimentel H, Kelley R, Salzberg SL. TopHat2: accurate alignment of transcriptomes in the presence of insertions, deletions and gene fusions. *Genome Biol* 2013; 14:R36; PMID:23618408; <http://dx.doi.org/10.1186/gb-2013-14-4-r36>
- [23] Koreny L, Field MC. Ancient eukaryotic origin and evolutionary plasticity of nuclear lamina. *Genome Biol Evol* 2016; PMID:27189989
- [24] Krüger A, Batsios P, Baumann O, Luckert E, Schwarz H, Stick R, Meyer I, Gräf R. Characterization of NE81, the first lamin-like nucleoskeleton protein in a unicellular organism. *Mol Biol Cell* 2012; 23:360-70; PMID:22090348
- [25] Kyte J, Doolittle RF. A simple method for displaying the hydropathic character of a protein. *J Mol Biol* 1982; 157:105-32; PMID:7108955; [http://dx.doi.org/10.1016/0022-2836\(82\)90515-0](http://dx.doi.org/10.1016/0022-2836(82)90515-0)
- [26] Letunic I, Bork P. Interactive Tree Of Life v2: online annotation and display of phylogenetic trees made easy. *Nucleic Acids Res* 2011; 39:W475-8; PMID:21470960; <http://dx.doi.org/10.1093/nar/gkr201>
- [27] Murphy SP, Simmons CR, Bass HW. Structure and expression of the maize (*Zea mays* L.) SUN-domain protein gene family: evidence for the existence of two divergent classes of SUN proteins in plants. *BMC Plant Biol* 2010; 10:269; PMID:21143845; <http://dx.doi.org/10.1186/1471-2229-10-269>
- [28] Oda Y, Fukuda H. Dynamics of Arabidopsis SUN proteins during mitosis and their involvement in nuclear shaping. *Plant J* 2011; 66:629-41; PMID:21294795; <http://dx.doi.org/10.1111/j.1365-313X.2011.04523.x>
- [29] Parry G. The plant nuclear envelope and regulation of gene expression. *J Exp Bot* 2015; 66:1673-85; PMID:25680795; <http://dx.doi.org/10.1093/jxb/erv023>

- [30] Pawar V, Poulet A, Détourné G, Tatout C, Vanrobays E, Evans DE, Graumann K. A novel family of plant nuclear envelope-associated proteins. *J Exp Bot* 2016; doi: 10.1093/jxb/erw332.
- [31] Price MN, Dehal PS, Arkin AP. FastTree 2 – Approximately Maximum-Likelihood Trees for Large Alignments. *PLoS One* 2010; 5:e9490; PMID:20224823; <http://dx.doi.org/10.1371/journal.pone.0009490>
- [32] Project AG, Albert VA, Barbazuk WB, dePamphilis CW, Der JP, Leebens-Mack J, Ma H, Palmer JD, Rounsley S, Sankoff D, et al. The Amborella Genome and the Evolution of Flowering Plants. *Science* 2013; 342:1241089; PMID:23403266
- [33] Soltis PS, Soltis DE. Ancient WGD events as drivers of key innovations in angiosperms. *Curr Opin Plant Biol* 2016; 30:159-165; PMID:27064530; <http://dx.doi.org/10.1016/j.pbi.2016.03.015>
- [34] Talavera G, Castresana J. Improvement of phylogenies after removing divergent and ambiguously aligned blocks from protein sequence alignments. *Syst Biol* 2007; 56:564-77; PMID:17654362; <http://dx.doi.org/10.1080/10635150701472164>
- [35] Tamura K, Goto C, Hara-Nishimura I. Recent advances in understanding plant nuclear envelope proteins involved in nuclear morphology. *J Exp Bot* 2015; 66:1641-7; PMID:25711706; <http://dx.doi.org/10.1093/jxb/erv036>
- [36] Varas J, Graumann K, Osman K, Pradillo M, Evans DE, Santos JL, Armstrong SJ. Absence of SUN1 and SUN2 proteins in *Arabidopsis thaliana* leads to a delay in meiotic progression and defects in synapsis and recombination. *Plant J* 2015; 81:329-46; PMID:25412930; <http://dx.doi.org/10.1111/tpj.12730>
- [37] Wang H, Dittmer TA, Richards EJ. *Arabidopsis* CROWDED NUCLEI (CRWN) proteins are required for nuclear size control and heterochromatin organization. *BMC Plant Biol* 2013; 13:200; PMID:24308514; <http://dx.doi.org/10.1186/1471-2229-13-200>
- [38] Xu XM, Meulia T, Meier I. Anchorage of plant RanGAP to the nuclear envelope involves novel nuclear-pore-associated proteins. *Curr Biol* 2007; 17:1157-63; PMID:17600715; <http://dx.doi.org/10.1016/j.cub.2007.05.076>
- [39] Zhao Q, Brkljacic J, Meier I. Two distinct interacting classes of nuclear envelope-associated coiled-coil proteins are required for the tissue-specific nuclear envelope targeting of *Arabidopsis* RanGAP. *Plant Cell* 2008; 20:1639-1651; PMID:18591351; <http://dx.doi.org/10.1105/tpc.108.059220>
- [40] Zhou X, Meier I. Efficient plant male fertility depends on vegetative nuclear movement mediated by two families of plant outer nuclear membrane proteins. *Proc Natl Acad Sci U S A* 2014; 111:11900-5; PMID:25074908; <http://dx.doi.org/10.1073/pnas.1323104111>
- [41] Zhou X, Graumann K, Evans DE, Meier I. Novel plant SUN-KASH bridges are involved in RanGAP anchoring and nuclear shape determination. *J Cell Biol* 2012; 196:203-11; PMID:22270916; <http://dx.doi.org/10.1083/jcb.201108098>
- [42] Zhou X, Graumann K, Wirthmueller L, Jones JDG, Meier I. Identification of unique SUN-interacting nuclear envelope proteins with diverse functions in plants. *J Cell Biol* 2014; 205:677-92; PMID:24891605; <http://dx.doi.org/10.1083/jcb.201401138>
- [43] Zhou X, Graumann K, Meier I. The plant nuclear envelope as a multifunctional platform LINCed by SUN and KASH. *J Exp Bot* 2015a; 66:1649-59; <http://dx.doi.org/10.1093/jxb/erv082>
- [44] Zhou X, Groves NR, Meier I. Plant nuclear shape is independently determined by the SUN-WIP-WIT2-myosin XI-i complex and CRWN1. *Nucl Austin Tex* 2015b; 6:144-53