



## Distinct types of short open reading frames are translated in plant cells

Igor Fesenko, Ilya Kirov, Andrey Knyazev, et al.

*Genome Res.* published online August 6, 2019

Access the most recent version at doi:[10.1101/gr.253302.119](https://doi.org/10.1101/gr.253302.119)

---

<b>P&lt;P</b>	Published online August 6, 2019 in advance of the print journal.
<b>Accepted Manuscript</b>	Peer-reviewed and accepted for publication but not copyedited or typeset; accepted manuscript is likely to differ from the final, published version.
<b>Creative Commons License</b>	This article is distributed exclusively by Cold Spring Harbor Laboratory Press for the first six months after the full-issue publication date (see <a href="http://genome.cshlp.org/site/misc/terms.xhtml">http://genome.cshlp.org/site/misc/terms.xhtml</a> ). After six months, it is available under a Creative Commons License (Attribution-NonCommercial 4.0 International), as described at <a href="http://creativecommons.org/licenses/by-nc/4.0/">http://creativecommons.org/licenses/by-nc/4.0/</a> .
<b>Email Alerting Service</b>	Receive free email alerts when new articles cite this article - sign up in the box at the top right corner of the article or <a href="#">click here</a> .

---

---

Advance online articles have been peer reviewed and accepted for publication but have not yet appeared in the paper journal (edited, typeset versions may be posted when available prior to final publication). Advance online articles are citable and establish publication priority; they are indexed by PubMed from initial publication. Citations to Advance online articles must include the digital object identifier (DOIs) and date of initial publication.

---

To subscribe to *Genome Research* go to:  
<https://genome.cshlp.org/subscriptions>

---

Published by Cold Spring Harbor Laboratory Press

# 1 **Distinct types of short open reading frames are translated in** 2 **plant cells**

3

4 Igor Fesenko<sup>1,\*</sup>, Ilya Kirov<sup>2</sup>, Andrey Kniazev<sup>1</sup>, Regina Khazigaleeva<sup>1</sup>, Vassili Lazarev<sup>3,4</sup>, Daria  
5 Kharlampieva<sup>3</sup>, Ekaterina Grafskaja<sup>3,4</sup>, Viktor Zgoda<sup>5</sup>, Ivan Butenko<sup>3</sup>, Georgy Arapidi<sup>1,3</sup>, Anna  
6 Mamaeva<sup>1</sup>, Vadim Ivanov<sup>1</sup>, Vadim Govorun<sup>3</sup>.

7 *<sup>1</sup>Laboratory of functional genomics and plant proteomics, Shemyakin and Ovchinnikov Institute of*  
8 *Bioorganic Chemistry, Moscow, Russian Federation; <sup>2</sup>Laboratory of marker-assisted and genomic*  
9 *selection of plants, All-Russian Research Institute of Agricultural Biotechnology, Moscow, Russian*  
10 *Federation; <sup>3</sup>Research Institute for Physico-Chemical Medicine, Moscow, Russian Federation; <sup>4</sup>Moscow*  
11 *Institute of Physics and Technology, Dolgoprudny, Moscow region, Russia; <sup>5</sup>Laboratory of System*  
12 *Biology, Institute of Biomedical Chemistry, Moscow, Russian Federation.*

13

14 **Corresponding author(s).**

15 \* Igor Fesenko, e-mail: [fesigor@gmail.com](mailto:fesigor@gmail.com)

16

17 **Running title:** Translation of sORFs in moss

18

19 **Keywords:** short open reading frames, plant peptides, LC-MS/MS, evolution, alternative splicing,  
20 lncRNA

21

22

## **ABSTRACT**

23 Genomes contain millions of short (<100 codons) open reading frames (sORFs), which are usually  
24 dismissed during gene annotation. Nevertheless, peptides encoded by such sORFs can play important  
25 biological roles, and their impact on cellular processes has long been underestimated. Here, we

26 analyzed approximately 70,000 transcribed sORFs in the model plant *Physcomitrella patens* (moss).  
27 Several distinct classes of sORFs that differ in terms of their position on transcripts and the level of  
28 evolutionary conservation are present in the moss genome. Over 5000 sORFs were conserved in at  
29 least one of ten plant species examined. Mass spectrometry analysis of proteomic and peptidomic  
30 datasets suggested that tens sORFs located on distinct parts of mRNAs and long non-coding RNAs  
31 (lncRNAs) are translated, including conserved sORFs. Translational analysis of the sORFs and main  
32 ORFs at a single locus suggested the existence of genes that code for multiple proteins and peptides  
33 with tissue-specific expression. Functional analysis of four lncRNA-encoded peptides showed that  
34 sORFs-encoded peptides are involved in regulation of growth and differentiation in moss. Knocking  
35 out lncRNA-encoded peptides resulted in a decrease of moss growth. By contrast, the overexpression  
36 of these peptides resulted in a diverse range of phenotypic effects. Our results thus open new  
37 avenues for discovering novel, biologically active peptides in the plant kingdom.

38

39

## INTRODUCTION

40

41 The genomes of nearly all organisms contain hundreds of thousands of short open reading  
42 frames (sORFs; <100 codons) whose coding potential has been the subject of recent reviews  
43 (Andrews and Rothnagel 2014; Couso 2015; Hellens et al. 2016; Couso and Patraquim 2017;  
44 Rothnagel and Menschaert 2018; Ruiz-Orera and Alba 2019). However, gene annotation algorithms  
45 are generally not suited for dealing with sORFs because short sequences are unable to obtain high  
46 conservation scores, which serve as an indicator of functionality (Ladoukakis et al. 2011).  
47 Nevertheless, using various bioinformatic approaches, sORFs with high coding potential have been  
48 identified in a range of organisms including fruit flies, mice, yeast and *Arabidopsis thaliana*  
49 (Ladoukakis et al. 2011; Hanada et al. 2013; Aspden et al. 2014; Bazzini et al. 2014). The first  
50 systematic study of sORFs was conducted on baker's yeast, where 299 previously non-annotated  
51 sORFs were identified and tested in genetic experiments (Kastenmayer et al. 2006). Subsequently,  
52 4561 conserved sORFs were identified in the genus *Drosophila*, 401 of which were postulated to be

53 functional, taking into account their syntenic positions, low  $K_A/K_S$  ( $<0.1$ ) values and transcriptional  
54 evidence (Ladoukakis et al. 2011). In a recent study, Mackowiak and colleagues predicted the  
55 presence of 2002 novel conserved sORFs (from 9 to 101 codons) in *Homo sapiens*, *Mus musculus*,  
56 *Danio rerio*, *Drosophila melanogaster* and *Caenorhabditis elegans* (Mackowiak et al. 2015). The first  
57 comprehensive study of sORFs in plants postulated the existence of thousands of sORFs with high  
58 coding potential in *Arabidopsis* (Lease and Walker 2006; Hanada et al. 2007; Hanada et al. 2013),  
59 including 49 that induced various morphological changes and had visible phenotypic effects.

60           Recent studies have pointed to the important roles of sORF-encoded peptides (SEPs) in  
61 cells (Magny et al. 2013; Nelson et al. 2016; D'Lima et al. 2017; Huang et al. 2017; Matsumoto et al.  
62 2017; Rubtsova et al. 2018). However, unraveling the roles of SEPs is a challenging task, as is their  
63 detection at the biochemical level. In animals, SEPs are known play important roles in a diverse  
64 range of cellular processes (Kondo et al. 2010; Magny et al. 2013). By contrast, only a few functional  
65 SEPs have been reported in plants, including POLARIS (PLS; 36 amino acids), EARLY NODULIN GENE  
66 40 (ENOD40; 12, 13, 24 or 27 amino acids), ROTUNDIFOLIA FOUR (ROT4; 53 amino acids), KISS OF  
67 DEATH (KOD; 25 amino acids), BRICK1 (BRK1; 84 amino acids), Zm-908p11 (97 amino acids) and  
68 Zm-401p10 (89 amino acids) (Andrews and Rothnagel 2014; Tavormina et al. 2015). These SEPs  
69 help modulate root growth and leaf vascular patterning (Chilley et al. 2006), symbiotic nodule  
70 development (Djordjevic et al. 2015), polar cell proliferation in lateral organs and leaf  
71 morphogenesis (Narita et al. 2004), and programmed cell death (apoptosis) (Blanvillain et al. 2011).

72           To date, functional sORFs have been found in a variety of transcripts, including  
73 untranslated regions of mRNA (5' leader and 3' trailer sequences), lncRNAs, and microRNA  
74 transcripts (pri-miRNAs) (Andrews and Rothnagel 2014; Laing et al. 2015; Laouressergues et al. 2015;  
75 Couso and Patraquim 2017; Brunet et al. 2019). Evidence for the transcription of potentially  
76 functional sORFs has been obtained in *Populus deltoides*, *Phaseolus vulgaris*, *Medicago truncatula*,  
77 *Glycine max* and *Lotus japonicus* (Guillen et al. 2013). The transcription of sORFs can be regulated by  
78 stress conditions and depends on the developmental stage of the plant (De Coninck et al. 2013;  
79 Hanada et al. 2013; Rasheed et al. 2016). Indeed, sORFs might represent an important source of  
80 advanced traits required under stress conditions. During stress, genomes undergo widespread

81 transcription to produce a diverse range of RNAs (Kim et al. 2010; Mazin et al. 2014); therefore, a  
82 large portion of sORFs becomes accessible to the translation machine for peptide production. Stress  
83 conditions can lead to the transcription of sORFs located in genomic regions that are usually non-  
84 coding (Giannakakis et al. 2015). Such sORFs appear to serve as raw materials for the birth and  
85 subsequent evolution of new protein-coding genes (Couso and Patraquim 2017; Ruiz-Orera and Alba  
86 2019).

87           The transcription of an sORF does not necessarily indicate that it fulfills any biological  
88 role, as opposed to being a component of the so-called translational noise (Guttman et al. 2013).  
89 According to ribosomal profiling data, thousands of lncRNAs display high ribosomal occupancy in  
90 regions containing sORFs in mammals (Ingolia et al. 2011; Aspden et al. 2014; Bazzini et al. 2014).  
91 However, lncRNAs can have the same ribosome profiling patterns as canonical non-coding RNAs  
92 (e.g., rRNA) that are known not to be translated, implying that these lncRNAs are unlikely to produce  
93 functional peptides (Guttman et al. 2013). In addition, identification of SEPs via mass spectrometry  
94 analyses has found many fewer peptides than predicted sORFs (Slavoff et al. 2013; Aspden et al.  
95 2014). Thus, the abundance, lifetime and other features of SEPs are generally unclear.

96           We performed a comprehensive analysis of the sORFs that have canonical AUG start  
97 codons and high coding potential in the *Physcomitrella patens* genome. The translation of tens of  
98 sORFs was confirmed by mass-spectrometry analysis. From these, candidate lncRNA-encoded  
99 peptides were selected for further analysis, which provided evidence for their biological functions.

100

## RESULTS

### 101 **Discovery and classification of potential coding sORFs in the moss genome**

102 Our approach is summarized in Fig. 1A. At the first stage of analysis, we used the sORF finder tool  
103 (Hanada et al. 2010) to identify single-exon sORFs starting with an AUG start codon and less than  
104 300 bp long. This approach resulted in the identification of 638,439 sORFs with coding potential (CI  
105 index) in all regions of the *P. patens* genome.

106 We selected 70,095 unique sORFs located on transcripts annotated in the moss genome  
107 (phytozome.jgi.doe.gov) and/or our dataset (Fesenko et al. 2015) for further analysis, as well as

108 those on lncRNAs from two databases - CANTATAdb (Szczesniak et al. 2016) and GreenNC (Paytuví  
109 Gallart et al. 2016); sORFs located in repetitive regions were discarded (Supplemental Table S1).  
110 These selected sORFs, which were 33 to 303 bp long, were located on 33,981 transcripts (22,969  
111 genes), with up to 28 sORFs per transcript (Supplemental Fig. S1A).

112 We then classified the sORFs based on their location on the transcript: 63,109 “genic-sORFs”  
113 (located on annotated transcripts, but not on lncRNA), 1241 “intergenic-sORFs” (located on  
114 transcripts from our dataset and not annotated in the current version of the genome) and 5745  
115 “lncRNA-sORFs” (located on lncRNAs from CANTATAdb (Szczesniak et al. 2016), GreenNC (Paytuví  
116 Gallart et al. 2016) or our data set (Fesenko et al. 2017); Fig. 1B). The genic-sORFs include 11,998  
117 upstream ORFs (uORFs; for 5'-UTR location), 9443 downstream ORFs (dORFs; for 3'-UTR location),  
118 36,732 coding sequence-sORFs (CDS-sORFs; sORFs overlapping with main ORFs (+1 frame) in non-  
119 canonical +2 and +3 reading frames) and 3485 interlaced-sORFs (overlapping with both the CDS and  
120 5'-UTR or CDS and 3'-UTR on the same transcript) (Fig. 1B, Supplemental Fig. S1B).

121 As expected based on the sORF finder search strategy (Hanada et al. 2010), the sORF set was  
122 enriched in CDS-sORFs (52%, Fisher's exact test,  $P$ -value  $< 10^{-16}$ ), whereas dORFs, uORFs and  
123 interlaced-sORFs were underrepresented (Fisher's exact test,  $P$ -value  $< 10^{-16}$ ) compared to a random  
124 exonic fragments set, which was used as a negative control. On average, CDS-sORFs (median size of  
125 22 codons) were shorter than uORFs (median size of 35 codons; Mann-Whitney  $U$  test  $P$ -value  $< 10^{-16}$ )  
126 and dORFs (median length 32 codons, Mann-Whitney  $U$  test  $P$ -value  $< 10^{-16}$ ). The median size of  
127 interlaced-sORFs was 49 codons, which is significantly longer than other genic-sORFs (Mann-  
128 Whitney  $U$  test  $P = 0.0021$ ) (Fig. 1C).

129 To estimate the number of conserved transcriptable sORFs, we performed a TBLASTN  
130 search (e-value cutoff 0.00001) of each sORF sequence against the reconstructed genomes of three *P.*  
131 *patens* ecotypes, Villersexel, Reute, and Kaskasia, as well as the transcriptomes of ten plant species  
132 (Supplemental Fig. S2). We found 5034 conserved sORFs with detectable homologous sequences in at  
133 least one species (Supplemental Table S1, Supplemental Fig. S3). A conservation analysis of the  
134 sORFs in the reconstructed *P. patens* ecotypes showed that 2.4% (1618) of the sORFs were lacking  
135 either the start or stop codons in at least one species. We then examined the differences in selection

136 pressure at the amino acid level between different major groups of conservative sORFs (CDS-sORFs,  
137 uORFs, dORFs, lncRNA-sORFs, interlaced-sORFs) using the criterion of  $K_A/K_S$ . Higher retention rates  
138 were observed for uORFs and dORFs, whereas CDS-sORFs and lncRNA-ORFs were under strong  
139 positive selection (Supplemental Fig. S4). These observations are in agreement with the fact that  
140 some type of sORFs (for example, uORFs) play a regulatory role instead of being translated (Barbosa  
141 et al. 2013).

142

### 143 **Experimental evidence for the translation of sORFs**

144 Obtaining evidence for the translation of sORFs is an important step towards identifying functional  
145 SEPs. We analyzed the Kozak consensus sequences (Kozak 1986) surrounding sORF start codons.  
146 Kozak consensus sequence plays an important role in translation initiation (Kozak 1997). Depending  
147 on the presence of the purine in position -3 and the G in position +4 (where +1 is “A” in the “AUG”  
148 codon) the Kozak was considered to be “strong” (both are present), “medium” (one is present) or  
149 “weak” (neither are present) (Kozak 1997). According to our results, 41816 (~60%) of the predicted  
150 sORFs were surrounded by “strong” and “medium” Kozak sequences. These values were significantly  
151 smaller than those of annotated protein coding ORFs (87%, Fisher’s exact test  $P$ -value  $< 2.2 \times 10^{-16}$ ).  
152 We then verified the translation of our predicted sORFs using mass-spectrometry (MS) analysis.  
153 Taking into account the shortage of proteomic methods for identifying small proteins or peptides, in  
154 the current study, we generated two datasets: the “peptidomic” dataset - endogenous peptides  
155 extracted from three types of moss cells: gametophores, protonemata and protoplasts and the  
156 “proteomic” dataset - tryptic peptides generated in a standard proteomic pipeline (Supplemental  
157 Table S2). All datasets were mapped with MaxQuant against a custom database containing our sORFs  
158 together with nuclear, chloroplast and mitochondrial moss protein sequences (see details in the  
159 Methods). PSMs (peptide spectrum matches) were identified at 1% FDR, and ambiguous peptides  
160 were filtered out. This resulted in 1177 PSMs corresponding to 296 distinct peptide sequences in the  
161 peptidomic dataset and 920 PSMs corresponding to 532 peptide sequences in the proteomic dataset.  
162 To generate a high-confidence sORF candidate set, we increased our acceptance threshold to a  
163 minimum Posterior Error Probability (PEP) of 0.01 and Andromeda score of higher than 60. The final

164 set underwent a manual inspection of spectra. As a result, we confirmed the translation of 46 sORFs:  
165 17 in gametophores, 29 in protonemata, and 14 in protoplasts (“confident sORFs”, Fig. 2A,  
166 Supplemental Table S3). The length of these small protein-coding sORFs ranged from 14 to 99 amino  
167 acids (aa), which were generally longer than untranslatable sORFs (Mann-Whitney *U* test *P*-value =  
168  $5.33 \times 10^{-6}$ ) (Fig. 2B). We observed that PSMs supporting SEP identifications had lower average  
169 quality than those mapped to the protein sequences (Supplemental Figs. S5A and S5B). This finding  
170 is in agreement with data obtained for the animal kingdom (Slavoff et al. 2013; Mackowiak et al.  
171 2015). The quality of spectra and the values of PSMs supporting the expression of SEPs were better  
172 in the “peptidomic” dataset (Supplemental Fig. S5C). Also, translatable sORFs were longer for those  
173 identified in the peptidomic dataset (Supplemental Fig. S5D). Approximately 63% of the translated  
174 sORFs (29 sORFs) contained “strong” and “medium” Kozak elements, which is similar to the results  
175 obtained for all predicted sORFs (~60%). This result suggests that translation initiation may differ  
176 for sORFs and protein coding ORFs.

177 The most prominent group of small protein-coding sORFs consisted of CDS-sORFs (19 sORFs, 41.3%)  
178 (Fig. 2C). Also, the translation of uORFs (6 sORFs, 13%) and dORFs (9 sORFs, 19.6%) was confirmed  
179 by our analysis. Based on our MS data, we identified seven loci with at least two translated ORFs  
180 (annotated as main ORF and sORF), including 5 CDS-sORFs, that represent putative multi-coding  
181 genes (Fig. 2D; Supplemental Table S4). Some of the putative multi-coding genes were translated  
182 simultaneously with protein-coding ORFs in the same type of moss cell (e.g. Pp3c11\_sORF461), while  
183 others showed different patterns of sORF and main ORF translation (e.g. Pp3c1\_sORF1909). These  
184 findings indicate that small protein-coding CDS-sORFs are expressed simultaneously with main ORFs  
185 and the translation of sORFs and proteins located together in the same locus might be regulated in a  
186 tissue-specific manner.

187 The translation of 9 sORFs located on lncRNAs was also detected by our analysis. The level of  
188 transcription of some lncRNAs (according to the previous data (Fesenko et al. 2017) and Phytozome  
189 12.0 expression atlas) and evidences of translation for the corresponding lncRNA-sORFs are shown  
190 in Fig. 2E. Three of these SEPs, Pp3c18\_sORF57 (40aa), Pp3c9\_sORF1544 (41aa) and  
191 Pp3c25\_sORF1000 (61aa), were common to all three cell types and were confirmed by several

192 unique endogenous peptides (Fig. 2E). These data may point to biological significance for the  
193 peptides translated from these sORFs rather than the sORFs having regulatory functions in the  
194 translation of the main ORF. To explore this notion, we investigated the functions of four SEPs  
195 encoded by lncRNAs (see below).

196

### 197 **Most small protein-coding sORFs are not evolutionarily conserved**

198 Analysis of the evolutionary conservation of sORFs is often a key step in revealing biologically active  
199 sORFs (Andrews and Rothnagel 2014). To investigate whether the trend in small protein-coding  
200 sORF evolution differs from that of the other sORFs, we estimated the number of species in which  
201 homologs can be found and the selection pressure ( $K_A/K_S$ ) on translatable sORFs on an evolutionary  
202 timescale using the transcriptomes of the ten abovementioned species. Overall, we found 5 sORFs  
203 had evidence of translation and conservation in at least one species, 4 of them were under negative  
204 selection ( $K_A/K_S \ll 1$ ). Thus, analysis of sORF sequence conservation showed that only 11% of our  
205 small protein-coding sORFs have signature of conservation between species.

206

### 207 **Alternative splicing regulates the number of sORFs in protein-coding transcripts**

208 Alternative splicing (AS) is a universal process among eukaryotic organisms and more than 50% of *P.*  
209 *patens* genes are alternatively spliced (Chang et al. 2014; Wu et al. 2014; Fesenko et al. 2017). AS  
210 events may lead to the specific gain, loss or truncation of different groups of sORFs located on the  
211 transcripts of the same gene. We found 6092 alternatively spliced sORFs (AS-sORFs) belonging to  
212 transcripts from 4389 genes. CDS-sORFs were significantly overrepresented (Supplemental Fig. S6)  
213 while interlaced-sORFs, uORFs and dORFs were significantly underrepresented among AS-sORFs  
214 compared to the control set of random exonic fragments. We found that about half of the entire set of  
215 AS-sORFs (48%, 2933) underwent complete excision (complete sORF removal from an isoform; Fig.  
216 3). The complete excision of sORFs occurred significantly more frequently in uORFs (57% of all AS-  
217 sORFs) than in the other AS-sORF groups (20–44% of all AS-sORFs, Fisher's exact test  $P$ -value  $< 10^{-6}$ ).  
218 Among small protein-coding AS-sORFs, we found 3 affected by stop codon excision. 2 of the

219 translatable AS-sORFs were affected by start codon excision and one undergone complete excision.  
220 We then randomly selected thirteen different AS-sORFs with/without evidence of translation and  
221 searched for the corresponding isoforms in the transcriptomes of three types of moss cells. RT-PCR  
222 analysis revealed the transcription of these isoforms, confirming that they could indeed be translated  
223 (Supplemental Fig. S7). Moreover, some sORFs contained isoforms showing tissue-specific  
224 transcription. These observations led to the hypothesis that the translation of sORFs is regulated by  
225 AS.

226 The formation of a premature termination codon (PTC) as a result of alternative splicing events,  
227 might lead to mRNA decay (Ge and Porse 2014; Karousis et al. 2016) and rapid nonsense-mediated  
228 decay (NMD)-coupled degradation of sORF-encoded peptides (Popp and Maquat 2013). Using  
229 recently published transcriptomic data from moss NMD-deficient mutants (Lloyd et al. 2018), we  
230 investigated whether our translatable sORFs were present on NMD-targeted transcripts. The only  
231 one CDS-sORF (Pp3c7\_sORF1583) was potentially present on such transcripts. Therefore, it is  
232 difficult to judge if AS-sORFs can trigger NMD-dependent transcript degradation.

233 Thus, our analysis demonstrated that AS might regulates the excision of sORFs from the  
234 transcriptome of *P. patens*, preventing AS-sORF translation by start or stop codon as well as complete  
235 sORF excision.

236

### 237 **The sequence similarity analysis reveals sORFs with high identity to coding genes**

238 Competitive inhibitors of protein-protein interactions (PPI) are referred to as MicroProteins (miPs)  
239 or small interfering peptides (siPEPs) and can be generated by alternative splicing or evolutionarily  
240 generated by domain loss (Seo et al. 2011; Staudt and Wenkel 2011; Eguen et al. 2015). Using  
241 BLASTP (e-value <  $10^{-6}$ ) similarity searches, we identified 363 sORFs resulting from AS events that  
242 partially overlapped with the main ORF, thereby generating truncated versions of the proteins (cis-  
243 sORFs; Supplemental Table S5). We found that 60 cis-sORFs harbored intrinsically disordered  
244 regions (IDRs, (van der Lee et al. 2014)), while 30 ones contained parts of 28 different domains  
245 (Supplemental Table S5). However, we did not identify small protein-coding cis-sORFs in our dataset.

246 It could be explained by a significant overlap with the protein sequences, whereas we filtered out the  
247 'ambiguous' PSMs.

248 We then identified 272 sORFs that shared similarity with annotated proteins but were  
249 located on other transcripts (trans-sORFs, see in Supplemental Table S5). Trans-sORFs may have  
250 originated through the divergence of ancient paralogous genes, which occurred after the paleo  
251 duplication of the moss genome (Rensing et al. 2007; Rensing et al. 2008). In fact, 159 (58.5%) trans-  
252 sORFs shared similarity to genes from at least one species. In addition, all of these trans-sORFs are  
253 under strong purifying selection ( $K_A/K_S \ll 1$ ).

254 Several distinct clusters with sORF-encoded peptides sharing similarity with more than four  
255 proteins from distinct genes were detected (Supplemental Fig. S8). Each cluster encompasses genes  
256 from different protein families, including one containing leucine-rich repeat and zinc-finger domains  
257 involved in protein-protein and protein-nucleic acid interactions, respectively. We examined the co-  
258 expression data and compared the distribution of correlation coefficient values between potential  
259 SEPs and their targets with those from randomly selected pairs (10 iterations) of genes. On average,  
260 these sORF-protein pairs had higher correlation coefficients than randomly selected gene pairs  
261 (Wilcoxon Rank Sum and Kolmogorov-Smirnov Tests  $P$ -value  $< 0.05$ ), implying that sORF-bearing  
262 and target genes are frequently co-expressed.

### 263 **SEPs regulate moss growth**

264 Despite the recent finding that 10% of overexpressed intergenic sORFs have clear phenotypes in  
265 *Arabidopsis* (Hanada et al. 2013), the functions of most sORFs and SEPs in plants are generally  
266 unknown. Known bioactive SEPs in plants are encoded by sORFs located on short non-protein-coding  
267 transcripts, which can be referred to as lncRNAs (Rohrig et al. 2002; Chilly et al. 2006). In this  
268 context, it would be important to determine how many plant lncRNAs encode peptides, as well as  
269 the biological functions of these SEPs. Our pipeline allowed us to identify translated sORFs, including  
270 those encoded by lncRNAs. Some of these lncRNA-sORFs showed tissue-specific transcription and  
271 translation patterns, while others were expressed in all types of moss cells (Fig. 2E). We reasoned  
272 that stably expressed lncRNA-sORFs can produce peptides that play fundamental roles in various  
273 cellular processes. To explore this hypothesis, we examined the impact of lncRNA-sORF

274 overexpression and knockout on moss morphology using four lncRNAs-sORFs: Pp3c9\_sORF1544,  
275 Pp3c25\_sORF1253, Pp3c25\_sORF1000 and Pp3c18\_sORF57 (Fig. 2E). The translation of these SEPs  
276 was confirmed by several unique peptides and they contained “strong” and “medium” Kozak  
277 elements. We obtained multiple independent mutant lines for each of these lncRNAs-sORFs  
278 (Supplemental Figs. S9, S10, S11 and S12). Both the overexpression and knockout of sORFs resulted  
279 in morphological changes, implying that these peptides play a role in growth and development of *P.*  
280 *patens* (Figs. 4 and 5, Supplemental Table S6).

281 Overexpression of a 41-aa peptide (*PSEP1*, *Physcomitrella patens* sORF encoded peptide 1) encoded  
282 by the lncRNA-sORF Pp3c9\_sORF1544 resulted in longer caulonema cells (filaments implicated in a  
283 rapid radial extension of the protonemal tissues) compared to the wild-type and *psep1* knockout  
284 lines (Fig. 4A-F and G, Supplemental Figs. S13A-F and S14A-D). Rapid growth in the *PSEP1*  
285 overexpressing lines (*OE*) was accompanied by earlier aging and cell death (Supplemental Fig. S15).  
286 By contrast, there was a small, but significant difference in growth rate between the wild-type and  
287 *psep1* mutant lines grown on solid media and in the liquid culture without glucose (Fig. 4H,  
288 Supplemental Fig. S13A-F).

289 The lines with a knockout in a 57-aa peptide (*psep3* KO) encoded by conservative lncRNA-sORF  
290 Pp3c25\_sORF1253 displayed a decrease in growth rate and altered filament branching (Fig. 4I-O). In  
291 the wild-type moss plants, a pale-green diffuse network of caulonemal filaments surrounded the  
292 central zone (principally chloronemata), while *psep3* KO mutant lines displayed short lateral  
293 filaments on medium without glucose and ammonium tartrate, which favors chloronemal growth  
294 (Supplemental Fig. S13G-I). Overexpression of *PSEP3* (*PSEP3* OE) resulted in a significant decrease in  
295 growth rate compared to the wild type (Fig. 4P). Moreover, much of the *PSEP3* OE protonemal tissue  
296 grown on medium without ammonium tartrate turned brown (Fig. 5K-N, Supplemental Figure S14E-  
297 L).

298 Similar to the results for the *psep3* knockout, knocking out a 61-aa peptide (*psep25* KO) encoded by  
299 conserved lncRNA-sORF Pp3c25\_sORF1000 also resulted in a decrease in growth rate and altered  
300 protonemal architecture on medium without glucose but supplemented with ammonium tartrate  
301 (Figs. 5A-G, Supplemental Fig. S13J-M). *PSEP25* - overexpressing mutant lines displayed a slight



**327 sORFs with high coding potential are not conserved among genomes**

328 Although analyzing the conservation of short amino acid sequences is not trivial (Moyers and Zhang  
329 2016), hundreds of conserved sORFs have recently been identified in plants, yeast and animals  
330 (Ladoukakis et al. 2011; Hanada et al. 2013; Mackowiak et al. 2015; Brunet et al. 2019). The number  
331 of sORFs conserved in the plant kingdom is undoubtedly underestimated due to the low sensitivity of  
332 tools used for conservation analysis and the limited number of available sequenced genomes from  
333 closely related species. Our pipeline allowed us to identify 5034 conserved sORFs among the  
334 transcriptomes of ten different plant species, 5 of which showed evidence of translation according to  
335 our MS data. Three of five conserved sORFs belonged to lncRNAs. These data are in line with a  
336 previously published study showing that a large fraction of small ORFs in mouse genome evolves  
337 neutrally (Ruiz-Orera et al. 2018). We also found that uORFs and dORFs were significantly  
338 underrepresented among the sORFs that are conserved in the closest related species. We even  
339 detected rapid inactivation of uORFs and dORFs in the reconstructed genomes of three *P. patens*  
340 ecotypes due to disruptions in the start or stop codons (47% of the total disrupted sORFs). As the  
341 occurrence of sORFs downstream or upstream of the main ORF can be deleterious to its translation  
342 or induce nonsense-mediated decay (NMD), we cannot rule out the possibility that this may cause  
343 strong selection pressure and the rapid elimination of uORFs and dORFs (Iacono et al. 2005; Neafsey  
344 and Galagan 2007; Johnstone et al. 2016; Ruiz-Orera and Alba 2019). Taken together, these findings  
345 suggest that sORFs located in untranslated regions of mRNAs are evolving rapidly and may play  
346 regulatory roles rather than encoding bioactive peptides.

347 In recent studies, thousands of alternative proteins were experimentally detected in human  
348 cell lines (Vanderperre et al. 2013; Samandi et al. 2017; Brunet et al. 2019). In *P. patens*, we found  
349 tens of thousands of sORFs (CDS-sORFs) that overlapped with the CDS of protein-coding genes. The  
350 evolution of CDS-sORFs is undoubtedly an expensive process for the cell, as these elements may be  
351 located in regions encoding protein domains and influence the structure and function of the protein  
352 encoded by the main ORF (Cherry 2010). We found both CDS-sORFs originated from regions  
353 associated with known protein domains and CDS-sORFs from disordered regions, with higher  
354 conservation for CDS-sORFs originated from protein domain-encoding regions. These results

355 indicate that the evolution of CDS-sORFs depends on their locations insight main CDS sequence.  
356 However, whether sORFs are preferentially generated in fast-evolving regions of proteins or whether  
357 the selective pressure on sORFs leads to changes in protein-coding sequences is still unknown.

### 358 **Analysis of sORF translation: approaches that makes sense**

359 It was recently suggested that sORFs are randomly generated in a genome(Couso and Patraquim  
360 2017). Assuming that the average length of a sORF is approximately 60 bp and that sORFs do not  
361 overlap, these elements occupy a substantial portion of the moss genome. This raises the question: to  
362 what extent are sORFs present in the transcriptome and the proteome of a cell? According to  
363 ribosome profiling data from a wide variety of species, sORFs translation appears to occur in a  
364 pervasive manner (Ingolia et al. 2011; Guttman et al. 2013; Bazzini et al. 2014; Couso and Patraquim  
365 2017). However, ribosome-profiling data alone are not sufficient to classify transcripts as coding or  
366 noncoding (Guttman et al. 2013). Mass-spectrometry studies have thus far confirmed the presence of  
367 a few dozen SEPs in the peptidomes of animal cells (Slavoff et al. 2013; Prabakaran et al. 2014;  
368 Mackowiak et al. 2015; Ma et al. 2016; Tharakan et al. 2019). Comparisons of ribosome profiling and  
369 mass spectrometry results have led to the conclusion that MS detects peptides arising from the most  
370 highly translated sORFs (Aspden et al. 2014; Bazzini et al. 2014). However, a recent study showed  
371 that there are no technical obstacles to the detection of sORF-encoded peptides by mass  
372 spectrometry (Verheggen et al. 2017).

373 In previous studies, only standard proteomics analysis was used to identify SEPs. We  
374 reasoned that analyzing endogenous peptide pools instead of tryptic peptides has several  
375 disadvantages in terms of SEP identification: 1) standard proteomic approaches are not suitable for  
376 the isolation and analysis of small and low-abundance peptide molecules; and 2) SEPs are shorter  
377 than standard proteins and it is unlikely that more than one tryptic fragment will be detected in a  
378 single proteomic experiment. Moreover, peptidomic approaches can theoretically be used to identify  
379 full-length SEPs in a cell. We did not observe any significant overlap between the sORFs detected  
380 using proteomic and peptidomic approaches. Thus, our study demonstrates the advantage of using  
381 complementary approaches for building a complete list of SEPs.

382           According to our MS data, the translation patterns of most small protein-coding sORFs tend  
383 to be tissue specific (Fig. 2A). We suggest that the slight overlap in tissue-specific expression among  
384 SEPs from various types of moss cells could be due to either specific SEP post-translational  
385 modification (PTM) patterns, tissue-specific transcription of sORFs, or the limitations of mass-  
386 spectrometry in detecting low-abundance or modified sORF-encoded peptides. According to our  
387 results, alternative splicing is an additional mechanism that control tissue-specific sORF expression  
388 in plant cells. Also, the number of sORFs that were commonly translated between two types of moss  
389 cells was higher for related cell types: protonemata and gametophores (two growth stages) as well  
390 as protonemata and protoplasts (protoplasts were generated from the protonemata). These  
391 observations indicate tissue-specific characteristics of SEPs translation and modification rather than  
392 technical limitation in detection.

### 393 **Functionality of SEPs**

394 We identified tens of small protein-coding sORFs representing multiple sORF types and suggested  
395 various functions for the types of sORFs. Clear evidence of transcription and translation points to a  
396 possible biological significance of the small protein-coding sORFs that we identified here. Based on  
397 our results (evolution, alternative splicing analysis), we suggest that the majority of uORFs play  
398 regulatory roles instead of having peptide-encoding functions.

399 By contrast, CDS- and lncRNA-sORFs have greater potential to encode bioactive peptides, as they are  
400 more highly conserved, frequently contain known protein domains and, according to the MS data,  
401 produce peptides. We identified 19 small protein-coding CDS-sORFs in our dataset, seven of which  
402 were translated simultaneously with previously annotated longer protein-coding ORF. This finding is  
403 in agreement with a recent study on mammals, reporting that a gene MIEF1 translational product is  
404 not the canonical protein but the small 70 amino acid alternative MiD51 protein is (Delcourt et al.  
405 2018).

406           One possible role for CDS-sORFs that are similar to known proteins is to mimic the similar  
407 protein to interfere with its function. MiPs (or siPEPs) are important modulators of protein–protein  
408 and protein–DNA interactions that, for example, prevent the formation of functional protein  
409 complexes (Seo et al. 2013; Graeff et al. 2016). We found that approximately 30% of cis-SEPs harbor

410 protein domains such as protein kinase domains and MYB-like DNA-binding domain or IDRs. Also,  
411 some sORFs with disordered regions might mediate protein–protein or protein–nucleic acid  
412 interactions, as suggested previously (Mackowiak et al. 2015). However, we failed to identify the  
413 translation of such sORFs using stringent identification criteria in our mass-spectrometry analysis.  
414 Therefore, this point requires further confirmation.

415         The transcription of the non-coding portions of the genome into lncRNAs is thought to give  
416 rise to the translation of sORFs located within them. Nevertheless, the functions of these peptides are  
417 unclear and require more detailed investigation. According to our results, knocking out the selected  
418 lncRNA-encoded peptides was not lethal in moss, but did influence moss growth and development.  
419 All SEP knockouts showed a decrease in growth rate compared to the wild-type plants. By contrast,  
420 we found that plants overexpressing lncRNA-encoded peptides showed more phenotypic differences  
421 compared to the wild-type plants and knockouts. We observed both a significant increase in growth  
422 rate (*PSEP1* OE) and in the number of leafy shoots (*PSEP25* OE) and a decrease in growth rate in  
423 *PSEP3* OE, *PSEP18* OE, and *PSEP25* OE lines. The differences between the wild-type and mutant lines  
424 often appeared only under certain growth conditions – solid or liquid media with/without glucose or  
425 tartrate ammonia. These data may point to a tight regulation of lncRNA-encoded peptide translation  
426 in cell. In light of these findings, we hypothesized that lncRNA-encoded peptides may not be vital but  
427 be important for the survival under certain conditions by serving as a raw material for the evolution.  
428 According to the recently proposed classification of small ORFs, lncRNA-sORFs used in our functional  
429 analysis may be referred to both lncORFs and short CDSs (Couso and Patraquim 2017). Both short  
430 CDSs and lncRNAs have a median size of 79 aa and 24 aa in animal genomes, respectively (Couso and  
431 Patraquim 2017). We suggest that the differences in types of predicted sORFs between plant and  
432 animal genomes require further investigation. Our results lay the groundwork for the systematic  
433 analysis of functional peptides encoded by sORFs.

434         The possible evolution of non-coding portions of the genome into protein-coding genes is  
435 also a subject of intensive debate (Carvunis et al. 2012; McLysaght and Guerzoni 2015; Couso and  
436 Patraquim 2017; Ruiz-Orera and Alba 2019). According to our data, putative homologous sORFs  
437 tended to differ in length in most cases. Thus, we suggest that most sORFs expanded during

438 evolution, providing support for the notion that they function as raw materials for selection;  
439 however, this point requires further confirmation.

## 440 **METHODS**

### 441 ***Physcomitrella patens* growth conditions**

442 *Physcomitrella patens* subsp. *patens* (“Gransden 2004”, Frieburg) protonemata were grown on BCD  
443 medium supplemented with 5 mM ammonium tartrate (BCDAT) or 0.5% glucose during a 16-h  
444 photoperiod at 25°C in 9-cm Petri dishes (Nishiyama et al. 2000). For all analyses, the protonemata  
445 were collected every 5 days. The gametophores were grown on free-ammonium tartrate BCD  
446 medium under the same conditions, and 8-week-old gametophores were used for analysis.  
447 Protoplast was prepared from protonemata as described previously (Fesenko et al. 2015).

448 For morphological analysis, protonemal tissue 2 mm in diameter were inoculated on BCD and  
449 BCDAT 9-cm Petri dishes. For growth rate measurements, photographs were taken at 7 d intervals  
450 over 42 days. Protonemal tissues and cells were photographed using a Microscope Digital Eyepiece  
451 DCM-510 attached to a Stemi 305 stereomicroscope or Olympus CKX41.

452

### 453 **Identification of coding sORFs in the *P. patens* genome**

454 To identify sORFs with high coding potential, the sORFinder (Hanada et al. 2010) tool was utilized.  
455 Intron sequences and CDS were used as negative and positive sets, respectively. Additional details  
456 are described in the Supplemental Methods. To select for sORFs that are transcribed, located in the  
457 exons of transcripts, and have introns, a BED file was generated using a Python script (GffParser.py)  
458 and intersected with exon positions extracted from a gff3 file of *P. patens* genome annotations. To  
459 identify intergenic-sORFs, the BED file was also intersected with transcribed regions determined  
460 based on our RNA-seq data (Fesenko et al. 2017). Using an R script, sORFs fully overlapping with  
461 exons were removed; 75,685 sORFs remained after this step. Identical sORFs were removed from the  
462 dataset. In addition, sORFs overlapping repetitive regions identified by RepeatMasker (Tempel  
463 2012), as well as sORFs comprising parts of annotated *P. patens* main and alternative protein

464 isoforms, were also removed from the dataset, resulting in a final dataset of sORFs comprising  
465 70,095 sequences.

466

#### 467 **sORF classification**

468 The step-by-step procedure performed for sORF classification is illustrated in Supplemental Fig. S17.  
469 In the first step, lncRNA-sORFs were identified by searching for identical sORFs in known lncRNA  
470 databases, including CANTATAdb (Szczesniak et al. 2016), GreenNC (Paytuyi Gallart et al. 2016) and  
471 our previously published moss dataset (Fesenko et al. 2017). After this sORF BED file was  
472 intersected with the latest moss genome annotation V3.3 (Lang et al. 2018). The locations of the  
473 sORFs on transcripts were determined, resulting in the further classification of genic-sORFs into  
474 uORFs, dORFs, CDS-sORFs and interlaced-sORFs. sORFs were denoted as upstream or downstream if  
475 they were fully separated from the longer protein-coding ORF as previously described (Calviello et al.  
476 2016; Samandi et al. 2017).

477 Because alternative splicing leads to inaccuracy in genome annotation, the locations of a  
478 subset of genic-sORFs cannot be unambiguously classified, as they can be located in different regions  
479 in different isoforms of the same gene. All sORFs located on transcripts that were not annotated in  
480 the *P. patens* genome V3.3 but were identified using our RNA-seq data were classified as intergenic-  
481 sORFs. To detect alternatively spliced sORFs (AS-sORFs), a BED file with sORF locations was  
482 intersected with a BED file containing intron coordinates for all isoforms. Those sORFs that  
483 overlapped for both exons (see above) and introns were classified as AS-sORFs.

484

#### 485 **Evolutionary conservation analysis**

486 The transcriptomes of nine plant species were downloaded from Phytozome v12: *Sphagnum fallax*  
487 (release 0.5), *Marchantia polymorpha* (release 3.1), *Selaginella moellendorffii* (release 1.0), *Spirodela*  
488 *polyrhiza* (release 2), *Arabidopsis thaliana* (TAIR 10), *Zea mays* (Ensembl-18), *Oryza sativa* (release  
489 7), *Volvox carteri* (release 2.1) and *Chlamydomonas reinhardtii* (release 5.5). The transcriptome of  
490 *Ceratodon purpureus* was *de novo* assembled using Trinity (Haas et al. (2013)). To identify

491 transcribed homologous sequences, TBLASTN (word size = 3) was performed using sORF peptide  
492 sequences as queries and the transcriptome sequences of the abovementioned species as subjects.  
493 The following cutoffs parameters were used to distinguish reliable alignments: e-value <  $10^{-6}$  and  
494 query coverage > 60%. Our e-value cutoff was obtained by applying a multiple comparison  
495 correction (Bonferroni correction) of 0.05, which is commonly used in biological experiments.

496         Pairwise  $K_A/K_S$  ratios were calculated using the codeml algorithm with PAML software (Yang  
497 2007). The calculation procedure, which was facilitated using a custom-made Python script  
498 (protein\_Ka\_Ks\_codeml.py), included alignment extraction from the TBLASTN output, PAL2NAL  
499 (Suyama et al. 2006) correction of the nucleotide alignment using the corresponding aligned protein  
500 sequences and calculation of  $K_A/K_S$  ratios using codeml. The script implements packages from  
501 biopython (Cock et al. 2009). To estimate homologous sORF lengths, a Python script  
502 (sORF\_completeness\_v2.0.py) was designed. Additional details are described in the Supplemental  
503 Methods.

504

#### 505 **Gene Ontology (GO) terms enrichment analysis**

506 GO enrichment analysis was performed using the topGO bioconductor R package using the Fisher's  
507 exact test in conjunction with the 'classic' algorithm (false discovery rate [FDR] < 0.05). GO terms  
508 assigned to *P. patens* genes were downloaded from Phytozome. Only GO terms containing >5 genes in  
509 a background dataset were considered in the enrichment analysis. Redundant GO terms were  
510 removed using the web-based tool REVIGO (Supek et al. 2011).

#### 511 **Peptide and protein extraction**

512 Endogenous peptide extraction was conducted as described previously (Fesenko et al. 2015). Proteins were  
513 extracted as described previously (Fesenko et al. 2016). Additional details are described in the  
514 Supplemental Methods.

#### 515 **Mass-spectrometry analysis and peptide identification**

516 Mass-spectrometry analysis was performed using three biological and three technical repeats for the  
517 proteomic and peptidomic datasets (Supplementary Table S2). Analysis was performed on two

518 different mass spectrometers: a TripleTOF 5600+ mass spectrometer with a NanoSpray III ion source  
519 (ABSciex,Canada) and a Q Exactive HF mass spectrometer (Q Exactive HF Hybrid Quadrupole-  
520 Orbitrap mass spectrometer, Thermo Fisher Scientific, USA). Additional details are described in the  
521 Supplemental Methods.

522 All datasets were searched individually with MaxQuant v1.5.8.3 (Tyanova et al. 2016) against  
523 a custom database containing 32926 proteins from annotated genes in the latest version of the moss  
524 genome (V3.3, (Lang et al. 2018)), 85 moss chloroplast proteins, 42 moss mitochondrial proteins and  
525 70052 predicted sORF peptides (Supplemental Code). MaxQuant's protein FDR filter was disabled,  
526 while 1% FDR was used to select high-confidence PSMs, and ambiguous peptides were filtered out.  
527 Moreover, any PSMs with Andromeda scores of less than 30 were discarded (to exclude poor MS/MS  
528 spectra). For the dataset of endogenous peptides (named "peptidomic", Supplementary Table S2), the  
529 parameter "Digestion Mode" was set to "unspecific" and modifications were not permitted. All other  
530 parameters were left as default values. For the dataset of tryptic peptides (named "proteomic") the  
531 parameter "Digestion Mode" was set to "specific" (the Trypsin/P), MaxQuant's protein FDR filter was  
532 disabled, and the peptide FDR remained at 1%. All other parameters were left as default values.  
533 Features of the PSMs (length, intensity, number of spectra, Andromeda score, intensity coverage and  
534 peak coverage) were extracted from MaxQuant's msms.txt files. Annotated spectra for identified  
535 sORFs were exported from MaxQuant (Supplemental Fig. S18) and manually inspected.

536 To filter out MS peptides that do not provide unambiguous evidence of sORF peptide  
537 expression, we assessed the number of times a peptide occurred in the whole moss genome by  
538 searching for exact matches to the MS peptides in the six-frame translated genome (see  
539 Supplemental Methods).

#### 540 **RT-PCR analysis of AS-sORFs**

541 Total RNA from gametophores, protonema and protoplasts was isolated as previously described  
542 (Cove et al. 2009). RNA quality and quantity were evaluated via electrophoresis in an agarose gel  
543 with ethidium bromide staining. The precise concentration of total RNA in each sample was  
544 measured using a Quant-iT™ RNA Assay Kit, 5–100 ng on a Qubit 3.0 (Invitrogen, US) fluorometer.  
545 The cDNA for RT-PCR was synthesized using an MMLV RT Kit (Evrogen, Russia) according to the

546 manufacturer's recommendations employing oligo(dT)<sub>17</sub> -primers from 2 µg total RNA after DNase  
547 treatment. The primers were designed using Primer-BLAST (Ye et al. 2012) (Supplemental Table  
548 S7). The minus reverse transcriptase control (-RT) contained RNA without reverse transcriptase  
549 treatment to confirm the absence of DNA in the samples. The RT-PCR products were resolved on an  
550 1.5% agarose gel and visualized using ethidium bromide staining.

#### 551 **Generation of overexpression and knockout lines**

552 To obtain PSEP1 (Pp3c9\_sORF1544), PSEP3 (Pp3c25\_sORF1253), PSEP25 (Pp3c25\_sORF1000) and  
553 PSPE18 (Pp3c18\_sORF57) overexpression lines, PCR was carried out using genomic DNA as a  
554 template and the PEP4f, PEP4r, pep3FXho, pep3RNhe, pep25FXho, pep25RNhe, pep18FXho and  
555 pep18RNhe primers, respectively (Supplemental Table S7). Amplicons were cloned into the pPLV27  
556 vector (GenBank JF909480) using the ligation-independent cloning (LIC) procedure (Aslanidis and  
557 de Jong 1990; De Rybel et al. 2011). The resulting plasmids were named pPLV-Hpa-4FR (*PSEP1*),  
558 pPLV-Hpa-3FR (*PSEP3*), pPLV-Hpa-25FR (*PSEP25*), pPLV-Hpa-18FR (*PSEP18*) and used for  
559 transformation. Additional details are described in the Supplemental Methods.

560 PSEP1 (sORF Pp3c9\_sORF1544), PSEP3 (Pp3c25\_sORF1253), PSEP25 (Pp3c25\_sORF1000),  
561 PSPE18 (sORF Pp3c18\_sORF57) knockout lines were created using the CRISPR/Cas9 system  
562 (Collonnier et al. 2017). The coding sequences were used to search for CRISPR RNA (crRNA)  
563 preceded by a *S. pyogenes* Cas9 PAM motif (NGG) using the web tool CRISPR DESIGN  
564 (<http://crispr.mit.edu/>). The crRNA closest to the translation start site (ATG) was selected for  
565 cloning (Supplemental Table S7).

566 Protoplasts were transformed using PEG transformation protocol (Schaefer and Zryd 1997).  
567 Additional details are described in the Supplemental Methods. The plasmids pACT-CAS9 (for CAS9  
568 expression) and pBNRF (resistance to G418) were kindly provided by Dr. Fabien Nogu e.  
569 Independent knockout and overexpression mutant lines have been obtained (Supplemental Figs. S9-  
570 12).

571 The ploidy level of the *PSEP1* overexpression and *psep1* knock-out lines were estimated using  
572 flow cytometry. Protoplasts were fixed in cold 70% methanol, washed in TBS with 0.1% Triton X-  
573 100, then washed with TBS and stained with 500 ng/ml DAPI. The fluorescence was analyzed with a

574 flow cytometer NovoCyte (ACEA Biosciences) and Novoexpress data software. Fluorescence was  
575 excited at 405 nm, and detection was at 445/45 nm.

576

#### 577 **DATA ACCESS**

578 All raw mass spectrometry data from this study have been deposited to the ProteomeXchange  
579 Consortium via the PRIDE (Vizcaino et al. 2016) partner repository with the dataset identifiers  
580 PXD007922, PXD007923, PXD007973.

581

#### 582 **SOFTWARE AVAILABILITY**

583 All data were analyzed using Python (<http://www.python.org>, v 3.5), and R (R Core Team 2017). All  
584 scripts are available at Zenodo (doi: 10.5281/zenodo.1160331) and are maintained in the GitHub  
585 code repository: [https://github.com/Kirovez/Scripts\\_sORFs\\_MS](https://github.com/Kirovez/Scripts_sORFs_MS) and in Supplemental Code.

586

#### 587 **ACKNOWLEDGEMENTS**

588 We thank Dr. James Lloyd for his assistance in NMD-targeted transcript analysis. This work was supported  
589 by the Russian Science Foundation (project No.17-14-01189). Some of mass spectrometric  
590 measurements were performed using the equipment of the “Human Proteome” Core Facility of the  
591 Orekhovich Institute of Biomedical Chemistry (Russia) which is supported by the Ministry of  
592 Education and Science of the Russian Federation.

#### 593 **Authors' contributions**

594 IF and IK conceived and designed experiments. AK, AM performed moss transformation experiments.  
595 RK, VL, DK, EG, VZ, IB and AM performed the proteomics analyses. IF, IK and GA performed the  
596 statistical and bioinformatics analyses. IF, IK, VI and VG wrote the manuscript with input from all  
597 authors. IF supervised the project. All authors read and approved the final manuscript.

#### 598 **DISCLOSURE DECLARATION**

599 The authors declare that they have no significant competing financial, professional, or personal  
600 interests that might have influenced the performance or presentation of the work described in this  
601 manuscript.

602

603 **REFERENCES**

- 604 Andrews SJ, Rothnagel JA. 2014. Emerging evidence for functional peptides encoded by short  
605 open reading frames. *Nat Rev Genet* **15**(3): 193-204.
- 606 Aslanidis C, de Jong PJ. 1990. Ligation-independent cloning of PCR products (LIC-PCR).  
607 *Nucleic acids research* **18**(20): 6069-6074.
- 608 Aspden JL, Eyre-Walker YC, Phillips RJ, Amin U, Mumtaz MA, Brocard M, Couso JP. 2014.  
609 Extensive translation of small Open Reading Frames revealed by Poly-Ribo-Seq. *eLife*  
610 **3**: e03528.
- 611 Barbosa C, Peixeiro I, Romao L. 2013. Gene expression regulation by upstream open reading  
612 frames and human disease. *PLoS genetics* **9**(8): e1003529.
- 613 Bazzini AA, Johnstone TG, Christiano R, Mackowiak SD, Obermayer B, Fleming ES, Vejnar CE,  
614 Lee MT, Rajewsky N, Walther TC et al. 2014. Identification of small ORFs in  
615 vertebrates using ribosome footprinting and evolutionary conservation. *The EMBO*  
616 *journal* **33**(9): 981-993.
- 617 Blanvillain R, Young B, Cai YM, Hecht V, Varoquaux F, Delorme V, Lancelin JM, Delseny M,  
618 Gallois P. 2011. The Arabidopsis peptide kiss of death is an inducer of programmed  
619 cell death. *The EMBO journal* **30**(6): 1173-1183.
- 620 Brunet MA, Brunelle M, Lucier JF, Delcourt V, Levesque M, Grenier F, Samandi S, Leblanc S,  
621 Aguilar JD, Dufour P et al. 2019. OpenProt: a more comprehensive guide to explore  
622 eukaryotic coding potential and proteomes. *Nucleic acids research* **47**(D1): D403-  
623 D410.
- 624 Calviello L, Mukherjee N, Wyler E, Zauber H, Hirsekorn A, Selbach M, Landthaler M,  
625 Obermayer B, Ohler U. 2016. Detecting actively translated open reading frames in  
626 ribosome profiling data. *Nature methods* **13**(2): 165-170.
- 627 Carvunis AR, Rolland T, Wapinski I, Calderwood MA, Yildirim MA, Simonis N, Charloteaux B,  
628 Hidalgo CA, Barbette J, Santhanam B et al. 2012. Proto-genes and de novo gene birth.  
629 *Nature* **487**(7407): 370-374.
- 630 Chang CY, Lin WD, Tu SL. 2014. Genome-Wide Analysis of Heat-Sensitive Alternative Splicing  
631 in *Physcomitrella patens*. *Plant Physiol* **165**(2): 826-840.
- 632 Cherry JL. 2010. Expression level, evolutionary rate, and the cost of expression. *Genome*  
633 *biology and evolution* **2**: 757-769.
- 634 Chilley PM, Casson SA, Tarkowski P, Hawkins N, Wang KL, Hussey PJ, Beale M, Ecker JR,  
635 Sandberg GK, Lindsey K. 2006. The POLARIS peptide of Arabidopsis regulates auxin  
636 transport and root growth via effects on ethylene signaling. *Plant Cell* **18**(11): 3058-  
637 3072.
- 638 Cock PJ, Antao T, Chang JT, Chapman BA, Cox CJ, Dalke A, Friedberg I, Hamelryck T, Kauff F,  
639 Wilczynski B et al. 2009. Biopython: freely available Python tools for computational  
640 molecular biology and bioinformatics. *Bioinformatics* **25**(11): 1422-1423.
- 641 Collonnier C, Epert A, Mara K, Maclot F, Guyon-Debast A, Charlot F, White C, Schaefer DG,  
642 Nogue F. 2017. CRISPR-Cas9-mediated efficient directed mutagenesis and RAD51-  
643 dependent and RAD51-independent gene targeting in the moss *Physcomitrella*  
644 *patens*. *Plant Biotechnol J* **15**(1): 122-131.
- 645 Couso JP. 2015. Finding smORFs: getting closer. *Genome Biol* **16**: 189.
- 646 Couso JP, Patraquim P. 2017. Classification and function of small open reading frames.  
647 *Nature reviews Molecular cell biology* **18**(9): 575-589.
- 648 Cove DJ, Perroud PF, Charron AJ, McDaniel SF, Khandelwal A, Quatrano RS. 2009. Isolation of  
649 DNA, RNA, and protein from the moss *Physcomitrella patens* gametophytes. *Cold*  
650 *Spring Harbor protocols* **2009**(2): pdb prot5146.

- 651 D'Lima NG, Ma J, Winkler L, Chu Q, Loh KH, Corpuz EO, Budnik BA, Lykke-Andersen J,  
652 Saghatelyan A, Slavoff SA. 2017. A human microprotein that interacts with the mRNA  
653 decapping complex. *Nat Chem Biol* **13**(2): 174-180.
- 654 De Coninck B, Carron D, Tavormina P, Willem L, Craik DJ, Vos C, Thevissen K, Mathys J,  
655 Cammue BP. 2013. Mining the genome of *Arabidopsis thaliana* as a basis for the  
656 identification of novel bioactive peptides involved in oxidative stress tolerance. *J Exp*  
657 *Bot* **64**(17): 5297-5307.
- 658 De Rybel B, van den Berg W, Lokerse A, Liao CY, van Mourik H, Moller B, Peris CL, Weijers D.  
659 2011. A versatile set of ligation-independent cloning vectors for functional studies in  
660 plants. *Plant Physiol* **156**(3): 1292-1299.
- 661 Delcourt V, Brunelle M, Roy AV, Jacques JF, Salzet M, Fournier I, Roucou X. 2018. The Protein  
662 Coded by a Short Open Reading Frame, Not by the Annotated Coding Sequence, Is the  
663 Main Gene Product of the Dual-Coding Gene MIEF1. *Mol Cell Proteomics* **17**(12):  
664 2402-2411.
- 665 Djordjevic MA, Mohd-Radzman NA, Imin N. 2015. Small-peptide signals that control root  
666 nodule number, development, and symbiosis. *J Exp Bot* **66**(17): 5171-5181.
- 667 Eguen T, Straub D, Graeff M, Wenkel S. 2015. MicroProteins: small size-big impact. *Trends*  
668 *Plant Sci* **20**(8): 477-482.
- 669 Fesenko I, Khazigaleeva R, Kirov I, Kniazhev A, Glushenko O, Babalyan K, Arapidi G, Shashkova  
670 T, Butenko I, Zgoda V et al. 2017. Alternative splicing shapes transcriptome but not  
671 proteome diversity in *Physcomitrella patens*. *Scientific reports* **7**(1): 2698.
- 672 Fesenko I, Seredina A, Arapidi G, Ptushenko V, Urban A, Butenko I, Kovalchuk S, Babalyan K,  
673 Knyazev A, Khazigaleeva R et al. 2016. The *Physcomitrella patens* Chloroplast  
674 Proteome Changes in Response to Protoplastation. *Front Plant Sci* **7**: 1661.
- 675 Fesenko IA, Arapidi GP, Skripnikov AY, Alexeev DG, Kostyukova ES, Manolov AI, Altukhov  
676 IA, Khazigaleeva RA, Seredina AV, Kovalchuk SI et al. 2015. Specific pools of  
677 endogenous peptides are present in gametophore, protonema, and protoplast cells of  
678 the moss *Physcomitrella patens*. *Bmc Plant Biol* **15**: 87.
- 679 Ge Y, Porse BT. 2014. The functional consequences of intron retention: alternative splicing  
680 coupled to NMD as a regulator of gene expression. *BioEssays : news and reviews in*  
681 *molecular, cellular and developmental biology* **36**(3): 236-243.
- 682 Giannakakis A, Zhang J, Jenjaroenpun P, Nama S, Zainolabidin N, Aau MY, Yarmishyn AA, Vaz  
683 C, Ivshina AV, Grinchuk OV et al. 2015. Contrasting expression patterns of coding and  
684 noncoding parts of the human genome upon oxidative stress. *Scientific reports* **5**:  
685 9737.
- 686 Graeff M, Straub D, Eguen T, Dolde U, Rodrigues V, Brandt R, Wenkel S. 2016. MicroProtein-  
687 Mediated Recruitment of CONSTANS into a TOPLESS Trimeric Complex Represses  
688 Flowering in *Arabidopsis*. *PLoS genetics* **12**(3): e1005959.
- 689 Guillen G, Diaz-Camino C, Loyola-Torres CA, Aparicio-Fabre R, Hernandez-Lopez A, Diaz-  
690 Sanchez M, Sanchez F. 2013. Detailed analysis of putative genes encoding small  
691 proteins in legume genomes. *Front Plant Sci* **4**.
- 692 Guttman M, Russell P, Ingolia NT, Weissman JS, Lander ES. 2013. Ribosome Profiling  
693 Provides Evidence that Large Noncoding RNAs Do Not Encode Proteins. *Cell* **154**(1):  
694 240-251.
- 695 Hanada K, Akiyama K, Sakurai T, Toyoda T, Shinozaki K, Shiu SH. 2010. sORF finder: a  
696 program package to identify small open reading frames with high coding potential.  
697 *Bioinformatics* **26**(3): 399-400.
- 698 Hanada K, Higuchi-Takeuchi M, Okamoto M, Yoshizumi T, Shimizu M, Nakaminami K, Nishi R,  
699 Ohashi C, Iida K, Tanaka M et al. 2013. Small open reading frames associated with  
700 morphogenesis are hidden in plant genomes. *P Natl Acad Sci USA* **110**(6): 2395-2400.

- 701 Hanada K, Zhang X, Borevitz JO, Li WH, Shiu SH. 2007. A large number of novel coding small  
702 open reading frames in the intergenic regions of the Arabidopsis thaliana genome are  
703 transcribed and/or under purifying selection. *Genome Res* **17**(5): 632-640.
- 704 Hellens RP, Brown CM, Chisnal MAW, Waterhouse PM, Macknight RC. 2016. The Emerging  
705 World of Small ORFs. *Trends Plant Sci* **21**(4): 317-328.
- 706 Huang JZ, Chen M, Chen, Gao XC, Zhu S, Huang H, Hu M, Zhu H, Yan GR. 2017. A Peptide  
707 Encoded by a Putative lncRNA HOXB-AS3 Suppresses Colon Cancer Growth. *Mol Cell*  
708 **68**(1): 171-184 e176.
- 709 Iacono M, Mignone F, Pesole G. 2005. uAUG and uORFs in human and rodent 5'untranslated  
710 mRNAs. *Gene* **349**: 97-105.
- 711 Ingolia NT, Lareau LF, Weissman JS. 2011. Ribosome Profiling of Mouse Embryonic Stem  
712 Cells Reveals the Complexity and Dynamics of Mammalian Proteomes. *Cell* **147**(4):  
713 789-802.
- 714 Johnstone TG, Bazzini AA, Giraldez AJ. 2016. Upstream ORFs are prevalent translational  
715 repressors in vertebrates. *The EMBO journal* **35**(7): 706-723.
- 716 Karousis ED, Nasif S, Muhlemann O. 2016. Nonsense-mediated mRNA decay: novel  
717 mechanistic insights and biological impact. *Wiley interdisciplinary reviews RNA* **7**(5):  
718 661-682.
- 719 Kastenmayer JP, Ni L, Chu A, Kitchen LE, Au WC, Yang H, Carter CD, Wheeler D, Davis RW,  
720 Boeke JD et al. 2006. Functional genomics of genes with small open reading frames  
721 (sORFs) in *S-cerevisiae*. *Genome Res* **16**(3): 365-373.
- 722 Kim TS, Liu CL, Yassour M, Holik J, Friedman N, Buratowski S, Rando OJ. 2010. RNA  
723 polymerase mapping during stress responses reveals widespread nonproductive  
724 transcription in yeast. *Genome Biol* **11**(7): R75.
- 725 Kondo T, Plaza S, Zanet J, Benrabah E, Valenti P, Hashimoto Y, Kobayashi S, Payre F,  
726 Kageyama Y. 2010. Small Peptides Switch the Transcriptional Activity of Shavenbaby  
727 During *Drosophila* Embryogenesis. *Science* **329**(5989): 336-339.
- 728 Kozak M. 1986. Point mutations define a sequence flanking the AUG initiator codon that  
729 modulates translation by eukaryotic ribosomes. *Cell* **44**(2): 283-292.
- 730 -. 1997. Recognition of AUG and alternative initiator codons is augmented by G in position +4  
731 but is not generally affected by the nucleotides in positions +5 and +6. *The EMBO*  
732 *journal* **16**(9): 2482-2492.
- 733 Ladoukakis E, Pereira V, Magny EG, Eyre-Walker A, Couso JP. 2011. Hundreds of putatively  
734 functional small open reading frames in *Drosophila*. *Genome Biol* **12**(11).
- 735 Laing WA, Martinez-Sanchez M, Wright MA, Bulley SM, Brewster D, Dare AP, Rassam M,  
736 Wang D, Storey R, Macknight RC et al. 2015. An Upstream Open Reading Frame Is  
737 Essential for Feedback Regulation of Ascorbate Biosynthesis in Arabidopsis. *Plant Cell*  
738 **27**(3): 772-786.
- 739 Lang D, Ullrich KK, Murat F, Fuchs J, Jenkins J, Haas FB, Piednoel M, Gundlach H, Van Bel M,  
740 Meyberg R et al. 2018. The *Physcomitrella patens* chromosome-scale assembly  
741 reveals moss genome structure and evolution. *Plant J* **93**(3): 515-533.
- 742 Laressergues D, Couzigou JM, Clemente HS, Martinez Y, Dunand C, Becard G, Combier JP.  
743 2015. Primary transcripts of microRNAs encode regulatory peptides. *Nature*  
744 **520**(7545): 90-93.
- 745 Lease KA, Walker JC. 2006. The Arabidopsis unannotated secreted peptide database, a  
746 resource for plant peptidomics. *Plant Physiol* **142**(3): 831-838.
- 747 Lloyd JPB, Lang D, Zimmer AD, Causier B, Reski R, Davies B. 2018. The loss of SMG1 causes  
748 defects in quality control pathways in *Physcomitrella patens*. *Nucleic acids research*  
749 **46**(11): 5822-5836.

- 750 Ma J, Diedrich JK, Jungreis I, Donaldson C, Vaughan J, Kellis M, Yates JR, 3rd, Saghatelian A.  
751 2016. Improved Identification and Analysis of Small Open Reading Frame Encoded  
752 Polypeptides. *Analytical chemistry* **88**(7): 3967-3975.
- 753 Mackowiak SD, Zauber H, Bielow C, Thiel D, Kutz K, Calviello L, Mastrobuoni G, Rajewsky N,  
754 Kempa S, Selbach M et al. 2015. Extensive identification and analysis of conserved  
755 small ORFs in animals. *Genome Biol* **16**.
- 756 Magny EG, Pueyo JL, Pearl FM, Cespedes MA, Niven JE, Bishop SA, Couso JP. 2013. Conserved  
757 regulation of cardiac calcium uptake by peptides encoded in small open reading  
758 frames. *Science* **341**(6150): 1116-1120.
- 759 Matsumoto A, Pasut A, Matsumoto M, Yamashita R, Fung J, Monteleone E, Saghatelian A,  
760 Nakayama KI, Clohessy JG, Pandolfi PP. 2017. mTORC1 and muscle regeneration are  
761 regulated by the LINC00961-encoded SPAR polypeptide. *Nature* **541**(7636): 228-  
762 232.
- 763 Mazin PV, Fisunov GY, Gorbachev AY, Kapitskaya KY, Altukhov IA, Semashko TA, Alexeev DG,  
764 Govorun VM. 2014. Transcriptome analysis reveals novel regulatory mechanisms in a  
765 genome-reduced bacterium. *Nucleic acids research* **42**(21): 13254-13268.
- 766 McLysaght A, Guerzoni D. 2015. New genes from non-coding sequence: the role of de novo  
767 protein-coding genes in eukaryotic evolutionary innovation. *Philosophical  
768 transactions of the Royal Society of London Series B, Biological sciences* **370**(1678):  
769 20140332.
- 770 Moyers BA, Zhang J. 2016. Evaluating Phylostratigraphic Evidence for Widespread De Novo  
771 Gene Birth in Genome Evolution. *Molecular biology and evolution* **33**(5): 1245-1256.
- 772 Narita NN, Moore S, Horiguchi G, Kubo M, Demura T, Fukuda H, Goodrich J, Tsukaya H. 2004.  
773 Overexpression of a novel small peptide ROTUNDIFOLIA4 decreases cell proliferation  
774 and alters leaf shape in *Arabidopsis thaliana*. *Plant J* **38**(4): 699-713.
- 775 Neafsey DE, Galagan JE. 2007. Dual modes of natural selection on upstream open reading  
776 frames. *Molecular biology and evolution* **24**(8): 1744-1751.
- 777 Nelson BR, Makarewich CA, Anderson DM, Winders BR, Troupes CD, Wu F, Reese AL,  
778 McAnally JR, Chen X, Kavalali ET et al. 2016. A peptide encoded by a transcript  
779 annotated as long noncoding RNA enhances SERCA activity in muscle. *Science*  
780 **351**(6270): 271-275.
- 781 Nishiyama T, Hiwatashi Y, Sakakibara I, Kato M, Hasebe M. 2000. Tagged mutagenesis and  
782 gene-trap in the moss, *Physcomitrella patens* by shuttle mutagenesis. *DNA research :  
783 an international journal for rapid publication of reports on genes and genomes* **7**(1): 9-  
784 17.
- 785 Paytuvi Gallart A, Hermoso Pulido A, Anzar Martinez de Lagran I, Sanseverino W, Aiese  
786 Cigliano R. 2016. GREENC: a Wiki-based database of plant lncRNAs. *Nucleic acids  
787 research* **44**(D1): D1161-1166.
- 788 Popp MW, Maquat LE. 2013. Organizing principles of mammalian nonsense-mediated mRNA  
789 decay. *Annual review of genetics* **47**: 139-165.
- 790 Prabakaran S, Hemberg M, Chauhan R, Winter D, Tweedie-Cullen RY, Dittrich C, Hong E,  
791 Gunawardena J, Steen H, Kreiman G et al. 2014. Quantitative profiling of peptides  
792 from RNAs classified as noncoding. *Nature communications* **5**: 5429.
- 793 R Core Team. 2017. R: A Language and Environment for Statistical Computing.
- 794 Rasheed S, Bashir K, Nakaminami K, Hanada K, Matsui A, Seki M. 2016. Drought stress  
795 differentially regulates the expression of small open reading frames (sORFs) in  
796 *Arabidopsis* roots and shoots. *Plant signaling & behavior* **11**(8): e1215792.
- 797 Rensing SA, Ick J, Fawcett JA, Lang D, Zimmer A, Van de Peer Y, Reski R. 2007. An ancient  
798 genome duplication contributed to the abundance of metabolic genes in the moss  
799 *Physcomitrella patens*. *BMC evolutionary biology* **7**: 130.

- 800 Rensing SA, Lang D, Zimmer AD, Terry A, Salamov A, Shapiro H, Nishiyama T, Perroud PF,  
801 Lindquist EA, Kamisugi Y et al. 2008. The Physcomitrella genome reveals  
802 evolutionary insights into the conquest of land by plants. *Science* **319**(5859): 64-69.
- 803 Rohrig H, Schmidt J, Miklashevichs E, Schell J, John M. 2002. Soybean ENOD40 encodes two  
804 peptides that bind to sucrose synthase. *Proc Natl Acad Sci U S A* **99**(4): 1915-1920.
- 805 Rothnagel J, Menschaert G. 2018. Short Open Reading Frames and Their Encoded Peptides.  
806 *Proteomics* **18**(10): e1700035.
- 807 Rubtsova M, Naraykina Y, Vasilkova D, Meerson M, Zvereva M, Prassolov V, Lazarev V,  
808 Manuvera V, Kovalchuk S, Anikanov N et al. 2018. Protein encoded in human  
809 telomerase RNA is involved in cell protective pathways. *Nucleic acids research*  
810 **46**(17): 8966-8977.
- 811 Ruiz-Orera J, Alba MM. 2019. Translation of Small Open Reading Frames: Roles in Regulation  
812 and Evolutionary Innovation. *Trends in genetics : TIG* **35**(3): 186-198.
- 813 Ruiz-Orera J, Verdaguer-Grau P, Villanueva-Canas JL, Messegueur X, Alba MM. 2018.  
814 Translation of neutrally evolving peptides provides a basis for de novo gene  
815 evolution. *Nature ecology & evolution* **2**(5): 890-896.
- 816 Samandi S, Roy AV, Delcourt V, Lucier JF, Gagnon J, Beaudoin MC, Vanderperre B, Breton MA,  
817 Motard J, Jacques JF et al. 2017. Deep transcriptome annotation enables the discovery  
818 and functional characterization of cryptic small proteins. *eLife* **6**.
- 819 Schaefer DG, Zryd JP. 1997. Efficient gene targeting in the moss *Physcomitrella patens*. *Plant*  
820 *J* **11**(6): 1195-1206.
- 821 Seo PJ, Hong SY, Kim SG, Park CM. 2011. Competitive inhibition of transcription factors by  
822 small interfering peptides. *Trends Plant Sci* **16**(10): 541-549.
- 823 Seo PJ, Park MJ, Park CM. 2013. Alternative splicing of transcription factors in plant  
824 responses to low temperature stress: mechanisms and functions. *Planta* **237**(6):  
825 1415-1424.
- 826 Slavoff SA, Mitchell AJ, Schwaid AG, Cabili MN, Ma J, Levin JZ, Karger AD, Budnik BA, Rinn JL,  
827 Saghatelian A. 2013. Peptidomic discovery of short open reading frame-encoded  
828 peptides in human cells. *Nat Chem Biol* **9**(1): 59-+.
- 829 Staudt AC, Wenkel S. 2011. Regulation of protein function by 'microProteins'. *EMBO reports*  
830 **12**(1): 35-42.
- 831 Supek F, Bosnjak M, Skunca N, Smuc T. 2011. REVIGO summarizes and visualizes long lists of  
832 gene ontology terms. *Plos One* **6**(7): e21800.
- 833 Suyama M, Torrents D, Bork P. 2006. PAL2NAL: robust conversion of protein sequence  
834 alignments into the corresponding codon alignments. *Nucleic acids research* **34**(Web  
835 Server issue): W609-612.
- 836 Szczesniak MW, Rosikiewicz W, Makalowska I. 2016. CANTATadb: A Collection of Plant Long  
837 Non-Coding RNAs. *Plant Cell Physiol* **57**(1): e8.
- 838 Tavormina P, De Coninck B, Nikonorova N, De Smet I, Cammue BP. 2015. The Plant  
839 Peptidome: An Expanding Repertoire of Structural Features and Biological Functions.  
840 *Plant Cell* **27**(8): 2095-2118.
- 841 Tempel S. 2012. Using and understanding RepeatMasker. *Methods Mol Biol* **859**: 29-51.
- 842 Tharakan R, Kreimer S, Ubaida-Mohien C, Lavoie J, Olexiouk V, Menschaert G, Ingolia NT,  
843 Cole RN, Ishizuka K, Sawa A et al. 2019. A methodology for discovering novel brain-  
844 relevant peptides: Combination of ribosome profiling and peptidomics. *Neuroscience*  
845 *research*.
- 846 Tyanova S, Temu T, Cox J. 2016. The MaxQuant computational platform for mass  
847 spectrometry-based shotgun proteomics. *Nat Protoc* **11**(12): 2301-2319.
- 848 van der Lee R, Buljan M, Lang B, Weatheritt RJ, Daughdrill GW, Dunker AK, Fuxreiter M,  
849 Gough J, Gsponer J, Jones DT et al. 2014. Classification of intrinsically disordered  
850 regions and proteins. *Chemical reviews* **114**(13): 6589-6631.

- 851 Vanderperre B, Lucier JF, Bissonnette C, Motard J, Tremblay G, Vanderperre S, Wisztorski M,  
 852 Salzet M, Boisvert FM, Roucou X. 2013. Direct Detection of Alternative Open Reading  
 853 Frames Translation Products in Human Significantly Expands the Proteome. *Plos One*  
 854 **8**(8).
- 855 Verheggen K, Volders PJ, Mestdagh P, Menschaert G, Van Damme P, Gevaert K, Martens L,  
 856 Vandesompele J. 2017. Noncoding after All: Biases in Proteomics Data Do Not Explain  
 857 Observed Absence of lncRNA Translation Products. *J Proteome Res* **16**(7): 2508-2515.
- 858 Vizcaino JA, Csordas A, Del-Toro N, Dianes JA, Griss J, Lavidas I, Mayer G, Perez-Riverol Y,  
 859 Reisinger F, Ternent T et al. 2016. 2016 update of the PRIDE database and its related  
 860 tools. *Nucleic acids research* **44**(22): 11033.
- 861 Wu HP, Su YS, Chen HC, Chen YR, Wu CC, Lin WD, Tu SL. 2014. Genome-wide analysis of  
 862 light-regulated alternative splicing mediated by photoreceptors in *Physcomitrella*  
 863 *patens*. *Genome Biol* **15**(1): R10.
- 864 Yang Z. 2007. PAML 4: phylogenetic analysis by maximum likelihood. *Molecular biology and*  
 865 *evolution* **24**(8): 1586-1591.
- 866 Ye J, Coulouris G, Zaretskaya I, Cutcutache I, Rozen S, Madden TL. 2012. Primer-BLAST: a tool  
 867 to design target-specific primers for polymerase chain reaction. *Bmc Bioinformatics*  
 868 **13**: 134.
- 869

870

871 **FIGURE LEGENDS**

872 **Fig. 1. Several distinct types of sORFs are present in the moss genome. A** – Pipeline used in this  
 873 study to identify coding sORFs; **B** – Proposed classification of sORFs according to the types of  
 874 encoding transcripts: upstream ORFs (uORFs) and downstream ORFs (dORFs) in the untranslated  
 875 regions (UTRs) of canonical mRNAs; CDS-sORFs, which overlap with protein-coding sequences in  
 876 alternative (+2 or +3) reading frames or are truncated versions of proteins generated by alternative  
 877 splicing; interlaced-sORFs, which overlap both the protein-coding sequence and UTR on the same  
 878 transcript; lncRNA-sORFs and intergenic sORFs, which are located on short non-protein coding  
 879 transcripts.; **C** – Boxplot of the length distribution of sORFs in different groups.

880 **Fig. 2. Moss contains tens of small protein-coding sORFs. A** – Venn diagram showing the  
 881 distribution of the identified translatable sORFs among three types of moss cells; **B** – Length  
 882 distribution of various groups of small protein-coding sORFs; **C** - Distribution of small protein-coding  
 883 sORFs based on the suggested classification; **D** – Binary heatmap showing evidence of translation for  
 884 sORFs and proteins in multicoding genes in three moss tissues. G, N and P correspond to  
 885 gametophores, protonemata and protoplasts, respectively; **E** – Heatmap showing expression levels  
 886 (log<sub>10</sub> (RPKM)) for the lncRNAs (left) carrying sORFs (lncRNA-sORFs) and binary heatmap showing

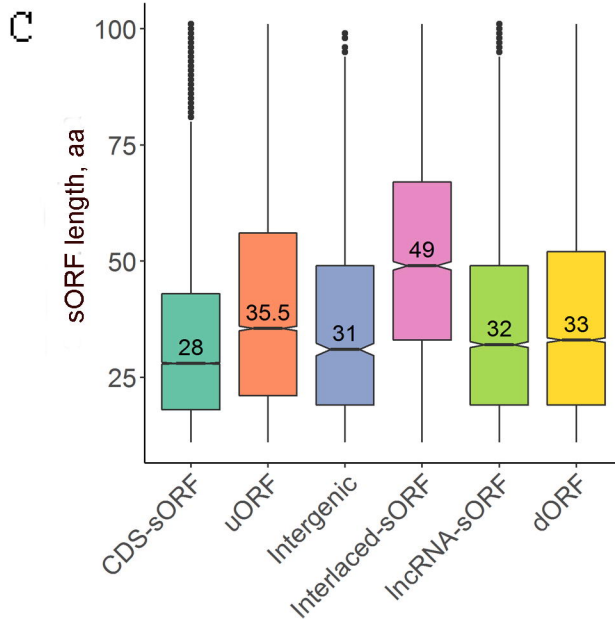
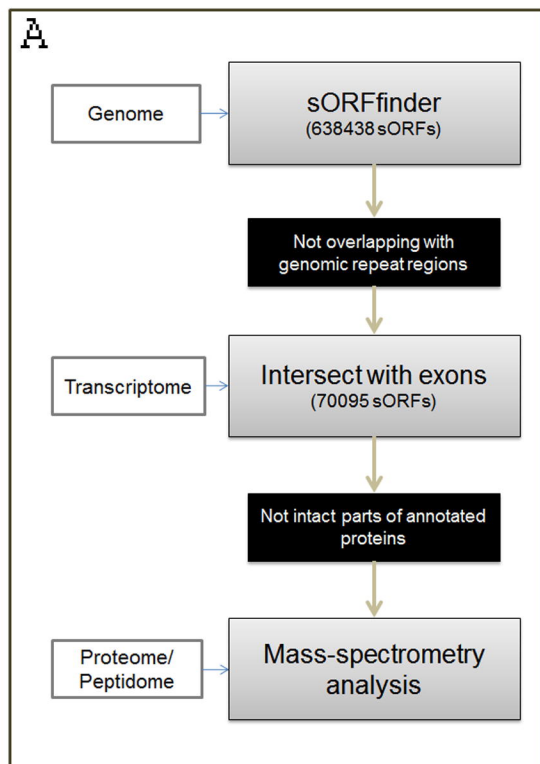
887 evidence of translation (determined as whether a peptide was identified (brown) or not (grey) in MS  
 888 data) for the corresponding lncRNA-sORFs (right) in three moss tissues, gametophores (G),  
 889 protonemata (N) and protoplasts (P).

890 **Fig. 3. Venn diagram showing the number of AS-sORFs influenced by different AS events.**

891 **Fig. 4. Morphology of wild type and sORF-encoded peptide mutant lines. The phenotypes of**  
 892 ***psep1* KO and *PSEP1* OE lines grown on BCD medium with 0.5% glucose:** A, D – wild type; B, E -  
 893 knockout of *PSEP1*; C, F – overexpression of *PSEP1*; G - the diameter of moss plants with  
 894 overexpression of *PSEP1* (Supplemental Figure S14A-D, Supplemental Table S6); H - the diameter of  
 895 moss plants with knockout of *PSEP1* (Supplemental Figure S13A-C, Supplemental Table S6). **The**  
 896 **phenotypes of *psep3* KO and *PSEP3* OE lines grown on BCD medium:** I, L - wild type; J, M -  
 897 knockout of *PSEP3*; K, N - overexpression of *PSEP3*; O - the diameter of moss plants with knockout of  
 898 *PSEP3* (Supplemental Figure S13G-I, Supplemental Table S6); P - the diameter of moss plants with  
 899 overexpression of the *PSEP3* (Supplemental Figure S14E-H, Supplemental Table S6). Scale bar: 0.5  
 900 mm. *P*-value was calculated by Student's unpaired *t*-test. \*\*\*\**P*-value < 0.0001; \*\*\**P*-value < 0.001  
 901 \*\**P*-value < 0.01; \**P*-value < 0.05.

902 **Fig. 5. Morphology of wild type and sORF-encoded peptide mutant lines. The phenotypes of**  
 903 ***psep25* KO and *PSEP25* OE lines grown on BCDAT medium:** A, D - wild type; B, E - knockout of  
 904 *PSEP25*; C, F – overexpression of *PSEP25*; G - the diameter of moss plants with knockout of *PSEP25*  
 905 (Supplemental Figure S13J-M, Supplemental Table S6); H - the diameter of moss plants with  
 906 overexpression of *PSEP25* (Supplemental Figure S14M-P, Supplemental Table S6); I, J, M - the  
 907 number of leafy shoots in wild type and three *psep25* KO lines. K, L, N - the number of leafy shoots in  
 908 wild type and two *PSEP25* OE lines on BCD medium. Arrows show young leafy gametophores. **The**  
 909 **phenotypes of *psep18* KO and *PSEP18* OE lines grown on BCD medium with 0.5% glucose:** O, R,  
 910 T – wild type; P - knockout of *PSEP18*; Q - the diameter of moss plants with knockout of *PSEP18*  
 911 (Supplemental Figure S16A-C, Supplemental Table S6); S, U– overexpression of *PSEP18*; V - the  
 912 diameter of moss plants with overexpression of *PSEP18* (Supplemental Figure S16D-F, Supplemental  
 913 Table S6).

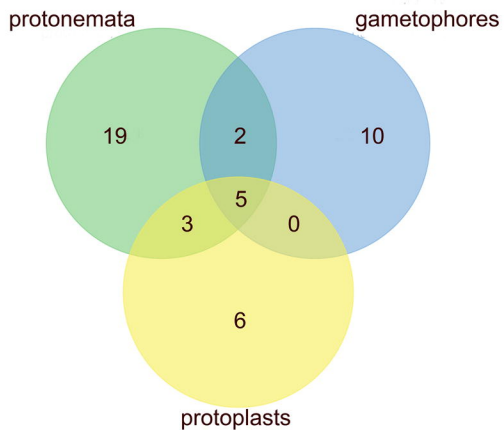
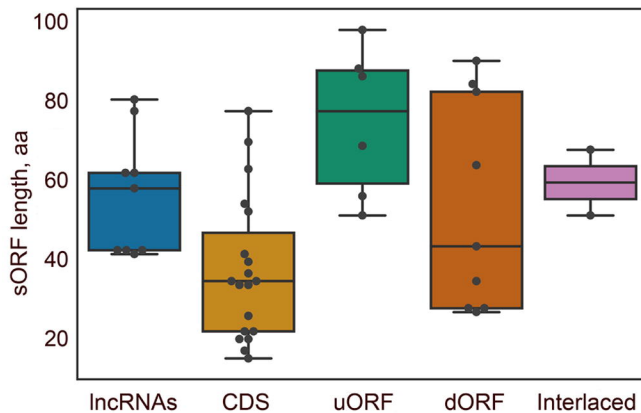
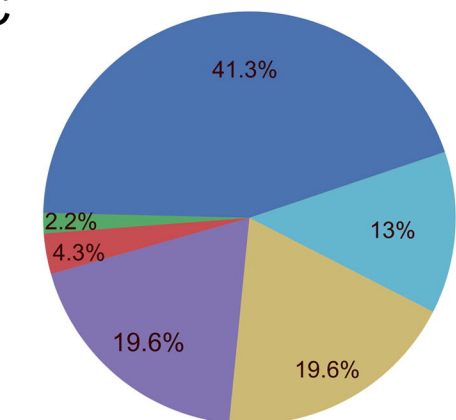
914 Scale bar: 0.5 mm. *P*-value was calculated by Student's unpaired *t*-test. \*\*\*\**P*-value < 0.0001; \*\*\**P*-  
915 value < 0.001 \*\**P*-value < 0.01; \**P*-value < 0.05.



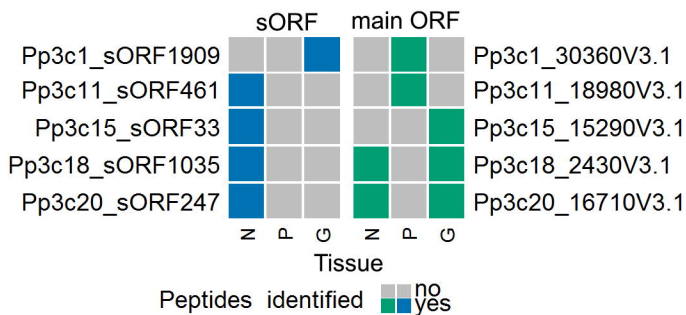
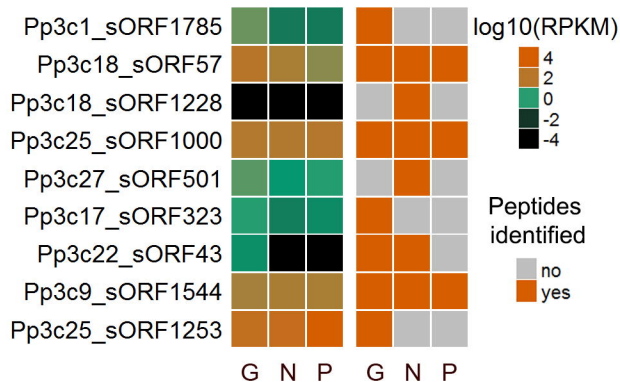
**B**

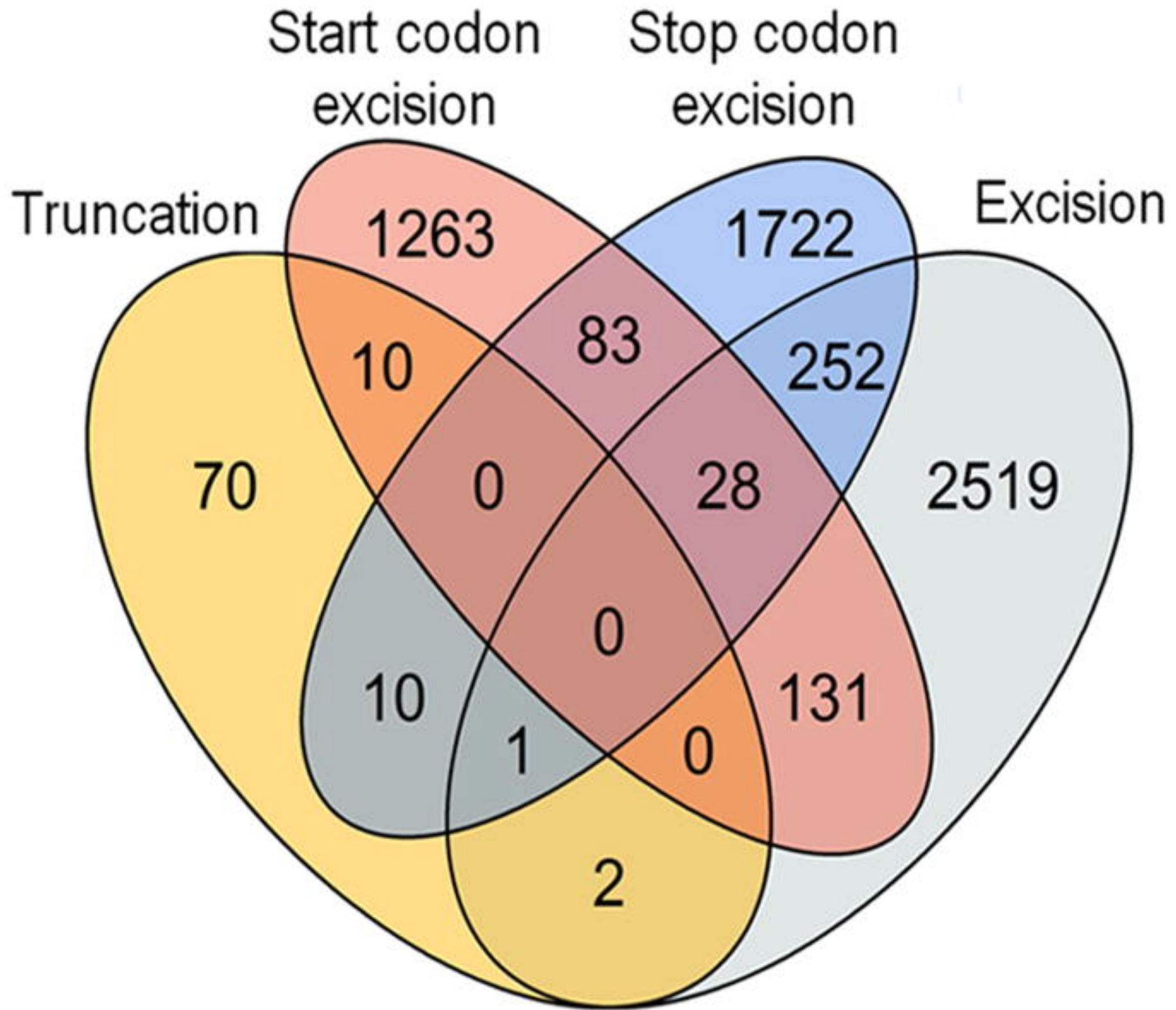
sORF class	RNA type	Evidence of transcription (RNA-seq data)	Evidence of translation (MS data)
Upstream sORFs (uORFs)		11998	6
Downstream sORFs (dORFs)		9444	9
Coding sequence-sORFs (CDS-sORFs)		36732	19
Interlaced-sORFs		3485	2
Intergenic/lncRNA-sORFs		1241/5745	1/9

main ORF   
  untranslated region   
  sORF   
  intron

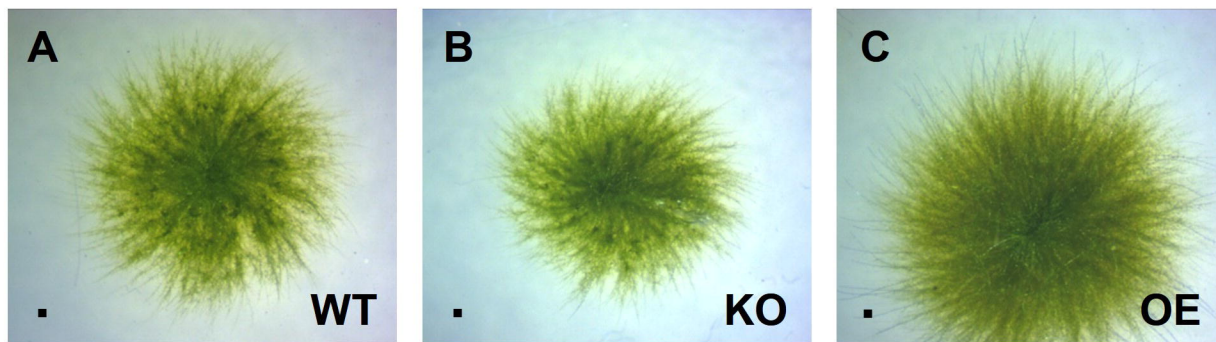
**A****B****C**

- CDS-sORFs
- uORFs
- dORFs
- lncRNA-sORFs
- Interlaced sORFs
- Intergenic sORFs

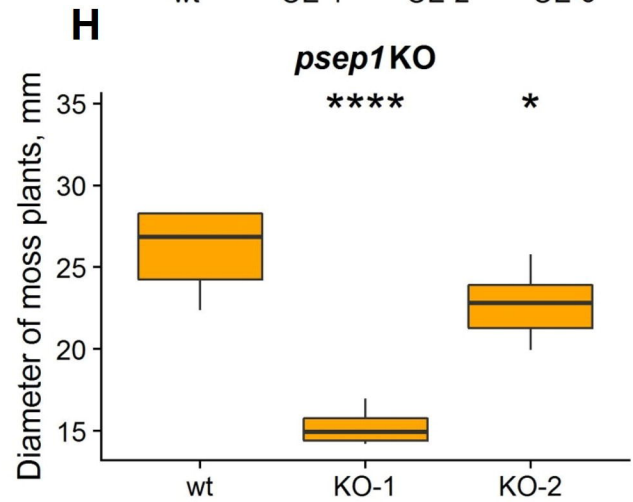
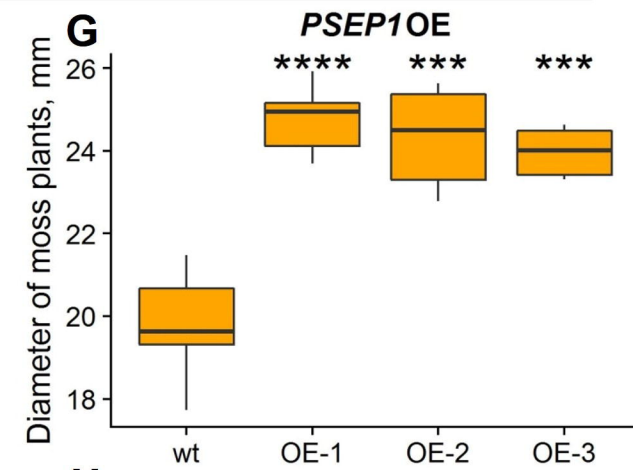
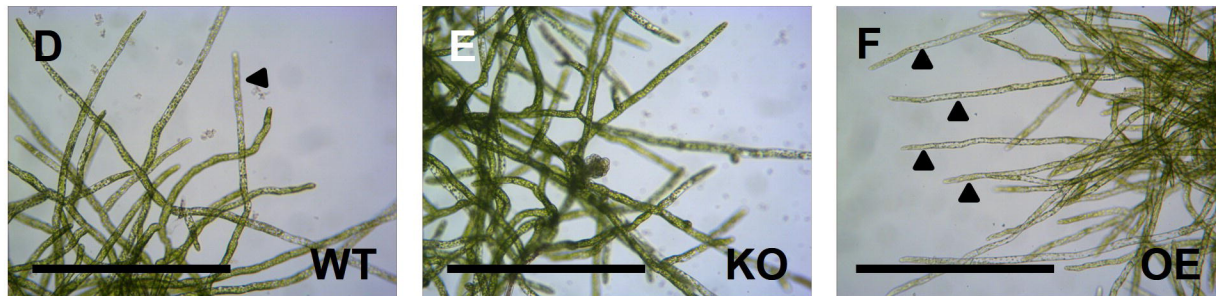
**D****E**



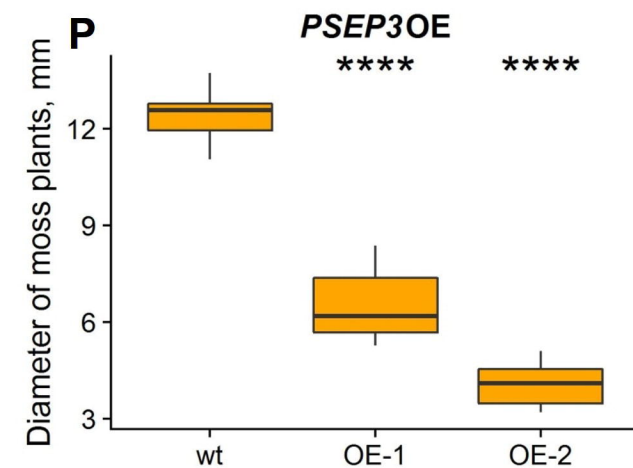
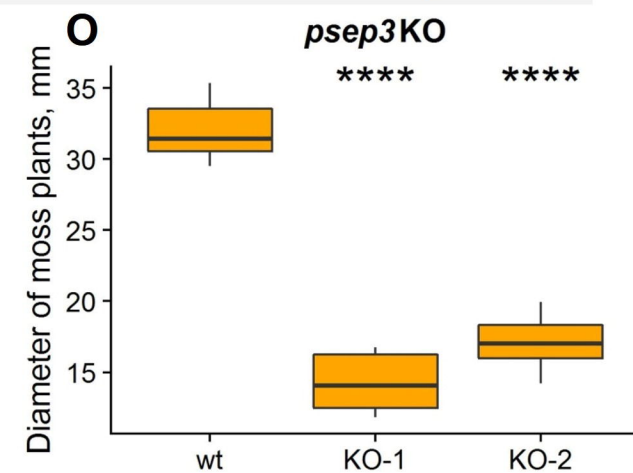
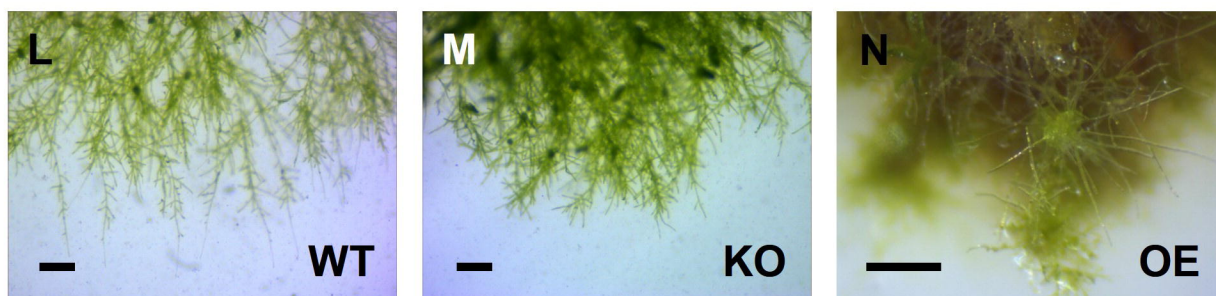
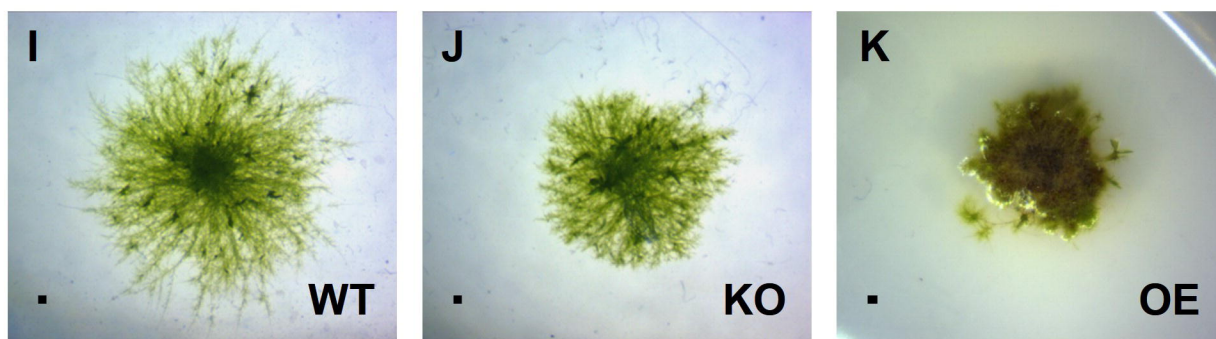
**A peptide PSEP1**  
MVQPLLARLASAAEFVALPGAILVAYFSTSRSTEPKRDHRK (41 aa)



Downloaded from [genome.cshlp.org](http://genome.cshlp.org) on June 13, 2026 . Published by Cold Spring Harbor Laboratory Press

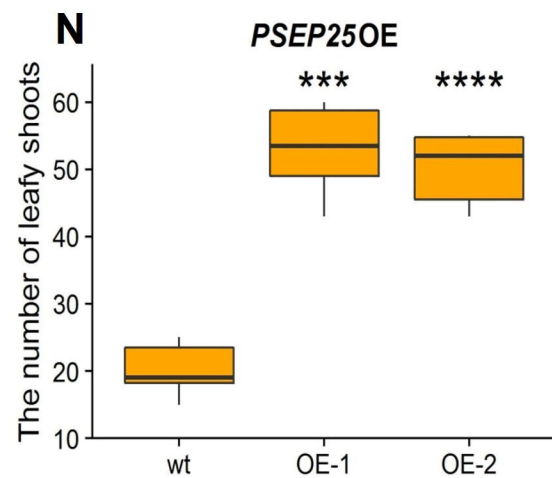
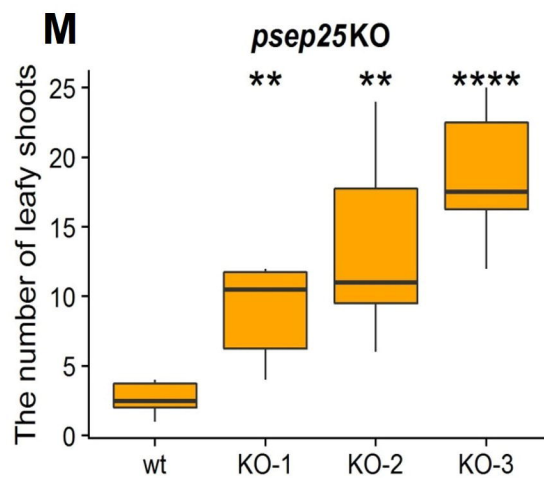
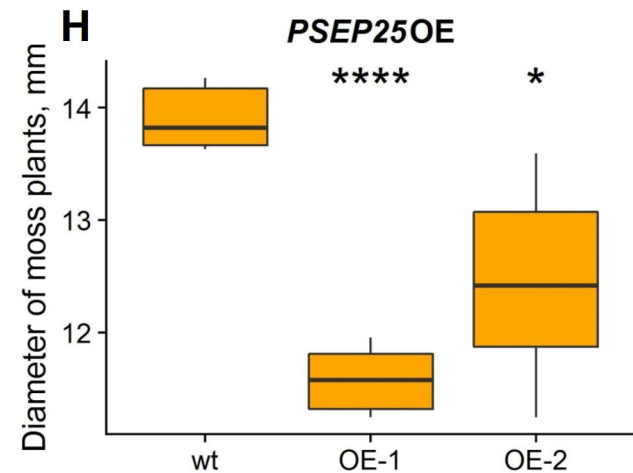
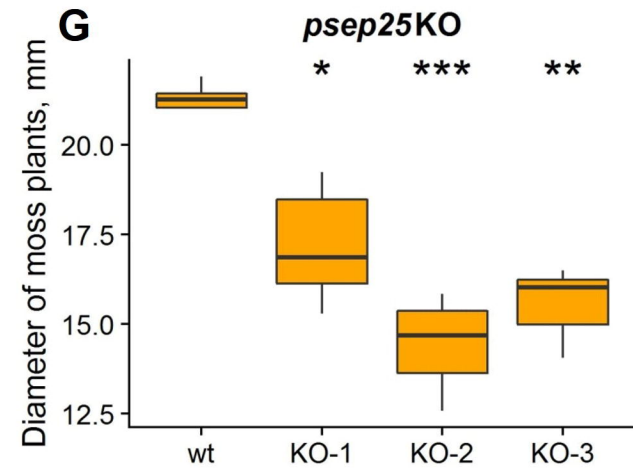
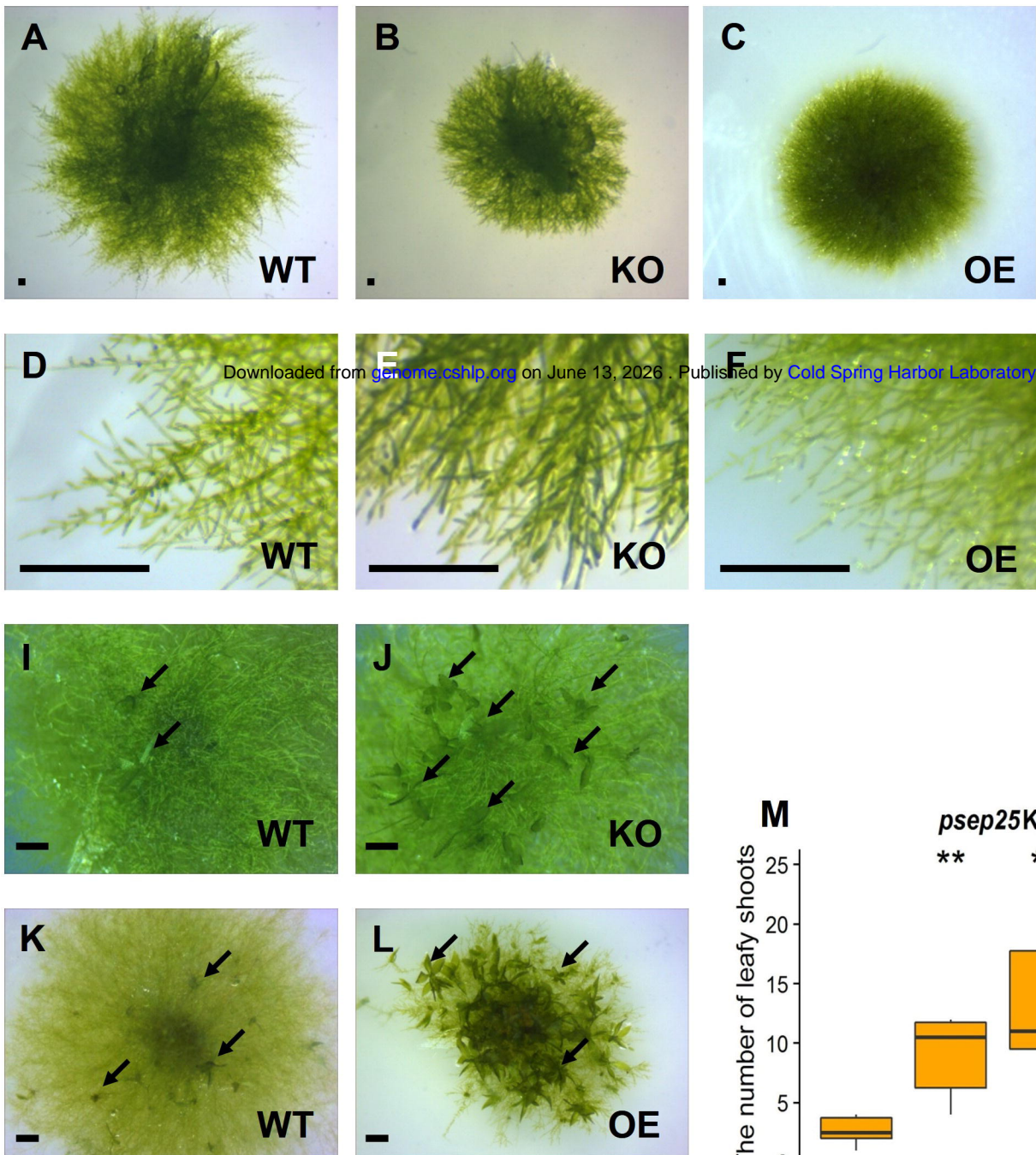


**A peptide PSEP3**  
MVHQDNSGSGLRFSFNHPNPPPNNNRPPSNPPVVRNPSSGRTPHPYPPPPHNYNGYPN (57 aa)



## A peptide PSEP25

MVQSKQGLSLLKFIPKVIRPQTSDVSSAVLWGTTAACGALWLVQPFDWIKEQITGPKEESK (61 aa)



## A peptide PSEP18

MQAFTDTQGYSSFNGPATTATTPPEVVGEGGKGWRPSS (40 aa)

