

RESEARCH ARTICLE

Open Access

# Newly evolved introns in human retrogenes provide novel insights into their evolutionary roles

Li-Fang Kang<sup>1</sup>, Zheng-Lin Zhu<sup>1\*</sup>, Qian Zhao<sup>1</sup>, Li-Yong Chen<sup>2\*</sup> and Ze Zhang<sup>1</sup>

## Abstract

**Background:** Retrogenes generally do not contain introns. However, in some instances, retrogenes may recruit internal exonic sequences as introns, which is known as intronization. A retrogene that undergoes intronization is a good model with which to investigate the origin of introns. Nevertheless, previously, only two cases in vertebrates have been reported.

**Results:** In this study, we systematically screened the human (*Homo sapiens*) genome for retrogenes that evolved introns and analyzed their patterns in structure, expression and origin. In total, we identified nine intron-containing retrogenes. Alignment of pairs of retrogenes and their parents indicated that, in addition to intronization (five cases), retrogenes also may have gained introns by insertion of external sequences into the genes (one case) or reversal of the orientation of transcription (three cases). Interestingly, many intronizations were promoted not by base substitutions but by cryptic splice sites, which were silent in the parental genes but active in the retrogenes. We also observed that the majority of introns generated by intronization did not involve frameshifts.

**Conclusions:** Intron gains in retrogenes are not as rare as previously thought. Furthermore, diverse mechanisms may lead to intron creation in retrogenes. The activation of cryptic splice sites in the intronization of retrogenes may be triggered by the change of gene structure after retroposition. A high percentage of non-frameshift introns in retrogenes may be because non-frameshift introns do not dramatically affect host proteins. Introns generated by intronization in human retrogenes are generally young, which is consistent with previous findings for *Caenorhabditis elegans*. Our results provide novel insights into the evolutionary role of introns.

## Background

Retroposition, or RNA-based duplication, is the process by which reverse-transcribed mRNAs are inserted into new genomic positions, which generates retrocopies [1]. Retrocopies are assumed not to carry the regulatory regions, but by chance they may obtain functions by recruiting new regulatory elements, and then become functional retrogenes [2-7]. These newly evolved genes may acquire introns in the untranslated regions by capture of nearby exons into a new genomic environment or fusion with host genes, which is chimerization based on intron gain [3-8]. Such retrogenes are usually

considered to be intronless because introns were not inherited from the parents. However, in some circumstances, retrogenes may recruit internal exonic sequences as introns [9,10], which is known as intronization [11].

Since intronization of retrogenes was first reported [9], this kind of evolutionary event has been commonly observed in plants. In *Arabidopsis* and *Populus*, 29 retrogenes have undergone intronization, which represent about 15.3 % of all known retrogenes [10]. In contrast, rare cases are reported in vertebrates [12,13]. Previously, only two retrogenes were found to be intronized in mammals [14]. This frequency is extremely low given the thousands of retrocopies in the human (*Homo sapiens*) genome [15-17]. How general retrogene intronization is remains unknown. In the present study, we scanned the human genome for intronized retrogenes

\* Correspondence: zhuzl@cqu.edu.cn; mzkcly@yahoo.com.cn

<sup>1</sup>College of Life Sciences, Chongqing University, Chongqing 400044, China

<sup>2</sup>Department of Anesthesiology, Research Institute of Surgery, Daping Hospital, Third Military Medical University, 10 Changjiang Zhilu, Chongqing 400042, China

53 and identified nine cases not reported previously. Our  
54 results provide new insights into the mechanism of in-  
55 tron gain and expression patterns of retrogenes.

## 56 **Methods**

### 57 **Scanning for intron gain in retrogenes**

58 The human genome data were downloaded from the  
59 UCSC Genome Browser database (release hg19) [18,19].  
60 Then, we used the approach of Zhu et al. [10] to search  
61 the data for retrocopies. First, we mapped human protein  
62 sequences onto the genome with tBLASTn [20] and used  
63 the Pseudopipe package [21] to process the raw align-  
64 ments with the default settings, including tBLASTn  
65 *e*-value cutoff (1e-10), coverage cutoff (70 %) and identity  
66 cutoff (40 %). Next, we retained candidates with more  
67 than three introns absent or only one or two introns ab-  
68 sent but with a small  $K_s$  (<2) or other RNA-based dupli-  
69 cation evidence, for example, a poly(A) track. Finally, as  
70 described previously [10], we set filters to discard possible  
71 DNA-based duplication cases. In brief, we discarded all  
72 retrocopies in which at least 50 % of the region overlapped  
73 with repeats or with flanking genes similar to the parental  
74 gene's flanking regions. We also discarded all retrocopies  
75 that aligned well with the introns of the parents. Ulti-  
76 mately, we identified 3436 retrocopies.

77 We wrote a series of PERL programs to look for  
78 intron-containing retrogenes on the basis of annotations  
79 from ENSEMBL (GRCh37) [22,23]. We identified 54  
80 candidates of intronized retrogenes for further study.

### 81 **Gene structure validation by transcription evidence**

82 We utilized the mRNA and EST annotations from the  
83 UCSC Genome Browser Database to search for transcrip-  
84 tion evidence of intron gain in retrogenes [18,19]. For  
85 each sample, we inspected the annotated intronic region  
86 to see whether there were transcripts that support its  
87 splicing. If transcripts were present, we mapped them on  
88 the human genome with BLAT [24] to check whether  
89 these transcripts uniquely correspond to the retroposed  
90 region. By this method, eight intron-containing retro-  
91 genes were validated (Additional files 1 and 2).

### 92 **$K_a$ and $K_s$ calculation**

93 We estimated the non-synonymous substitution rate ( $K_a$ ),  
94 synonymous substitution rate ( $K_s$ ) and  $K_a/K_s$  values be-  
95 tween the intronic regions of retrogenes and their parental  
96 copies, by implementing the codeml program in the  
97 PAML package following the Nei-Gojobori method  
98 [25,26] and analyzed the results with the likelihood ratio  
99 test. We did  $K_a/K_s$  estimation between the exonic regions  
100 of retrogenes and their parental copies in the same way.

### RT-PCR

101 In order to validate the structure of the retrogenes, we  
102 collected samples of 16 human tissues from Daping  
103 Hospital, Chongqing, for experiments (Additional file  
104 3). Following the manufacturer's instructions, we used  
105 TRIzol Reagent (Invitrogen, Carlsbad, CA) to isolate  
106 RNA and digested the contaminating genomic DNA  
107 with RNase-free DNase I (Promega, Madison, WI).  
108 cDNAs were synthesized with Moloney murine  
109 leukemia virus reverse transcriptase (Promega). We per-  
110 formed PCR in a 25  $\mu$ l reaction volume, and 5  $\mu$ l of the  
111 PCR products were electrophoresed on a 1.2 % agarose  
112 gel. To validate whether the smaller-sized bands repre-  
113 sented the retrogenes, we cloned and sequenced those  
114 PCR products. Ultimately, we identified two samples in  
115 which the sequences of the smaller-sized bands  
116 belonged to retrocopies and the larger bands to the par-  
117 ental genes (Additional file 4).  
118

### Peptide support for intronized retrogenes

119 To identify whether one retrogene was expressed at the  
120 protein level, we sought peptide evidence in the Pepti-  
121 deAtlas [27-29] and PRIDE [30,31] databases using the  
122 gene name. Each search result displayed experimental  
123 details including the fractionation and sequencing (by  
124 mass spectroscopy or other methods) of short peptides.  
125 Among the results, we extracted peptides that matched  
126 the protein sequence of the intronized retrogene. Given  
127 that one peptide may match many proteins, we also used  
128 BLASTp [32,33] to ensure that the peptide specifically  
129 mapped to the gene we targeted. We only retained pep-  
130 tides for which the best hit was a targeted protein.  
131

### Age estimation of the retrogenes

132 We examined the presence and absence of orthologs in  
133 the phylogenetic tree for vertebrates and used the estab-  
134 lished origination times of all human genes [34] to infer  
135 the times of origin of the retrogenes. For comparison we  
136 used the same method to estimate the time of origin of  
137 27 retrogenes that recruited introns by chimerization  
138 [8]. We mapped the results on the vertebrate phylogeny  
139 (Additional files 5, 6 and 7). The timeline and divergence  
140 time of species in the phylogeny were reconstructed  
141 based on data from the UCSC Genome Browser data-  
142 base and other sources [19,34-40].  
143

### Detection of splicing signals

144 We detected splicing signals of new introns with SROO-  
145 GLE [41]. For an intron X, if its upstream exon is Y and  
146 downstream exon is Z, we used X and Y to detect signals  
147 of the 5' splice site (SS) and X and Z for that of the  
148 branch site (BS), polypyrimidine tract (PPT), and 3' SS.  
149 We performed two detections for each intron; one was  
150 performed on the parental gene and the other was done  
151

152 for the retroposed sequence. The former and latter were  
 153 considered to represent the status before retroposition  
 154 and the current status, respectively. Finally, for each de-  
 155 tection, we recorded the percentile score for constitutive  
 156 introns, which was obtained from a data set composed  
 157 of >50,000 constitutive introns [41], because all introns  
 158 in our data set showed no evidence for alternative spli-  
 159 cing (Additional file 8).

## 160 Results

### 161 Identification of intron gain in retrogenes

162 We focused on identifying retrogenes that contain  
 163 introns and scanned the human genome using a pub-  
 164 lished pipeline [10]. We mapped all human proteins  
 165 onto the genome with tBLASTn [20] and extracted all  
 166 possible candidates of retrocopies from among the  
 167 results with PseudoPipe [21]. Then, we set filters to ex-  
 168 clude cases that did not fulfill the properties of retropo-  
 169 sition and obtained 3436 retrocopies. Finally, we  
 170 determined that 54 of the 3436 retrocopies contained  
 171 introns on the basis of gene structure annotations from  
 172 ENSEMBL [22,23].

173 We used two methods to validate the existence of retro-  
 174 regene introns. First, we collected information from the  
 175 UCSC Genome Browser database [18,19] and found  
 176 eight cases with confident transcriptional evidence (Addi-  
 177 tional file 2). Next, we performed experiments to valid-  
 178 ate the existence of the retrogene introns. Given the  
 179 high similarity in the flanking regions of new introns for  
 180 most retro-parental alignments, we designed pairs of prim-  
 181 ers whose products (Additional file 9) spanned the in-  
 182 tronic regions for both the retrogenes and their parental  
 183 genes. Theoretically, the amplified segments from the  
 184 retrogenes (without the intronic sequences) would be  
 185 smaller than those of the parental genes (with the in-  
 186 tronic sequences). By this method, we confirmed that  
 187 two retrogenes contained introns (Additional file 4), one  
 188 of which was one of the eight retrogenes mentioned  
 189 above. In total, we identified nine retrogenes that  
 T1 190 evolved introns in the retroposed regions (Table 1). Our  
 191 data did not include RNF113B and DCAF12, which were  
 192 reported in a previous study [14], because the parents of  
 193 these two retrogenes were lost after the divergence of  
 194 mammals from vertebrates, whereas our pipeline used  
 195 parental protein sequences as queries to search for retro-  
 196 copies. In addition, we discarded POM121L2 and  
 197 ARPM1, which were suggested to be intronized retro-  
 198 genes previously [8], because the alignment identities of  
 199 these retrogenes and their respective parents did not ful-  
 200 fill the criteria set in our pipeline (>40 % identity).

### 201 Mechanisms of intron gain in retrogenes

202 To clarify the intron-gain mechanisms of these retro-  
 203 genes, we produced protein and nucleotide sequence

**Table 1** Nine human retrogenes that gained introns investigated in this study

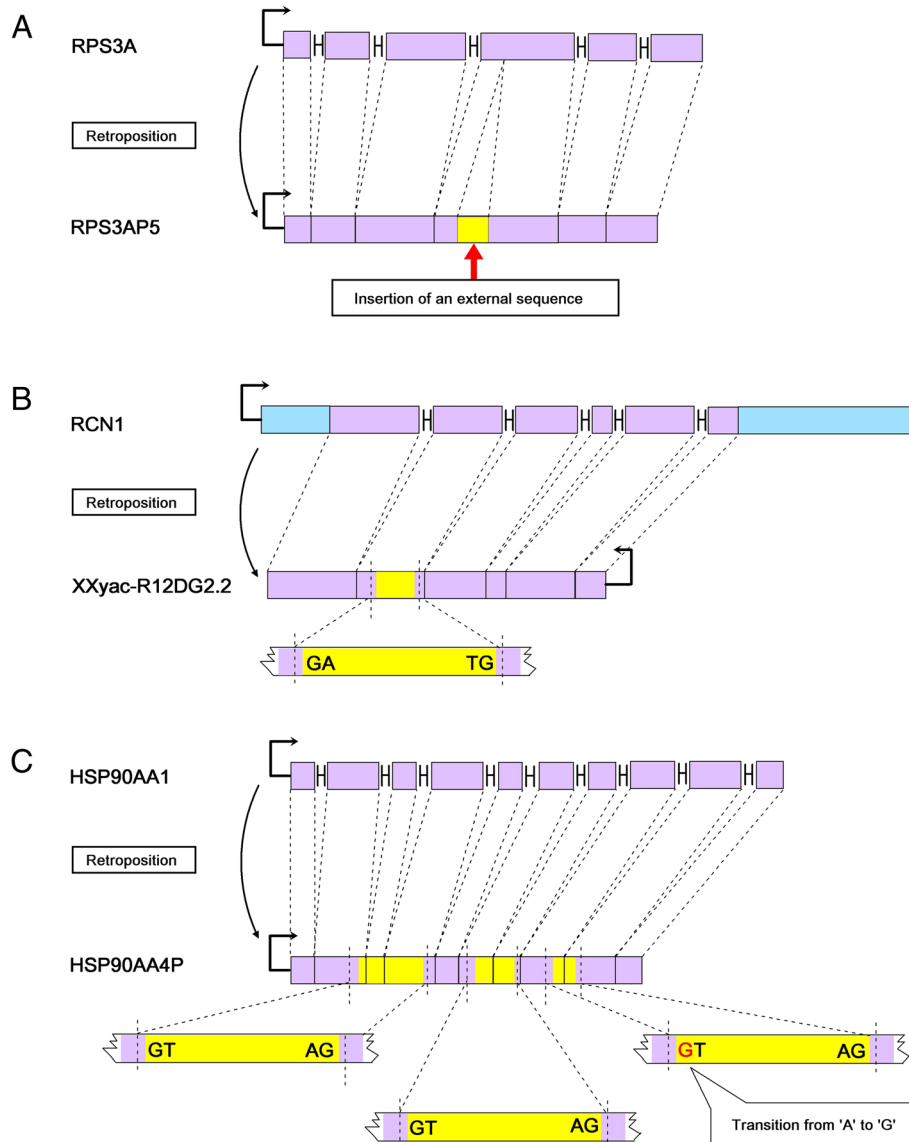
Retrogene	Parent	Movement	Intron (+)	Intron (-)	Evidence	
TMEM14D	TMEM14B	10 ← 6	1	4	A	t1.1 t1.2 t1.3 t1.4
RPS3AP5	RPS3A	10 ← 4	1	5	B	t1.5
XXyac-R12DG2.2	RCN1	13 ← 11	2*	5	B	t1.6
HSP90B2P	HSP90B1	15 ← 12	2	16	B	t1.7
HSP90AA4P	HSP90AA1	4 ← 14	3	9	A,B	t1.8
HSP90AA5P	HSP90AA1	3 ← 14	2	7	B	t1.9
CSMD3	RPL18	8 ← 19	1	5	B	t1.10
WBP2NL	SLC25A5	22 ← X	1	3	B	t1.11
AC019016.1	CSNK1A1	15 ← 5	2*	8	B	t1.12

In the column 'Movement', '10 ← 6' means a new gene on chromosome 10 is retroposed from a gene on chromosome 6, for example. 'Intron (-)' and 'Intron (+)' are the numbers of intron losses and intron gains in retrocopies, respectively. For 'Evidence', 'A', confirmed by RT-PCR; 'B', supported by convincing transcription evidence. "\*" means that the newly evolved intronic regions of XXyac-R12DG2.2 and AC019016.1 could be spliced in two patterns, respectively.

alignments for the retrogenes and their respective parental genes (Additional files 10 and 11). For RPS3AP5, we observed that its intronic region did not have counterparts in the parental gene. This result indicated that this retrogene did not gain the intron by intronization, but rather by insertion of an external sequence (Figure 1A). Using the inserted sequence as a query for a BLAT [24] search against the human genome, we identified more than five paralogous sequences with identity >95 % and coverage >70 %. The new intron may be derived from one of these paralogs. By checking the genome annotations in the UCSC Genome Browser database [18,19], we found that none of these paralogs were annotated as introns. Thus, the new intron may not have originated by 'reverse splicing', the process by which a spliced-out intronic RNA is inserted into a novel site of one RNA gene transcript by reversal of the splicing reaction [12,42]. The intron may have been created by a mechanism not reported previously.

We observed that three retrogenes (XXyac-R12DG2.2, CSMD3 and WBP2NL) were transcribed in the reverse direction relative to that of their parents. For XXyac-R12DG2.2 there are 10 annotated transcription patterns and introns appeared in four of the 10 patterns (Additional file 2). Taking ENST00000379050 as an example, the retrocopy contained a 170 bp intron, and its splicing donor and acceptor sites ('GT' and 'AG') had reverse counterparts ('AC' and 'AT') in the parental gene (Figure 1B, Additional files 10 and 11). Thus, transcription in the reverse orientation led to the origin of the intron splicing sites. For the remaining three transcription patterns (ENST00000522673, ENST00000519494 and ENST00000330825), the newly evolved intron was

F1



**Figure 1 Mechanisms of intron gain in retrogenes.** In the parental gene, rectangles represent exons, 'H'-like tags represent introns, the retroposed regions are indicated in purple, and other regions are indicated in blue. In the retrogene, the retroposed region is indicated in purple and the newly evolved intronic regions are indicated in yellow. Semi-rectangle lines with arrows indicate the direction of transcription. (A) The retrogene RPS3AP5 gained an intron by insertion of an external sequence; (B) the retrogene XXyac-R12DG2.2 evolved a new intron after transcription in the opposite orientation compared to the parent; (C) in retrogene HSP90AA4P three new introns were generated by intronization. There is no mutation at the splice sites in the two introns near the 5' terminus, whereas one transition from 'A' to 'G' (indicated in red) at the splice sites occurred in the intron near the 3' terminus.

237 shorter (127 bp) and the retroposed sequence was  
 238 located near the 3' end. In addition, the retrocopy is  
 239 inserted near the 3' end of a ncRNA gene candidate  
 240 (LOC 100190939, Additional file 12).

241 In CSMD3, the retroposed region was located at the  
 242 5' untranslated region (UTR) of the mRNA. Some part  
 243 of the retrocopy had changed into an intergenic se-  
 244 quence, and some part acted as a portion of an intron  
 245 (Additional files 2 and 12). The retrogene was located in

the first intron of WBP2NL (Additional file 12). Never-  
 246 theless, the retrocopy might be transcribed at least some  
 247 of the time, because an mRNA sequence, BC03789, sup-  
 248 ports the transcription of this retrogene (Additional file  
 249 1 and 2). We did not find evidence for protein-level ex-  
 250 pression of the three retrogenes that gained an intron  
 251 after transcription in the reverse orientation. The new  
 252 introns in these three retrogenes were annotated to be  
 253 in non-coding regions.  
 254

255 The remaining five retrogenes had gained introns  
256 through intronization, which generated 10 new introns.  
257 Taking HSP90AA4P as an example, three exonic  
258 sequences were changed into introns (Figure 1C). Eight  
259 of the 10 introns had the canonical splicing boundaries  
260 'GT-AG'. 80 % (8/10) of the introns arose in ORF and  
261 20 % (2/10) in UTRs.

262 In total, we observed three mechanisms of intron gain  
263 for these retrogenes. In addition to intronization, retro-  
264 genes may gain introns after insertion of external  
265 sequences or transcription in the opposite orientation  
266 compared to the parent (Figure 1).

#### 267 Non-frameshift introns generated by intronization had 268 greater evolutionary success

269 For the five retrogenes that underwent intronization, we  
270 examined the alignments of retrocopies and their corre-  
271 sponding parental sequences to assess whether these  
272 introns had disturbed the frame of putative translation  
273 inherited from the parental genes (Additional file 11). If  
274 one intron disturbed the frame, we termed it a frame-  
275 shift intron, otherwise it was considered to be a non-  
276 frameshift intron. The lengths of the corresponding  
277 sequences of the five retrogenes (70 %) were in multiples  
278 of three bases. We performed a manual check for each  
279 retrogene. At the location 100 bp upstream of the sec-  
280 ond intron of HSP90AA4P (from 5' to 3', HSP90AA4P-  
281 2), we observed an insertion of 23 bases. The length of  
282 HSP90AA4P-2 was 83 bp. Thus, compared with the par-  
283 ent, the intron and insertion led to an overall loss of 60  
284 bases (divisible by three) in the transcript. Similarly, for  
285 HSP90AA5P we observed an insertion of 22 bases  
286 located 1 bp upstream of the intron near the 5' terminus  
287 (HSP90AA5P-1) and a deletion of four bases located  
288 2 bp upstream of the intron near the 3' terminus  
289 (HSP90AA5P-2). The lengths of these two introns were  
290 439 and 254 bp, respectively. As in HSP90AA4P-2, both  
291 the indels and intronization shortened the coding  
292 sequences by 417 and 258 bp in HSP90AA5P-1 and  
293 HSP90AA5P-2, respectively (both numbers are divisible  
294 by three). Both were classified as non-frameshift introns.  
295 The two alternative spliced introns of AC019016.1 were  
296 annotated to be UTR-region introns according to the  
297 UCSC database [18,19] and Ensembl [22,23].

298 In total, eight of the 10 introns created by intronization  
299 were non-frameshift introns. This proportion (80 %) is sig-  
300 nificantly higher than the percentage of frameshift introns  
301 generated by chimerization based on intron-gain retro-  
302 genes (29.8 %, 16/49) ( $P$ -value = 0.017) [8]. From searches  
303 of PeptideAtlas [27-29] and PRIDE [30,31], we found that  
304 the predicted proteins of HSP90B2P, HSP90AA4P and  
305 HSP90AA5P had respective unique matching peptides  
T2 306 (Table 2), which indicated the true protein-coding activity  
307 of these transcripts. Consistent with findings for

*Caenorhabditis elegans* [11], our observations showed that 308  
non-frameshift introns had greater evolutionary success. 309

#### 310 Retrogenes underwent intronization by cryptic splicing 311 sites

312 Previous studies showed that most intronizations were  
313 caused by base substitutions at the 5' and 3' SS [10,11].  
314 However, we observed only four such cases (40 % of all  
315 cases) in our data set. By inspecting the EST annotations  
316 for the corresponding parental regions of all newly intro-  
317 nized introns, we found that none of these intronized  
318 introns was created by inheriting alternative splicing  
319 sites from the parental gene. What led to the creation of  
320 the other six retrogene introns? Since a retrogene does  
321 not contain introns compared with its parental gene, we  
322 proposed that the new introns were created by cryptic  
323 splice sites in the exonic regions of the parents. That is,  
324 cryptic splice sites were silent in the parents, but were  
325 activated in the retrogenes after retroposition and the  
326 new introns were generated. To test our hypothesis, we  
327 used SROOGLE [41] to detect the splicing signals (5'  
328 SS, 3' SS, the PPT located upstream of the 3' SS, and  
329 the BS located upstream of the PPT) of the retrogene  
330 introns and their respective corresponding regions in the  
331 parental genes. The splicing signals of introns in four of  
332 the six retrogenes were increased, except for those of  
333 TMEM14D and HSP90B2P (Figure 2, Table 3). For the  
334 latter two retrogenes, in the parental gene the corre-  
335 sponding regions of the retrogene introns had lower  
336 splicing signals compared with those of neighboring  
337 introns (Additional file 13). It is likely that these cryptic  
338 intronic regions were oppressed in the parental genomes  
339 and the oppression was released after retroposition. The  
340 splice sites of these six new introns pre-existed but were  
341 cryptic in the parental genes. After retroposition, the  
342 splice sites were activated in the novel genomic environ-  
343 ments. In addition, for the four introns that showed base  
344 substitutions at their splice sites, the splicing signals  
345 increased not only at the 5' SS and 3' SS but also at the  
346 BS and PPT (Table 3). In addition to point mutation, the  
347 change in gene structure after retroposition might also  
348 contribute to the evolution of new introns.

#### 349 Intronization tended to occur in young retrogenes

350 In *C. elegans*, intronization is reported to be a major con-  
351 tributor to intron creation and most introns generated by  
352 this mechanism are young [11]. In our data set, 66.7 % of  
353 retrogene introns (10/15) were created by intronization.  
354 This finding is consistent with previous studies [11]. We  
355 used the established origination times of all human genes  
356 to trace the time of origin of intronized retrogenes [34]  
357 and examined the presence and absence of the corre-  
358 sponding orthologs in the vertebrates phylogeny (Add-  
359 itional file 6). We found that 80 % (4/5) of the intronized

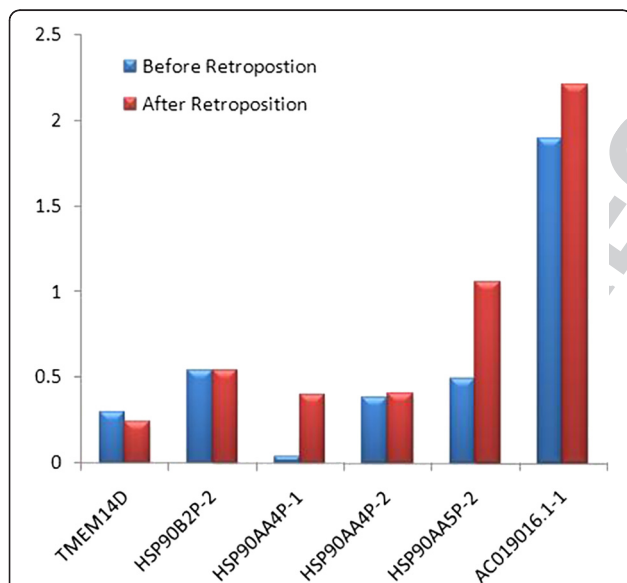
t2.1 **Table 2 Peptide support for intronized retrogenes**

t2.2	Gene name	Peptide match	Peptide database reference <sup>a</sup>	Location in protein seq	BLASTP hits <sup>b</sup>
t2.3	HSP90B2P	NLNFVKGVVDSGGLSLNVSCETLQQHK	PRIDE: 8670	86	Self (4e-19, 100 %)
t2.4		IEKAMVSQCLTESLALVASQYGWSGNMER	PRIDE: 8670	270	Self (4e-24, 100 %)
t2.5		AMVSQCLTESLALVASQYGWSGNMER	PRIDE: 8671; 8668	273	Self (7e-21, 100 %)
t2.6		MAETIQEVEDEYKAFCK	PRIDE: 8672	1	Self (9e-11, 100 %)
t2.7		CVFITDDFRDTPMK	PRIDE: 8669	72	Self (7e-08, 100 %)
t2.8	HSP90AA4P	HNNDEQYAWESSLR	PeptideAtlas: PAp00393519	93	Self (1e-07, 100 %)
t2.9		ADLNNLGTITK	PeptideAtlas: PAp01587648	20	Self (8e-04, 100 %)
t2.10		DQVANSTIVQR	PeptideAtlas: PAp00565957	207	Self (0.005, 100 %)
t2.11	HSP90AA5P	IKEIVKKHSQFIGYPITLFEVKKR	PeptideAtlas: PAp00040955; PAp00423980	33	Self (2e-17, 100 %)
t2.12		HGLEVIYMIELIDKYCVQQLK	PeptideAtlas: PAp00040711	199	Self (2e-15, 100 %)

t2.13 <sup>a</sup>, Database name and experiment numbers or identifiers. <sup>b</sup>, BLASTP search against the GenBank non-redundant protein database (e-value and maximum  
 t2.14 identity of the match are shown in parentheses [32,33]).

360 retrogenes were primate specific. We also recalculated the  
 361 ages of 27 chimerizations based on intronized retrogenes  
 362 with the same method [8] (Additional file 7) and found

that only 18.5 % of intronized retrogenes (5/27) were pri- 363  
 mate specific. This finding indicated that intronization 364  
 tended to occur in young retrogenes (proportion test, 365  
 $P = 0.023$ ). Furthermore, in our data set, no intronized retro- 366  
 gene (0/5) was retroposed from chromosome X ('out-of- 367  
 X'). The retrogenes from chromosome X were mostly old 368  
 and evolved after the divergence of eutherian mammals 369  
 (human or mouse) and marsupials (opossum) [34]. For 370  
 retrogenes that underwent intron gains by chimerization, 371  
 the proportion of 'out-of-X' retrogenes was 37 % (Addi- 372  
 tional file 7). Therefore, the comparison of 0 % and 37 % 373  
 reinforced the conclusion that intronization tended to 374  
 occur in young retrogenes. 375



**Figure 2 Comparison of splicing signals of retrogene introns (after retroposition) and their corresponding regions in the parental gene (before retroposition).** The y-axis is the sum of the percentile scores of four different signals, comprising the branch site (BS), polypyrimidine tract (PPT), and 5' and 3' splice sites (Table 3). The higher the score, the stronger the splicing signal. Scores for BS and PPT were calculated with algorithm 'K' [47] and those for the splice sites were calculated with 'M' [48] by SROOGLE [41]. Six different retrogene introns are plotted on the x-axis. If a retrogene evolved only one intron (TMEM14D), we used the gene name to represent the intron, or we marked different introns in one retrogene in the format of the gene name plus a serial number following the hyphen. For example, 'HSP90AA4P-1' represents the first intron (in the direction from 5' to 3') in the retrogene HSP90AA4.

#### Evolutionary rates of intronized retrogenes

To evaluate the evolutionary rates of retrogenes, we calculated  $K_a$ ,  $K_s$ , and  $K_a/K_s$  values between the intronic regions of retrogenes and their parental copies as well as between the exonic regions of retrogenes and their parental copies. The  $K_a$  values in the intronic regions were higher than those in the exonic regions (Mean<sub>intronic</sub> = 0.207, Mean<sub>exonic</sub> = 0.111, Wilcoxon two-sample test,  $P$ -value = 0.098; Table 4). Similarly,  $K_s$  values in the intronic regions were higher than those in the exonic regions (Mean<sub>intronic</sub> = 0.263, Mean<sub>exonic</sub> = 0.151, Wilcoxon two-sample test,  $P$ -value = 0.194). These findings are consistent with the conclusion that introns evolved faster than exons.

In addition, the exonic regions of most intronized retrogenes had  $K_a/K_s$  values smaller than 1 ( $P$ -value < 0.1), which suggested that the corresponding regions were under negative selection. By checking for evidence of expression, we found that three of the five intronized retrogenes showed evidence for expression at the protein level and the additional two retrogenes showed transcription evidence at the RNA level. This result

T4

t3.1 **Table 3 Percentile scores [40] of splicing signals of retrogene introns (after retroposition) and their corresponding**  
 t3.2 **regions in the parental gene (before retroposition)**

t3.3 t3.4 t3.5 Intron symbol	Splice sites	After retroposition				Before retroposition			
		BS (K)	PPT (K)	5' SS (M)	3' SS (M)	BS (K)	PPT (K)	5' SS (M)	3' SS (M)
t3.6 TMEM14D	GC-AG	0.14	0.06	0.02	0.02	0.14	0.04	0.01	0.11
t3.7 (HSP90B2P-1)	GT-AG	0.39	0.39	0	0.01	0.39	0.24	0	0
t3.8 HSP90B2P-2	GT-AG	0.5	0.03	0	0.01	0.5	0.02	0	0.02
t3.9 HSP90AA4P-1	GT-AG	0.21	0.03	0.04	0.12	0	0	0.04	0
t3.10 HSP90AA4P-2	GT-AG	0.03	0.06	0.04	0.28	0.03	0.06	0.04	0.25
t3.11 (HSP90AA4P-3)	GT-AG	0.56	0.33	0	0.03	0.56	0.22	0	0.02
t3.12 (HSP90AA5P-1)	TT-AG	0.25	0.27	0	0.54	0	0	0	0
t3.13 HSP90AA5P-2	GT-AG	0.45	0.45	0.01	0.15	0.12	0.15	0.01	0.21
t3.14 AC019016.1-1	GT-AG	0.91	0.35	0.11	0.84	0.61	0.35	0.11	0.82
t3.15 (AC019016.1-2)	GT-AG	0.91	0.35	0.47	0.84	0.61	0.35	0	0.82

t3.16 The higher the score, the stronger the splicing signal is. The scores for BS and PPT were calculated with the 'K' algorithms [47], and those for 5' SS and 3' SS were  
 t3.17 calculated with 'M' [48] by SROOGLE [41]. The intron symbol is in the format of the gene name plus a serial number following the hyphen. For example,  
 t3.18 'HSP90B2P-1' indicates the first intron (in the direction from 5' to 3') in HSP90B2P. If a retrogene evolved only one intron (TMEM14D), the intron is represented by  
 t3.19 the gene name. In the column 'Intron symbol', parentheses indicate that the splice sites underwent base substitution.

397 indicated that most intronized retrogenes were func-  
 398 tional and should be under negative selection.

399 With regard to the three retrogenes that gained  
 400 introns after transcription in the opposite orientation  
 401 compared with the parent, they were annotated to be in  
 402 the non-coding regions of other genes. We observed that  
 403 CSMD3 and WBP2NL evolved faster than the other re-  
 404 trogenes (Table 4). This finding is consistent with the  
 405 conclusion that non-coding regions such as UTR regions  
 406 are under less functional constraint than coding regions.  
 407 However, XXyac-R12DG2.2 evolved slowly relative to  
 408 that of CSMD3 and WBP2NL. Thus, XXyac-R12DG2.2  
 409 is likely to be under functional constraint.

## Discussion

In this study, we systematically searched the human ge-  
 411 nome for retrogenes that underwent intron gain in the  
 412 coding region and in total identified 15 retrogene  
 413 introns. These newly generated introns evolved at a fas-  
 414 ter rate than neighboring exons. In contrast to the find-  
 415 ings in plants [10], we found that intron gain events in  
 416 retrogenes were rare in humans. In spite of this rarity,  
 417 the mechanisms of intron creation in these retrogenes  
 418 are diverse. We found that retrogenes could gain introns  
 419 in three ways: insertion from an external sequence, tran-  
 420 scription in the opposite direction compared with the  
 421 parent, and intronization. For the latter method, in  
 422

t4.1 **Table 4 Substitution rates between the intronic and exonic regions of retrogenes and their corresponding regions of**  
 t4.2 **parental genes**

t4.3 t4.4 t4.5 t4.6 t4.7 t4.8 t4.9 t4.10 t4.11 t4.12 t4.13 Retrogene	Intronic region					Exonic region				
	$K_a$	$K_s$	$K_a/K_s$	<i>P</i> -value	Length	$K_a$	$K_s$	$K_a/K_s$	<i>P</i> -value	Length
t4.5 TMEM14D <sup>c</sup>	0.062	0.058	1.074	0.936	105	0.006	0.014	0.440	0.570	237
t4.6 RPS3AP5 <sup>a</sup>	NA	NA	NA	NA	NA	0.017	0.014	1.210	0.172	780
t4.7 XXyac-R12DG2.2 <sup>b</sup>	0.024	0.029	0.830	0.892	129	0.008	0.012	0.643	0.631	813
t4.8 HSP90B2Pa <sup>c*</sup>	0.823	0.597	1.379	0.526	144	0.045	0.067	0.678	0.091	2163
t4.9 HSP90AA4P <sup>c*</sup>	0.104	0.277	0.374	0.000	744	0.055	0.085	0.656	0.000	1374
t4.10 HSP90AA5P <sup>c*</sup>	0.087	0.215	0.406	0.001	672	0.088	0.221	0.400	0.082	897
t4.11 CSMD3 <sup>b</sup>	0.313	0.575	0.544	0.051	291	0.186	0.282	0.659	0.310	225
t4.12 WBP2NL <sup>b</sup>	0.033	0.088	0.373	0.192	177	0.385	0.377	1.021	0.919	684
t4.13 AC019016.1 <sup>c</sup>	0.083	0.082	1.010	0.978	636	0.081	0.175	0.466	0.084	273

t4.14  $K_a$  represents the non-synonymous substitution rate and  $K_s$  indicates the synonymous substitution rate. The *P*-value was calculated with the likelihood ratio test  
 t4.15 and the null hypothesis was  $K_a/K_s = 1$ . NA: not available (the corresponding parental sequence of the new intron in retrogene RPS3AP5 did not exist, because the  
 t4.16 intron was created by insertion of an external sequence). 'a', The retrogene gained introns by insertion of an external sequence. 'b', The retrogene gained introns  
 t4.17 after transcription in the opposite orientation compared to the parent. 'c', The retrogene gained introns by intronization. '\*', Evidence at the protein level for  
 t4.18 transcription of the retrogene was obtained.

423 addition to base substitution, retrogenes also may create  
424 introns in exonic regions via cryptic splice sites, which  
425 might be activated by the new gene structure after retro-  
426 position. Consistent with the findings in *C. elegans* [11],  
427 retrogene introns generated by intronization in humans  
428 are generally young and are mostly located in the coding  
429 region of the new gene. The retrogenes that underwent  
430 intronization in coding regions all retained the parental  
431 frames of translation and most showed expression evi-  
432 dence at the protein level. The significantly higher per-  
433 centage of non-frameshift introns implied that this kind  
434 of intron possessed a higher likelihood of persistence  
435 after intronization. The reason for this may be that  
436 frameshift introns mostly have a major effect on the pro-  
437 teins. Thus, non-frameshift introns are more likely to  
438 survive. However, non-frameshift introns may be neutral  
439 in effect, as proposed previously [43,44]. Furthermore,  
440 previous studies have shown that the rate of intron loss  
441 is much larger than that of intron gain in mammals  
442 [12,13,45]. Consequently, the older the retrogene is, the  
443 more probable the retrogene will lose the intronized  
444 exon, and this may explain why such introns are mainly  
445 observed in young retrogenes.

446 Some questions arise from careful examination of our  
447 observations. For example, for the retrogene RPS3AP5,  
448 in which the new intron was created by insertion of an  
449 external sequence, the process by which the new intron  
450 was created is unknown. In addition, in searches of  
451 UCSC [18,19], Ensembl [22,23], PeptideAtlas [27-29]  
452 and PRIDE [30,31], we did not obtain evidence of  
453 protein-level expression for the three retrogenes that  
454 gained introns after transcription in the reverse orienta-  
455 tion compared with their parents. The new introns in  
456 these three retrogenes were annotated to be in non-  
457 coding regions and appeared to be parts of existing  
458 intron-containing genes, as described previously [7].  
459 Thus, these retrogenes generally evolved faster than  
460 intronized retrogenes (Table 4).

461 For the eight non-frameshift introns generated by  
462 intronization, we examined whether they are under nat-  
463 ural selection by checking their genetic variation in dif-  
464 ferent human populations with the 1000 Genomes  
465 Browser [46]. However, we did not find insertions,  
466 deletions or mutations in splice sites in seven of these  
467 retrogenes (file 14), which implied that they are nearly  
468 fixed in all populations and may be under negative se-  
469 lection. In addition, there is a possibility that this pat-  
470 tern observed was caused by genetic drift because  
471 generation of new introns may be neutral. Finally,  
472 what is the importance of producing a shorter protein  
473 than the protein from the parent gene? This question  
474 may be answered by comparing the functions of the  
475 original proteins and that encoded by the retrogenes  
476 in the future.

## 477 Conclusions

478 Our results showed that retrogenes may gain introns in  
479 three ways: insertion from an external sequence, tran-  
480 scription in the reverse direction compared to that in  
481 the parent, and intronization. In addition to base substi-  
482 tution, intronization also may be promoted by cryptic  
483 splice sites. For introns generated by intronization, non-  
484 frameshift introns might have greater evolutionary suc-  
485 cess than frameshift introns, because non-frameshift  
486 introns have only a small effect on the host proteins or  
487 are neutral. Furthermore, intronization tended to occur  
488 in young retrogenes.

## 489 Additional files

490  
491  
492 **Additional file 1: Transcripts uniquely mapped to retrogenes.** This  
493 file lists transcripts that spanned the introns of their mapped retrogenes.

494 **Additional file 2: Evidence for transcription of retrogene introns**  
495 **(from the UCSC Genome Browser database).** This file contains  
496 snapshots from the UCSC Genome Browser database that displays the  
497 transcription of retrogenes that gained introns.

498 **Additional file 3: List of human tissues sampled for the**  
499 **experiments.** This file lists the human tissues that we used for the  
500 experiments to validate the existence of retrogene introns.

501 **Additional file 4: Experimental validation of retrogene introns in**  
502 **TMEM14D and HSP90AA4P.** This file shows the experimental results for  
503 validating the existence of retrogene introns.

504 **Additional file 5: Phylogenetic tree for vertebrates.** A diagram of the  
505 phylogenetic tree for vertebrates.

506 **Additional file 6: Chromosome and time of origin of intronized**  
507 **retrogenes.** This file shows the origination times of intronized  
retrogenes.

508 **Additional file 7: Chromosome and time of origin of retrogenes**  
509 **that gained introns by chimerization.** This file shows the origination  
510 times of retrogenes that gained introns by chimerization.

511 **Additional file 8: Transcription annotations (from the UCSC**  
512 **Genome Browser database) of retrogene introns in the parental**  
513 **gene.** This file contains snapshots from the UCSC Genome Browser  
514 database displaying transcription annotations of retrogene introns in the  
515 parental gene.

516 **Additional file 9: Sequences of primer pairs used to amplify the**  
517 **retrogenes and their parents.** A table that lists primer pairs we used to  
518 amplify the retrogenes and their parents.

519 **Additional file 10: Protein-level alignments of intron-gain**  
520 **retrogenes ("Sbjct") and their parents ("Query") by GeneWise.** This  
521 file contains alignments of intron-gain retrogenes and their parents in  
522 protein level.

523 **Additional file 11: Nucleotide-level alignments of retrogene introns**  
524 **('Sbjct', blue and red, splice sites) and parental genes ('Query',**  
525 **program NCBI-BLAST).** This file contains alignments of intron-gain  
526 retrogenes and their parents in DNA level.

527 **Additional file 12: Positions of three retrogenes (XXyac-R12DG2.2,**  
528 **CSMD3 and WBP2NL) in the human genome (from the UCSC**  
529 **Genome Browser database).** This file contains snapshots from the  
530 UCSC Genome Browser Database displaying the positions of three  
retrogenes.

531 **Additional file 13: Comparison of splicing signals (percentile score)**  
532 **in the corresponding region of the new intron in the parental gene**  
533 **and neighboring introns.** This file shows the results for the comparison  
534 of splicing signals in the corresponding region of the new intron in the  
535 parental gene and neighboring introns.



536 **Additional file 14: Genetic variation of four retrogenes in different**  
537 **human populations.** This file displays alignments of genomes of  
538 different human populations in the region of four retrogenes.

#### 539 Abbreviations

540 BS: branch site; PPT: polypyrimidine tract; SS: splice site.

#### 541 Authors' contributions

542 LFK and ZLZ together carried out the identification of intronized retrogenes  
543 and data analysis, and performed the statistical analyses. LFK performed the  
544 PCR analysis and helped to draft the manuscript. ZLZ conceived the study,  
545 participated in its design and analysis, and drafted the manuscript. QZ  
546 helped to perform the data analysis and statistical analyses, participated in  
547 the design of the study and helped to draft the manuscript. LYC provided  
548 the materials for experiments. ZZ participated in the design of the study and  
549 helped to draft the manuscript. All authors read and approved the final  
550 manuscript.

#### 551 Acknowledgements

552 Many thanks to Prof. Yong Zhang (Institute of Zoology, Chinese Academy of  
553 Sciences) and Prof. Tao Sang (Institute of Botany, Chinese Academy of  
554 Sciences) for invaluable suggestions and comments, and Dr. Quan-You Yu  
555 for help with the experiments. This work was supported by the Fundamental  
556 Research Funds for the Central Universities (No. CDJZR11290002) and by  
557 Natural Science Foundation Project of CQ CSTC (cstc2012jjB80007).

558 Received: 30 March 2012 Accepted: 28 July 2012

559 Published: 28 July 2012

#### 560 References

- 561 1. Brosius J: **Retroposons-seeds of evolution.** *Science* 1991, **251**(4995):753.
- 562 2. Betran E, Thornton K, Long M: **Retroposed new genes out of the X in**  
563 ***Drosophila*.** *Genome Res* 2002, **12**(12):1854–1859.
- 564 3. Wang W, Brunet FG, Nevo E, Long M: **Origin of sphinx, a young**  
565 **chimeric RNA gene in *Drosophila melanogaster*.** *Proc Natl Acad Sci USA*  
566 2002, **99**(7):4448–4453.
- 567 4. Nisole S, Lynch C, Stoye JP, Yap MW: **A Trim5-cyclophilin A fusion protein**  
568 **found in owl monkey kidney cells can restrict HIV-1.** *Proc Natl Acad Sci*  
569 *USA* 2004, **101**(36):13324–13328.
- 570 5. Sayah DM, Sokolskaja E, Berthoux L, Luban J: **Cyclophilin A**  
571 **retrotransposition into TRIM5 explains owl monkey resistance to HIV-1.**  
572 *Nature* 2004, **430**(6999):569–573.
- 573 6. Zhang J, Dean AM, Brunet F, Long M: **Evolving protein functional**  
574 **diversity in new genes of *Drosophila*.** *Proc Natl Acad Sci USA* 2004,  
575 **101**(46):16246–16250.
- 576 7. Baertsch R, Diekhans M, Kent WJ, Haussler D, Brosius J: **Retropcopy**  
577 **contributions to the evolution of the human genome.** *BMC Genomics*  
578 2008, **9**(1):466.
- 579 8. Fablet M, Bueno M, Potrzebowski L, Kaessmann H: **Evolutionary origin and**  
580 **functions of retrogene introns.** *Mol Biol Evol* 2009, **26**(9):2147–2156.
- 581 9. Lahn BT, Page DC: **Retroposition of autosomal mRNA yielded testis-specific**  
582 **gene family on human Y chromosome.** *Nat Genet* 1999, **21**(4):429–433.
- 583 10. Zhu Z, Zhang Y, Long M: **Extensive structural renovation of retrogenes in**  
584 **the evolution of the *Populus* genome.** *Plant Physiol* 2009, **151**(4):1943–1951.
- 585 11. Irimia M, Rukov JL, Penny D, Vinther J, Garcia-Fernandez J, Roy SW: **Origin**  
586 **of introns by 'intronization' of exonic sequences.** *Trends Genet*  
587 2008, **24**(8):378–381.
- 588 12. Roy SW, Fedorov A, Gilbert W: **Large-scale comparison of intron positions**  
589 **in mammalian genes shows intron loss but no gain.** *Proc Natl Acad Sci*  
590 *USA* 2003, **100**(12):7158–7162.
- 591 13. Coulombe-Huntington J, Majewski J: **Characterization of intron loss events**  
592 **in mammals.** *Genome Res* 2007, **17**(1):23–32.
- 593 14. Szczesniak MW, Ciomborowska J, Nowak W, Rogozin IB, Makalowska I: **Primate**  
594 **and rodent specific intron gains and the origin of retrogenes**  
595 **with splice variants.** *Mol Biol Evol* 2011, **28**(1):33–37.
- 596 15. Emerson JJ, Kaessmann H, Betran E, Long M: **Extensive gene traffic on the**  
597 **mammalian X chromosome.** *Science* 2004, **303**(5657):537–540.
- 598 16. Marques AC, Dupanloup I, Vinckenbosch N, Reymond A, Kaessmann H: **Emergence**  
599 **of young human genes after a burst of retroposition in**  
600 **primates.** *PLoS Biol* 2005, **3**(11):e357.
17. Vinckenbosch N, Dupanloup I, Kaessmann H: **Evolutionary fate of**  
601 **retroposed gene copies in the human genome.** *Proc Natl Acad Sci USA*  
602 2006, **103**(9):3220–3225.
- 603 18. Karolchik D, Hinrichs AS, Furey TS, Roskin KM, Sugnet CW, Haussler D,  
604 Kent WJ: **The UCSC Table Browser data retrieval tool.** *Nucleic Acids Res*  
605 2004, **32**(Database issue):D493–496.
- 606 19. Kuhn RM, Karolchik D, Zweig AS, Wang T, Smith KE, Rosenbloom KR, Rhead  
607 B, Raney BJ, Pohl A, Pheasant M, Meyer L, Hsu F, Hinrichs AS, Harte RA,  
608 Giardine B, Fujita P, Diekhans M, Dreszter T, Clawson H, Barber GP, Haussler  
609 D, Kent WJ: **The UCSC Genome Browser Database: update 2009.** *Nucleic*  
610 *Acids Res* 2009, **37**(Database issue):D755–761.
- 611 20. Altschul SF, Madden TL, Schaffer AA, Zhang J, Zhang Z, Miller W, Lipman DJ: **Gapped**  
612 **BLAST and PSI-BLAST: a new generation of protein database**  
613 **search programs.** *Nucleic Acids Res* 1997, **25**(17):3389–3402.
- 614 21. Zhang Z, Carriero N, Zheng D, Karro J, Harrison PM, Gerstein M: **PseudoPipe: an**  
615 **automated pseudogene identification pipeline.** *Bioinformatics* 2006,  
616 **22**(12):1437–1439.
- 617 22. Hubbard T, Barker D, Birney E, Cameron G, Chen Y, Clark L, Cox T, Cuff J,  
618 Curwen V, Down T, Durbin R, Eyras E, Gilbert J, Hammond M, Huminiecki L,  
619 Kasprzyk A, Lehvaslaiho H, Lijnzaad P, Melsopp C, Mongin E, Pettett R,  
620 Pocock M, Potter S, Rust A, Schmidt E, Searle S, Slater G, Smith J, Spooner  
621 W, Stabenau A, Stalker J, Stupka E, Ureta-Vidal A, Vastrik I, Clamp M: **The**  
622 **Ensembl genome database project.** *Nucleic Acids Res* 2002, **30**(1):38–41.
- 623 23. Hubbard TJ, Aken BL, Ayling S, Ballester B, Beal K, Bragin E, Brent S, Chen  
624 Y, Clapham P, Clarke L, Coates G, Fairley S, Fitzgerald S, Fernandez-Banet  
625 J, Gordon L, Graf S, Haider S, Hammond M, Holland R, Howe K,  
626 Jenkinson A, Johnson N, Kahari A, Keefe D, Keenan S, Kinsella R,  
627 Kokoćinski F, Kulesha E, Lawson D, Longden I, Megy K, Meidl P, Overduin  
628 B, Parker A, Pritchard B, Rios D, Schuster M, Slater G, Smedley D, Spooner  
629 W, Spudich G, Trevanion S, Vilella A, Vogel J, White S, Wilder S, Zadissa A,  
630 Birney E, Cunningham F, Curwen V, Durbin R, Fernandez-Suarez XM,  
631 Herrero J, Kasprzyk A, Proctor G, Smith J, Searle S, Flicek P: **Ensembl**  
632 **2009.** *Nucleic Acids Res* 2009, **37**(suppl 1):D690–697.
- 633 24. Kent WJ: **BLAT—the BLAST-like alignment tool.** *Genome Res* 2002,  
634 **12**(4):656–664.
- 635 25. Nei M, Gojobori T: **Simple methods for estimating the numbers of**  
636 **synonymous and nonsynonymous nucleotide substitutions.** *Mol Biol Evol*  
637 1986, **3**(5):418–426.
- 638 26. Yang Z: **PAML: a program package for phylogenetic analysis by**  
639 **maximum likelihood.** *Comput Appl Biosci* 1997, **13**(5):555–556.
- 640 27. Desiere F, Deutsch EW, Nesvizhskii AI, Mallick P, King NL, Eng JK, Aderem A,  
641 Boyle R, Brunner E, Donohoe S, Fausto N, Hafen E, Hood L, Katze MG,  
642 Kennedy KA, Kregenow F, Lee H, Lin B, Martin D, Ranish JA, Rawlings DJ,  
643 Samelson LE, Shih Y, Watts JD, Wollscheid B, Wright ME, Yan W, Yang L, Yi  
644 EC, Zhang H, Aebersold R: **Integration with the human genome of**  
645 **peptide sequences obtained by high-throughput mass spectrometry.**  
646 *Genome Biol* 2005, **6**(1):R9.
- 647 28. Deutsch EW, Lam H, Aebersold R: **PeptideAtlas: a resource for target**  
648 **selection for emerging targeted proteomics workflows.** *EMBO Rep* 2008,  
649 **9**(5):429–434.
- 650 29. Farrah T, Deutsch EW, Omenn GS, Campbell DS, Sun Z, Bletz JA, Mallick P,  
651 Katz JE, Malmström J, Ossola R, Watts JD, Lin B, Zhang H, Moritz RL,  
652 Aebersold R: **A high-confidence human plasma proteome reference set**  
653 **with estimated concentrations in PeptideAtlas.** *Mol Cell Proteomics* 2011,  
654 **10**(9):M110.006353.
- 655 30. Jones P, Cote RG, Martens L, Quinn AF, Taylor CF, Derache W, Hermjakob  
656 H, Apweiler R: **PRIDE: a public repository of protein and peptide**  
657 **identifications for the proteomics community.** *Nucleic Acids Res* 2006,  
658 **34**(Database issue):D659–663.
- 659 31. Vizcaino JA, Cote R, Reisinger F, Foster JM, Mueller M, Rameseder J,  
660 Hermjakob H, Martens L: **A guide to the Proteomics Identifications**  
661 **Database proteomics data repository.** *Proteomics* 2009, **9**(18):4276–4283.
- 662 32. Jenuth JP: **The NCBI. Publicly available tools and resources on the Web.**  
663 *Methods Mol Biol* 2000, **132**:301–312.
- 664 33. Johnson M, Zaretskaya I, Raytselis Y, Merezuk Y, McGinnis S, Madden TL:  
665 **NCBI BLAST: a better web interface.** *Nucleic Acids Res* 2008, **36**(Web Server  
666 issue):W5–9.
- 667 34. Zhang YE, Vibranovski MD, Landback P, Marais GA, Long M: **Chromosomal**  
668 **redistribution of male-biased genes in mammalian evolution with two**  
669 **bursts of gene gain on the X chromosome.** *PLoS Biol* 2010, **8**(10):  
670 e1000494.
- 671

- 672 35. Wu NW, Jalkanen S, Streeter PR, Butcher EC: Evolutionary conservation of  
673 tissue-specific lymphocyte-endothelial cell recognition mechanisms  
674 involved in lymphocyte homing. *J Cell Biol* 1988, **107**(5):1845–1851.
- 675 36. Trusov YA, Dear PH: A molecular clock based on the expansion of gene  
676 families. *Nucleic Acids Res* 1996, **24**(6):995–999.
- 677 37. Thomas JW, Touchman JW: Vertebrate genome sequencing: building a  
678 backbone for comparative genomics. *Trends Genet* 2002, **18**(2):104–108.
- 679 38. Zhao S, Shetty J, Hou L, Delcher A, Zhu B, Osoegawa K, de Jong P,  
680 Nierman WC, Strausberg RL, Fraser CM: Human, mouse, and rat genome  
681 large-scale rearrangements: stability versus speciation. *Genome Res*  
682 2004, **14**(10A):1851–1860.
- 683 39. Falkowski PG, Katz ME, Milligan AJ, Fennel K, Cramer BS, Aubry MP,  
684 Berner RA, Novacek MJ, Zapol WM: The rise of oxygen over the past  
685 205 million years and the evolution of large placental mammals.  
686 *Science* 2005, **309**(5744):2202–2204.
- 687 40. Waters PD, Delbridge ML, Deakin JE, El-Mogharbel N, Kirby PJ,  
688 Carvalho-Silva DR, Graves JA: Autosomal location of genes from the  
689 conserved mammalian X in the platypus (*Ornithorhynchus anatinus*):  
690 implications for mammalian sex chromosome evolution. *Chromosome*  
691 *Res* 2005, **13**(4):401–410.
- 692 41. Schwartz S, Hall E, Ast G: SROOGLE: webserver for integrative, user-friendly  
693 visualization of splicing signals. *Nucleic Acids Res* 2009, **37**(Web Server issue):  
694 W189–192.
- 695 42. Cavalier-Smith T: Selfish DNA and the origin of introns. *Nature* 1985,  
696 **315**:283–284.
- 697 43. Castillo-Davis CI, Bedford TBC, Hart DL: Accelerated rates of intron gain/  
698 loss and protein evolution in duplicate genes in human and mouse  
699 malaria parasites. *Mol Biol Evol* 2004, **21**(7):1422–1427.
- 700 44. Li W, Tucker AE, Sung W, Thomas WK, Lynch M: Extensive, recent intron  
701 gains in *Daphnia* populations. *Science* 2009, **326**(5957):1260–1262.
- 702 45. Roy SW, Gilbert W: Rates of intron loss and gain: implications for early  
703 eukaryotic evolution. *Proc Natl Acad Sci USA* 2005, **102**:5773–5778.
- 704 46. The 1000 Genomes Project Consortium: A map of human genome variation  
705 from population-scale sequencing. *Nature* 2010, **467**(7319):1061–1073.
- 706 47. Kol G, Lev-Maor G, Ast G: Human-mouse comparative analysis reveals that  
707 branch-site plasticity contributes to splicing regulation. *Hum Mol Genet*  
708 2005, **14**(11):1559–1568.
- 709 48. Schwartz SH, Silva J, Burstein D, Pupko T, Eyras E, Ast G: Large-scale  
710 comparative analysis of splicing signals and their corresponding splicing  
711 factors in eukaryotes. *Genome Res* 2008, **18**(1):88–103.

712 doi:10.1186/1471-2148-12-128

713 Cite this article as: Kang et al.: Newly evolved introns in human  
714 retrogenes provide novel insights into their  
715 evolutionary roles. *BMC Evolutionary Biology* 2012 **12**:128.

Submit your next manuscript to BioMed Central  
and take full advantage of:

- Convenient online submission
- Thorough peer review
- No space constraints or color figure charges
- Immediate publication on acceptance
- Inclusion in PubMed, CAS, Scopus and Google Scholar
- Research which is freely available for redistribution

Submit your manuscript at  
www.biomedcentral.com/submit

