



**HAL**  
open science

## **P-GRe: an efficient pipeline to maximised pseudogene prediction in plants/eucaryotes**

Sébastien Cabanac, Christophe Dunand, Catherine Mathé

► **To cite this version:**

Sébastien Cabanac, Christophe Dunand, Catherine Mathé. P-GRe: an efficient pipeline to maximised pseudogene prediction in plants/eucaryotes. 2024. hal-04483163

**HAL Id: hal-04483163**

**<https://ut3-toulouseinp.hal.science/hal-04483163>**

Preprint submitted on 29 Feb 2024

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution - NonCommercial 4.0 International License

# 1 **P-GRe : an efficient pipeline to maximised pseudogene prediction in** 2 **plants/eucaryotes**

3  
4

5 Sébastien Cabanac<sup>1</sup>, Christophe Dunand<sup>1</sup> and Catherine Mathé<sup>1</sup>

6 <sup>1</sup> Laboratoire de Recherche en Sciences Végétales, Université de Toulouse, CNRS, Université Paul  
7 Sabatier Toulouse 3, Toulouse INP, Auzeville-Tolosane, France

8 Corresponding authors.

9 E-mail addresses: [catherine.mathe-dehais@univ-tlse3.fr](mailto:catherine.mathe-dehais@univ-tlse3.fr) (C. Mathé) and [christophe.dunand@univ-tlse3.fr](mailto:christophe.dunand@univ-tlse3.fr)  
10 [christophe.dunand@univ-tlse3.fr](mailto:christophe.dunand@univ-tlse3.fr) (C. Dunand)

11  
12 **ABSTRACT**

13 Formerly considered as part of "junk DNA", pseudogenes are nowadays known for their role in the  
14 post-transcriptional regulation of functional genes. In addition, their identification allows a better  
15 understanding of gene evolution in the frame of multigenic families. Despite this, there is, to our  
16 knowledge, no fully automatic user-friendly software allowing the annotation of pseudogenes on a  
17 whole genome. Here, we present Pseudo-Gene Retriever (P-GRe), a fully automated pseudogene  
18 prediction software requiring only a genome sequence and its corresponding GFF annotation file.  
19 P-GRe detects the sequences of the pseudogenes on a whole genome and returns to the user all their  
20 genomic sequences and their pseudo-coding sequences. The ability of P-GRe to finely reconstruct the  
21 structure of pseudogenes also allow to obtain a set of proteins virtually encoded by the predicted  
22 pseudogenes. We show here that in 70% of the cases, virtual proteins constructed by P-GRe from  
23 *Arabidopsis thaliana* proteome and genome aligned better to their parent protein than their annotated  
24 counterpart.

25 **INTRODUCTION**

26 Pseudogenes are genomic sequences with homology to functional genes but that harbor deleterious  
27 mutations, such as loss of the start codon, loss of coding sequence part, gain of stop or frame-shifts.  
28 No longer coding for a functional protein, pseudogenes are rarely transcribed and are often described  
29 as having no function. It is known that part of the pseudogenes are transcribed, this part representing  
30 for example 15% of all the pseudogenes in mice (Sisu *et al.*, 2020). It has also been shown that these  
31 transcripts, originating from pseudogenes, could form duplexes with the mRNAs of homologous  
32 functional genes and thus participate in post-transcriptional regulation through the RNAi pathway  
33 (Tam *et al.*, 2008; Watanabe *et al.*, 2008; Guo *et al.*, 2009). Moreover, the exhaustive prediction of  
34 pseudogenes could allow a better understanding of dynamic of gene evolution in multigenic families  
35 often subjected to duplication and pseudogenization events.

36 Most of the current pseudogene prediction software relies on the homology between the known  
37 protein sequences of an organism and the sequences of the pseudogenes to predict the positions of  
38 these pseudogenes on the genome. Briefly, one or more local alignments of each protein sequence are  
39 performed in order to find an approximate position of the pseudogenes. The protein sequence with  
40 the highest similarity with the sequence found for each local alignment is defined as being encoded by  
41 the parent gene, as it is assumed that many pseudogenes are derived from the duplication of functional  
42 genes. A finer alignment is then carried out between the hits obtained at the end of the first local  
43 alignment and the associated parent sequences. Shiu's pipeline (Zou *et al.*, 2009) and PseudoPipe  
44 (Zhang *et al.*, 2006), two software commonly used for pseudogene prediction, compares protein

45 sequences to DNA sequence database in both DNA orientations using the tfasty local alignment  
46 algorithm (Pearson & Lipman, 1988).

47 If a large number of pseudogene prediction software exist, many are specific to a type of organism,  
48 such as Pseudofinder (Syberg-Olsen *et al.*, 2022) and Psi-Phi (Lerat & Ochman, 2004) which are  
49 dedicated to prokaryotes. Others target specific types of pseudogenes, such as PPFINDER (van Baren  
50 & Brent, 2006) and PΨFinder (Abrahamsson *et al.*, 2022) which are made to predict pseudogenes  
51 originating from transcript retrotransposition events. Software able to work on any type of organisms  
52 and pseudogenes are rarer, and often produce very different outputs: some can approximate a protein  
53 or a peptide virtually encoded by the pseudogenes, while others simply return the sequence of the  
54 pseudogenes without worrying about their pseudo-coding structure. Most generate their own results  
55 file in TSV format and surprisingly, to our knowledge, none of the software returns its results in the  
56 GFF or GTF reference formats mainly used for structural genome annotation. Last, these software are  
57 often not very user-friendly, requiring a tedious preparation of the data and the organization of the  
58 working directories according to precise instructions.

59 Pseudo-Gene Retriever (P-GRe) allows the prediction of the positions of pseudogenes on a whole  
60 genome as well as the precise reconstruction of their structures in pseudo-coding exon. P-GRe only  
61 requires the genome sequence and the corresponding structural annotation in GFF format. It is  
62 therefore optional to provide all the protein sequences, as P-GRe can generate them from genome  
63 and annotation files. P-GRe returns the set of genomic sequences of the pseudogenes, and also the set  
64 of pseudo-coding sequences. The positions of pseudogenes and their features on the genome are  
65 returned in GFF format, including pseudo-exons, start and stop codons as well as frame-shifts. Precise  
66 reconstruction of the pseudogenes' structure and detection of frame-shift events also enable P-GRe  
67 to return the list of proteins and peptides virtually encoded by the predicted pseudogenes. The P-GRe  
68 process integrates a set of new approaches, including the detection of frame-shifts by "chimera  
69 generation", as described in the Materials & Methods section, and the replacement of the local  
70 alignments, during the step of fine alignment, by global alignments. These alignments are corrected  
71 and refined by an algorithm inspired by Lindley's process (Lindley, 1952) and by a search for canonical  
72 splice sites, both described in the Materials & Methods section.

73

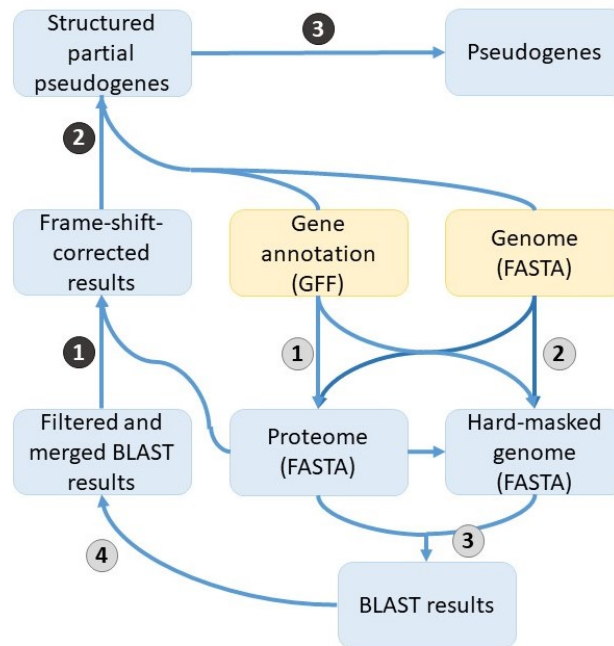
## 74 MATERIALS & METHODS

### 75 P-GRe

76 Like most pseudogene prediction software, P-GRe proceeds in two steps, the first consisting in finding  
77 the approximate position of the pseudogenes and the second in refining their structure (Figure 1). For  
78 the first step, P-GRe uses the GFFread software (Pertea & Pertea, 2020) to generate the set of proteins  
79 from the genome and the annotation file provided by the user. These same files are used for hard  
80 masking annotated gene sequences on the genome using maskfasta from the BEDTools suite (Quinlan  
81 & Hall, 2010). The generated proteins are then locally aligned to the masked genome using the tblastn  
82 algorithm (Altschul *et al.*, 1990), with the options -seg 'yes', -soft\_masking True, -db\_soft\_mask dust, -  
83 outfmt 6, -evalue 0.01, -word\_size 3, -gapextend 2 and -max\_intron\_length 10000. The local  
84 alignments obtained by tblastn are then selected on their percentage of identity. The latter must be  
85 greater than a threshold to be selected, this threshold being dependent on the alignment length and  
86 calculated by the following formula:

$$87 \quad T = \max\{75 - (l - 20) * 0.6875, 20\}$$

88 where  $T$  is the identity threshold and  $l$  is the length of the alignment. If several local alignments have  
89 been obtained with the same protein and are separated by less than 10 kb, all these local alignments  
90 are selected. In the end, all non-selected results are filtered out. Remaining overlapping alignments  
91 obtained from different proteins are filtered in order to keep only the longest alignment.

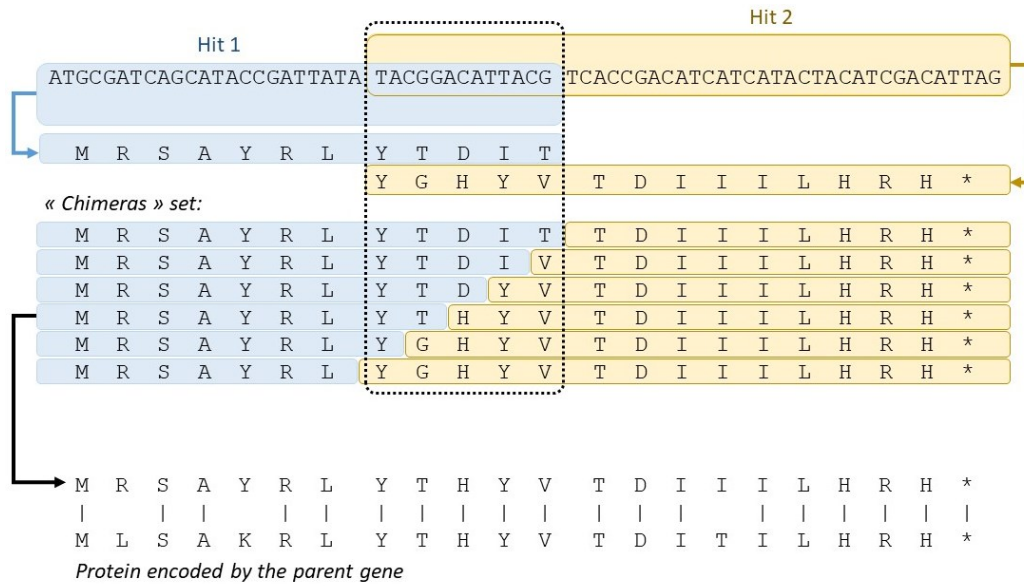


92

93 Figure 1. Diagram representing the operation of P-GRe. Mandatory inputs are colored in yellow. The  
94 first step of P-GRe, *i.e.* the detection of pseudogenes, is made up of four sub-steps numbered in gray  
95 circles: 1. Generation of the proteome, 2. Hard-masking of genes on the genome. 3. Local alignments  
96 of proteins on the masked genome and 4. Filtering and merging of alignment results. The second step  
97 of P-GRe, *i.e.* the reconstruction of the pseudogenes, consists of three sub-steps numbered on black  
98 circles: 1. Detection of the frame-shifts, 2. Determination of the structure of the pseudogenes and 3.  
99 Reconstruction of the terminal parts.

100 The second step, *i.e.* the prediction of the structure of each pseudogene, is itself divided into four  
101 stages.

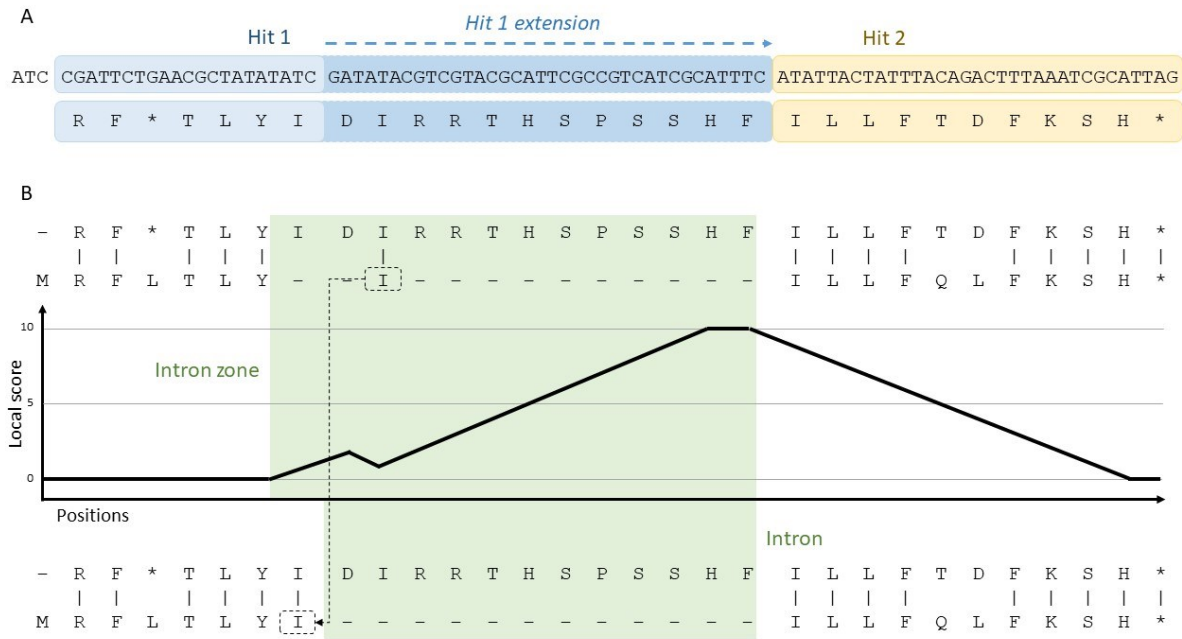
102 First, frame-shifts are identified from overlapping local alignments obtained by the same protein,  
103 where the length of the overlap is not divisible by three. The position of the frame-shift is precisely  
104 determined after translation in the different reading frames of the local alignments' area (Figure 2). All  
105 of the peptides obtained, called chimeras, are aligned locally with blastp algorithm with the protein  
106 sequence coded by the parent gene of the pseudogene. The chimera that aligns best, *i.e.* that achieves  
107 the lowest E-value, is used to determine the correct frame and the position of the frame-shift. For this  
108 step, the following parameters are set in the blastp algorithm: -evalue 0.01, -word\_size 3, -gapextend  
109 2, -max\_target\_seqs 1.



110

111 Figure 2. Example of frame-shift detection via "chimera construction". When two local alignments  
 112 originating from the same protein overlap, the peptides encoded by the sequences corresponding to  
 113 these alignments are generated. The two peptides are concatenated, with the amino acids encoded  
 114 by the overlapping positions of the N-ter side peptide (in blue) being gradually replaced by the  
 115 overlapping amino acid of the C-ter side peptide (in orange). All the chimeras thus generated are  
 116 aligned locally on the protein encoded by the parent gene of the pseudogene. The chimera returning  
 117 the best alignment (the one with the lowest E-value) is then used to determine the position of the  
 118 frame shift.

119 Secondly, because local alignments obtained in the first step generally do not cover the entire pseudo-  
 120 coding exon, a few bases at the ends or the beginning of the pseudo-coding exon may be missing. To  
 121 overcome this problem, P-GRe extends the positions of the alignments until the next alignment (Figure  
 122 3A). The expanded regions are then translated into a single amino acid sequence that is globally aligned  
 123 with the parent protein sequence. The extended gaps thus generally correspond to the pseudo-introns  
 124 on the pseudogene. Certain misalignments can falsify this structure, and are corrected by a process  
 125 inspired by the Lindley process, which is therefore called hereafter "*a priori* Lindley-inspired process"  
 126 (Figure 3B).



127

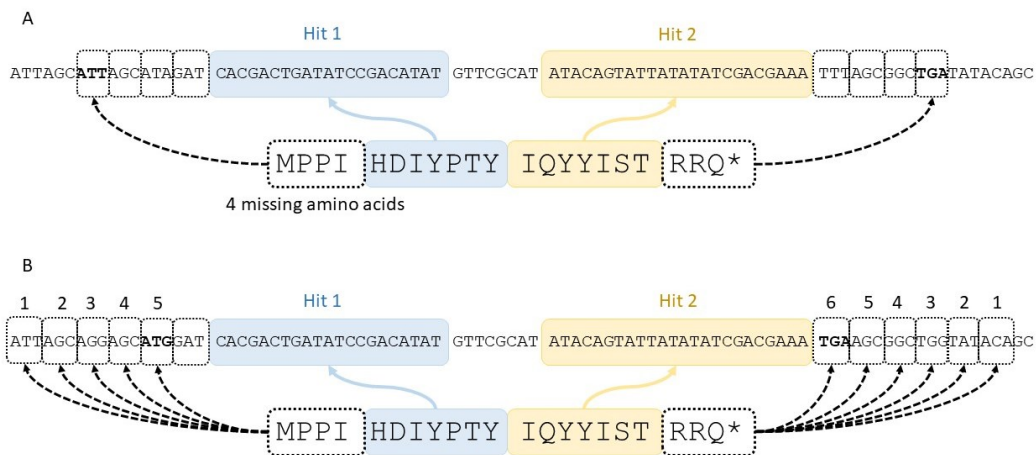
128 Figure 3. Example of alignment correction performed by P-GRe with a Lindley-inspired process. A. Each  
 129 local alignment obtained from a protein (hit) is extended to the next hit. B. The resulting peptide is  
 130 aligned with the protein encoded by the parent gene of the pseudogene. A process inspired by the  
 131 Lindley process assigns each position a local score, between 0 and 10, equal to the local score of the  
 132 previous position, +1 if a gap is present on the aligned parent sequence and -1 otherwise. Zones of  
 133 introns are identified by a sequence of positions with a non-zero score, with the last position having  
 134 the maximum score. Finally, the positions within the regions of introns to which amino acids are  
 135 aligned are realigned, the amino acids of the protein encoded by the parent gene being reassembled  
 136 with the amino acids encoded by the same coding sequence.

137 Briefly, the alignment between the supposed amino acid sequence and the parent protein sequence is  
 138 scanned from left to right. Each amino acid coded by the pseudogene is assigned a local score equal to  
 139 the score of the preceding amino acid, +1 if the latter is aligned with a gap, -1 otherwise. The local  
 140 score of an amino acid can never be negative and never exceed 10. Along the alignment, local score  
 141 "spikes" can thus be obtained and are used to define zones of introns, *i.e.* the longest series of positions  
 142 with non-zero scores whose first position has a score equal to 1 and last position has the highest score  
 143 of the series, this score being greater than or equal to 5. Furthermore, thanks to the information  
 144 present in the GFF file, the *a priori* Lindley-like process is activated at the approximate positions where  
 145 an intron is expected, *i.e.* five amino acids before each pseudo-coding exon change. The local score  
 146 can only vary when the process is activated, and the process is deactivated when the local score  
 147 reaches 0 after defining an intron zone, or if no gaps are found within 10 positions after an activation.  
 148 For this step, the alignments are carried out using the pairwise2 module of the BioPython tool suite if  
 149 the two sequences to be aligned are less than 300 amino acids (Cock *et al.*, 2009), and using stretcher  
 150 from the EMBOSS tool suite (Rice *et al.*, 2000) otherwise to be faster. As large gaps are expected in low  
 151 numbers, corresponding to introns, the alignments are performed with a strong gap opening penalty  
 152 (5) and no extension penalty. The BLOSUM62 matrix is used for substitution scoring. Amino acids  
 153 encoded by the pseudogene aligning with those of the parent protein within an intron zone are then  
 154 considered as misalignments.

155 Thanks to the information present in the GFF file, P-GRe can then reattach a misaligned amino acid  
 156 from the parent protein sequence to the amino acids encoded by the same coding exon, correcting

157 the alignment. Finally, to refine the structure of the pseudogenes, a search for canonical GT/AT splicing  
 158 sites is carried out at more or less 9 bases from each start and end of introns.

159 The third stage consists in reconstructing the N-terminal and C-terminal parts of the proteins virtually  
 160 coded by the pseudogenes. Thanks to the local alignments carried out during the search for the  
 161 pseudogenes, P-GRe determines for each pseudogene the positions at which a start codon and a stop  
 162 codon should be found (Figure 4A). Concerning the start codon, P-GRe will search for an ATG codon or  
 163 a degenerate ATG codon, *i.e.* with one substitution allowed. If no start codon is found at the expected  
 164 position, P-GRe will search for ATG codons from the expected position towards the start of the first  
 165 pseudo-coding exon. P-GRe also searches for a stop codon from the expected position towards the  
 166 end of the last pseudo-exon (Figure 4B).

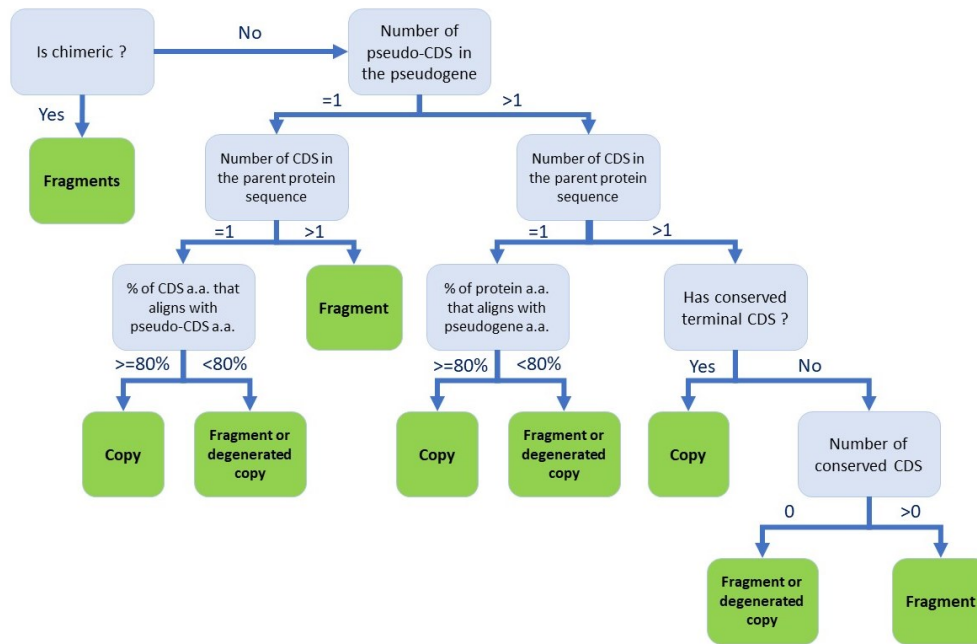


167  
 168 Figure 4. Examples of reconstructions of pseudogene ends by searching for start and stop codons. A.  
 169 Case where a degenerated start codon is found (ATG codon with a substitution) at the expected  
 170 position, and where a stop codon is also found at the expected position. B. If no start or degenerate  
 171 start codon is found, P-GRe will look for a start codon a bit further and gradually look for one, base by  
 172 base, towards the start of the first exon. For the stop codon, this search is done base by base towards  
 173 the end of the last exon.

174 Fourth, once the set of pseudogenes has been predicted, the pseudogenes separated by less than  
 175 2.5 kb and with no terminal stop codon are merged. This step makes it possible in particular to  
 176 reconstitute the pseudogenes whose different exons initially matched to different parent proteins  
 177 highly similar. This also allows to consider the rare cases of chimeric pseudogenes, *i.e.* pseudogenes  
 178 consisting of fragments of sequences originating from different genes.

### 179 Pseudogenes categorization

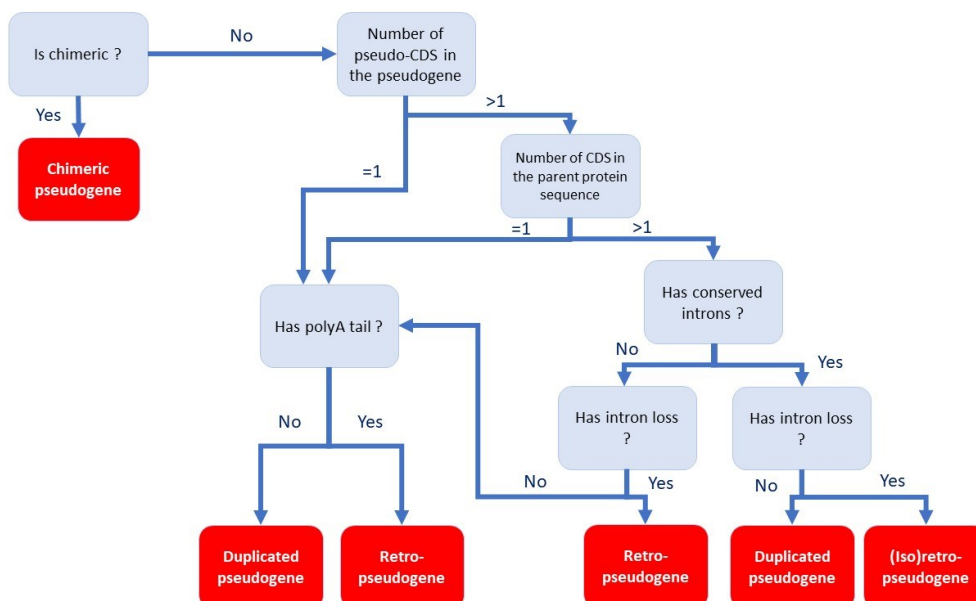
180 The calculations carried out by P-GRe during the different steps are used to categorize the different  
 181 predicted pseudogenes. Pseudogenes are categorized according to two characteristics: their  
 182 completeness and their type. Completeness can be noted as "Copy" if the pseudogene appears to be  
 183 a full copy of the parent gene or "Fragment" if it appears to be only a copy of part of it. This  
 184 characteristic is determined using a decision tree (Figure 5).



185

186 Figure 5. Decision tree to categorize the completeness of pseudogenes predicted by P-GRe. Conserved  
 187 coding sequence (CDS) are defined as pseudogene (pseudo-)CDS encoding an amino acid sequence  
 188 producing an alignment containing less than 20% gap with the CDS of a parent protein. These  
 189 percentages are calculated at the time of alignment corrected by the Lindley-inspired process.

190 The pseudogene type is determined from intron losses or retentions and the presence or absence of a  
 191 poly(A) tail. It may be noted "Duplicated pseudogene" if the pseudogene appears to be a duplication  
 192 of the parent gene, "Retro-pseudogene" if it appears to be the result of a reverse transcription of the  
 193 parent gene, or "(Iso)retro-pseudogene" if it appears to be derived of a reverse transcription of an  
 194 isoform of the parent gene. This characteristic is also determined using a decision tree (Figure 6). Note  
 195 that for this step, a poly(A) tail search can be carried out by P-GRe in an area of (maximum) 500 bases  
 196 following the last position on the 3' side of a pseudogene. A poly(A) tail is considered present if a series  
 197 of 20 bases containing at least 15 A is found in that area. Conversely, it is considered absent if no series  
 198 of 20 bases with more than 7 A is found in this area.



199



200 Figure 6. Decision tree to categorize the type of pseudogenes predicted by P-GRe. Intron loss or  
201 retention is analyzed at the time of alignments corrected by the Lindley-inspired process.  
202 Furthermore, losses or retention of introns are only sought between two conserved (pseudo-)CDSs,  
203 *i.e.*, CDSs of pseudogenes which produce an alignment with a CDS of a parent protein containing less  
204 than 20% gap.

205

## 206 **Quality of predictions and sensibility**

207 The sensitivity of P-GRe was tested using annotation, proteome and genome from *A. thaliana* available  
208 on the Ensembl database (TAIR10). Among the 924 known *A. thaliana* pseudogenes, those that  
209 overlapped pseudogenes predicted by P-GRe over at least 60% of their length were considered found.  
210 Note that this sensitivity metric and the data used for prediction are similar to those used by Xiao *et al.*  
211 (Xiao *et al.*, 2016), which achieved 78.9% and 81.3% sensitivity with Shiu's pipeline and PseudoPipe  
212 respectively. A similar method was used to identify the predicted pseudogenes that corresponded to  
213 transposable elements (TEs, annotated "transposable\_element\_gene" or "transposable\_element" on  
214 the TAIR10 annotation). In addition, the sequences corresponding to the elements annotated as  
215 "transposable\_element\_gene" were translated and the sequences virtually encoded by the  
216 pseudogenes predicted by P-GRe were locally aligned with these TE sequences using the blastp  
217 algorithm.

218 In addition to sensitivity, the quality of the predictions was also evaluated. All pseudogenes found were  
219 aligned locally against the entire proteome of *A. thaliana* using the blastp algorithm, keeping for each  
220 pseudogene only the best alignment. From these local alignments, percentage of identity, hit coverage,  
221 and total coverage (from the start of the first hit to the end of the last hit) were retained. Each protein  
222 sequence that lead to a best alignment was then semi-globally aligned with their matching pseudogene  
223 protein sequence using the pairwise2 module, with gap opening and gap extension penalties of 0.5  
224 and 0.1, respectively, and with the BLOSUM62 substitution matrix. From these semi-global alignments,  
225 the alignment scores were retained. In the rare cases where a pseudogene predicted by P-GRe  
226 overlapped several annotated pseudogenes, the score for each of the TAIR10-annotated pseudogenes  
227 was set to the sum of their scores. The same kind of semi-global alignment was used to compare  
228 pseudogenes predicted by P-GRe that overlapped known pseudogenes with those that did not overlap  
229 and were not associated with TEs.

230 Note that the prediction quality measurement, described in the previous paragraph, uses both local  
231 and semi-global alignments because it is common for pseudogenes to be a ends-truncated copy of a  
232 parent gene, or to be highly degenerated, where only a few short preserved and scattered regions  
233 remain.

## 234 **RESULTS & DISCUSSION**

### 235 **Comparison with known pseudogenes**

236 779 (84.31%) of the 924 annotated pseudogenes in *A. thaliana* were predicted by P-GRe, for a total of  
237 8459 predicted pseudogenes. The quality of the predictions of the 779 pseudogenes predicted by  
238 P-GRe was measured and compared to the corresponding annotated pseudogenes. Overall, P-GRe  
239 pseudogenes achieved better alignments with their parent proteins than annotated pseudogenes,  
240 which is particularly marked by an average alignment score higher by 159.31 points. In concrete terms,  
241 over 70% of the pseudogenes predicted by P-GRe match one of the *A. thaliana* proteins with an  
242 alignment score higher than that of their known counterparts annotated in TAIR10 (Table 1).

	Identity	Average hit coverage	Average total coverage	Average semi-global alignment score
P-GRe	0.64	0.56	0.81	1018.48
TAIR10	0.65	0.42	0.71	859.17

243 Table 1. Measurements of the quality of predictions made by P-GRe and comparison with TAIR10  
 244 annotated pseudogenes. The quality of the amino acid sequences virtually encoded by the set of  
 245 pseudogenes common to P-GRe and the TAIR10 annotation was measured by aligning each sequence  
 246 locally against the entire *A. thaliana* proteome. The identity, coverage rate of the subject protein by  
 247 hits and total coverage rate (from the first hit to the last) were measured. A semi-global alignment  
 248 between the query and subject sequences was then performed and the alignment score was recorded.

#### 249 TE-related pseudogenes

250 Of the pseudogenes predicted by P-GRe, 4505 aligned locally with the sequences of known and  
 251 annotated TEs in TAIR10. Among these TE-associated predictions, the positions of 2186 of them  
 252 overlap the positions of an annotated transposable element over at least 60% of its length. With regard  
 253 to the remaining 2319 predicted pseudogenes, it should be noted that their local alignments with the  
 254 TE genes obtain a low E-value of 7.73e-05 in average. For comparison, pseudogenes that overlap  
 255 annotated TEs produce alignments with a slightly higher average E-value of 8.77e-05 (Table 2). We  
 256 therefore assume that, rather than being false positives, a large proportion of these pseudogenes  
 257 predicted by P-GRe could correspond to still unknown TE genes or TE pseudogenes.

	Predicted pseudogenes overlapping known TE	TE-similar predicted pseudogenes not overlapping any known TE
Hits average identity	0.62	0.63
Hits average length	49.44	53.91
Hits average E-value	8.77e-05	7.73e-05

258 Table 2. Comparison of local alignments used to generate pseudogenes similar to transposable  
 259 elements (TEs). Predicted pseudogenes were considered to overlap a known TE if their positions  
 260 overlapped at least 60% of a known TE. The right column contains data for predicted pseudogenes that  
 261 did not overlap a known TE but still produced local alignments with annotated TE sequences.

#### 262 Non-TE-related pseudogenes

263 Among the pseudogenes predicted by P-GRe, 2913 do not cover any genes from TAIR10 annotation,  
 264 or any transposable elements or pseudogenes, nor do these pseudogenes produce any significant  
 265 alignment with known transposable elements in *A. thaliana*. These pseudogenes were constructed by  
 266 P-GRe with hits having an average E-value of 4.94e-04 and an average identity of 64%, which is similar  
 267 to those corresponding to annotated pseudogenes, as they were constructed from hits with an average  
 268 E-value of 1.05e-04 and an average identity of 63%. The only notable difference between these two  
 269 situations is the average length of the hits that enabled them to be generated. Indeed, pseudogenes  
 270 predicted by P-GRe and not corresponding to any known pseudogene were generated from hits with  
 271 an average length of 47.77 amino acids, whereas those corresponding to known pseudogenes were  
 272 generated from hits with an average length of 78.65 amino acids. As a consequence, the average sizes  
 273 of the complete pseudo-proteins are 100 amino acids and 454 amino acids respectively for previously  
 274 unknown and known pseudogenes. Nevertheless, when aligning these pseudo-proteins to their  
 275 respective parent proteins, an average relative alignment score of 3.76 per amino acid is reached, while  
 276 it is 3.70 per amino acid for the annotated ones (Table 3). Overall, this suggests that P-GRe may be able  
 277 to detect short pseudogene fragments that were missed in current annotation.

	Predicted overlapping pseudogenes	pseudogenes known	Predicted overlapping pseudogenes or TEs	pseudogenes not known
Hits average identity	0.63		0.64	
Hits average length	78.65		47.77	
Hits average E-value	1.05e-04		4.94e-04	
Average length	454		100	
Average relative alignment score	3.70		3.76	

278 Table 3. Comparison between pseudogenes predicted by P-GRe corresponding to annotated  
 279 pseudogenes and the other non-TE-related pseudogenes. The first three lines correspond to average  
 280 values concerning the local alignments used for generating the pseudogenes. The average length is  
 281 that of the peptides and proteins virtually encoded by the two categories of pseudogenes, and the  
 282 average relative alignment score is the average score of the pseudogenes obtained by dividing the  
 283 score of the semi-global alignment between a pseudogene and its parent gene by the length of the  
 284 pseudogene.

### 285 Pseudogene type and sequence conservation

286

287 The classification of pseudogenes carried out by P-GRe was compared between the three categories  
 288 of pseudogenes, *i.e.* predicted pseudogenes which overlap known pseudogenes, predicted  
 289 pseudogenes which overlap known TEs, and predicted pseudogenes which overlap neither and whose  
 290 sequences do not align with the sequences of the known TEs. Striking differences were found between  
 291 these categories (Table 4). For pseudogenes overlapping known pseudogenes, 19% were noted as full  
 292 copies, compared to 9% and 7%, respectively, for the other two pseudogene categories. 13% of these  
 293 pseudogenes were also noted as coming from duplication event, compared to 5% and 4% for the other  
 294 pseudogene categories. These data seem to demonstrate that the annotations of the pseudogenes of  
 295 *A. thaliana* are mainly focused on pseudogenes with relatively conserved sequences, close to their  
 296 parent gene sequences, which could explain the large number of new predicted pseudogenes. In  
 297 contrast, predicted pseudogenes that were not annotated appeared more heavily degraded, with 85%  
 298 of them classified as fragments or degraded copies, compared to 51% and 40% for pseudogenes  
 299 overlapping annotated pseudogenes or annotated TEs, respectively. In addition, the type of 80% of  
 300 them could not be determined, compared to 48% and 37% for the pseudogenes overlapping known  
 301 TEs and those not overlapping any known TEs or annotated pseudogenes, respectively.

302 Interestingly, the proportion of pseudogene noted as chimeric was much higher in predicted  
 303 pseudogenes overlapping TEs (50%), than in pseudogenes overlapping annotated pseudogenes (29%)  
 304 or in non-annotated pseudogenes (8%). This result was expected, as TEs have a role in the formation  
 305 of this type of pseudogene. However, it should be noted that the number of pseudogenes noted as  
 306 chimeric may be exaggerated due to pseudogenes having diverged so much from their parent that  
 307 they have similarities with several different genes. Chimeric pseudogenes aside, we also found that  
 308 10% of retropseudogenes originated from a retrotransposition event on an alternative mRNA (iso-  
 309 retrotransposon).

	Predicted overlapping pseudogenes	pseudogenes known	Predicted overlapping pseudogenes or TEs	pseudogenes not known
Completeness				

Copy	19%	9%	7%
Fragment	35%	23%	56%
Chimeric fragments	29%	50%	8%
Fragment or highly degenerated copy	16%	17%	29%
Type			
Chimeric pseudogene	29%	50%	8%
Duplicated pseudogene	13%	5%	4%
(Iso)retropseudogene	3%	1%	0%
Retropseudogene	7%	6%	7%
Undetermined	48%	37%	80%

310 Table 4. Type and completeness of the different categories of predicted pseudogenes.

311

## 312 Acknowledgements

313 The authors are thankful to the Paul Sabatier-Toulouse 3 University and to the Centre National de la  
314 Recherche Scientifique (CNRS) for granting their work. SC is the recipient of a fellowship from the  
315 “École Universitaire de Recherche (EUR)” TULIP-GS (ANR-18-EURE-0019). This study is set within the  
316 framework of the “Laboratoires d’Excellences (LABEX)” TULIP (ANR-10-LABX-41).

317

## 318 References

- 319 **Abrahamsson S, Eiengård F, Rohlin A, Dávila López M. 2022.** PΨFinder: a practical tool for the  
320 identification and visualization of novel pseudogenes in DNA sequencing data. *BMC bioinformatics*  
321 **23**: 59.
- 322 **Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ. 1990.** Basic local alignment search tool. *Journal*  
323 *of Molecular Biology* **215**: 403–410.
- 324 **van Baren MJ, Brent MR. 2006.** Iterative gene prediction and pseudogene removal improves genome  
325 annotation. *Genome Research* **16**: 678–685.
- 326 **Cock PJA, Antao T, Chang JT, Chapman BA, Cox CJ, Dalke A, Friedberg I, Hamelryck T, Kauff F,**  
327 **Wilczynski B, et al. 2009.** Biopython: freely available Python tools for computational molecular  
328 biology and bioinformatics. *Bioinformatics (Oxford, England)* **25**: 1422–1423.
- 329 **Guo X, Zhang Z, Gerstein MB, Zheng D. 2009.** Small RNAs Originated from Pseudogenes: cis- or trans-  
330 Acting? *PLOS Computational Biology* **5**: e1000449.
- 331 **Lerat E, Ochman H. 2004.** Psi-Phi: exploring the outer limits of bacterial pseudogenes. *Genome*  
332 *Research* **14**: 2273–2278.
- 333 **Lindley DV. 1952.** The theory of queues with a single server. *Mathematical Proceedings of the*  
334 *Cambridge Philosophical Society* **48**: 277–289.
- 335 **Pearson WR, Lipman DJ. 1988.** Improved tools for biological sequence comparison. - PMC.
- 336 **Pertea G, Pertea M. 2020.** GFF Utilities: GffRead and GffCompare. *F1000Research* **9**: ISCB Comm J-  
337 304.

- 338 **Quinlan AR, Hall IM. 2010.** BEDTools: a flexible suite of utilities for comparing genomic features.  
339 *Bioinformatics (Oxford, England)* **26**: 841–842.
- 340 **Rice P, Longden I, Bleasby A. 2000.** EMBOSS: the European Molecular Biology Open Software Suite.  
341 *Trends in genetics: TIG* **16**: 276–277.
- 342 **Sisu C, Muir P, Frankish A, Fiddes I, Diekhans M, Thybert D, Odom DT, Flicek P, Keane TM, Hubbard**  
343 **T, et al. 2020.** Transcriptional activity and strain-specific history of mouse pseudogenes. *Nature*  
344 *Communications* **11**: 3695.
- 345 **Syberg-Olsen MJ, Garber AI, Keeling PJ, McCutcheon JP, Husnik F. 2022.** Pseudofinder: Detection of  
346 Pseudogenes in Prokaryotic Genomes. *Molecular Biology and Evolution* **39**: msac153.
- 347 **Tam OH, Aravin AA, Stein P, Girard A, Murchison EP, Cheloufi S, Hodges E, Anger M,**  
348 **Sachidanandam R, Schultz RM, et al. 2008.** Pseudogene-derived small interfering RNAs regulate  
349 gene expression in mouse oocytes. *Nature* **453**: 534–538.
- 350 **Watanabe T, Totoki Y, Toyoda A, Kaneda M, Kuramochi-Miyagawa S, Obata Y, Chiba H, Kohara Y,**  
351 **Kono T, Nakano T, et al. 2008.** Endogenous siRNAs from naturally formed dsRNAs regulate  
352 transcripts in mouse oocytes. *Nature* **453**: 539–543.
- 353 **Xiao J, Sekhwal MK, Li P, Ragupathy R, Cloutier S, Wang X, You FM. 2016.** Pseudogenes and Their  
354 Genome-Wide Prediction in Plants. *International Journal of Molecular Sciences* **17**: 1991.
- 355 **Zhang Z, Carriero N, Zheng D, Karro J, Harrison PM, Gerstein M. 2006.** PseudoPipe: an automated  
356 pseudogene identification pipeline. *Bioinformatics (Oxford, England)* **22**: 1437–1439.
- 357 **Zou C, Lehti-Shiu MD, Thibaud-Nissen F, Prakash T, Buell CR, Shiu S-H. 2009.** Evolutionary and  
358 Expression Signatures of Pseudogenes in Arabidopsis and Rice. *Plant Physiology* **151**: 3–15.
- 359