

## Characterization of a novel alternatively-spliced 5' exon in the human insulin-like growth factor I (IGF-I) gene, expressed in liver and some cancers

Article (Accepted Version)

Wallis, Michael (2019) Characterization of a novel alternatively-spliced 5' exon in the human insulin-like growth factor I (IGF-I) gene, expressed in liver and some cancers. *Growth Hormone and IGF Research*, 46-47. pp. 36-43. ISSN 1096-6374

This version is available from Sussex Research Online: <http://sro.sussex.ac.uk/id/eprint/84510/>

This document is made available in accordance with publisher policies and may differ from the published version or from the version of record. If you wish to cite this item you are advised to consult the publisher's version. Please see the URL above for details on accessing the published version.

### **Copyright and reuse:**

Sussex Research Online is a digital repository of the research output of the University.

Copyright and all moral rights to the version of the paper presented here belong to the individual author(s) and/or other copyright owners. To the extent reasonable and practicable, the material made available in SRO has been checked for eligibility before being made available.

Copies of full text items generally can be reproduced, displayed or performed and given to third parties in any format or medium for personal research or study, educational, or not-for-profit purposes without prior permission or charge, provided that the authors, title and full bibliographic details are credited, a hyperlink and/or URL is given for the original metadata page and the content is not changed in any way.

1     **Characterization of a novel alternatively-spliced 5' exon in the human insulin-**  
2           **like growth factor I (IGF-I) gene, expressed in liver and some cancers**

3

4

Michael Wallis

5

6

Biochemistry and Biomedicine Group, School of Life Sciences, University of Sussex, Brighton BN1

7

9QG, UK

8

9

email: m.wallis@sussex.ac.uk

10

11

Running title: **A novel 5' exon in the human IGF-I gene**

12

13

This is the author's accepted version of an article published in

14

Growth Hormone & IGF Research, Volumes 46–47, June–August 2019, Pages 36-43

15

<https://doi.org/10.1016/j.ghir.2019.06.002>

16

17 **Abstract**

18

19 In mammals, the large IGF-I gene comprises 6 exons, which are subject to alternative splicing. All  
20 transcripts contain exons 3 and 4, encoding mature IGF-I, but the other exons are included in various  
21 combinations, giving at least 6 possible mature mRNAs. At the 5' end, exons 1 and 2 are spliced  
22 alternatively to exon 3, giving different leader/signal sequences. It is shown in this study that in human  
23 an additional exon (designated exon 0) is present, upstream of exon 1. This can be spliced directly to  
24 exon 3 or, less frequently, into exon 1. Exon 0 is utilised in liver, in about 24% of IGF-I transcripts, to  
25 a minor extent in prostate and endometrium (<1% of transcripts), but not in any of 29 other normal  
26 human tissues examined. The exon 0 sequence includes an in-frame ATG/AUG, potentially providing a  
27 translation start point giving an IGF-I precursor with a very long signal peptide. However, this ATG is  
28 very close to the 5' end, and may not be included in all transcripts; an in-frame ATG in exon 3 could  
29 provide an alternative start point. Utilization of exon 0 was detected in other apes, and to a small extent  
30 in Old World monkeys, but not in New World monkeys, prosimians or various non-primate mammals.  
31 Exon 0 was not expressed in most human tumours, but was utilized in many prostate tumours, at levels  
32 much greater than seen in normal prostate, and in liver tumours, at a lower level than in normal liver.

33

34 **Keywords:** IGF-I, exon 0, alternative splicing, liver, primates, prostate cancer

35

36

## 37 **1. Introduction**

38

39 Insulin-like growth factor I (IGF-I) is a small protein formed from a larger precursor, encoded by a  
40 very large gene extending over ~85 kb in human. In mammals this gene comprises at least 6 exons,  
41 with the sequence encoding mature IGF-I included in exons 3 and 4. Alternative splicing allows these  
42 two exons to be combined with various combinations of the others, leading to production of at least 6  
43 different IGF-I precursors [1-3] all of which include the sequence of mature IGF-I, which is generally  
44 strongly conserved [4-6]. At the 5' end the alternative splicing leads to production of mRNAs and IGF-  
45 I precursors with differing leader/signal sequences while at the 3' end varying amino acid (aa)  
46 sequences are produced, which can be processed to a number of peptides with potential biological  
47 actions [2-5]. IGF-I is produced in many tissues, and has both endocrine and paracrine actions. It  
48 mediates the actions of pituitary growth hormone, forming the basis of the somatomedin hypothesis  
49 [7,8]. Initially this mediating role was thought to involve circulating IGF-I, but more recent evidence  
50 suggests that paracrine actions may be more important [8]. Nevertheless, most IGF-I in the body is  
51 found in the circulation, associated with IGF binding proteins [1]. The source of this circulating IGF-I  
52 appears to be largely the liver [9,10] though its function is unclear [8].

53

54 Alternative splicing at the 5' end of the IGF-I gene has been widely studied, especially in rodents and  
55 human, though its significance remains poorly understood. In mammals (but not other vertebrates,  
56 where there is no equivalent of mammalian exon 2 [11]) exon 3 can be spliced to either exon 1 or exon  
57 2, giving mRNA variants that encode IGF-I precursors with signal peptides of respectively 48 or 32  
58 amino acid residues (aas) preceding the mature IGF-I sequence. Studies using cell-free systems  
59 showed that with rat mRNAs, signal peptides of these lengths are indeed produced [12,13] although  
60 initiation of translation can also occur at an ATG/AUG within exon 3, giving a signal peptide of 22 aas  
61 [13]. Translational efficiency appeared to be much greater for the 32 and 22 aa signal peptide forms  
62 than for the 48 aa variant [13]. Further complications are introduced because transcription of both exon  
63 1 and exon 2 is initiated at multiple sites, giving mRNAs with untranslated regions (utrs) of variable  
64 length [14,15]; for the shortest of these in exon 1, the potential initiating AUG is missing, so translation  
65 would presumably start in exon 3. Yet another complication in rat is the occurrence of a deletion of 186  
66 nt within the 5' utr in some transcripts, apparently due to another splicing event [16]; this deletion

67 removes some of the potential transcription initiation sites. The significance of all this variation at the  
68 5' end of the IGF-I mRNA (and potentially the N-terminus of the IGF-I precursor) is unclear. Use of  
69 alternative leader sequences, with different promoter sequences may allow differential regulation of  
70 expression, for example between tissues or developmental stages [17-20] or in response to growth  
71 hormone [21-23]. Initiation of translation at different points can occur with differing efficiency, as  
72 discussed. Use of different signal sequences could result in varying intracellular localization of IGF-I  
73 precursors and/or post-translational modification [24]. The physiological importance of these various  
74 possibilities remains to be determined.

75

76 In light of this, the variation of the alternative splicing observed for IGF-I between different tissues,  
77 developmental stages and species is clearly of interest. In primates, variation in splicing across tissues  
78 has been assessed to some extent in human, but has been little studied in other primate species. The  
79 availability of large amounts of data on gene transcription in the ncbi SRA database makes possible the  
80 assessment of IGF-I splicing in a number of primate species. Such a study is reported here, focussing  
81 particularly on splicing at the 5' end of the gene. It was found that splicing does vary considerably  
82 among different tissues and species. Most notably, a novel exon, encoding a third 5' leader sequence,  
83 was found to be used in human and some other primates, expressed primarily in liver and some  
84 tumours, and associated with a substantial proportion of IGF-I transcripts.

85

## 86 **2. Methods**

87

### 88 *2.1. Sequence and transcriptomic data*

89 The sequences of IGF-I genes for a range of primates and other mammals were obtained using the  
90 Ensembl genome browser ([www.ensembl.org](http://www.ensembl.org)). Releases used were as follows: human (*Homo sapiens*)  
91 GRCH38, chromosome 12; chimpanzee (*Pan troglodytes*) Pan\_tro\_3.0, chromosome 12; gorilla  
92 (*Gorilla gorilla*) gorGor4, chromosome 12; orangutan (*Pongo pygmaeus*) PPYG2, chromosome 12;  
93 gibbon (*Nomascus leucogenys*) Nleu\_3.0, chromosome 10; rhesus macaque (*Macaca mulatta*)  
94 Mmul\_8.0.1, chromosome 11; African green monkey (vervet, *Chlorocebus sabaeus*) ChlSabl.1,  
95 chromosome 11; marmoset (*Callithrix jacchus*) ASM275486v1, contig NTIC01000002.1; squirrel  
96 monkey (*Saimiri boliviensis*) SaiBol1.0, scaffold JH378139; tarsier (*Tarsier syrichta* now renamed

97 *Carlito syrichta*) Tarsius\_syrichta-2.0.1, scaffold KE941348.1; mouse lemur (*Microcebus murinus*)  
98 Mmur\_3.0, chromosome 7; galago (*Otolemur garnettii*) OtoGar3, scaffold GL873526.1; tree shrew  
99 (*Tupaia belangeri*) tupBell, scaffold GeneScaffold\_475; rat (*Rattus norvegicus*) Rnor\_6.0,  
100 chromosome 7; rabbit (*Oryctolagus cuniculus*) OryCun2.0, chromosome 4; dog (*Canis familiaris*)  
101 CanFam3.1, chromosome 15. Exons in these genes were identified by comparison with the human and  
102 rat genes; the gene structure for IGF-I is known to be well conserved across mammals [4].

103

104 Transcriptomes were accessed via the SRA database ([trace.ncbi.nlm.nih.gov/Traces/sra](http://trace.ncbi.nlm.nih.gov/Traces/sra)). The SRA  
105 projects used in this study are summarized in Supplementary Table 1.

106

### 107 2.2. Transcriptome analysis

108 Each transcriptome database (experiment) was subjected to BLAST searches [25] on the SRA website,  
109 using IGF-I sequences from the appropriate species. For each transcriptome, sequences with exon 3  
110 linked at the 5' end to exon 0, exon 1 and exon 2 respectively were used to determine the relative usage  
111 of these three exons. Full sequences of these Blast Queries are given in Supplementary Table 2.

112 Expression of IGF-I in various tissues was determined by BLAST analysis with the sequence encoding  
113 mature IGF-I (since this is included in all IGF-I precursor transcripts) and expressed as BLAST  
114 hits/million reads.

115

### 116 2.3. Exon profiling

117 To obtain a graphical visualization of the relative usage of the exons at the 5' end of the IGF-I gene, an  
118 in-house perl programme was used. The perl script is given in Supplementary Method 1. Input for this  
119 programme was the Blast alignment output produced by searching the corresponding transcriptomic  
120 database with a section of the IGF-I gene extending from about 5000 nt upstream of exon 1 to about  
121 2000 nt downstream of exon 3. The output from the programme gave the number of representations in  
122 the alignment for each nucleotide, which was converted to graphical form using Microsoft Excel.

123 Identification of repetitive sequence elements was performed using RepeatMasker [26].

124

### 125 2.4. Statistical analysis

126 Statistical analysis was carried out using the R package ([www.R-project.org](http://www.R-project.org)).

127

128

129 **3. Results**

130

131 *3.1. A novel 5' utr and leader sequence for human IGF-I (exon 0)*

132 Investigation of sequences spliced to the 5' end of IGF-I exon 3 in human liver transcriptomes (initially  
133 in SRA project ERP003613) revealed not only sequences corresponding to exons 1 and 2, but  
134 additional sequence corresponding to a region about 1380 nucleotides (nt) upstream of the 3' end of  
135 exon 1. BLAST analysis of this region identified a new IGF-I leader exon extending over about 145 nt  
136 (Fig.1), though the majority of reads corresponded to a shorter sequence of 110 nt. This novel exon is  
137 referred to as exon 0. In the majority of transcripts (~90%) it is spliced to the 5' end of exon 3, but in a  
138 minority it is spliced into a site 82 nt upstream of the 3' end of exon 1. Exon 0 contains an in-frame  
139 ATG/AUG which would allow production of a very long signal peptide of 65 aas (Fig.1). However,  
140 this ATG is close to the apparent (putative) 5' end of the majority of transcripts, and has weak  
141 surrounding sequence for an initiating AUG [27], so it is doubtful whether initiation starts at this point.  
142 If not, initiation would be expected to start at the first in-frame ATG/AUG in exon 3, giving a signal  
143 peptide of 25 aas; the equivalent signal sequence in rat is known to be functional [13] though a little  
144 shorter (22 aas).

145

146 Exon profiling for the 5' end of the human IGF-I gene, including exons 0-3 and ~4750 nt upstream of  
147 exon 0 confirmed the expression of exon 0 in human liver, as well as exons 1 and 2 (Fig. 2). Minor  
148 peaks between exons 0 and 1, exons 1 and 2 and exons 2 and 3 may reflect traces of unspliced or  
149 undegraded introns; none of these showed evidence for splicing to exons. The small peak ~3800 nt  
150 upstream of exon 0 was identified by RepeatMasker as an L1-like repetitive element, with again no  
151 evidence for splicing to an exon. All the exon 3 sequences in this transcriptome that ran up from the 5'  
152 splice site were spliced to exon 0, exon 1 or exon 2. On the basis of the profile of Fig. 2, exon 1 is the  
153 major 5' leader sequence, but exon 0 appears to be more highly represented than exon 2.

154

155 Expression of exon 0 was observed in all normal human liver transcriptomes examined (from 11  
156 individuals in 5 separate SRA projects). Exon 1 to exon 3 splices predominated (71% of all such

157 splices), with exon 0-exon 3 splices at 22%, exon 2-exon 3 6% and exon 0-exon 1-exon 3 2% (Fig. 3).  
158 However, the range of values was large (Fig.3), probably reflecting the fact that liver samples were  
159 obtained from a wide range of individuals with no controlling for sex, age or physiological state. In  
160 many cases such information was not available. Where replicate values for individuals were available,  
161 agreement was close.

162

### 163 *3.2. Expression of IGF-I and alternative splicing in normal human tissues*

164 Any one of three upstream exons can be used in the expression of the human IGF-I gene, with exon 0,  
165 1 or 2 spliced to the 5' of exon 3 and in addition, a minor variant in which exon 0 is spliced into exon 1  
166 (Fig. 1). Each of these variants gives a different leader sequence on the derived mRNA, and potentially  
167 a different signal peptide on the IGF-I precursor. The patterns of expression of these variants in a wide  
168 variety of human tissues was determined using transcriptomic data from SRA projects ERP003613 and  
169 ERP006650 [28] and is summarized in Fig. 4. Fig. 4A shows the expression level of IGF-I for each  
170 tissue, determined by a BLAST search using the coding sequence for mature IGF-I, and expressed as  
171 BLAST hits/ $10^6$  reads. Expression levels vary markedly between tissues, and for some tissues  
172 (particularly endometrium, liver and smooth muscle) between individuals. Replicates for individuals  
173 generally showed good agreement. The level seen for liver was similar to the mean for the 11  
174 individuals of Fig. 3 ( $4.20 \pm 0.92$  BLAST hits/ $10^6$  reads). Several tissues show an expression  
175 level similar to that in liver. All 32 tissues studied showed some expression of IGF-I, though  
176 levels in skin, bone marrow, skeletal muscle and pancreas were very low.

177

178 Fig. 4B indicates the relative proportions of splices of exons 0, 1 and 2 to exon 3 for each tissue. Use of  
179 exon 0 was seen in liver and to a very minor extent (less than 1%) in prostate gland and endometrium,  
180 but not in any other tissue. Use of exon 1 was predominant in all tissues. The percent utilization of  
181 exon 2 varied between tissues, and this variation was shown to be significant by one-way analysis of  
182 variance ( $p < 0.001$ ). High variability was seen for some tissues; this, and the failure to detect any exon  
183 2-exon 3 splices in some tissues may have been due to the low number of exon 2-exon 3 BLAST hits  
184 obtained for tissues with low overall expression of IGF-I.

185

### 186 *3.3. Expression and alternative splicing of IGF-I in non-human primates*



187 An alignment of exon 0-like sequences derived from genomic sequences of various primates and a few  
188 non-primate mammals is shown in Fig. 5. Strong similarity is seen between primate sequences, but is  
189 less clear cut for non-primates. Rat sequence was not included because a convincing alignment could  
190 not be produced. The positions of two potential initiating ATGs is shown. After splicing to exon 3 the  
191 more upstream of these (ATG-1) is in-frame with the IGF-I coding sequence for the human sequence.  
192 The chimpanzee sequence is identical to that of human. Gorilla, orangutan and gibbon sequences are  
193 very similar to human, but insertion of an additional base (A) in the long A-tract changes the reading  
194 frame so that the downstream ATG (ATG-2) rather than ATG-1 is in-frame. Sequences for two Old-  
195 World monkeys (macaque and vervet) differ more from human; insertion of 6 bases leaves ATG-1 in-  
196 frame. Sequences of two New-World monkeys (marmoset and squirrel monkey) show a 9 base  
197 deletion, compared with human, in the region of the long A-tract; this would potentially leave the  
198 reading frame unchanged, but mutation of ATG-1 to CTG would prevent translation initiation from this  
199 site. Sequences from prosimians (tarsier, mouse lemur and galago) diverge further from human, and  
200 include a number of insertions/deletions (indels). Sequences of non-primates are even more divergent,  
201 with convincing alignment difficult to establish at the 5' end.

202

203 Expression of IGF-I was seen in livers of three apes (chimpanzee, gorilla and orangutan; no liver  
204 transcriptomic data is available in the SRA for gibbon) and in the various other primates and non-  
205 primates studied (Fig. 6A). Notably the expression level in human liver was lower than that in any  
206 other species except chimpanzee (though for gorilla, orangutan and squirrel monkey data from only  
207 two individuals were available, so this conclusion must be tentative). The level of IGF-I mRNA in rat  
208 liver was about 25-fold greater than that in human.

209

210 Utilization of exon 0 was seen in liver transcriptomes from the apes for which data is available  
211 (chimpanzee, gorilla and orangutan) though in all cases it was lower than in human. Exon 0 was also  
212 utilized, but to a very minor extent (<1%) in the Old-World monkeys (Fig. 6B). In none of the other  
213 species examined was there any utilization of an exon corresponding to human exon 0, with the 5' of  
214 exon 3 being spliced only to exon 1 or exon 2, except for the mouse lemur where a sequence between  
215 exons 1 and 2 was utilized to a minor extent (~3%).

216

217 The relative importance of exons 1 and 2 varies markedly between species, with exon 1 predominating  
218 in some species, including the apes and mouse lemur, but close to equal usage of exons 1 and 2 in some  
219 cases (e.g. galago and rat) (Fig. 6B). Notably, in the two New-World monkeys use of exon 2  
220 predominates (Fig. 6B); the IGF-I system is unusual in this group in several other ways, including  
221 accelerated evolution of IGF-I and of its receptor [29].

222

#### 223 *3.4. IGF-I alternative splicing in human tumours*

224 As indicated, in normal human tissues, utilization of exon 0 is confined largely to liver. The possibility  
225 that this exon is expressed in tumours was explored by examining SRA transcriptomic databases for a  
226 number of human cancers. In a number of tumour types examined, including breast, colorectal and  
227 uterine cancer and Ewing sarcoma (SRA projects SRP63460, SRP042620, SRP029880, SRP060016)  
228 utilization of exon 0 was not detected. In endometrial cancer (SRP040442) very low expression of  
229 exon 1 was detected in 2 of the 9 individuals studied, a level similar to that seen in normal  
230 endometrium. Low expression of exon 1 was also detected infrequently in lung cancer (in 2 of 40  
231 individuals examined; project ERP001058).

232

233 The situation in prostate cancer is markedly different, as seen in three SRA projects (ERP000550,  
234 ERP017433, ERP006077). Study ERP006077 [30] included the largest number of cases, and results  
235 derived from this are shown in Fig. 7; similar results were obtained for the other two studies.  
236 Expression levels of IGF-I (Fig. 7A) in most tumours were similar to those in normal prostate (Fig. 4),  
237 but in some they were much lower and in one case (T23) markedly higher. Many of these tumours  
238 showed utilization of exon 0, with 17/25 showing such utilization at >1% of all 5' splices to exon 3,  
239 6/25 utilization at >10%, and in one case >25% (Fig. 7B). Utilization of exon 0 in normal prostate is  
240 very low (<1% in the study of Fig. 4; no expression detected in two other studies, ERP076277,  
241 SRP163292). For 5 of the patients included in the project illustrated in Fig. 7 apparently normal tissue  
242 adjacent to the prostate tumour was also investigated; utilization of exon 0 was very low (<1%) in 3 of  
243 these cases, but it was higher (1.7% and 9.0%) in the other 2.

244

245 Many liver tumours express IGF-I at levels similar to those seen in normal liver (SRA studies  
246 SRP030040, SRP064431, SRP007560). Utilization of exon 0 in these tumours was variable, but

247 generally rather lower than that seen in normal liver (Fig. 3) or, in most cases, apparently normal tissue  
248 adjacent to the tumours. Fig. 8 summarises data derived from study SRP064431.

249

250

#### 251 **4. Discussion**

252

253 Alternative splicing of the IGF-I gene has been recognized for many years, with, in mammals, the 6  
254 exons of this complex gene differentially spliced at both the 5' and 3' ends, resulting in at least 6  
255 different mRNAs and IGF-I precursors, all including the mature IGF-I sequence [1-3,6]. At the 5' end  
256 of the gene, two exons have been identified, exons 1 and 2 (Fig. 1), either of which can be spliced to  
257 exon 3. The relative usage of these 5' exons in various tissues, developmental stages and disease  
258 conditions has been widely studied, but their significance remains unclear. Usage of exon 1 is generally  
259 greater than that of exon 2, although the ratio between the two varies and relative exon 2 usage is  
260 increased by growth hormone treatment [19]. For both these exons the transcription start points can  
261 vary, giving mRNAs with 5' utrs of differing length [14,15]. Initiation of translation can occur within  
262 both exon 1 and exon 2, giving signal peptides of 48 or 32 aas respectively, of which the C-terminal 27  
263 residues are encoded by exon 3 and therefore identical in the two signal peptides [13]. Initiation can  
264 also occur within exon 3, giving a signal peptide of 22 aas in rat (25 aas in human); exons 1 and 2  
265 would then be untranslated. Initiation within exons 2 or 3 is more efficient than within exon 1 [13], so  
266 the higher usage of exon 1 may not necessarily mean that it is of more importance in overall IGF-I  
267 production. The 48-residue signal peptide encoded when initiation occurs within exon 1 is  
268 exceptionally long, and may have a specific role, such as directing the protein to a particular  
269 subcellular compartment.

270

271 It is shown here that in human the situation is further complicated by utilization of an additional  
272 upstream exon (exon 0) in the IGF-I gene. Alternative splicing of this gives rise to two novel leader  
273 sequences; in the predominant one exon 0 is spliced directly to exon 3, while less frequently exon 0 is  
274 spliced into exon 1. Utilization of exon 0 in human IGF-I has not been generally recognized previously  
275 and is not shown in the GTEx database (<https://gtexportal.org/home/gene/IGF1>), though it is predicted  
276 in three entries in the ncbi nucleotide database ([www.ncbi.nlm.nih.gov/nucleotide/](http://www.ncbi.nlm.nih.gov/nucleotide/)), produced by

277 automated computer analysis (Variant X1, XM\_017019259; Variant X3, XM\_017019262; Variant X4,  
278 XM\_017019263). In normal human tissues utilization of exon 0 is confined mostly to liver, where it  
279 comprises about 24% of splices to exon 3 (Fig. 3, Fig. 4). An AUG/ATG codon 114 nt upstream of the  
280 splice site would be in-frame with the IGF-I mature sequence when exon 0 is spliced to exon 3, giving  
281 a very long putative signal peptide of 65 aas. Whether translation does in fact initiate at this site is not  
282 clear. Analysis of transcriptomic data cannot define precisely the transcription start point, but few  
283 transcripts extend up 114 nt from the splice site, and those that do mostly extend only a short distance  
284 upstream of this point (Fig. 5). There does not appear to be a TATA box or CCAAT element upstream  
285 of the exon 0 sequence, suggesting that as for exon 1 multiple transcription start points may apply, so  
286 many of the transcripts may be too short to allow translation initiation at the in-frame AUG/ATG.  
287 Translation would then probably start at the AUG 7 nt from the 5' end of exon 3-derived sequence,  
288 giving a signal peptide of 25 aas. In the exon 0-exon 1-exon 3 variant, the AUG/ATG 114 nt upstream  
289 of the exon 0 splice site would not be in-frame, but an AUG/ATG 22 nt downstream of this would,  
290 potentially allowing translation of a precursor with a very long signal peptide of 85 aas.

291

292 The origin of exon 0 has been clarified by analysing sequences from various primates. A sequence  
293 upstream of exon 1 similar to that of human exon 0 was identifiable in all primates, and in several non-  
294 primate mammals, but not rat, where homologous sequence may be absent or just unrecognizable due  
295 to divergent evolution. Utilization of exon 0, judged by its appearance in liver transcripts, is only seen  
296 in Old World monkeys and apes, and in the former is very low (< 1% of all splices to exon 3; Fig. 6).  
297 It would thus appear that exon 0 arose by exaption of an existing sequence present upstream of the  
298 IGF-I gene, possibly only becoming fully functional in apes. Also notable from the comparative study  
299 of Fig. 6 are (1) the relatively low expression of IGF-I in human liver compared with many other  
300 species, the contrast being particularly marked in the case of rat, and (2) the rather low utilization of  
301 exon 2 in human and other apes, compared with most other primates and non-primate mammals.

302

303 Expression of exon 0 was largely confined to liver and not detected in 29 other human tissues  
304 examined, though very low expression was detected in prostate and endometrium. Notably, for many of  
305 these tissues, expression of IGF-I was similar to, or higher than, that in liver, though variation between  
306 individuals was high. Unlike exon 0, expression of exon 2 was seen in most human tissues, but always

307 less than that of exon 1, though the ratio between the two varied considerably. ANOVA showed that  
308 such variation was significant, indicating independent regulation of the expression of these two exons.  
309

310 The potential involvement of IGF-I in cancer has been discussed extensively [31-37], and the  
311 possibility of aberrant expression of exon 0 in cancer was investigated. In most cancers examined  
312 expression of exon 0 was absent (breast, colorectal, uterus, Ewing sarcoma) or very low (endometrium,  
313 lung). A marked exception was prostate cancer, where a high proportion of patients showed substantial  
314 expression of exon 0. The significance of this is unclear, but the abnormal expression of a potent  
315 growth factor could well play a part in tumour growth and progression. Exon 0 expression was also  
316 seen in liver cancer (Fig. 8), but at a rather lower level than seen in normal liver tissue.

317  
318 In conclusion, the study described here shows that in human and apes, the IGF-I gene includes a novel  
319 exon, at the 5' end. Alternative splicing allows either exon 0, exon 1 or exon 2 to be spliced to exon 3,  
320 potentially giving 3 IGF-I precursors with different signal peptides. A further minor variant can also be  
321 produced, resulting from splicing of exon 0 into exon 1. Since alternative splicing also occurs at the 3'  
322 end of the IGF-I gene the total number of potential variants of the IGF-I precursor becomes quite large.  
323 The significance of the alternative leader sequences for human IGF-I remains unclear, but it is notable  
324 that unlike exons 1 or 2 utilization of exon 0 is very tissue specific, for normal tissues being confined  
325 largely to liver. In mice prevention of utilization of exon 2 has no apparent effect on development,  
326 growth or circulating IGF-I levels [38]. Whether the same might be true for exon 0 in humans is not  
327 clear, but its utilization in many prostate tumours does suggest that it may play a role in human disease.

328

329

## 330 **5. Funding**

331

332 This research did not receive any specific grant from funding agencies in the public, commercial, or  
333 not-for-profit sectors.

334

## 335 **6. Declaration of interest**

336

337 None.

338

339 **Appendix A. Supplementary data**

340

341 Supplementary data associated with this article can be found in the online version, at

342

343

344 **References**

345

346 [1] W.H. Daughaday, P. Rotwein, Insulin-like growth factors I and II. Peptide, messenger ribonucleic  
347 acid and gene structures, serum, and tissue concentrations, *Endocr. Rev.* 10 (1989) 68-91.

348

349 [2] A. Philippou, M. Maridaki, S. Pneumaticos, M. Koutsilieris, The complexity of the *IGF1* gene  
350 splicing, posttranslational modification and bioactivity, *Mol. Med.* 20 (2014) 202-214.

351

352 [3] G. Vassilakos, E.R. Barton, Insulin-like growth factor 1 regulation and its actions in skeletal  
353 muscle, *Compr. Physiol.* 9 (2019) 413-438.

354

355 [4] P. Rotwein, Diversification of the insulin-like growth factor 1 gene in mammals, *PLoS ONE* 12  
356 (2017) e0189642.

357

358 [5] P. Rotwein, Variation in the insulin-like growth factor 1 gene in primates, *Endocrinology* 158  
359 (2017) 804-814.

360

361 [6] M. Wallis, New insulin-like growth factor (IGF)-precursor sequences from mammalian genomes:  
362 the molecular evolution of IGFs and associated peptides in primates, *Growth Horm. IGF Res.* 19  
363 (2009) 12-23.

364

365 [7] W.D. Salmon, W.H. Daughaday, A hormonally controlled serum factor which stimulates sulfate  
366 incorporation by cartilage in vitro, *J. Lab. Clin. Med.* 49 (1957) 825–836.

367

368 [8] D. LeRoith, C. Bondy, S. Yakar, J.-L. Liu, A. Butler, The somatomedin hypothesis: 2001,  
369 *Endocr. Rev.* 22 (2001) 53-74.

370

371 [9] K. Sjögren, J.-L. Liu, K. Blad, S. Skrtic, O. Vidal, V. Wallenius, D. LeRoith, J. Törnell, O.G.  
372 Isaksson, J.-O. Jansson, C. Ohlsson, Liver-derived insulin-like growth factor I (IGF-I) is the  
373 principal source of IGF-I in blood but is not required for postnatal body growth in mice, *Proc.*  
374 *Natl. Acad. Sci. U.S.A.* 96 (1999) 7088–7092.

375

376 [10] S. Yakar, J.-L. Liu, B. Stannard, A. Butler, D. Accili, B. Sauer, D. LeRoith, Normal growth and  
377 development in the absence of hepatic insulin-like growth factor I, *Proc. Natl. Acad. Sci. U.S.A.*  
378 96 (1999) 7324–7329.

379

380 [11] P. Rotwein, Insulinlike growth factor 1 gene variation in vertebrates, *Endocrinology* 159 (2018)  
381 2288-2305.

382

383 [12] P. Rotwein, R.J. Folz, J.I. Gordon, Biosynthesis of human insulin-like growth factor I (IGF-I).  
384 The primary translation product of IGF-I mRNA contains an unusual 48-amino acid signal  
385 peptide, *J. Biol. Chem.* 262 (1987) 11807-11812.

386

387 [13] H. Yang, M.L. Adamo, A.P. Koval, M.C. McGuinness, H. Ben-Hur, Y. Yang, D. LeRoith, C.T.  
388 Roberts, Alternative leader sequences in insulin-like growth factor I mRNAs modulate  
389 translational efficiency and encode multiple signal peptides, *Mol. Endocrinol.* 9 (1995) 1380-  
390 1395.

391

392 [14] M.L. Adamo, H. Ben-Hur, D. LeRoith, C.T. Roberts, Transcription initiation in the two leader  
393 exons of the rat IGF-I gene occurs from disperse versus localized sites, *Biochem. Biophys. Res.*  
394 *Commun.* 176 (1991) 887-893.

395

396 [15] E. Jansen, P.H. Steenbergh, D. LeRoith, C.T. Roberts, J.S. Sussenbach, Identification of multiple

- 397 transcription start sites in the human insulin-like growth factor-I gene, *Mol. Cell. Endocrinol.* 78  
398 (1991) 115-125.
- 399
- 400 [16] A. Shimatsu, P. Rotwein, Sequence of two insulin-like growth factor I mRNAs differing within  
401 the 5' untranslated region, *Nuc. Acids Res.* 15 (1987) 7196.
- 402
- 403 [17] M.L. Adamo, H. Ben-Hur, C.T. Roberts, D. LeRoith, Regulation of start site usage in the leader  
404 exons of the rat insulin-like growth factor-I gene by development, fasting, and diabetes, *Mol.*  
405 *Endocrinol.* 5 (1991) 1677-1686.
- 406
- 407 [18] E. Jansen, P.H. Steenbergh, F.M.A. van Schaik, J.S. Sussenbach, The human IGF-I gene contains  
408 two cell type-specifically regulated promoters, *Biochem. Biophys. Res. Commun.* 187 (1992)  
409 1219-1226.
- 410
- 411 [19] W.L. Lowe, C.T. Roberts, S.R. Lasky, D. LeRoith, Differential expression of alternative 5'  
412 untranslated regions in mRNAs encoding rat insulin-like growth factor I, *Proc. Natl. Acad. Sci.*  
413 *U.S.A.* 84 (1987) 8946-8950.
- 414
- 415 [20] C.A. West, T.R. Arnett, S.M. Farrow, Expression of insulin-like growth factor I (IGF-I) mRNA  
416 variants in rat bone, *Bone* 19 (1996) 41-46.
- 417
- 418 [21] D.J. Chia, J.J. Young, A.R. Mertens, P. Rotwein, Distinct alterations in chromatin organization of  
419 the two IGF-I promoters precede growth hormone-induced activation of IGF-I gene transcription,  
420 *Mol. Endocrinol.* 24 (2010) 779-789.
- 421
- 422 [22] H.L. Foyt, F. Lanau, M. Woloschak, D. LeRoith, C.T. Roberts, Effect of growth hormone on  
423 levels of differentially processed insulin-like growth factor I mRNAs in total and polysomal  
424 mRNA populations, *Mol. Endocrinol.* 6 (1992) 1881-1888.
- 425
- 426 [23] J.M. Pell, J.C. Saunders, R.S. Gilmour, Differential regulation of transcription initiation from



- 427 insulin-like growth factor-I (IGF-I) leader exons and of tissue IGF-I expression in response to  
428 changed growth hormone and nutritional status in sheep, *Endocrinology* 132 (1993) 1797-1807.  
429
- 430 [24] J.G. Simmons, J.J. van Wyk, E.C. Hoyt, P.K.Lund, Multiple transcription start sites in the rat  
431 insulin-like growth factor-I gene give rise to IGF-I messenger-RNAs that encode different IGF-I  
432 precursors and are processed differently in-vitro, *Growth Factors* 9 (1993) 205-221.  
433
- 434 [25] S.F. Altschul, W. Gish, W. Miller, E.W. Myers, D.J. Lipman, Basic Local Alignment Tool, *J. Mol.*  
435 *Biol.* 215 (1990) 403-410.  
436
- 437 [26] A.F.A. Smit, R. Hubley, P. Green, *RepeatMasker Open-4.0*, (2013-2015)  
438 <<http://www.repeatmasker.org>>.  
439
- 440 [27] M. Kozak, At least six nucleotides preceding the AUG initiator codon enhance translation in  
441 mammalian cells, *J. Mol. Biol* 196 (1987) 947-950.  
442
- 443 [28] L. Fagerberg, B.M. Hallström, P. Oksvold, C. Kampf, D. Djureinovic, J. Odeberg, M. Habuka, S.  
444 Tahmasebpour, A. Danielsson, K. Edlund, A. Asplund, E. Sjöstedt, E. Lundberg, C.A-K.  
445 Szigyarto, M. Skogs, J.O. Takanen, H. Berling, H. Tegel, J. Mulder, P. Nilsson, J.M. Schwenk, C.  
446 Lindskog, F. Danielsson, A. Mardinoglu, Å. Sivertsson, K. von Feilitzen, M. Forsberg, M.  
447 Zwahlen, I. Olsson, S. Navani, M. Huss, J. Nielsen, F. Ponten, M. Uhlén, Analysis of the human  
448 tissue-specific expression by genome-wide integration of transcriptomics and antibody-based  
449 proteomics, *Mol. Cell. Proteomics* 13 (2014) 397-406.  
450
- 451 [29] M. Wallis, Coevolution of insulin-like growth factors, insulin and their receptors and binding  
452 proteins in New World Monkeys, *Growth Horm. IGF Res.* 25 (2015) 158-167.  
453
- 454 [30] A.W. Wyatt, F. Mo, K. Wang, B. McConeghy, S. Brahmabhatt, L. Jong, D.M. Mitchell, R.L.  
455 Johnston, A. Haegert, E. Li, J. Liew, J. Yeung, R. Shrestha, A.V. Lapuk, A. McPherson, R.  
456 Shukin, R.H. Bell, S. Anderson, J. Bishop, A. Hurtado-Coll, H. Xiao, A.M. Chinnaiyan, R.

- 457 Mehra, D. Lin, Y. Wang, L. Fazli, M.E. Gleave, S.V. Volik, C.C. Collins, Heterogeneity in the  
458 inter-tumor transcriptome of high risk prostate cancer, *Genome Biol.* 15 (2014) art. no. 426.  
459
- 460 [31] P.F. Christopoulos, P. Msaouel, M. Koutsilieris, The role of the insulin-like growth factor-1  
461 system in breast cancer, *Mol. Cancer* 14 (2015) art. no. 43.  
462
- 463 [32] E.J. Gallagher, D. LeRoith, The proliferating role of insulin and insulin-like growth factors in  
464 cancer, *Trends Endocrinol. Metab.* 21 (2010) 610-618.  
465
- 466 [33] E.J. Gallagher, D. LeRoith, Minireview:IGF, insulin and cancer, *Endocrinology* 152 (2011) 2546-  
467 2551.  
468
- 469 [34] C. Gennigens, C. Menetrier-Caux, J.P. Droz, Insulin-like growth factor (IGF) family and prostate  
470 cancer, *Crit. Rev. Oncol. Hematol.* 58 (2006) 124-145.  
471
- 472 [35] C. Kang, D. LeRoith, E.J. Gallagher, Diabetes, obesity, and breast cancer, *Endocrinology* 159  
473 (2018) 3801-3812.  
474
- 475 [36] A. Kasprzak, W. Kwasniewski, A. Adamek, A. Gozdzicka-Jozefiak, Insulin-like growth factor  
476 (IGF) axis in cancerogenesis, *Mut. Res.* 772 (2017) 78-104.  
477
- 478 [37] P.G. Vigneri, E.Tirro, M.S. Pennisi, M. Massimino, S. Stella, C. Romano, L. Manzella, The  
479 insulin/IGF system in colorectal cancer development and resistance to therapy, *Front. Oncol.* 5  
480 (2015) art. no. 230.  
481
- 482 [38] L. Temmerman, E. Slonimsky, N. Rosenthal, Class 2 IGF-1 isoforms are dispensible for viability,  
483 growth and maintenance of IGF-1 serum levels, *Growth Horm. IGF Res.* 20 (2010) 255-263.  
484  
485

486 **Legends for Figures**

487

488 **Fig. 1.** The human IGF-I gene. A. Organization of the gene showing the 6 previously-recognized exons  
 489 (solid rectangles), the novel exon 0 (open rectangle), and the possible alternative splices at the 5' end,  
 490 the minor exon 0-exon 1-exon 3 variant being shown with dotted lines. Introns 3 and 5 are shown  
 491 truncated; their full lengths are approx. 55.9 kb (intron 3) and 14.9 kb (intron 5). B. Sequence of exon  
 492 0. The start of the exon at the 5' end is poorly defined and there may be multiple transcription start  
 493 points; the region under the dashes is represented in few transcripts. Flanking sequences are shown in  
 494 light grey. C. The potential signal peptide derived from exon 0 spliced to exon 3. Solid arrows indicate  
 495 possible translation initiation sites; which of these is actually used is not clear (see text). The N-  
 496 terminal end of mature IGF-I is also included; the broken arrow indicates the cleavage site for the  
 497 signal peptide.

498

499 **Fig. 2.** Exon profile for the 5' end of IGF-I gene in human liver transcriptome (SRA project  
 500 ERP003613; 5 experiments combined). Arrows indicate positions of in-frame translation initiation sites  
 501 (ATG/AUG).

502

503 **Fig. 3.** Utilization of exons 0, 1 and 2 in human liver. Values are means  $\pm$  SEM for 11 individuals from  
 504 5 SRA projects.

505

506 **Fig. 4.** Expression level and 5' exon utilization for IGF-I in human tissues. Based on data in SRA  
 507 projects ERP003613 and ERP006650. A. Expression level (BLAST hits/ $10^6$  reads), based on BLAST  
 508 searching with the sequence encoding mature human IGF-I. B. Utilization of exons 0, 1 and 2 in each  
 509 tissue, expressed as a percentage of all splices to exon 3. The exon 0-exon 3 and exon 0-exon 1-exon 3  
 510 splices were combined. Error bars show S.E.M. for 2-7 individuals.

511

512 **Fig. 5.** Sequences of IGF-I exon 0 region for a number of primates and other mammals. The human  
 513 sequence is given in full. For other sequences, a dot (.) indicates identity to human and a dash (-)  
 514 indicates a gap. The exon 0 region is included in the large box, which is open at the left hand side  
 515 because the 5' end is not defined. ATG-1 and ATG-2 indicate the positions of two potential translation  
 516 start sites, in-frame in some sequences (see text). a and b indicate the 5' limits seen in the human

517 normal liver transcriptomes, b being reached in all 5 projects studied (including that illustrated in Fig.  
 518 2), a in only two projects. Abbreviations of species names are as follows: Hsa, *Homo sapiens* (human),  
 519 Ptr, *Pan troglodytes* (chimpanzee), Ggo, *Gorilla gorilla* (gorilla), Ppy, *Pongo pygmaeus* (orangutan),  
 520 Nle, *Nomascus leucogenys* (gibbon), Mmu, *Macaca mulatta* (rhesus macaque), Csa, *Chlorocebus*  
 521 *sabaeus* (African green monkey), Cja, *Callithrix jacchus* (marmoset), Sbo, *Saimiri boliviensis* (squirrel  
 522 monkey), Tsy, *Tarsius syrichta* (tarsier), Mmur, *Microcebus murinus* (mouse lemur), Oga, *Otolemur*  
 523 *garnettii* (galago), Tbe, *Tupaia belangeri* (tree shrew), Ocu, *Oryctolagus cuniculus* (rabbit), Cfa, *Canis*  
 524 *familiaris* (dog).

525

526 **Fig. 6.** Expression level and 5' exon utilization for IGF-I in liver of various primates plus tree shrew,  
 527 rat, rabbit and dog. A. Expression level (BLAST hits/ $10^6$  reads), based on BLAST searching with the  
 528 sequence encoding mature IGF-I for the appropriate species. Numbers in brackets indicate the number  
 529 of individuals in the sample. B. Utilization of exons 0 (white), 1 (grey) and 2 (black) in each tissue,  
 530 expressed as a percentage of all splices to exon 3. The exon 0-exon 3 and exon 0-exon 1-exon 3 spliced  
 531 forms were combined. Error bars show  $\pm$  S.E.M. except where  $n < 3$ , where individual values are shown  
 532 (x). Species abbreviations as for Fig. 5 plus Rno, *Rattus norvegicus* (rat).

533

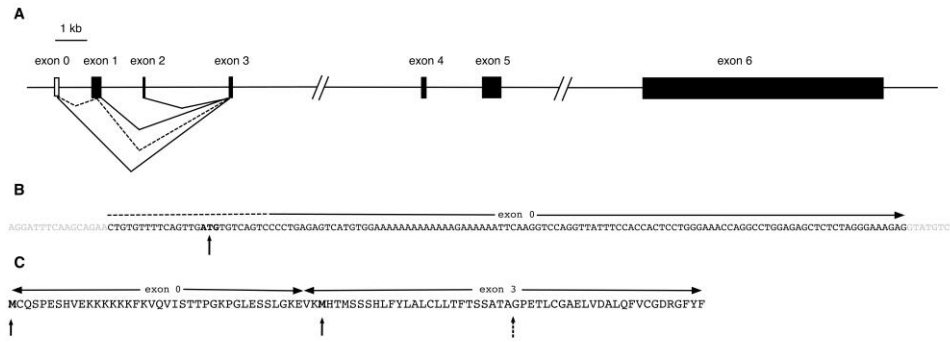
534 **Fig. 7.** Expression level and 5' exon utilization for IGF-I in prostate cancers from 25 individual  
 535 patients. Based on data in SRA project ERP006077. A. Expression level (BLAST hits/ $10^6$  reads), based  
 536 on BLAST searching with the sequence encoding mature IGF-I. B. Utilization of exons 0 (light grey),  
 537 1 (grey) and 2 (black) in each tumour, expressed as a percentage of all splices to exon 3. The exon 0-  
 538 exon 3 and exon 0-exon 1-exon 3 spliced forms were combined. Error bars show  $\pm$  S.E.M. except  
 539 where  $n < 3$ .  $n = 3$ , except where individual points are shown (x,  $n = 1$  or 2).

540

541 **Fig 8.** Expression level and 5' exon utilization for IGF-I in liver cancers from 4 individual patients.  
 542 Based on data in SRA project SRP064431. A. Expression level (BLAST hits/ $10^6$  reads), based on  
 543 BLAST searching with the sequence encoding mature human IGF-I for each tumour (T) or adjacent  
 544 apparently normal tissue (N). B. Utilization of exons 0 (white), 1 (grey) and 2 (black) in each tumour  
 545 (T) or adjacent apparently normal tissue (N), expressed as a percentage of all splices to exon 3. The  
 546 exon 0-exon 3 and exon 0-exon 1-exon 3 spliced forms were combined.

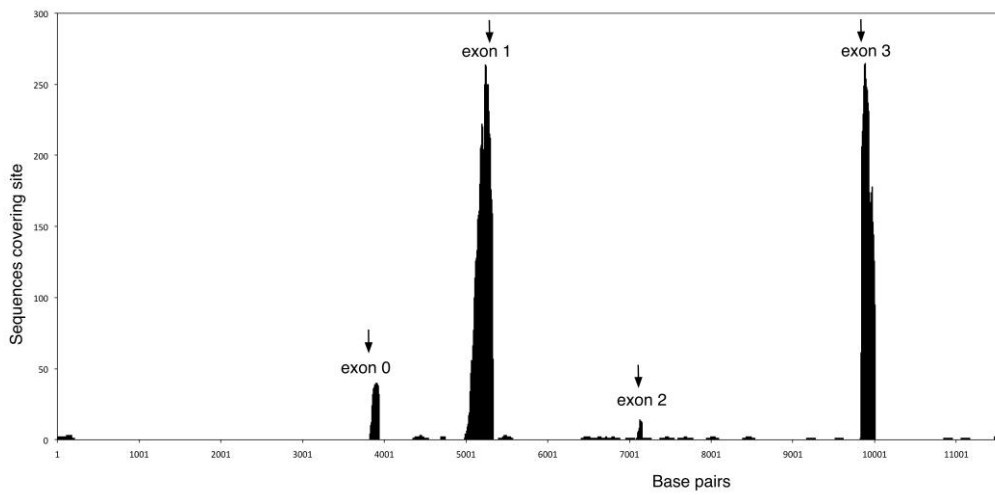
547

548 Fig. 1.



549

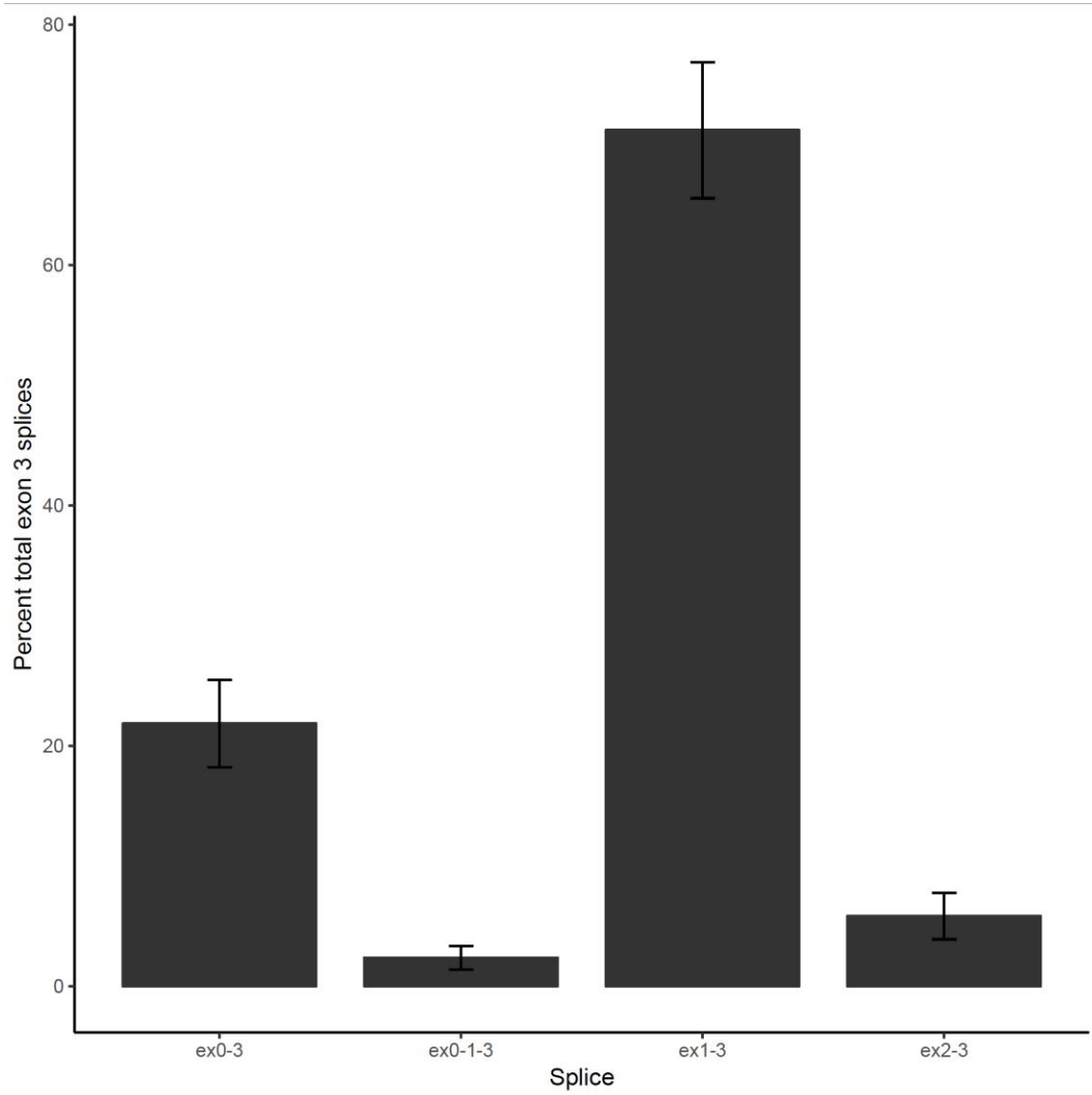
550 Fig. 2.



551

552

553 Fig. 3.

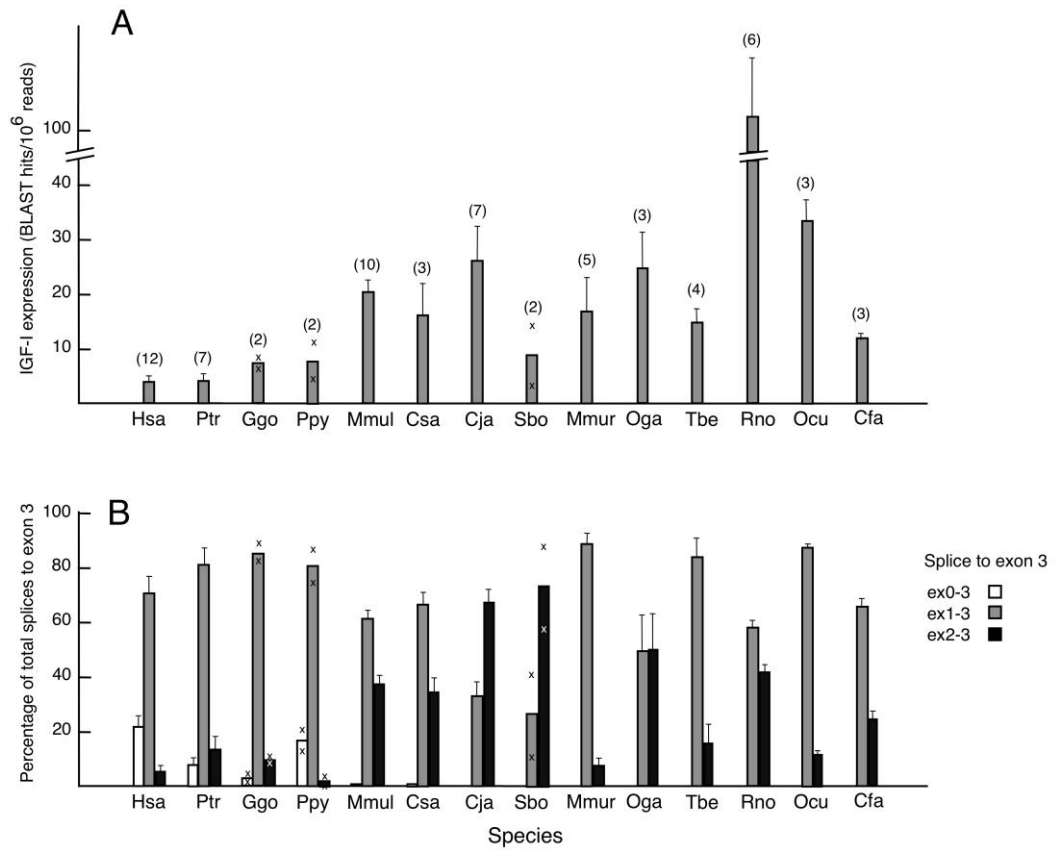


554

555

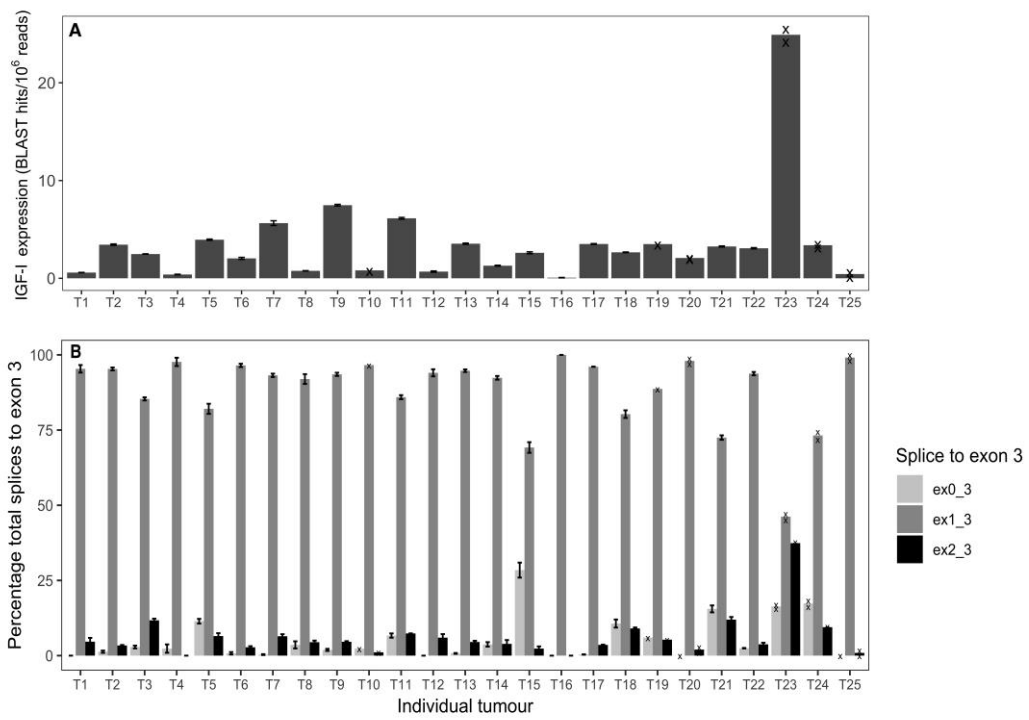


561 Fig. 6.



562

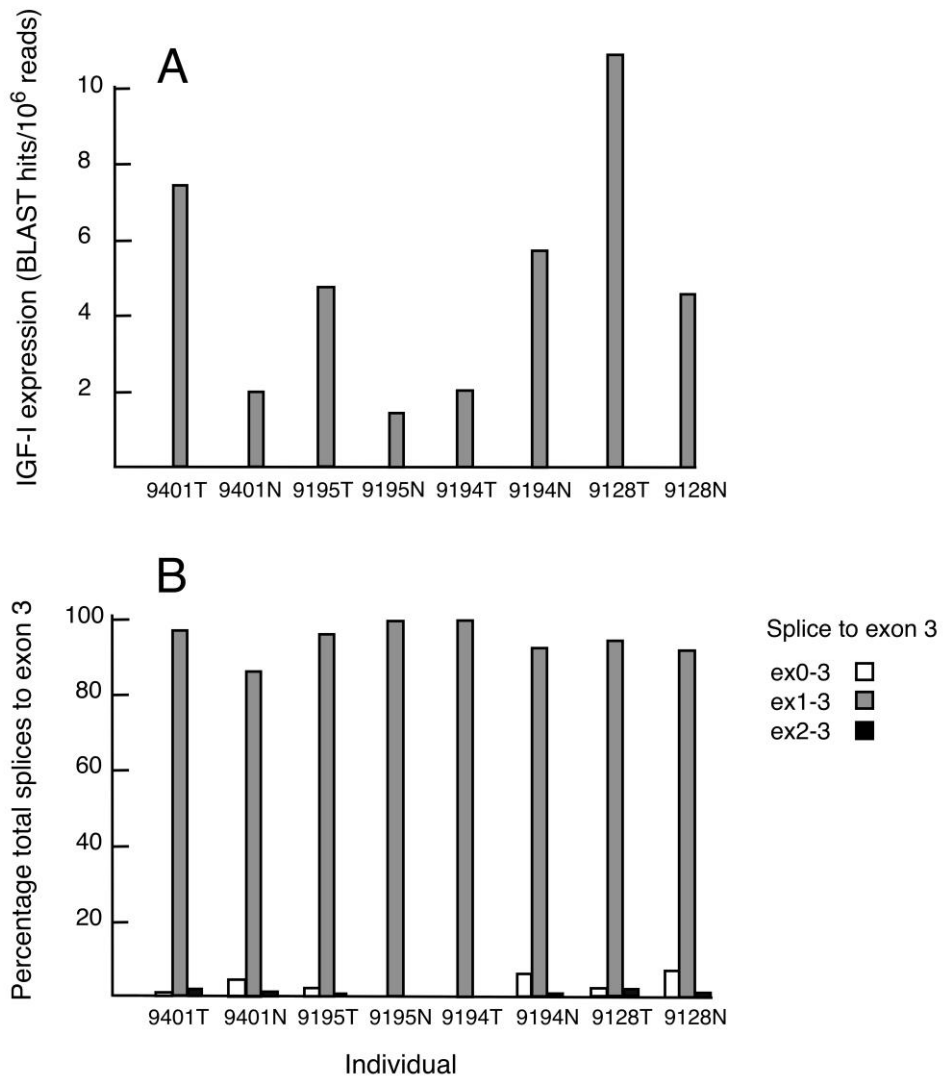
563 Fig. 7.



564



565 Fig. 8.



567 Fig. 8.

568