

1 **DOI: [10.1016/j.bbagrm.2016.03.014](https://doi.org/10.1016/j.bbagrm.2016.03.014)**

2  
3  
4  
5 **REVISED MANUSCRIPT (text with changes Marked)**

6  
7  
8  
9 **MIR retroposon exonization promotes evolutionary variability and generates species-specific**  
10 **expression of IGF-1 splice variants**

11  
12  
13 Giosuè Annibalini<sup>1\*†</sup>, Pamela Bielli<sup>2†</sup>, Mauro De Santi<sup>1</sup>, Deborah Agostini<sup>1</sup>, Michele Guescini<sup>1</sup>, Davide Sisti<sup>1</sup>,  
14 Serena Contarelli<sup>1</sup>, Giorgio Brandi<sup>1</sup>, ~~Claudio Sette~~<sup>2</sup>, Anna Villarini<sup>3</sup>, Vilberto Stocchi<sup>1</sup>, ~~Claudio Sette~~<sup>2</sup> and  
15 Elena Barbieri<sup>1^</sup>

16  
17 <sup>1</sup>Department of Biomolecular Sciences, University of Urbino Carlo Bo, 61029 Urbino, Italy.

18 <sup>2</sup>Department of Biomedicine and Prevention, University of Rome Tor Vergata, 00133 Rome, Italy; Laboratory  
19 of Neuroembryology, Fondazione Santa Lucia, 00133 Rome, Italy.

20 <sup>3</sup>Department of Preventive & Predictive Medicine, Fondazione IRCCS Istituto Nazionale dei Tumori, 20133  
21 Milan, Italy.

22 <sup>^</sup>IIM, Interuniversity Institute of Myology

23  
24 \* To whom correspondence should be addressed. Tel:+39 0722-303418 Fax:+39 0722-303401; Email:  
25 giosue.annibalini@uniurb.it

26 †Giosuè Annibalini and Pamela Bielli equally contributed to this work.

27

28

1 **Abstract**

2 Insulin-like growth factor-1 (IGF-1) is a pleiotropic hormone exerting mitogenic and anti-apoptotic effects.  
3 Inclusion or exclusion of exon 5 into the IGF-1 mRNA gives rise to three transcripts, IGF-1Ea, IGF-1Eb and  
4 IGF-1Ec, which yield three different C-terminal extensions called Ea, Eb and Ec peptides. The biological  
5 significance of **the** IGF-1 splice variants and how the E-peptides affect the actions of mature IGF-1 are  
6 largely unknown. In this study we investigated the origin and conservation of the IGF-1 E-peptides and we  
7 compared the **pattern of expression of the** IGF-1 isoforms *in vivo*, in nine mammalian species, and *in vitro*  
8 using human and mouse IGF-1 minigenes.

9 Our analysis showed that only IGF-1Ea is conserved among all vertebrates, whereas IGF-1Eb and IGF-1Ec  
10 are an evolutionary novelty originated from the exonization of a mammalian interspersed repetitive-b (MIR-b)  
11 element. Both IGF-1Eb and IGF-1Ec mRNAs were constitutively expressed in all mammalian species  
12 analyzed but their expression ratio varies greatly **among** species. Using IGF-1 minigenes we demonstrated  
13 that divergences in *cis*-acting regulatory elements between human and mouse conferred species-specific  
14 features to the exon 5 region. Finally, the protein-coding **sequences** of ~~the~~ exon 5 showed low rate of  
15 synonymous mutations and contain disorder-promoting amino acids, suggesting a regulatory role for these  
16 domains.

17 In conclusion, ~~the~~ exonization of **a** MIR-b element in the *IGF-1* gene determined **gain of the** exon 5 **gain**  
18 during mammalian evolution. Alternative splicing of this novel exon added new regulatory elements at the  
19 mRNA and protein level potentially able to regulate the mature IGF-1 across tissues and species.

20

21

22 **Keywords**

23 IGF-1 isoforms, alternative splicing, retroposon exonization, synonymous sites, intrinsically disordered  
24 regions

25

26

27 **1. Introduction**

28 The mammalian insulin-like growth factor (*IGF*) -1 gene is a single copy gene composed of six exons and  
29 five introns, which gives rise to an immature IGF-1 peptide (IGF-1 pro-peptide) containing a signal peptide at  
30 the 5' end of the gene, a core region and an E-peptide at the 3' end (**Fig.1 Figs. 1A and 1B**). Both signal  
31 peptide and E-peptide are then removed by protease cleavage to form the 70 amino acid-long mature IGF-1  
32 peptide (**IGF-1 core**), which displays growth-promoting and metabolic functions.

33 Four of the six *IGF-1* exons are subjected to alternative splicing (**Fig.1 Figs. 1A and 1B**) [1]. **Exon 1 and exon**  
34 **2 are mutually exclusive first exons and generate different signal peptides. Transcripts containing exon 1 are**  
35 **referred to as Class I transcripts whereas those containing exon 2 are referred to as Class II transcripts. The**

1 Class II IGF-1 knockout mice indicated that class II isoforms are dispensable for fetal and postnatal growth,  
2 and the significance of the alternate signal peptides encoded by the first two IGF-1 exons was discussed in  
3 [2-3]. At the 3' end of the gene, alternative splicing ~~the 3'-end of the pre-mRNA~~ yields three different mRNA  
4 transcripts, each encoding distinct carboxyl-terminal portions of E-peptide followed by the 3'-untranslated  
5 region (3'UTR) (~~Fig. 4~~ Figs. 1A and 1B). Thus, although alternative splicing generates different precursor  
6 peptides, it does not alter the ~~sequence structure~~ of the mature IGF-1 peptide.

7 Splicing of exon 4 with exon 6 yields the most common IGF-1Ea variant, which encodes the 35 amino acid-  
8 long Ea peptide. The first 16 amino acids of the Ea peptide are encoded by exon 4 and are common in all ~~of~~  
9 ~~the three different~~ E-peptides, whereas the remaining 19 amino acids are encoded by exon 6 ~~and are unique~~  
10 ~~to this isoform~~. The IGF-1Eb variant is produced when exon 4 is spliced with exon 5 and encodes the Eb  
11 peptide, which contains the 16 common amino acids and 61 additional amino acids encoded by exon 5. ~~IGF-~~  
12 ~~1Eb transcript terminates in exon 5 and excludes exon 6 from the mRNA, hence it has a completely different~~  
13 ~~3'UTR compared to IGF-1Ea and IGF-1Ec~~. The third variant, named IGF-1Ec in humans, is generated by  
14 usage of a cryptic 5' splice site (c5'ss) (named IGF633) [4] present in exon 5, which is in turn spliced with  
15 exon 6. The ~~cryptic 5'ss c5'ss~~ of *IGF-1* exon 5 deviates from the vertebrate consensus and is commonly  
16 used in rodents and rabbits [5-7] but rarely and in a tissue-specific manner in human [1, 4]. Notably, although  
17 this variant is named IGF-1Eb in rodents, for clarity we will use ~~throughout this manuscript~~ the human IGF-  
18 1Ec nomenclature ~~throughout this manuscript~~, regardless of the species. The human Ec peptide has a  
19 predicted length of 40 amino acids, with the common 16 amino acids deriving from the exon 4, 16 from exon  
20 5 and the last 8 amino acids from exon 6. Since expression of IGF-1Ec was linked to mechanical stimuli in  
21 muscle and other mechanosensitive cells, this variant is sometime referred to as Mechano Growth Factor  
22 (MGF) [1].

23 Few studies have addressed the regulation of *IGF-1* alternative splicing. In human cells, splicing of exon 5 is  
24 regulated by the antagonistic activities of the serine-arginine protein splicing factor-1 (SRSF1) and  
25 heterogeneous nuclear ribonucleoproteins A1 (hnRNPA1) [8-9]. SRSF1 was proposed to increase splicing of  
26 the IGF-1Eb variant by binding to a purine-rich exonic splicing enhancer (ESE) in exon 5 and preventing the  
27 recruitment of hnRNPA1, which functions as a splicing repressor [8-9]. Notably, neither splicing factor was  
28 capable of regulating the IGF-1Ec variant in human IGF-1 minigenes [8-9]. The reason of the lack of usage  
29 of the ~~cryptic 5'ss c5'ss~~, and hence for the lack of IGF-1Ec isoform production, ~~reported~~ in these studies, has  
30 yet to be clarified. Similarly, whether or not the mouse ~~cryptic 5'ss c5'ss~~, which is closer to the canonical  
31 splicing donor consensus, can be recognized in these experimental settings has not been tested.

32 Furthermore, it is currently unknown whether SRSF1 and hnRNP A1 function specifically ~~on~~ with the human  
33 IGF-1 pre-mRNA or whether they also play a role in *IGF-1* alternative splicing in other species.

34 To date, the fate and the biological functions of ~~the~~ E-peptides are not entirely clear [1, 10-12]. Interestingly,  
35 gene structure comparison showed that ~~the~~ *IGF-1* exon 6 is conserved in all vertebrates, whereas exon 5 is  
36 conserved only among mammals [13-15]. In addition, comparative studies on mammalian *IGF-1* show that  
37 both the IGF-1 core and Ea peptide are subjected to strong purifying selection, whereas sequences of the Eb

1 and Ec peptides are more variable [16-17]. In particular it was shown that the ratio of non-synonymous (dN;  
2 amino-acid altering) to synonymous (dS; silent) substitutions is about 10-fold higher for Eb and Ec peptides  
3 compared to Ea peptide, suggesting that ~~the Eb and Ec peptides might have little~~ they have no specific roles  
4 or, at most, species-specific functions [17]. Indeed, it is usual to view poor protein sequence conservation,  
5 i.e. evolving under neutral or nearly neutral conditions, as evidence of reduced functional importance [17-18].  
6 However, Lin and colleagues [19] recently demonstrated that the “dual-coding” DNA sequences of *IGF-1*  
7 exon 6, which is translated in two alternative reading frame to give rise to Ea and Ec peptides, showed a  
8 very low dS rate suggesting additional sequence constraints beyond those dictated by the amino acid  
9 sequence of the Ea and Ec peptides. Notably, the picture of largely neutral evolution of synonymous  
10 substitutions in mammals has been recently challenged [20], showing that selection may constrain not only  
11 dN<sub>r</sub> to preserve amino acid sequence, but also dS<sub>r</sub> to preserve regulatory elements in nucleotide sequences,  
12 such as ESE, RNA secondary structures and microRNA target sites [19-21]. However, whether the dS drop  
13 of the dual-coding exon 6 has contributed to the observed local increase in dN/dS ratio of E-peptides and  
14 hence to an inaccurate estimate of evolutionary rate is still unclear.

15 In this study, we have analyzed the conservation and the mRNA expression pattern of IGF-1 isoforms in  
16 different mammalian species and evaluated the evolutionary pressure within the different coding regions of  
17 the mammalian *IGF-1* gene. Moreover, we have used human and mouse minigene systems to compare the  
18 mechanisms regulating the *IGF-1* exon 5 splicing between these two species.

## 21 **2. Materials and Methods**

### 22 **2.1 Sequences and databases**

23 Orthologous of the human *IGF-1* gene were obtained from the UCSC Genome Browser Database  
24 (<http://genome.ucsc.edu>; Feb. 2009 GRCh37/hg19) and Ensembl website (<http://www.ensembl.org>). The  
25 “CDS FASTA alignment from multiple alignments” data, derived from the “multiz100way” alignment data  
26 prepared from 100 vertebrate genomes [22], were downloaded using the Table Browser tool of the UCSC  
27 Genome Browser. Sequences were subsequently realigned using MUSCLE [23] and protein coding  
28 sequences from 27 mammalian species were extracted from these alignment datasets. Supplementary file  
29 S1 contains the IGF-1 sequence of 27 mammalian species in FASTA format and the neighbor-joining tree  
30 used in the present study. Percent nucleotide and amino acid identities between sequences were computed  
31 using the CLUSTALW procedure [24]. The splice site scores were obtained for each 5' and 3' ss sequence  
32 using the Maximum Entropy scores [25]. For Transposable elements analysis, we used RepeatMasker  
33 (<http://www.repeatmasker.org>) and Repbase annotations [26-28].

### 34 **2.2 Human and animal tissues**

35 Freshly frozen normal human liver (2 males, mean age 54 +/- 9 years), adipose (2 males and 2 females,  
36 mean age 52 +/- 12 years) and muscle (3 males, mean age 25 +/- 6 years) samples were provided by the

1 complex structure of biomarkers (DOSMM) of National Cancer Institute of Milan. The macaque (*Macaca*  
2 *mulatta*) autoptic specimens were kindly provided by Elena Borra (Department of Neuroscience, University of  
3 Parma, Parma, Italy). Three 4-year olds macaques were tested in this study (2 female and 1 male). One-  
4 month old CD1 female mice (n=3) (*Mus Musculus*), 2-month old Young Sprague Dawley male rat (n=3)  
5 (*Rattus norvegicus*) and three 8 month-old New Zealand male rabbits (*Oryctolagus cuniculus*) (Charles River  
6 Laboratories, Milan, Italy) were housed with a 12-h light/dark cycle and free access to standard laboratory  
7 chow and water. Care and handling were in accordance with the *Guide for the Care and Use of Laboratory*  
8 *Animals* by Ministero della Sanità D.L. 116 (1992) and approved by the university committee for animal  
9 experiments. Animals were sacrificed with an overdose of anaesthetics (ketamine in combination with  
10 xylazine). Other tissues used in this study were collected at local slaughter during routine meat inspection:  
11 pig (*Sus scrofa*) (3 males), cow (*Bos taurus*) (3 females), sheep (*Ovis aries*) (3 female) and goat (*Capra*  
12 *hircus*) (2 males and 1 female) were tested in this study. The age of the animals ranged from 2 to 5 years. All  
13 tissues (about 30 mg) were immediately submerged in RNAlater stabilization solution (Qiagen, Milan, Italy)  
14 and left at least 10 min at room temperature. Then, tissues were stored at -80°C until RNA extraction.

15

### 16 **2.3 RNA extraction and cDNA synthesis**

17 The tissues were removed from RNAlater and transferred into a clean Eppendorf tube where the total RNA  
18 was extracted and purified using the Omega Bio-Tek E.Z.N.A.TM Total RNA kit (VWR International s.r.l.,  
19 Milan, Italy) according to the manufacturer's instructions. The amount and quality of RNA were assessed  
20 with DU-640 UV Spectrophotometer (Beckman Coulter). After DNA digestion with DNase I enzyme (Qiagen,  
21 Milan, Italy) complementary DNA was synthesized from 1 µg of total RNA using Omniscript RT (Qiagen,  
22 Milan, Italy) and random hexamers or anchored oligo-dT primers where specified.

23

### 24 **2.4 Real time PCR quantification of human, macaque and mouse IGF-1 isoforms**

25 Serial dilution (1:4) of three recombinant plasmids containing the IGF-1Ea, IGF-1Eb and IGF-1Ec sequences  
26 of human and mouse were prepared in order to generate standard curves for plotting CT values against  
27 number of molecules. Molecules of each plasmid were calculated using the concentration of each plasmid,  
28 Avogadro's constant, the molecular weight of double-stranded DNA and the size of the target amplicon [29].  
29 The copy number of genes was calculated in individual samples using a corresponding reference plasmid  
30 cDNA clone at known concentration. The amount of target transcripts was normalized to glyceraldehyde-3-  
31 phosphate dehydrogenase (GAPDH). Percentage of IGF-1 isoforms was calculated as [(mRNA copy number  
32 of single isoform/(mRNA copy numbers of IGF1Ea+IGF1Eb+IGF1Ec)\*100]. Expression of IGF-1 isoforms  
33 ~~was were~~ compared using 2-ways analysis of variance with interactions: species and IGF-1 isoforms were  
34 used as predictive factors. In order to meet assumption of homoscedasticity percentage were arcsin radq  
35 transformed. Post-hoc analysis was performed using Bonferroni correction.

36 Real-time quantitative PCR was performed with two microliters of cDNA and 300 nM of of each primer in an  
37 Applied Biosystems StepOnePlus™ Real Time PCR System using SYBR Select Master Mix (Applied

1 Biosystems, Monza, Italy). The real-time PCR conditions were: 50°C for 2 min, 95°C for 2 min followed by 40  
2 cycles of three-steps at 95°C for 15 sec, 60°C for 15 sec and 72°C for 30 sec. The specificity of the  
3 amplification products was confirmed by examining thermal denaturation plots and by sample separation in a  
4 4% DNA agarose gel. The sequences and the annealing positions of the primers used to quantify the IGF-1  
5 isoforms were shown in Supplementary Table S1.

## 6 7 **2.5 3' RACE-PCR**

8 The mRNA splicing pattern of IGF-1 was assayed by 3-RACE-PCR using an anchored oligo-dT primer (5'-  
9 AAGCAGTGGTATCAACGCAGAGTACT<sub>(30)</sub>NV-3') as reverse transcription primer. 1 µg of total liver RNA  
10 was reverse transcribed into cDNA. To amplify the 3'-end IGF-1 variants a reverse primer corresponding to  
11 the anchor sequence of the RT primers was used in combination with the following forward primers: human,  
12 macaque, rabbit, (5'-CCTCCTCGCATCTCTTCTACCTG-3'); mouse and rat (5'-  
13 GCTATGGCTCCAGCATTCG-3'); and pig (5'-CGTGGATGAGTGCTGCTTC-3'). The PCR reaction were as  
14 follows: 95° for 10 min; 35 cycles of three-steps at 95°C for 30 sec, 60°C for 30 sec, 72°C for 1 min and 30  
15 sec followed by a final elongation cycle (72°C for 5 min). The PCR products were loaded on 4.0% agarose  
16 gel and DNA fragments were eluted from gel, subcloned and sequenced.

## 17 18 **2.6 Conventional RT-PCR**

19 RT-PCR was performed in 50 µl of reaction volume with 4 µl of cDNA, 800 nM of primers and 25 µl of 2X  
20 HotStartTaq mix (Qiagen, Milan, Italy). The sequences and the annealing positions of the primers used for  
21 RT-PCR were shown in Supplementary Table S1. RT-PCR conditions involved an initial denaturation step at  
22 95°C for 10 min, followed by 35 cycles with denaturation step at 95°C for 30 sec, annealing at 60°C for 30  
23 sec and extension at 72°C for 30 sec. PCR products were size-fractionated by electrophoresis on 4.0%  
24 agarose gels and visualized by ethidium bromide staining under UV light. Amplification products were  
25 purified with QIAquick gel extraction kit (Qiagen), cloned and sequenced.

## 26 27 **2.7 Plasmid constructs**

28 The 5' and 3'end of IGF-1 mouse minigene were amplified using primers #(1,2) and #(3,4), respectively, from  
29 C57 mice genomic DNA. After enzymatic digestion, the PCR products were cloned in KpnI/Sall and Sall/NotI  
30 restriction sites of pCI vector (Promega). The 5' and 3'end of the IGF-1 human minigene were amplified  
31 using primers #(5,6) and #(7,8), respectively, from HeLa cell genomic DNA and cloned in EcoRI/Sall and  
32 Sall/NotI restriction sites of pCI vector. The mGAvsTG and hTGvsGA IGF-1 mutant minigenes were  
33 constructed using the mega-primer strategy [30]. The mouse and human IGF-1 5' mutant ends were  
34 generated using primers #(1,9,2) or #(5,10,6) respectively. All oligonucleotide sequences are listed in  
35 Supplementary Table S1. After enzymatic digestion, the PCR products were subcloned in the corresponding  
36 wild-type minigene. PCR reactions were performed using Phusion Hot Start High-Fidelity DNA polymerase  
37 (Thermo Scientific) according to manufacturer's instruction. All plasmids were sequenced and validated.

1 **2.8 Cell cultures and transfections and cell extract preparation**

2 Cell cultures, transfections and sample preparation were carried out by standard methods as previously  
3 described [31]. Briefly, ~~HEK293T and primary mouse hepatocytes~~ human and mouse cells lines were  
4 transfected with various combination of vectors as indicated using Lipofectamine 2000 (Invitrogen). Twenty-  
5 four hours after transfection, cells were collected for RNA and protein isolation, as previously described [31].  
6 For RNAi, cell were transfected twice with 60nM siRNAs using Lipofectamine RNAi Max according to  
7 manufacturer's instruction (Invitrogen). Sequences of siRNAs are listed in the Supplemental Table 1.

8

9 **2.9 Splicing assay and PCR analyses**

10 Twenty-four hours after transfection, cells were collected for RNA extraction using TRIzol (Ambion, Life  
11 Technologies) according to the manufacturer's instructions. After DNase digestion (Roche), 1 µg of total  
12 RNA was retrotranscribed using M-MLV reverse transcriptase (Promega). cDNA was used as template for  
13 conventional PCR reactions (GoTaq G2, Promega) in presence of specific primers and PCR products were  
14 analysed on agarose gel. In details, primers #(11,13) were used to amplify mIGF-1Eb isoform; primers  
15 #(11,12) to amplify mouse and human IGF-1Ea and IGF-1Ec isoforms; primers #(11,14) to amplify hIGF-  
16 1Eb; primers #(11,15) and #(11,16) to amplify total mouse and human IGF-1, respectively. In splicing assay  
17 human IGF-1Ea and IGF-1Eb were amplified using primers #(11,12,14) in the same PCR reaction. All  
18 ~~oligonucleotide sequences~~ The sequences and the annealing positions of the primers are showed listed in  
19 Supplementary Table S1.

20

21 **2.10 IGF-1 evolutionary analysis**

22 The MG94xREV\_3x4 codon model with transition/transversion bias correction and nine base frequency  
23 parameters was fitted to sequence alignments by using the HyPhy 2.2.3 package [32]. A neighbor-joining  
24 tree was built using the IGF-1 core alignment of 27 mammalian species and the tree topology along with the  
25 alignment were used as input for IGF-1Ea, IGF-1Eb and IGF-1Ec sequence analysis (Supplementary file  
26 S1). For every branch in a phylogenetic tree, the expected number of synonymous substitutions per  
27 synonymous site (dS) and its counterpart for nonsynonymous substitutions (dN) were estimated by using  
28 maximum likelihood [33] allowing for independent dN and dS values for each branch (the local parameters  
29 option). Mean and 95% confidence interval for each dS, dN and dN-dS values are calculated using bootstrap  
30 approach with n=1000 resampling [34].

31 The web application of DisCons [35] was used to predict the protein disorder on the amino acid level, using  
32 IUPred long disorder prediction parameter [36].

33

34

35

36

## 1 3 Results

2

### 3 3.1 IGF-1 exon 5 is conserved in Mammalia but not in other vertebrates

4 Previous studies demonstrated that the IGF-1Ea splice variant, which skips exon 5, is the most common E-  
5 peptide isoform expressed among vertebrates and is the predominant form in amphibians and birds [13-15].  
6 To trace the evolutionary history of the IGF-1 isoforms that include exon 5 (i.e. IGF-1Eb and IGF-1Ec), we  
7 performed BLAST searches in mammalian and non-mammalian draft genomes using the Ensembl database  
8 and the UCSC vertebrate genome alignment provided at the UCSC Genome Browser site (Fe. 2009  
9 GRCh37/hg19). ClustalX alignment of *IGF-1* sequences was performed along the exon 6 5 and exon 5 6  
10 regions among 23 representative amniota vertebrates (Supplementary Fig. S1). We retrieved orthologs of  
11 the exon 5 sequence from all placental mammals (Eutheria), except Guinea pig (*Cavia aperea porcellus*),  
12 probably because of insufficient sequencing coverage. In addition, we identified exon 5 sequences from the  
13 marsupial short-tailed opossums (*Didelphidae*) as well as from the monotremata egg-laying platypus  
14 (*Prototheria*). By contrast, we failed to identify exon 5 orthologs in all examined non-mammalian vertebrates  
15 (fish, amphibians, reptiles and birds), including turkey (*Meleagris gallopavo*), chicken (*Gallus gallus*), zebra  
16 finch (*Taeniopygia guttata*) and green anole lizard (*Anolis carolinensis*) (Supplementary Fig. S1). These data  
17 show that only exon 6 is conserved in all vertebrates and suggest that the Ea peptide represents the  
18 ancestral IGF-1 E-peptide, whereas the exon 5-encoded Eb and Ec peptides have been acquired during  
19 evolution in the mammalian lineage.

20

### 21 3.2 MIR-b retroposon exonization in the *IGF-1* gene explains the origin of the Mammalia alternatively 22 spliced exon 5 in Mammalia

23 ~~Recent studies suggest that~~ Transposon exonization is a major source of new exons in higher eukaryotes  
24 [37]. This process occurs when transposon-derived intronic sequences are recognized by the spliceosome  
25 and are exonized. Using the RepeatMasker web server, we searched for transposon elements in the  
26 genomic region encompassing *IGF-1* exon 5. No repetitive sequences were detected using the search tools  
27 “abblast”, “rmbblast” or “cross-match” whereas the recently introduced “nhmmer” search engine ~~of the~~  
28 ~~RepeatMasker web server~~, which uses the Dfam database, [27-28], highlighted a transposon belonging to the  
29 Mammalian Interspersed Repetitive-b (MIR-b) element (bit score 13.6; E-value 0.0031). The *IGF-1* exon 5  
30 displays about 60% nucleotide identity to 103 nt of MIR-b consensus sequence from Repbase (Fig. 1C-2).  
31 The exonized MIR-b element is inserted in the antisense orientation relative to the *IGF-1* gene. Both the  
32 polypyrimidine tract and the *IGF-1* exon 5 splice sites are embedded within the MIR-b sequence (Fig. 1C-2).  
33 The 3'ss of *IGF-1* exon 5 is in close proximity to the ~~65-nt~~ 65 nt that form the highly conserved core region of  
34 MIR-b while the cryptic 5'ss lies inside within this region (Fig. 1C-2). This configuration is similar to those  
35 described for others antisense MIR-b exons [35-36] and ~~ef for disease-linked~~ genes containing MIR  
36 pseudoexons [38].

37

### 1 3.3 Comparison of the MIR-derived IGF-1 exon 5 sequence within mammalian species

2 Alignment of the IGF-1 exon 5 ~~sequence nucleotides~~—originated from exonization of MIR-b element among  
3 ~~human, macaque, mouse, rat, rabbit, pig, cow, sheep and goat~~ mammalian orthologous showed a level of  
4 nucleotide identity ~~ranging from 56 to 62%~~ between MIR-b consensus sequence and IGF-1 exon 5 coding for  
5 Ec peptide (49 bp in human and 52 bp in other species) ~~ranging from 56 to 62%~~ (Fig. 2A 3A). This analysis  
6 indicates that, after exonization, the exon 5 sequence coding for the Ec peptides has been fairly well  
7 conserved among species, even though the percent of identity of nucleotide sequence was lower in rodents  
8 and logomorpha than ~~in~~ others species (Fig. 2A-3A). Notably, the 3'ss of exon 5 is conserved and relatively  
9 strong, whereas the ~~cryptic 5'ss c5'ss~~ has only a limited match to the consensus, except in mouse and rat  
10 (Figs. 2A-3A and 2B-3B). Interestingly, humans and murids showed a marked difference in ~~c5'ss strength of~~  
11 ~~the cryptic 5'ss of exon 5~~. This splice site is very weak in humans; by contrast it appears strong in mouse  
12 and rat (Fig. 2B-3B). Further analysis of ~~the strength of~~ the IGF-1 c5'ss strength ~~exon 5'ss~~ among 58  
13 mammalian genomic sequence datasets available at the UCSC genome browser showed that most  
14 mammals (79%) ~~have a cryptic 5'ss with~~ display medium strength, while 16% have ~~high strength a strong~~  
15 ~~5'ss~~, including all members of the two rodent families Cricetidae: prairie vole (*Microtus ochrogaster*), chinese  
16 hamster (*Cricetulus barabensis griseus*), golden hamster (*Mesocricetus auratus*) and Muridae: mouse (*Mus*  
17 *musculus*) and rat (*Rattus norvegicus*) (Supplementary Fig. S2). These analyses suggest that processing of  
18 both IGF-1Eb and IGF-1Ec mRNA variants are expected in most mammals, whereas the weak nature of the  
19 ~~cryptic 5'ss c5'ss~~ of human exon 5 will likely prevent its recognition by the spliceosome, yielding  
20 predominantly the IGF-1Eb isoform. By contrast, the relatively strong c5'ss of exon 5 in the Muridae Family  
21 will favour processing of the IGF-1Ec. Therefore, our analyses suggest that, under normal conditions, the  
22 IGF-1 alternative splice variants are likely expressed in a species-specific manner.

### 23 3.4 Expression pattern of IGF-1 mRNA isoforms ~~in mammals within mammalian species~~

24 To test this hypothesis, we quantified by real-time PCR the IGF-1Ea, IGF-1Eb and IGF-1Ec mRNA variants  
25 in skeletal muscle, adipose tissues and liver of human (*Homo sapiens*), macaque (*Macaca mulatta*) and  
26 mouse (*Mus musculus*). These species represent examples of weak, intermediate and strong ~~exon 5 c5'ss~~  
27 ~~cryptic 5'ss of exon 5~~, respectively (Fig. 2B-3B). As predicted by our analysis, the expression profile of these  
28 IGF-1 mRNA variants ~~changed varies~~ among species (Anova test;  $p < 0.01$ ) (Fig. 3A-4A). In particular the  
29 IGF-1Ea and IGF-1Eb isoforms are predominant in human skeletal muscle, adipose tissue and liver,  
30 whereas the IGF-1Ec isoform accounts only for 1-5% of total IGF-1 mRNAs (Fig. 3A-4A). Conversely, IGF-  
31 1Ea and IGF-1Ec isoforms predominate in mouse tissues, whereas the IGF-1Eb isoform was barely  
32 detectable ( $\leq 1\%$  of total IGF-1 mRNAs) (Fig. 3A-4A), as also reported previously [39]. In macaque (*Macaca*  
33 *mulatta*), the expression of IGF-1Eb was higher compared to IGF-1Ec in skeletal muscle and adipose tissue,  
34 whereas these splice variants were expressed at approximately equivalent level in the liver (Fig. 3A-4A).  
35 Noteworthy, the same PCR primers have been used to amplify human and macaque IGF-1 mRNA variants,  
36 ruling out differences in amplification efficiencies between species.

1 To further confirm the species-specific expression pattern of IGF-1 isoforms, we performed a 3' RACE-PCRs  
2 using liver-derived cDNAs from human, macaque and or mouse. In addition, we also analyzed liver tissues  
3 from rat (*Rattus norvegicus*; strong c5'ss), rabbit (*Oryctolagus cuniculus*; intermediate c5'ss strength) and pig  
4 (*Sus scrofa*; intermediate c5'ss strength). IGF-1Ea, IGF-1Eb and IGF-1Ec mRNA products were amplified in  
5 macaque, rabbit and pig liver, which represent examples of intermediate cryptic c5'ss- (Fig. 3B-4B and  
6 Supplementary Figs. S3A, S3B and S3C for the complete nucleotide sequences and putative termination  
7 sites). On the contrary, the IGF-1Ec and IGF-1Eb splice variants were not detected in human and Muridae  
8 liver, respectively, probably because their expression level is too low for the sensitivity of the 3' RACE-PCR  
9 analysis. Notably, human, mouse, rabbit and pig IGF-1 isoforms show multiple termination sites (Fig. 3B and  
10 Supplementary Figs. S3A, S3B and S3C) (Supplementary Fig. S3A). Isoform-specific PCR analysis showed  
11 that both IGF-1Eb and IGF-1Ec mRNA variants were expressed in other species characterized by an  
12 intermediate c5'ss, such as cow (*Bos taurus*), sheep (*Ovis aries*) and goat (*Capra hircus*) (Supplementary  
13 Fig. S3D S3B).

14 Collectively, these results document that the IGF-1Eb isoform is more widely expressed than previously  
15 reported and point out a sharp difference in the expression of IGF-1 isoforms expression between human  
16 and mouse.

17

### 18 3.5 The mouse- and human-specific IGF-1 isoforms pattern is not dependent on the cell context

19 In order to analyze the species-specific splicing of exon 5 in human and mouse, we constructed minigene  
20 systems consisting of exon 4, intron 4, exon 5, part of intron 5 and exon 6 of the corresponding *IGF-1* genes  
21 (Fig. 4A-5A). Minigene splicing assay performed in four human HEK293T cells (Fig. 5B) and three in primary  
22 mouse cell lines hepatocytes (Fig. 5C) showed that alternative splicing of the human minigene (hIGF-1) did  
23 not generate the IGF-1Ec variant (Figs. 5B and 5C-Fig. 4B, left panel and Supplementary Fig. S4). This  
24 result is These results are in line with the barely detectable amounts of this mRNA found by quantitative real-  
25 time PCR in human tissues (Fig. 3A-4A). By contrast, the mouse minigene (mIGF-1) yielded mainly the IGF-  
26 1Ea and IGF-1Ec variants in all cell types analyzed (Figs. 5B and 5C-Fig. 4B, right panel and Supplementary  
27 Fig. S4), similarly to what observed in mouse tissues (Fig. 3A-4A), with splicing in primary mouse  
28 hepatocytes mostly shifted toward the IGF-1Ec isoform in mouse cell types (Fig. 5C). We suppose that these  
29 little differences might be due to the lack of part of intron 5 and/or of chromatin context, or to the presence of  
30 a different promoter in the minigene [40-41]. Collectively, these results indicate that the mouse and human  
31 IGF-1 minigenes mostly recapitulate the splicing pattern observed in tissues.

32 Since the mhIGF-1 and hmIGF-1 splicing patterns were mostly preserved regardless of the cell context, it is  
33 likely that splicing of exon 5 is not regulated by a specific human or mouse trans-acting environment. To  
34 validate this hypothesis, we forced the expression of selected splicing factors to determine their effect on the  
35 mouse and human IGF-1 minigenes. Human *IGF-1* exon 5 alternative splicing is regulated by competition  
36 between SRSF1 and hnRNPA1 [8-9]. Thus, we tested whether or not the differences in human and mouse  
37 exon 5 alternative splicing could be influenced by modulation of the expression of SRSF1, hnRNP A1 and

1 other similar splicing factors (Fig. 5 6A and 6B and Supplementary Fig. S5). As expected [8-9],  
2 overexpression of hnRNP A1 promoted splicing of the IGF-1Ea variant by the hIGF-1 minigene, whereas  
3 overexpression of SRSF1 favored splicing of the IGF-1Eb variant (Fig. 5B, left panel 6A). However,  
4 overexpression of neither splicing factor was capable of inducing the IGF-1Ec variant from the hIGF-1  
5 minigene (Fig. 5B, left panel 6A). Moreover, splicing assays in presence of several other hnRNPs (A2, F, G,  
6 H, I and K), SR (SRSF3 and 7) and SR-like (TRA2 $\alpha$  and 2 $\beta$ ) proteins showed that none of them was capable  
7 of inducing the IGF-1Ec variant from the hIGF-1 minigene in HEK293T cells (Fig. 5B, left panel 6A), even  
8 though most of these splicing factors influenced the IGF-1Eb/IGF-1Ea ratio (Fig. 5B, left panel 6A). Similarly,  
9 depletion of hnRNPA1, SRSF1 or TRA2 $\beta$  did not promote splicing of IGF-1Ec isoform by the endogenous  
10 human gene (Fig. 5D, left panel and Supplementary Fig. S5B) and hIGF-1 minigene (Supplementary Fig.  
11 S5C). Since, splicing of IGF-1Ec is promoted by local muscle tissues damage in rodent [7], we tested  
12 whether oxidative stress condition in a human cell line (LNCaP) may promote splicing of this isoform.  
13 Notably, we observed modulation of IGF-1Eb/IGF-1Ea ratio but not splicing of IGF-1Ec (Supplementary Fig.  
14 S5D). These results strongly suggest that lack or low expression of the IGF-1Ec variant is unlikely due to  
15 deficient expression of trans-acting factors in human cells. ~~Similar conclusions could be drawn by~~  
16 ~~experiments performed with the mIGF-1 minigene (Fig. 6B). In this case,~~ In the case of the mIGF-1  
17 minigene, overexpression (Fig. 5B, right panel 5C) or depletion (Supplementary Fig. S5C) of most splicing  
18 factors modulated the IGF-1Ec/IGF-1Ea ratio between the two canonical variants ~~expressed by the mIGF-1~~  
19 ~~minigene~~, but none of them affected the low levels of basal expression of the IGF-1Eb variant produced by  
20 ~~this the mIGF-1 minigene.~~ Accordingly, depletion of hnRNPA1, SRSF1 or TRA2 $\beta$  did not promote splicing of  
21 IGF-1Eb isoform from the endogenous mouse gene (Fig. 5D, right panel and Supplementary Fig. S5B).  
22 These results argue against a key role played by trans-acting factors in the alternative processing of *IGF-1*  
23 exon 5 observed in both human and mouse.

24

### 25 **3.6 The relative strength of the cryptic 5'ss of *IGF-1* exon 5 contributes to the species-specific** 26 **production of the IGF-1Ec isoform in human and mouse**

27 Given the marked difference in strength between the mouse and human cryptic c5'ss (Fig. 2B 3B), we  
28 hypothesized that evolutionary modifications of this splice site explains its different usage in the two species.  
29 To verify this possibility, we replaced by site-directed mutagenesis the weak human c5'ss with the strong  
30 mouse splice site (hMUT: TGvsGA) (Fig. 6A 7A). Notably, strengthening of human c5'ss partially recovered  
31 splicing of IGF-1Ec (Fig. 6B 7B). By contrast, weakening of the mouse c5'ss by reverting it to the human  
32 sequence (mMUT: GAvsTG) (Fig. 6A 7A) prevented its usage and favored selection of another cryptic 5'ss  
33 located downstream the canonical one (IGF-1Ec\* in Figs. 6A 7A and 6C 7C). These results strongly suggest  
34 that the strength of the cryptic c5'ss in exon 5 is a crucial factor in determining the species-specific  
35 expression of the IGF-1Ec splice variant.

36 We noticed that humanization of the mouse c5'ss (mMUT: GAvsTG) was not sufficient to enhance splicing of  
37 the IGF-1Eb variant from the mouse minigene (Fig. 6C 7C). Moreover, splicing was not affected by

1 overexpression of SRSF1, which promotes the IGF-1Eb variant from the human minigene (Supplementary  
2 Fig. ~~S6A S4B and S4C~~). Similarly, unlike with the mouse minigene, overexpression of SRSF1 did not  
3 decrease the IGF-1Ec/IGF-1Ea ratio in presence of hMUT (TGvsGA) minigene (Supplementary Fig. ~~S6A~~  
4 ~~S4B~~). Altogether, these results suggest that other sequence elements, in addition to the ~~eryptie~~ c5'ss,  
5 contribute to species-specific variations of *IGF-1* alternative splicing, possibly by conferring sensitivity to  
6 trans-acting factors.

7

### 8 **3.7 Molecular evolution of mammalian IGF-1 isoforms and prediction of IGF-1 domain structural** 9 **properties**

10 The selective pressure acting on IGF-1 domains is currently largely unknown [17, 19]. The higher variability  
11 of E-peptide sequences compared to core region sequences among mammals might reflect a less stringent  
12 functional requirement for this portion of the IGF-1 protein [17]. On the other hand, the extremely low  
13 synonymous mutation rates found on the dual-coding exon 6 across 29 mammalian species [19] suggests an  
14 enrichment of conserved regulatory elements in this region of the gene [18].

15 To better estimate the evolutionary pressure acting on IGF-1 domains we compared the non-synonymous  
16 ( $dN$ ) and synonymous ( $dS$ ) substitutions among 27 mammalian protein-coding sequences of *IGF-1* [32]. As  
17 shown in Table 1, there was substantial variation in the substitution rate within the gene regions encoding  
18 the different domains. Specifically, the Ea, Eb and Ec peptides showed a marked reduction of  $dS$  and only a  
19 slight increase of  $dN$  (significant only for Eb peptide) compared to the core region. Thus, the increase of the  
20 E-peptides  $dN/dS$  ratio was only marginally caused by an increase of  $dN$ , rather, it is likely due to a  
21 significantly smaller  $dS$ .

22 This result confirms and extends ~~the~~ previous findings ~~of Lin M.F. et al.~~ [19], ~~as since~~ we demonstrated  
23 strong reduction in the rate of on synonymous substitution both in the dual-coding region of exon 6 and in the  
24 single-coding region of exon 5. The strong evolutionary constraint on these exons is consistent with previous  
25 observations on alternatively spliced exons and is a signature of enrichment of functional elements, such as  
26 splicing enhancer or ~~exclusion silencer~~ elements [18, 21]. Accordingly, ~~the~~ exon 5 contains binding sites for  
27 SRSF1 and hnRNP A1 [8-9] and ~~the~~ exon 5 alternative splicing is regulated by multiple splicing factors (Fig.  
28 ~~5 6~~). Intriguingly, recent studies demonstrated that protein-coding regions with low synonymous mutation  
29 rate are significantly enriched in Intrinsically Disordered Regions (IDRs) [42-43]. IDRs are polypeptide chains  
30 lacking a stable tertiary structure and the relaxed protein structural constraint provide an advantage when a  
31 protein-coding region simultaneously contains additional regulatory sites such as splicing enhancer and  
32 repressor sequences [43]. In order to analyze the structural properties of IGF-1 domains, we used DisCons  
33 [35], a sequence analysis tool able to identify protein disorder segments in an evolutionary context. As  
34 expected, among 27 mammalian species analyzed, the IGF-1 core region was predicted to be structured for  
35 the most part (Fig. 7 8). Conversely all the E-peptides were predicted to be almost completely “constrained  
36 disorder” (Fig. 7 8). Namely, they were enriched in disorder-promoting residues and there was conservation  
37 of the property of disorder across mammals. These results do not depend on the disorder predictor algorithm  
38 because core results were qualitatively replicated using PONDR® VSL2 instead of IUPred (data not shown).

1 Thus, during the course of mammalian evolution, the structure of IGF-1 protein ~~has been~~ ~~was~~ strongly  
2 conserved to maintain both the folded structure of IGF-1 core and its intrinsically disordered E-peptide tails.  
3  
4

#### 5 **4. Discussion**

6 In the present work we provided evolutionary evidence supporting the emergence of a functional alternative  
7 exon 5 by ~~an~~ exonization of ~~a~~ MIR-b element in mammalian *IGF-1* intron 4 (Fig. 1C). The exon 5 gain leads  
8 to the generation of two new IGF-1 E-peptides in the mammalian lineage: Eb and Ec. Our work indicates that  
9 the Ea peptide represents the ancestral E-peptide, common to all vertebrates, whereas Eb and Ec peptides  
10 are an evolutionary novelty appeared in Mammalia after the exonization of MIR-b transposon by exaptation  
11 [44]. The MIR element ~~transposon~~ belongs to a family of transposable elements, which were actively  
12 propagating prior to ~~the~~ mammalian radiation [45]. The MIR family, with a conservation rate between 60 and  
13 70%, is one of the most diverged transposable element families identifiable in the human genome and only  
14 recent molecular approaches have improved their chances to be identified [27, 46]. For example, MIRs have  
15 a 95% chance of being missed by RepeatMasker if they are 50 bp in length, and 50% chance of being  
16 missed if 100 bp long [46]. In our study, we ~~were~~ successfully ~~identified in-identifying~~ and ~~aligned~~ ~~aligning~~ the  
17 MIR-b element to 103 nucleotides of the human *IGF-1* exon 5 sequence only by using the recently  
18 introduced “hmmer” search engine of RepeatMasker, which uses the hidden Markov model search tool  
19 nhmmer and the Dfam library [27-28].

20 Several genome-wide studies have focused on the exaptation and exonization of MIR elements and found  
21 over 100 exonized MIR elements in the human genome, located both in the UTRs and the CDSs of  
22 annotated genes [47-48]. These studies showed that exonization of retroposed sequences can occur at any  
23 time following their insertion [49-50]. Notably, both DNA sequence analyses performed in our study and  
24 mRNA analyses conducted in birds [13], fish [15] and reptiles [14] showed that *IGF-1* exon 5 is absent in  
25 these vertebrates. Therefore, the emergence of *IGF-1* exon 5 may have taken place around the mammalian  
26 diversification, which coincides with the timing of integration of the MIR elements [51].

27 The fact that MIR-b integrated in the *IGF-1* exon 5 in its antisense orientation and contributed ~~the~~ ~~to~~ splice  
28 sites and the oligopyrimidine tract may have facilitated its exonization immediately after integration [50].  
29 Moreover, the 3'ss of *IGF-1* exon 5 is located near the highly conserved core region of MIR-b, a feature that  
30 is present in many of the exonized-MIR described so far [38, 50] and that may have also contributed to its  
31 early exonization [50]. The conservation and the relative strength of the 3'ss of *IGF-1* exon 5, together with  
32 the maintenance of open reading frame in all ~~the major~~ ~~main~~ mammalian branches, suggest that during  
33 evolution the creation of a functional 3'ss inside the MIR-b integrated in the *IGF-1* intron 4 was essential for  
34 successful exonization of this element [52].

35 Exonization of MIR-b in the mammalian *IGF-1* gene also created a ~~cryptic~~ c5'ss that allows splicing of the  
36 IGF-1Ec variant in some species. Our analyses highlighted that, in contrast to the 3'ss, the mammalian  
37 ~~cryptic~~ c5'ss shows a considerable variation in term of the number of matching nucleotides to the vertebrate

1 5'ss consensus, with the most marked difference observed between humans and Muridae (Figs. ~~2A, 3A~~ and  
2 ~~2B-3B~~ and Supplementary Fig. S2). This led us to hypothesize that the spliceosome ability to recognize the  
3 ~~cryptic c5'ss~~ present in ~~IGF-1~~ exon 5 may vary among mammalian species. Hence, the relative quantity of  
4 the IGF-1Eb and IGF-1Ec isoforms might ~~also differ among between humans, Muridae and others~~ mammals.  
5 Accordingly, we observed a differential expression pattern of IGF-1 isoforms across 9 different mammalian  
6 species analyzed (Fig. ~~3-4A-4B~~ and Supplementary Fig. ~~S3 S3A and S3B~~). In particular, we found a marked  
7 difference in ~~IGF-1~~ splicing between humans and Muridae, with a prevalence of IGF-1Eb in humans and  
8 IGF-1Ec in Muridae. In other mammalian species ~~analysed~~, including the primate macaque, ~~the~~ expression  
9 of ~~Eb the IGF-1Eb~~ isoform was slightly higher ~~than~~ or equal to the ~~IGF1-Ec~~ IGF-1Ec. Hence, the expression  
10 of IGF-1Eb is not restricted to humans, and this splice variant should no longer be considered human-  
11 specific [1, 39].

12 Divergence of alternative splicing represents one of the major driving forces to shape phenotypic differences  
13 across species [53-57]. Such changes could arise from ~~the divergences~~ in *cis*-regulatory elements and/or  
14 *trans*-acting splicing factors [58]. Therefore, we next focused on the mechanism underlying the diversification  
15 of ~~IGF-1~~ alternative splicing between human and mouse. For this purpose we constructed minigene-based  
16 systems of the ~~mouse and human IGF-1~~ alternatively spliced regions. Both minigenes recapitulated the  
17 splicing pattern observed in tissues, ~~as indeed~~ the human and mouse minigenes produced very low levels of  
18 IGF-1Ec and IGF-1Eb isoforms, respectively (~~Figs. 5A and 5B~~ Fig. 4 and Supplementary Fig. S4).  
19 Nevertheless, our minigene splicing assays also showed some small differences with respect to mouse  
20 tissues. ~~For instance, alternative splicing of the mouse minigene in mouse primary hepatocytes is mostly~~  
21 ~~shifted toward the IGF-1Ec isoform. We suppose that these differences might be due to alterations of the~~  
22 ~~genomic and epigenetic structure of the IGF-1 gene with respect to the minigene.~~ In particular, mouse  
23 minigene produces a small amount of the IGF-1Eb variant, whereas this isoform is minimally expressed in  
24 mouse tissues (Figs. 3A and 4B and Supplementary Fig. S4). Moreover, alternative splicing in the mouse  
25 cell lines is mostly shifted toward the IGF-1Ec isoform (Fig. 4B and Supplementary Fig. S4). We suppose  
26 that these differences might be due to alterations of the genomic and epigenetic structure of the *IGF-1* gene  
27 with respect to the minigene ~~and/or to the different promoter present in the minigenes, as all these features~~  
28 ~~have been shown to modulate the outcome of pre-mRNA splicing [40-41]. For instance, the minigene lacks~~  
29 ~~part of intron 5 and it might miss elements that modulate splicing of the endogenous gene. Moreover, it is~~  
30 ~~known that the chromatin context, epigenetic modifications within the transcription unit and the presence of~~  
31 ~~different promoters can influence alternative splicing events. Thus, it is likely that the small differences~~  
32 ~~observed are due to alterations of the chromatin context in the minigene.~~ Nevertheless, our splicing assays  
33 indicate that the IGF-1 minigenes are reliable tools for the analysis of the splicing regulation of the  
34 corresponding endogenous genes. By using these minigenes, we tested whether *cis*-acting elements and/or  
35 *trans*-acting factors were responsible for the different splicing patterns observed *in vivo* between mouse and  
36 human IGF-1 mRNA. Alternative splicing outcome is determined by competition between splice sites in  
37 different exons and is influenced by *trans*-acting splicing factors, which in turn influence splice site

1 recognition by the spliceosome [59]. As a consequence, both changes in the sequence of *cis*-acting  
2 elements and in the expression levels of splicing factors can modulate usage of specific exons [58, 60]. Our  
3 results argue against changes in expression of specific splicing factors as a key determinant between the  
4 splicing differences across species. Indeed, we found that splicing of the mouse and human IGF-1  
5 minigenes is mostly preserved across different cell type contexts (~~Figs. 5A, 5B and 5C~~ Fig. 4B and  
6 ~~Supplementary Figure S4~~). Moreover, by ~~forcing~~ **modulating** the expression of specific splicing factors it was  
7 not possible to promote splicing of the IGF-1Ec variant from the human minigene (~~hIGF-1~~) (Fig. 5B, *left panel*  
8 ~~and Supplementary Fig. S5C~~) or endogenous human gene (Fig. 5D, *left panel* and Supplementary Fig. S5B).  
9 Similarly, neither the overexpression (Fig. 5B ~~5C~~, *right panel*) nor the siRNA-mediated depletion (Fig. 5D,  
10 *right panel* and Supplementary Figs. S5B and S5C) of splicing factors affected the expression of mouse IGF-  
11 ~~1Eb isoform. nor that of IGF-1Eb from the mouse minigene (mIGF-1)~~. We tested several hnRNPs and SR  
12 proteins, including SRSF1 and hnRNP A1 that were previously shown to ~~modulate alternative regulate~~  
13 splicing of human *IGF-1* [8-9]. Most of these *trans*-acting factors modulated the **human IGF-1Eb/IGF-1Ea**  
14 **and mouse IGF-1Ec/IGF-1Ea** ratio ~~between the variants produced by the minigene~~, but none of them  
15 promoted the variant typical of the other species. Thus, our results suggest that the mechanism regulating  
16 this process does not rely on differential expression of species-specific *trans*-acting factors.

17 An alternative possibility is the presence of *cis*-acting elements conferring species-specific features to the  
18 exon 5 region. This hypothesis is in line with a recent study showing that most vertebrate specie-specific  
19 splicing patterns are primarily *cis*-directed [54, 58]. Accordingly, our analysis revealed that the consensus  
20 bases within the **cryptic c5'ss** are highly divergent among mammalian species (Fig. 2A–3A and  
21 Supplementary Fig. S2). This observation led us to hypothesize that variation in the **cryptic c5'ss** strength  
22 might represent a potential source of species-specific splicing pattern. In support of this hypothesis, ~~we~~  
23 ~~found that~~ strengthening the human **c5'ss** by replacing it with the mouse **one 5'ss** (hMUT: TGvsGA) was  
24 sufficient to promote usage of this splice site and to yield IGF-1Ec variant (Figs. 6A and ~~6B-7A and 7B~~). On  
25 the other hand, weakening the mouse **c5'ss** by swapping it with the human sequence (mMUT: GAvsTG)  
26 completely prevented its usage (Figs. 6A and ~~6C-7A and 7C~~). Nevertheless, our results also indicate that  
27 other *cis*-acting elements likely contribute to the splicing pattern typically observed in tissues. First, we found  
28 that “murinization” of the human **cryptic c5'ss** only partially promoted splicing of the IGF-1Ec variant, but did  
29 not restore the IGF-1Ec/IGF-1Ea ratio observed with the mouse minigene and in mouse tissues (Figs. 3A  
30 ~~and 6B-7B~~). ~~This result suggests that additional sequence elements differing between human and mouse~~  
31 ~~exon 5 might also contribute to its alternative splicing, perhaps by recruiting more or less efficiently specific~~  
32 ~~splicing factors. In support of the presence of additional cis-acting elements is also the analysis of the mutant~~  
33 ~~mouse IGF-1 minigene. Furthermore, “humanization” of the mouse cryptic c5'ss (mMUT: GAvsTG)~~  
34 prevented splicing of the IGF-1Ec variant, but did not promote IGF-1Eb variant (Figs. 6A and ~~6C-7A and 7C~~).  
35 By contrast, we observed the selection of another cryptic splice site (\*5'ss) located 38 nt downstream of the  
36 canonical one (Figs. 6A and ~~6C-7A and 7C~~). Notably, this cryptic splice site is not conserved in the human  
37 gene, thus explaining the generation of the unusual IGF-1 mRNA variant instead of IGF-1Eb from the mMUT

1 (GAVsTG) minigene. Moreover, overexpression of SRSF1 did not promote splicing of the IGF-1Eb variant  
2 from the mMUT (GAVsTG) minigene (Supplementary Fig. S6A S4B). SRSF1-dependent splicing of human  
3 IGF-1Eb requires its binding to a purine-rich enhancer located in exon 5 [8-9]. Sequence alignment of human  
4 and mouse exon 5 revealed that the SRSF1 binding site is not conserved in the mouse (Supplementary Fig.  
5 S6B S4D), possibly explaining the different behavior of SRSF1 in presence of the mouse minigene. These  
6 results strongly indicate that the strength of the **cryptic c5'ss** in exon 5 is a prerequisite for *IGF-1* alternative  
7 splicing diversification between species, but also highlight that the co-evolution of additional ~~cis-acting~~ **exonic**  
8 **and intronic cis-regulatory** elements contribute to such diversification [40-41]. Notably, our phylogenetic  
9 analysis of the cryptic 5'ss indicates that this splice site is relative strong not only in mouse but also in rat  
10 and all members of Cricetidae family (Fig. 2B-3B and Supplementary Fig. S2). It is possible that  
11 strengthening of this **cryptic c5'ss** in a common ancestor of Muridae and Cricetidae, as well as the  
12 emergence of other splicing regulatory motifs, have determined **the marked shift toward production of the**  
13 **IGF-1Ec isoform in these species and possibly a the-acquisition-of** rodent-specific **gain** of functional  
14 properties. ~~Accordingly~~ **Interestingly**, a recent study comparing the IGF-1 isoform expression in both mouse  
15 and human muscle samples at different ages, showed that the age-related change of IGF-1 splice variants is  
16 species-dependent, **and, unlike the mouse, only the human IGF-1Eb isoform was regulated during ageing in**  
17 **skeletal muscles** [61]. Thus, the results obtained with murine models on IGF-1Ec/mechano growth factor  
18 regulation must be interpreted with ~~more-than-usual~~ caution and additional studies comparing IGF-1  
19 expression levels across species are needed to clarify the functional role of the IGF-1 isoforms.  
20 ~~The regulation of IGF-1 is a critically important factor in several conditions such as the pathophysiology of~~  
21 ~~several cancers, or the mitogenic and myogenic processes during muscle development, regeneration or~~  
22 ~~hypertrophy~~. Compared with mature IGF-1 relatively little is known about the mechanism of action of the  
23 ~~different E peptides [1, 10-12]. How the E-peptides affect the actions of IGF-1 is still under debate.~~ From an  
24 evolutionary point of view the unchanged persistence over long evolutionary periods of MIRb-derived *IGF-1*  
25 exon 5 implies its functional relevance [52, 62]. Moreover, E-peptides are protein-coding regions in which  
26 synonymous mutation rates are extremely low compared to IGF-1 core (Table 1), indicating additional  
27 sequence constraints beyond those dictated by the structure and function of the proteins. These additional  
28 constraints probably stem from the demands of regulatory sites involved in transcript splicing (Fig. 5) and [8-  
29 9, 18], nevertheless the presence of other regulatory sites such as specific RNA secondary structures and  
30 microRNA targets [20-21] cannot be excluded. The systematic analysis of the role of synonymous variants  
31 and the comparative splicing evaluation of mammalian sequence divergences [63-64] will help to  
32 characterize the functional roles of these regulatory elements.  
33 Analysis of structural properties of IGF-1 domains adds a new layer of complexity to the function of E-  
34 peptides. Indeed, our work represents the first evidence that E-peptides contain disorder-promoting amino  
35 acids and that there is substantial evolutionary pressure to keep the different E-peptides as intrinsically  
36 disordered regions (IDRs) (Fig. 7-8). There is a growing interest on IDRs since they are usually enriched in  
37 posttranslational modification sites and may exert a number of regulatory functions on their "host" protein

1 [65-66]. Intriguingly, we and others recently demonstrated that the IGF-1 protein retaining C-terminal E  
2 peptides are the predominant forms produced intracellularly, instead of mature IGF-1, and are subjected to  
3 extensive post-translational modifications [67-68], further hinting to their functional relevance. ~~Intriguingly, in  
4 murine skeletal muscle, the IGF-1 protein retaining the C-terminal Ea peptide was the predominant form,  
5 instead of mature IGF-1, and exists in both non-glycosylated and glycosylated forms. Moreover  
6 overexpression of IGF-1 lacking any E-peptide did not promote muscle hypertrophy in young mice, further  
7 hinting to its functional relevance.~~

8 Exonization of previously non-coding sequences, together with creation of novel domain combinations, has  
9 been directly related to the increase of organismal complexity [66, 69-71]. Accordingly, we propose that MIR-  
10 b exonization during mammalian evolution determined the IGF-1 exon 5 gain and hence the addition of two  
11 new disordered tails to IGF-1: the Eb and Ec-tails. Thus, novel exon and E-peptide combinations may have  
12 created new layers of regulation to mature IGF-1 in mammalian species. Targeting these regulatory  
13 elements may represent a new strategy to control IGF-1 bioavailability in different physiological/pathological  
14 conditions, with particular attention to possible differences between species.

## 19 **Acknowledgement**

20 We wish to thank Dr. Franco Berrino of the Department of Predictive & Preventive Medicine, National  
21 Institute of Cancer, Milan and Dr. Giorgio Arnaldi, Department of Internal Medicine, Polytechnic University of  
22 Marche Region, Ancona, Italy for a critical reading of the manuscript. This research was supported by the  
23 Italian Ministry of Health, "Ricerca finalizzata 2009" (grant number: RF-2009-1532789) and "Ricerca  
24 Finalizzata 2011 (grant number: GR-2011-02348423); by the Associazione Italiana Ricerca sul Cancro  
25 (AIRC IG14581) and by Telethon (CGP14095).

## References

- [1] E.R. Barton, The ABCs of IGF-I isoforms: impact on muscle hypertrophy and implications for repair, *Appl Physiol Nutr Metab* 31 (2006) 791-797.
- [2] L. Temmerman, E. Slonimsky, N. Rosenthal, Class 2 IGF-1 isoforms are dispensable for viability, growth and maintenance of IGF-1 serum levels, *Growth Horm IGF Res* 20 (2010) 255-263.
- [3] A.M. Oberbauer, The Regulation of IGF-1 Gene Transcription and Splicing during Development and Aging, *Front Endocrinol (Lausanne)* 4 (2013) 39.
- [4] S.L. Chew, P. Lavender, A.J. Clark, R.J. Ross, An alternatively spliced human insulin-like growth factor-I transcript with hepatic tissue expression that diverts away from the mitogenic IBE1 peptide, *Endocrinology* 136 (1995) 1939-1944.
- [5] C.T. Roberts, Jr., S.R. Lasky, W.L. Lowe, Jr., W.T. Seaman, D. LeRoith, Molecular cloning of rat insulin-like growth factor I complementary deoxyribonucleic acids: differential messenger ribonucleic acid processing and regulation by growth hormone in extrahepatic tissues, *Mol Endocrinol* 1 (1987) 243-248.
- [6] S. Yang, M. Alnaqeeb, H. Simpson, G. Goldspink, Cloning and characterization of an IGF-1 isoform expressed in skeletal muscle subjected to stretch, *J Muscle Res Cell Motil* 17 (1996) 487-495.
- [7] M. Hill, G. Goldspink, Expression and splicing of the insulin-like growth factor gene in rodent muscle is associated with muscle satellite (stem) cell activation following local tissue damage, *J Physiol* 549 (2003) 409-418.
- [8] P.J. Smith, E.L. Spurrell, J. Coakley, C.J. Hinds, R.J. Ross, A.R. Krainer, S.L. Chew, An exonic splicing enhancer in human IGF-I pre-mRNA mediates recognition of alternative exon 5 by the serine-arginine protein splicing factor-2/alternative splicing factor, *Endocrinology* 143 (2002) 146-154.
- [9] A. Clery, R. Sinha, O. Anczukow, A. Corriero, A. Moursy, G.M. Daubner, J. Valcarcel, A.R. Krainer, F.H. Allain, Isolated pseudo-RNA-recognition motifs of SR proteins can regulate splicing using a noncanonical mode of RNA recognition, *Proc Natl Acad Sci U S A* 110 (2013) E2802-2811.
- [10] P. Rotwein, Editorial: the fall of mechanogrowth factor?, *Mol Endocrinol* 28 (2014) 155-156.
- [11] A. Armakolas, M. Kaparelou, A. Dimakakos, E. Papageorgiou, N. Armakolas, A. Antonopoulos, C. Petraki, M. Lekarakou, P. Lelovas, M. Stathaki, C. Psarros, I. Donta, P.S. Galanos, P. Msaouel, V.G. Gorgoulis, M. Koutsilieris, Oncogenic Role of the Ec Peptide of the IGF-1Ec Isoform in Prostate Cancer, *Mol Med* 21 (2015) 167-179.
- [12] R.W. Matheny, Jr., B.C. Nindl, M.L. Adamo, Minireview: Mechano-growth factor: a putative product of IGF-I gene expression involved in tissue repair and regeneration, *Endocrinology* 151 (2010) 865-875.
- [13] Y. Kajimoto, P. Rotwein, Structure of the chicken insulin-like growth factor I gene reveals conserved promoter elements, *J Biol Chem* 266 (1991) 9724-9731.
- [14] A.M. Sparkman, T.S. Schwartz, J.A. Madden, S.E. Boyken, N.B. Ford, J.M. Serb, A.M. Bronikowski, Rates of molecular evolution vary in vertebrates for insulin-like growth factor-1 (IGF-1), a pleiotropic locus that regulates life history traits, *Gen Comp Endocrinol* 178 (2012) 164-173.
- [15] D.M. Tiago, V. Laize, M.L. Cancela, Alternatively spliced transcripts of *Sparus aurata* insulin-like growth factor 1 are differentially expressed in adult tissues and during early development, *Gen Comp Endocrinol* 157 (2008) 107-115.
- [16] M. Reinecke, C. Collet, The phylogeny of the insulin-like growth factors, *Int Rev Cytol* 183 (1998) 1-94.

- 1 [17] M. Wallis, New insulin-like growth factor (IGF)-precursor sequences from mammalian  
2 genomes: the molecular evolution of IGFs and associated peptides in primates, *Growth*  
3 *Horm IGF Res* 19 (2009) 12-23.
- 4 [18] Y. Xing, C. Lee, Evidence of functional selection pressure for alternative splicing events  
5 that accelerate evolution of protein subsequences, *Proc Natl Acad Sci U S A* 102 (2005)  
6 13526-13531.
- 7 [19] M.F. Lin, P. Kheradpour, S. Washietl, B.J. Parker, J.S. Pedersen, M. Kellis, Locating  
8 protein-coding sequences under selection for additional, overlapping functions in 29  
9 mammalian genomes, *Genome Res* 21 (2011) 1916-1928.
- 10 [20] J.V. Chamary, J.L. Parmley, L.D. Hurst, Hearing silence: non-neutral evolution at  
11 synonymous sites in mammals, *Nat Rev Genet* 7 (2006) 98-108.
- 12 [21] Y. Xing, C. Lee, Can RNA selection pressure distort the measurement of Ka/Ks?, *Gene*  
13 370 (2006) 1-5.
- 14 [22] M. Blanchette, W.J. Kent, C. Riemer, L. Elnitski, A.F. Smit, K.M. Roskin, R. Baertsch, K.  
15 Rosenbloom, H. Clawson, E.D. Green, D. Haussler, W. Miller, Aligning multiple genomic  
16 sequences with the threaded blockset aligner, *Genome Res* 14 (2004) 708-715.
- 17 [23] R.C. Edgar, MUSCLE: multiple sequence alignment with high accuracy and high  
18 throughput, *Nucleic Acids Res* 32 (2004) 1792-1797.
- 19 [24] J.D. Thompson, T.J. Gibson, D.G. Higgins, Multiple sequence alignment using ClustalW  
20 and ClustalX, *Curr Protoc Bioinformatics Chapter 2* (2002) Unit 2 3.
- 21 [25] G. Yeo, C.B. Burge, Maximum entropy modeling of short sequence motifs with applications  
22 to RNA splicing signals, *J Comput Biol* 11 (2004) 377-394.
- 23 [26] J. Jurka, Repbase update: a database and an electronic journal of repetitive elements,  
24 *Trends Genet* 16 (2000) 418-420.
- 25 [27] T.J. Wheeler, J. Clements, S.R. Eddy, R. Hubley, T.A. Jones, J. Jurka, A.F. Smit, R.D.  
26 Finn, Dfam: a database of repetitive DNA based on profile hidden Markov models, *Nucleic*  
27 *Acids Res* 41 (2013) D70-82.
- 28 [28] T.J. Wheeler, S.R. Eddy, nhmmer: DNA homology search with profile HMMs,  
29 *Bioinformatics* 29 (2013) 2487-2489.
- 30 [29] G. Annibaldi, M. Guescini, D. Agostini, R.D. Matteis, P. Sestili, P. Tibollo, M. Mantuano, C.  
31 Martinelli, V. Stocchi, The expression analysis of mouse interleukin-6 splice variants argued  
32 against their biological relevance, *BMB Rep* 45 (2012) 32-37.
- 33 [30] S. Barik, Site-directed mutagenesis in vitro by megaprimer PCR, *Methods Mol Biol* 57  
34 (1996) 203-215.
- 35 [31] P. Bielli, R. Busa, S.M. Di Stasi, M.J. Munoz, F. Botti, A.R. Kornblihtt, C. Sette, The  
36 transcription factor FBI-1 inhibits SAM68-mediated BCL-X alternative splicing and  
37 apoptosis, *EMBO Rep* 15 (2014) 419-427.
- 38 [32] S.L. Pond, S.D. Frost, S.V. Muse, HyPhy: hypothesis testing using phylogenies,  
39 *Bioinformatics* 21 (2005) 676-679.
- 40 [33] S.V. Muse, Estimating synonymous and nonsynonymous substitution rates, *Mol Biol Evol*  
41 13 (1996) 105-114.
- 42 [34] R. Szklarczyk, J. Heringa, S.K. Pond, A. Nekrutenko, Rapid asymmetric evolution of a dual-  
43 coding tumor suppressor INK4a/ARF locus contradicts its function, *Proc Natl Acad Sci U S*  
44 *A* 104 (2007) 12807-12812.
- 45 [35] M. Varadi, M. Guharoy, F. Zsolyomi, P. Tompa, DisCons: a novel tool to quantify and  
46 classify evolutionary conservation of intrinsic protein disorder, *BMC Bioinformatics* 16  
47 (2015) 153.
- 48 [36] Z. Dosztanyi, V. Csizmok, P. Tompa, I. Simon, IUPred: web server for the prediction of  
49 intrinsically unstructured regions of proteins based on estimated energy content,  
50 *Bioinformatics* 21 (2005) 3433-3434.

- 1 [37] R. Sorek, The birth of new exons: mechanisms and evolutionary consequences, *RNA* 13  
2 (2007) 1603-1608.
- 3 [38] I. Vorechovsky, Transposable elements in disease-associated cryptic exons, *Hum Genet*  
4 127 (2010) 135-154.
- 5 [39] A.M. Oberbauer, J.M. Belanger, G. Rincon, A. Canovas, A. Islas-Trejo, R. Gularte-Merida,  
6 M.G. Thomas, J.F. Medrano, Bovine and murine tissue expression of insulin like growth  
7 factor-I, *Gene* 535 (2014) 101-105.
- 8 [40] S. Naftelberg, I.E. Schor, G. Ast, A.R. Kornblihtt, Regulation of alternative splicing through  
9 coupling with transcription and chromatin structure, *Annu Rev Biochem* 84 (2015) 165-198.
- 10 [41] N.N. Singh, M.N. Lawler, E.W. Ottesen, D. Upreti, J.R. Kaczynski, R.N. Singh, An intronic  
11 structure enabled by a long-distance interaction serves as a novel target for splicing  
12 correction in spinal muscular atrophy, *Nucleic Acids Res* 41 (2013) 8144-8165.
- 13 [42] E. Kovacs, P. Tompa, K. Liliom, L. Kalmar, Dual coding in alternative reading frames  
14 correlates with intrinsic protein disorder, *Proc Natl Acad Sci U S A* 107 (2010) 5429-5434.
- 15 [43] M. Macossay-Castillo, S. Kosol, P. Tompa, R. Pancsa, Synonymous constraint elements  
16 show a tendency to encode intrinsically disordered protein segments, *PLoS Comput Biol* 10  
17 (2014) e1003607.
- 18 [44] J. Brosius, S.J. Gould, On "genomenclature": a comprehensive (and respectful) taxonomy  
19 for pseudogenes and other "junk DNA", *Proc Natl Acad Sci U S A* 89 (1992) 10706-10710.
- 20 [45] L. Lin, P. Jiang, S. Shen, S. Sato, B.L. Davidson, Y. Xing, Large-scale analysis of exonized  
21 mammalian-wide interspersed repeats in primate genomes, *Hum Mol Genet* 18 (2009)  
22 2204-2214.
- 23 [46] A.P. de Koning, W. Gu, T.A. Castoe, M.A. Batzer, D.D. Pollock, Repetitive elements may  
24 comprise over two-thirds of the human genome, *PLoS Genet* 7 (2011) e1002384.
- 25 [47] B. Mersch, N. Sela, G. Ast, S. Suhai, A. Hotz-Wagenblatt, SERpredict: detection of tissue-  
26 or tumor-specific isoforms generated through exonization of transposable elements, *BMC*  
27 *Genet* 8 (2007) 78.
- 28 [48] N. Sela, B. Mersch, N. Gal-Mark, G. Lev-Maor, A. Hotz-Wagenblatt, G. Ast, Comparative  
29 analysis of transposed element insertion within human and mouse genomes reveals Alu's  
30 unique role in shaping the human transcriptome, *Genome Biol* 8 (2007) R127.
- 31 [49] R. Baertsch, M. Diekhans, W.J. Kent, D. Haussler, J. Brosius, Retrocopy contributions to  
32 the evolution of the human genome, *BMC Genomics* 9 (2008) 466.
- 33 [50] M. Krull, M. Petrusma, W. Makalowski, J. Brosius, J. Schmitz, Functional persistence of  
34 exonized mammalian-wide interspersed repeat elements (MIRs), *Genome Res* 17 (2007)  
35 1139-1145.
- 36 [51] F.S. de Souza, L.F. Franchini, M. Rubinstein, Exaptation of transposable elements into  
37 novel cis-regulatory elements: is the evidence always strong?, *Mol Biol Evol* 30 (2013)  
38 1239-1251.
- 39 [52] J. Schmitz, J. Brosius, Exonization of transposed elements: A challenge and opportunity for  
40 evolution, *Biochimie* 93 (2011) 1928-1934.
- 41 [53] H. Keren, G. Lev-Maor, G. Ast, Alternative splicing and evolution: diversification, exon  
42 definition and function, *Nat Rev Genet* 11 (2010) 345-355.
- 43 [54] N.L. Barbosa-Morais, M. Irimia, Q. Pan, H.Y. Xiong, S. Gueroussov, L.J. Lee, V.  
44 Slobodeniuc, C. Kutter, S. Watt, R. Colak, T. Kim, C.M. Misquitta-Ali, M.D. Wilson, P.M.  
45 Kim, D.T. Odom, B.J. Frey, B.J. Blencowe, The evolutionary landscape of alternative  
46 splicing in vertebrate species, *Science* 338 (2012) 1587-1593.
- 47 [55] J. Merkin, C. Russell, P. Chen, C.B. Burge, Evolutionary dynamics of gene and isoform  
48 regulation in Mammalian tissues, *Science* 338 (2012) 1593-1599.
- 49 [56] T. Lappalainen, M. Sammeth, M.R. Friedlander, P.A. t Hoen, J. Monlong, M.A. Rivas, M.  
50 Gonzalez-Porta, N. Kurbatova, T. Griebel, P.G. Ferreira, M. Barann, T. Wieland, L. Greger,  
51 M. van Iterson, J. Almlof, P. Ribeca, I. Pulyakhina, D. Esser, T. Giger, A. Tikhonov, M.

- 1 Sultan, G. Bertier, D.G. MacArthur, M. Lek, E. Lizano, H.P. Buermans, I. Padioleau, T.  
2 Schwarzmayr, O. Karlberg, H. Ongen, H. Kilpinen, S. Beltran, M. Gut, K. Kahlem, V.  
3 Amstislavskiy, O. Stegle, M. Pirinen, S.B. Montgomery, P. Donnelly, M.I. McCarthy, P.  
4 Flicek, T.M. Strom, H. Lehrach, S. Schreiber, R. Sudbrak, A. Carracedo, S.E. Antonarakis,  
5 R. Hasler, A.C. Syvanen, G.J. van Ommen, A. Brazma, T. Meitinger, P. Rosenstiel, R.  
6 Guigo, I.G. Gut, X. Estivill, E.T. Dermitzakis, Transcriptome and genome sequencing  
7 uncovers functional variation in humans, *Nature* 501 (2013) 506-511.
- 8 [57] A. Necsulea, H. Kaessmann, Evolutionary dynamics of coding and non-coding  
9 transcriptomes, *Nat Rev Genet* 15 (2014) 734-748.
- 10 [58] Q. Gao, W. Sun, M. Ballegeer, C. Libert, W. Chen, Predominant contribution of cis-  
11 regulatory divergence in the evolution of mouse alternative splicing, *Mol Syst Biol* 11 (2015)  
12 816.
- 13 [59] D.L. Black, Mechanisms of alternative pre-messenger RNA splicing, *Annu Rev Biochem* 72  
14 (2003) 291-336.
- 15 [60] M. Chen, J.L. Manley, Mechanisms of alternative splicing regulation: insights from  
16 molecular and genomics approaches, *Nat Rev Mol Cell Biol* 10 (2009) 741-754.
- 17 [61] M. Sandri, L. Barberi, A.Y. Bijlsma, B. Blaauw, K.A. Dyar, G. Milan, C. Mammucari, C.G.  
18 Meskers, G. Pallafacchina, A. Paoli, D. Pion, M. Roceri, V. Romanello, A.L. Serrano, L.  
19 Toniolo, L. Larsson, A.B. Maier, P. Munoz-Canoves, A. Musaro, M. Pende, C. Reggiani, R.  
20 Rizzuto, S. Schiaffino, Signalling pathways regulating muscle mass in ageing skeletal  
21 muscle: the role of the IGF1-Akt-mTOR-FoxO pathway, *Biogerontology* 14 (2013) 303-323.
- 22 [62] V. Gotea, W. Makalowski, Do transposable elements really contribute to proteomes?,  
23 *Trends Genet* 22 (2006) 260-267.
- 24 [63] F. Pagani, M. Raponi, F.E. Baralle, Synonymous mutations in CFTR exon 12 affect splicing  
25 and are not neutral in evolution, *Proc Natl Acad Sci U S A* 102 (2005) 6368-6372.
- 26 [64] Z.E. Sauna, C. Kimchi-Sarfaty, Understanding the contribution of synonymous mutations to  
27 human disease, *Nat Rev Genet* 12 (2011) 683-691.
- 28 [65] M. Buljan, G. Chalancon, S. Eustermann, G.P. Wagner, M. Fuxreiter, A. Bateman, M.M.  
29 Babu, Tissue-specific splicing of disordered segments that embed binding motifs rewires  
30 protein interaction networks, *Mol Cell* 46 (2012) 871-883.
- 31 [66] G. Thiulin-Pardo, L. Avilan, M. Kojadinovic, B. Gontero, Fairy "tails": flexibility and function  
32 of intrinsically disordered extensions in the photosynthetic world, *Front Mol Biosci* 2 (2015)  
33 23.
- 34 [67] J. Durzynska, A. Philippou, B.K. Brisson, M. Nguyen-McCarty, E.R. Barton, The pro-forms  
35 of insulin-like growth factor I (IGF-I) are predominant in skeletal muscle and alter IGF-I  
36 receptor activation, *Endocrinology* 154 (2013) 1215-1224.
- 37 [68] M. De Santi, G. Annibalini, E. Barbieri, A. Villarini, L. Vallorani, S. Contarelli, F. Berrino, V.  
38 Stocchi, G. Brandi, Human IGF1 pro-forms induce breast cancer cell proliferation via the  
39 IGF1 receptor, *Cell Oncol (Dordr)* (2015).
- 40 [69] M. Buljan, A. Frankish, A. Bateman, Quantifying the mechanisms of domain gain in animal  
41 proteins, *Genome Biol* 11 (2010) R74.
- 42 [70] P.R. Romero, S. Zaidi, Y.Y. Fang, V.N. Uversky, P. Radivojac, C.J. Oldfield, M.S. Cortese,  
43 M. Sickmeier, T. LeGall, Z. Obradovic, A.K. Dunker, Alternative splicing in concert with  
44 protein intrinsic disorder enables increased functional diversity in multicellular organisms,  
45 *Proc Natl Acad Sci U S A* 103 (2006) 8390-8395.
- 46 [71] V.N. Uversky, The most important thing is the tail: multitudinous functionalities of  
47 intrinsically disordered protein termini, *FEBS Lett* 587 (2013) 1891-1901.
- 48  
49  
50

## Figure captions

**Figure 1. Schematic representation of the *IGF-1* gene (A), its splice variants (B) and the exonization of MIR-b in the *IGF-1* gene (C).** (A) Map of the *IGF-1* gene showing exons (boxes), introns (solid lines), splicing options (dashed lines), cryptic 5' splice site (c5'ss) in exon 5 and poly(A) sites (pA). The graph is not drawn to scale. (B) Splice variants of the *IGF-1* gene; human *IGF-1* NCBI RefSeq transcripts: class I IGF-1Ea (NM\_000618), class I IGF-1Eb (NM\_001111285), class I IGF-1Ec (NM\_001111283) and class II IGF-1Ea (NM\_001111284). (C) Pairwise alignment of the human *IGF-1* gene sequence and the MIR-b inverse consensus sequence created by RepeatMasker using "nhmmer" search engine (intronic nucleotides in lower-case letters, exonic nucleotides in upper-case letters). The polypyrimidine tract and exon 5 splice sites are shown in bold in *IGF-1* sequence; the bold region of MIR-b sequence corresponds to the 65-nt conserved central domain of MIR-b sequence. The middle line shows exact matches (spaces), gaps (dashes) and mismatches (i=transition; v=transversion).

~~**Figure 1. Schematic representation of the *IGF-1* gene and its splice variants.** Exons (boxes), introns (lines), cryptic 5' splice site (c5'ss) in exon 5 (upper panel) and splice variants of the *IGF-1* gene (lower panel) are indicated.~~

~~**Figure 2. *IGF-1* exon 5 is an exapted retroposon of the MIR-b family.** Pairwise alignment of the human *IGF-1* exon 5 sequence and the MIR-b inverse consensus sequence created by RepeatMasker using "nhmmer" search engine (intronic nucleotides in lower-case letters, exonic nucleotides in upper-case letters). The polypyrimidine tract and exon 5 splice sites are shown in bold in *IGF-1* sequence; the bold region of MIR-b sequence correspond to the 65-nt conserved central domain of MIR-b sequence. The middle line shows exact matches (spaces), gaps (dashes) and mismatches (i=transition; v=transversion).~~

~~**Figure 3. (A) Multiple alignment of MIR-derived exon 5 of *IGF-1* gene among mammals. Figure 2. Comparative analyses of exon 5 splice sites among mammals.** (A) Multiple alignment of MIR-derived exon 5 of *IGF-1* gene among mammals. Partial intron 4 and exon 5 sequences is shown (intronic nucleotides in lower-case letters, exonic nucleotides in upper-case letters). The nucleotides that did not match the 3'ss and 5'ss consensus sequences are highlighted in grey (y=purine; r=pyrimidine). Nucleotides conservation is marked at the lower edge with asterisks indicating full conservation relative to the MIR-b consensus sequence (lower row). The percentage of nucleotide identity between mammalian sequences and MIR-b consensus sequence is indicated under "MIR-b identity". (B) Splice site score of the 3'ss and c5'ss of *IGF-1* exon 5 was calculated using the MaxEntScan algorithm [25].~~

~~**Figure 3. Expression pattern of IGF-1 mRNA isoforms within mammalian species.** (A) Quantification of IGF-1Ea, IGF-1Eb and IGF-1Ec isoforms by real time RT-PCR on skeletal muscle, adipose tissue and liver of Human, Macaque and Mouse. Percentage mean (+/-SD) of IGF-1 isoforms was calculated as: [(mRNA of single isoform/(mRNAs of IGF1Ea + IGF1Eb + IGF1Ec)\*100]. n=2-4. (B) Agarose gel of 3' RACE-PCR analysis performed in liver tissues from different mammalian species. PCR products were eluted from gel, subcloned and sequenced (See Supplementary Figs. S3A, S3B and S3C). A schematic representation of IGF-1 gene and the relative position of forward primers used in the 3'RACE-PCR is also shown. See Materials and Methods section for universal reverse primer information. Diagrammatic representation of splice variants is given on the left of the gel; their sizes (bp) are indicated on the right.~~

~~**Figure 4. Expression pattern of IGF-1 mRNA isoforms within mammalian species.** (A) Quantification of IGF-1Ea, IGF-1Eb and IGF-1Ec isoforms by real time RT-PCR on skeletal muscle, adipose tissue and liver of Human, Macaque and Mouse. Percentage of IGF-1 isoforms was calculated as: [(mRNA of single isoform/(mRNAs of IGF1Ea + IGF1Eb + IGF1Ec)\*100]. The percentage mean (+/- S.D.) of IGF-1Ea is shown. (B) The 3' RACE-PCR products obtained from liver tissues of mammalian species. The 3' RACE-PCR products were loaded on 4.0% agarose gel and DNA fragments were eluted from gel, subcloned and sequenced. Dashes mark: 1, PCR fragment correspondent to IGF-1Ea isoform; 2, PCR fragment correspondent to IGF-1Eb isoform; and 3, PCR fragment correspondent to IGF-1Ec isoform.~~

~~**Figure 4. The expression pattern of human and mouse IGF-1 minigenes.** (A) Schematic representation of IGF-1 minigenes. Exons (boxes), introns (lines), and cryptic 5' splice site (c5'ss) in exon 5 are indicated (left panel). A diagram of primers position used for RT-PCR analysis and the expected PCR products are also shown (right panel). A minigene-specific forward primer (P1) was used to amplify minigene-derived IGF1 isoforms. Primers P1-P4 were used to amplify a constant region of minigene-derived IGF-1 mRNA; primers P1-P3 were used to amplify IGF-1Ea and IGF-1Ec isoforms;~~

1 primers P1-P2 were used to amplify IGF-1Eb isoform. All oligonucleotide sequences are listed in Supplementary Table  
2 S1. (B) RT-PCR analysis of splicing assay performed in human HEK293T cells (*left panel*) and primary mouse  
3 hepatocytes (*right panel*) in presence of human (hIGF-1) or mouse (mIGF-1) minigenes. Untransfected cells (-) or cells  
4 transfected with the empty vector (vect.) were used as a PCR control. L34 was used as loading control.

5 ~~**Figure 5. The expression pattern of human and mouse IGF-1 minigenes.**~~ (A) Schematic representation of IGF-1  
6 minigenes. Exons (*boxes*), introns (*lines*), and cryptic 5' splice site (c5'ss) in exon 5 are indicated (*upper panel*). A  
7 diagram of primers position used for RT-PCR analysis is also shown (*lower panel*): primers P1-P3 were used to amplify  
8 mouse and human IGF-1Ea and IGF-1Ec isoforms; primers P1-P2 were used to amplify human and mouse IGF-1Eb  
9 isoform; primers P1-P2-P3 were used to amplify human IGF-1Ea and IGF-1Eb isoforms. All oligonucleotide sequences  
10 are listed in Supplementary Table S1. (B and C) RT-PCR analysis of *in vivo* splicing assay performed in human  
11 HEK293T cells (B) and primary mouse hepatocytes (C) in presence of human (hIGF-1) or mouse (mIGF-1) minigenes.

12  
13  
14 **Figure 5. The IGF-1 splicing pattern of human and mouse minigenes after the splicing factors overexpression**  
15 **(A,B) and effect of siRNA-mediated depletion of splicing factors on endogenous IGF-1 isoforms expression in**  
16 **human and mouse cells (C,D).** (A) Diagram of primers position used to amplify minigene-derived IGF1 isoforms. (B)  
17 RT-PCR analysis of splicing assay performed in HEK293T cell line in presence of human (*left panel*) or mouse (*right*  
18 *panel*) IGF-1 minigenes and the indicated splicing factors. The bar graph shows the ratio of densitometric analysis  
19 between IGF-1Eb and IGF-1Ea (*left panel*) or IGF-1Ec and IGF-1Ea (*right panel*) isoforms (mean  $\pm$  SD, n=3). C)  
20 Diagram of primers position used to amplify endogenous gene-derived IGF-1 isoforms. (D) RT-PCR analysis of splicing  
21 assay of endogenous IGF-1 gene performed in human LNCaP cells (*left panel*) and in mouse NIH/3T3 cells (*left panel*)  
22 depleted with the indicated splicing factors.

23 ~~**Figure 6. The IGF-1 splicing pattern of human and mouse minigenes after the splicing factors overexpression.**~~  
24 A,B) RT-PCR analysis of splicing assay of human (A) or mouse (B) IGF-1 minigenes performed in HEK293T cell line in  
25 presence of indicated splicing factors. Primers P1-P3 were used to amplify mouse and human IGF-1Ea and IGF-1Ec  
26 isoforms (A,B). Primers P1-P2-P3 were used to amplify human IGF-1Ea and IGF-1Eb isoforms (A); primers P1-P2 were  
27 used to amplify mouse IGF-1Eb isoform (B); primers P1-P4 were used to amplify total human and mouse *IGF-1*. The  
28 primer positions are indicated in Fig. 5. The bar graph shows the ratio of densitometric analysis between IGF-1Eb and  
29 IGF-1Ea (A) or IGF-1Ec and IGF-1Ea (B) isoforms (mean  $\pm$  SD, n=3).

30  
31  
32 **Figure 7. Figure 6. Comparative analysis of human and mouse IGF-1 alternative splicing on wild-type and**  
33 **mutated minigenes.** (A) Scheme of human (hMUT: TGvsGA) and mouse (mMUT: GAvsTG) mutant minigenes. The  
34 mutated bases are underlined. Position of the canonical (c5'ss) and non-canonical (\*5'ss) cryptic 5' splice site used in  
35 mMUT minigene are also shown. (B and C) RT-PCR analysis of splicing assay performed in human HEK293T cells in  
36 presence of human (B) and mouse (C) IGF-1 mutant minigenes. RT-PCR was performed using primers as indicated in  
37 Fig. 4.

38  
39  
40 **Figure 7 8. The sliding-window plot of the disorder conservation score of IGF-1 domains.** The x axis represents  
41 the codon positions and the y axis shows disorder conservation profiles for IGF-1 core (black line), Ea (dashed line), Eb  
42 (dotted line) and Ec (gray line) peptides. The disorder conservation score for IGF-1 domains was calculated using the  
43 Disorder Conservation (DisCons) software with default parameters [35]. Low values correspond to order, high values to  
44 disorder.

## Supplementary data captions

**Supplementary Figure S1.** Partial sequence alignment of IGF-1 genomic sequences among 23 amniota vertebrata retrieved from Ensembl database (23 amniota vertebrates Pecan). (A) The nucleotide sequences surrounding the alternatively spliced exon 6 were shown. The exon 6 sequences are in upper case letters and highlighted in gray, the intron sequences are in lower case letters. The 3' splice site of exon 6 is indicated by bold letters. Sequence gaps are indicated by dashed lines. Human chromosomal location:GRCh38:12:102537063:102537363. (B) The nucleotide sequences surrounding the alternatively spliced exon 5 were shown. The exon 5 sequences are in upper case letters and highlighted in black and gray: the black nucleotides designated a part of exon 5 common to IGF-1Eb and IGF-1Ec variants while the grey region is included only in the IGF-1Eb isoform; the intron sequences are in lower case letters. The 3' splice site and the cryptic 5' splice site (c5'ss) of exon 5 are shown in bold. Sequence gaps are indicated by dashed lines. Human chromosomal location:GRCh38: 12:102458404:102458924. Ensembl Genomes: Human (*Homo sapiens*) GRCh38.p5; Mouse (*Mus musculus*) GRCh38.p4; Rat (*Rattus norvegicus*) Rnor\_6.0; Rabbit (*Oryctolagus cuniculus*) OryCun2.0; Chimpanzee (*Pan troglodytes*) CHIMP2.1.4 ; Gorilla (*Gorilla gorilla gorilla*) gorGor3.1; Orangutan (*Pongo abelii*) PPYG2; Macaque (*Macaca mulatta*) MMUL 1.0 ; Marmoset (*Callithrix jacchus*) C\_jacchus3.2.1; Cat (*Felis catus*) Felis\_catus\_6.2; Dog (*Canis lupus familiaris*) CanFam3.1; Horse (*Equus caballus*) Equ Cab 2; Cow (*Bos taurus*) UMD3.1; Sheep (*Ovis aries*) Oar\_v3.1; Pig (*Sus scrofa*) Sscrofa10.2; Opossum (*Monodelphis domestica*) monDom5 ; Platypus (*Ornithorhynchus anatinus*) OANA5; Turkey (*Meleagris gallopavo*) Turkey\_2.01; Chicken (*Gallus gallus*) Galgal4; Zebra Finch (*Taeniopygia guttata*) taeGut3.2.4; Anole lizard (*Anolis carolinensis*) AnoCar2.0.

**Supplementary Figure S1. ClustalX alignment of IGF-1 sequences along the IGF-1 exon 6 and exon 5 regions among 23 representative amniota vertebrates.** > Alignment of IGF-1 sequences along the exon 6 region of 23 amniota vertebrata (chromosome:GRCh38:12:102537063:102537363:-1) > Alignment of IGF-1 sequences along the exon 5 region of 23 amniota vertebrata (chromosome:GRCh38:12:102458404:102458924:-1)

**Supplementary Figure S2.** Analyses of the 3' and the 5' splice-site signals of the IGF-1 exon 5 among 58 mammalian genomic sequence datasets available at the UCSC genome browser. Evolutionary tree of 58 organisms and the reconstructed ancestral state was adapted from the UCSC Genome Browser. Strength of 3' and 5'ss in exon5 obtained using the Maximum Entropy scores is shown for each species [Yeo G. and Burge CB. *Comput. Biol.* (2004);11:377-394]. Species with relative strong exon 5 5'ss are boxed. n.d.= not determined.

**Supplementary Figure S2. The strength of the IGF-1 exon 5 splice sites among 58 mammalian genomic sequence datasets available at the UCSC genome browser.** The phylogenetic tree for the 58 organisms was adapted from the UCSC Genome Browser. The splice site scores were obtained for each 5' and 3' splice site sequence using the Maximum Entropy scores [Yeo G. and Burge CB. *Comput. Biol.* (2004);11:377-394]. Species with relative strong exon 5 5'ss are in highlighted bold in the table and boxed in the phylogenetic tree.

**Supplementary Figure S3.** Species-Specific differences in expression pattern of the IGF-1 isoforms among mammals. (A, B and C) 3' RACE-PCR sequences obtained from liver tissues of mammalian species (see Figure 3B). 3' RACE-PCR products were sequenced and the nucleotide sequences corresponding to IGF-1Ea (A), IGF-1-1Eb (B) and IGF-1Ec (C) were aligned using Clustal W [Thompson J.D. et al. *Curr Protoc Bioinformatics* Chapter 2 (2002) Unit 23]. Poly(A) Signal Miner (<http://dnafsminer.bic.nus.edu.sg/PolyA.html>) was used to predict poly(A) signal and, if detected in the sequences, highlighted in bold. PCR product sizes (nucleotides) for each sequence are indicated in parentheses. The nucleotide sequences matched to the following NCBI Reference Sequence (RefSeq): IGF-1Ea (Pig 350 and 283: NM\_214256; Rabbit 537: XM\_008256720/ XM\_008256719; Macaque 537: NM\_001260726; Human 627 and 534: NM\_000618/ NM\_001111284; Mouse 391: NM\_001314010/ NM\_001111275/ NM\_001111276; Rat 381: NM\_178866/ NM\_001082479); IGF-1Eb (Pig 548 and 490: XM\_005664196/ XM\_005664195; Rabbit 725 and 678: XM\_008256717/ XM\_008256716; Macaque 685: NC\_027903; Human 734, 688 and 491: NM\_001111285); IGF-1Ec (Rabbit 589: XM\_008256718; Macaque 589: XM\_015152532/ XM\_015152534; Pig 401: XM\_005664197; Rat 466: NM\_001082477/ NM\_001082478; Mouse 539 and 450: NM\_010512/ NM\_001111274). (D) Amplification of IGF-1Eb and IGF-1Ec isoforms in liver tissue of goat, sheep, pig and cow with isoform-specific PCR strategy. Map of IGF-1 gene with the relative position of primers used to separately amplify IGF-1 isoforms is shown on the left. The IGF-1F forward primer is common to all mammals (5'-CGTGGATGAGTGCTGCTTC-3'); the IGF-1R1 reverse primer is specific for IGF-1Eb amplification in goat, sheep and cow (5'-TCCTTCTGTTCCCTCCTGG-3'); the IGF-1R2 reverse primer is specific for

1 pig IGF-1Eb isoform (5'-CCCTCCTGGGTGTTTCTTTG-3'); the IGF-1R3 reverse primer is specific for IGF-1Ec  
2 amplification in all mammals (5'-CTTCAAATGTACTTCCTTTCC-3'). The conventional PCR was carried out as described  
3 in Material and Methods. The PCR products were loaded on 4.0% agarose gel and DNA fragments were eluted from gel,  
4 subcloned and sequenced. The nucleotide sequences matched to the following NCBI Reference Sequence (RefSeq):  
5 IGF-1Eb (Goat: XM\_013963970/ XM\_013963971/ XM\_005680531/ XM\_005680532/ XM\_005680533/ XM\_005680534/  
6 XM\_005680535; Sheep: NC\_019460; Cow; XM\_015471061/ XM\_015471062/ XM\_015471063/ XM\_015471064; Pig:  
7 XM\_005664196/XM\_005664195); IGF-1Ec (Goat: NC\_022297; Sheep: NC\_019460; Cow: XM\_005206490; Pig:  
8 XM\_005664197).

9 **Supplementary Figure S3. The IGF-1 isoform expression pattern of several mammalian species analysed by 3'**  
10 **RACE-PCR (Figure S3A) and isoform-specific PCR strategy (Figure S3B). Supplementary Figure S3B.** Expression  
11 analysis of IGF-1Eb and IGF-1Ec splice variants in liver tissue of goat, sheep, pig and cow. The schematic  
12 representation of the exon/intron structure of IGF-1 gene and position of primers used to separately amplify IGF-1  
13 isoforms were shown for IGF-1Eb (left) and IGF-1Ec (right). The PCR products were loaded on 4.0% agarose gel and  
14 DNA fragments were eluted from gel, subcloned and sequenced. The following primers were used to amplified the IGF-  
15 1Eb and IGF-1Ec isoforms in liver tissues of different mammalian species: IGF-1F (forward primer common to all  
16 mammals): 5'-CGTGGATGAGTGCTGCTTC-3'; and the following reverse primers for IGF-1Eb amplification: IGF-1R4  
17 (goat, sheep and cow): 5'-TCCTTCTGTTCCCTCCTGG-3' and IGF-1R2 (pig): 5'-CCCTCCTGGGTGTTTCTTTG-3'; and  
18 for IGF-1Ec amplification: IGF-1R3 (common to all mammals): 5'-CTTCAAATGTACTTCCTTTCC-3'.

19  
20  
21 **Supplementary Figure S4.** The expression pattern of human and mouse IGF-1 minigenes in different human and  
22 mouse cell lines. RT-PCR analysis of splicing assay of human (hIGF-1; *upper panel*) and mouse (mIGF-1, *lower panel*)  
23 IGF-1 minigenes performed in the indicated human and mouse cell lines. Untransfected cells (-) or cells transfected with  
24 the empty vector (vect.) were used as a PCR control. L34 was used as loading control. A diagram of primers position  
25 used for RT-PCR analysis is also shown.

26  
27  
28 **Supplementary Figure S5.** Cis-acting elements contribute to species-specific variation of IGF-1 alternative splicing.  
29 (A,B) Representative Western-blot and RT-PCR analysis showing the expression level of hnRNP, SR and SR-like  
30 proteins in the splicing assay experiments.  $\beta$ -actin was used as loading control. C) RT-PCR analysis of splicing assay of  
31 human (hIGF-1) and mouse (mIGF-1) IGF-1 minigenes performed in HEK293T cells silenced with the indicated siRNAs.  
32 A representative western-blot showing the expression level of the silenced proteins is also shown.  $\beta$ -actin was used as  
33 loading control. D) RT-PCR analysis of splicing assay of endogenous human IGF-1 gene performed in LNCaP treated for  
34 3hrs with 100 $\mu$ M Sodium Arsenite ( $\text{Na}_3\text{AsO}_3$ ).

35  
36  
37 **Supplementary Figure S6.** (A) RT-PCR analysis of splicing assay performed in HEK293T cell line of mouse and human  
38 wild-type (WT), hMUT (TGvsGA) (*left panel*) and mMUT (GAvsTG) (*right panel*) IGF-1 minigenes in presence or not of  
39 SRSF1. The ratio between IGF-1Ec/IGF-1Ea (*left panel*) and IGF-1Eb/IGF-1tot (*right panel*) is also indicated. n.d.= not  
40 determined. A diagram of primers position used for RT-PCR analysis is also shown. (B) Representation of SRSF1  
41 binding site (black box) in IGF-1 exon 5 identified by Smith P.J. et al. [Endocrinology (2002) 143:146-154] and Cléry A. et  
42 al. [PNAS (2013) 110:E2802-E2811]. Canonical cryptic 5' splice site (c5'ss) is also indicated.

43  
44  
45 **Supplementary file S1.** IGF-1 isoform nucleotide sequences of the 27 mammalian species used for IGF-1 evolutionary  
46 analysis. The "CDS FASTA alignment from multiple alignments" data, derived from the "multiz100way" alignment data  
47 prepared from 100 vertebrate genomes [Blanchette M. et al. Genome Res (2004), 14, 708-715], were downloaded using  
48 the Table Browser tool of the UCSC Genome Browser. Sequences were subsequently realigned using MUSCLE [Edgar  
49 R.C. Nucleic Acids Res (2004), 32, 1792-1797] and protein coding sequences from 27 mammalian species were  
50 extracted from these alignment datasets. The phylogenetic tree in Newick format was constructed using the IGF-1 core  
51 alignment of 27 mammalian species and is shown at the end of supplementary file S1.

1 **Table 1. Substitution rates of IGF-1 domains among 27 mammalian species**

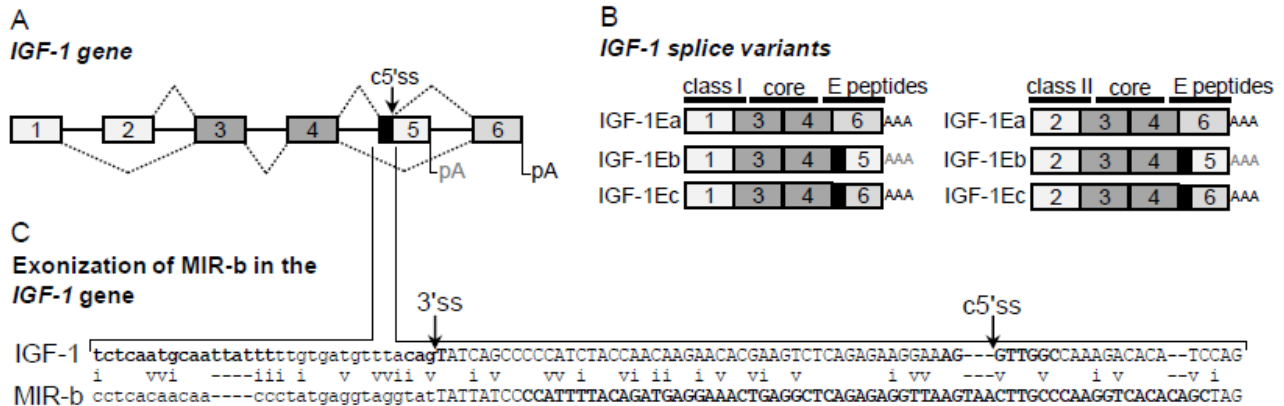
IGF-1 domain	Average length (codons)	dS	dN	dN-dS
IGF-1 core	70.0	0.33 (0.23;0.44) <sup>a</sup>	0.01 (0.01;0.02) <sup>a</sup>	-0.31 (-0.42;-0.22) <sup>a*</sup>
Ea peptide	35.0	0.13 (0.09;0.18) <sup>b</sup>	0.02 (0.01;0.02) <sup>a</sup>	-0.12 (-0.16;-0.08) <sup>b*</sup>
Eb peptide	63.3	0.09 (0.06;0.12) <sup>b</sup>	0.05 (0.03;0.06) <sup>b</sup>	-0.04 (-0.06;-0.02) <sup>c*</sup>
Ec peptide	41.1	0.06 (0.04;0.09) <sup>b</sup>	0.02 (0.02;0.03) <sup>a</sup>	-0.04 (-0.07;-0.02) <sup>c*</sup>

2 All rates are estimated by maximum likelihood using the MG94xREV\_3x4 (local) codon substitution model. Mean (95%  
 3 confidence intervals) and *P* values for the neutrality tests (dN-dS expected value=0) are estimated by using parametric  
 4 bootstrap based on 1000 replicates. \**P*<0.05. Multiple comparisons among IGF-1 domains were performed with Analysis  
 5 of variance bootstrap based. Values in the same column bearing the same letter are not significantly different. The  
 6 difference dN-dS was used in place of a more common ratio dN/dS, to avoid numerical issues when dS is zero, which is  
 7 possible because HyPhy permits synonymous substitution rates to vary from site to site.

8  
 9

1  
2  
3

FIGURE 1



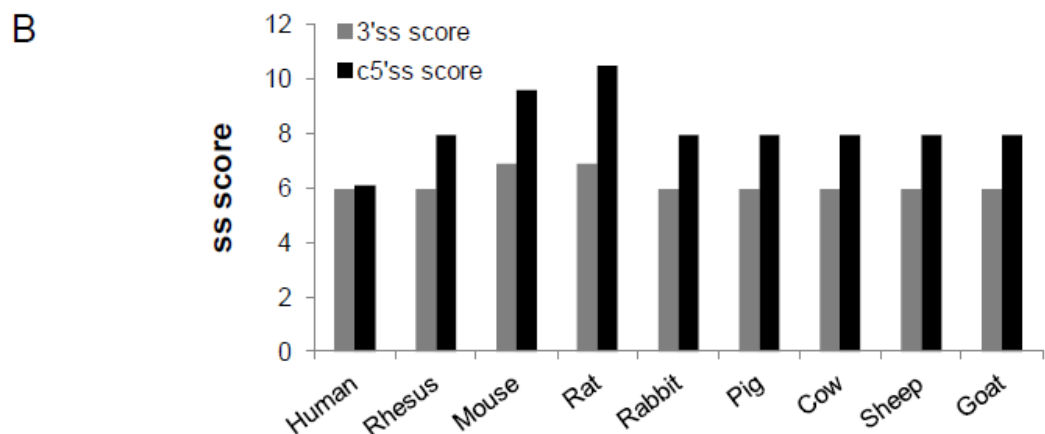
4  
5  
6  
7  
8  
9

FIGURE 2

**A**

	Sequence	MIR-b identity (%)
Human	cagTATCAGCCCCCATCTACCAACAAGAACACGAAGTCTCAGAG---AAGGAAAGGTGGGC	59.2
Macaque	cagTATCAGCCCCCATCTACCAACAAGAACACGAAGTCTCAGAGGAGAAGGAAAGGTGGGC	60.2
Mouse	cagTCCCGTCCCTATCGACAAAACAAGAAAACGAAGCTGCAAAGGAGAAGGAAAGGTGAGC	56.3
Rat	cagTCCAGCCCCATCTCGACACACAAGAAAAGGAAGCTGCAAAGGAGAAGGAAAGGTGAGT	58.3
Rabbit	cagTATCAGCCTCCATCTACCAACAAGAAAATGAAGTCTCAGAGGAGAAGGAAAGGTGGGC	56.3
Pig	cagTATCAGCCCCCATCTACCAACAAGAAAACGAAGTCTCAGAGGCGAAGGAAAGGTGGGC	60.2
Cow	cagTATCAGCCCCCATCTACCAACAAGAAAATGAAGTCTCAGAGGAGAAGGAAAGGTGGGC	62.1
Sheep	cagTATCAGTCCCATCTACCAACAAGAAAATGAAGTCTCAGAGGAGA-GGAAAGGTGGGC	61.2
Goat	cagTATCAGTCCCATCTACCAACAAGAAAATGAAGTCTCAGAGGAGAAGGAAAGGTGGGC	61.2
MIR-b	tatTATTATCCCATTTTACAGATGAGGAAACTGAGGCTCAGAGAGGTTAAGTAACCTGGCC	

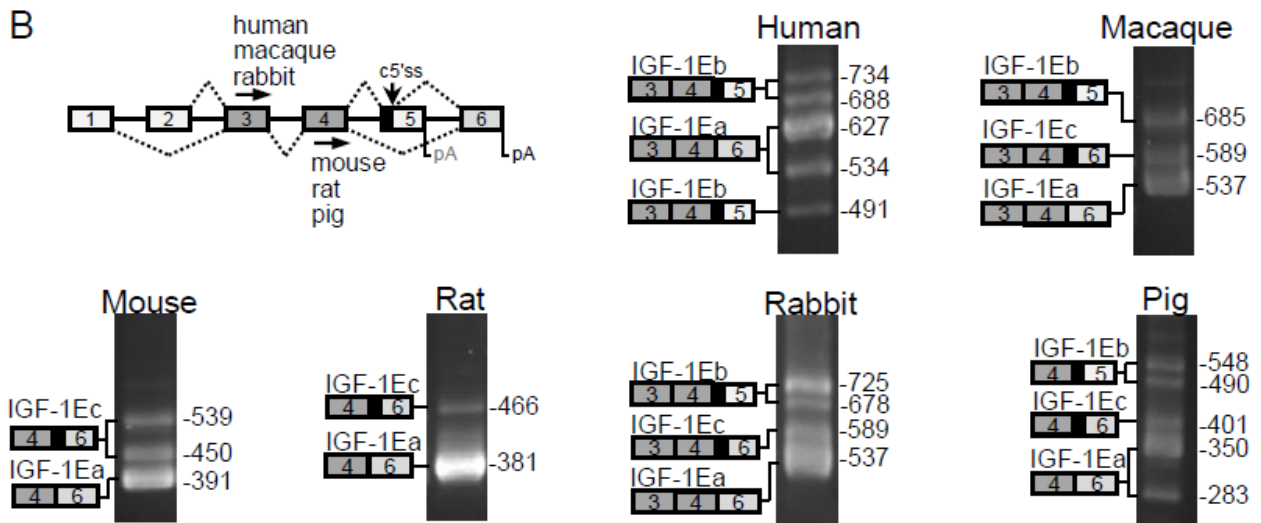
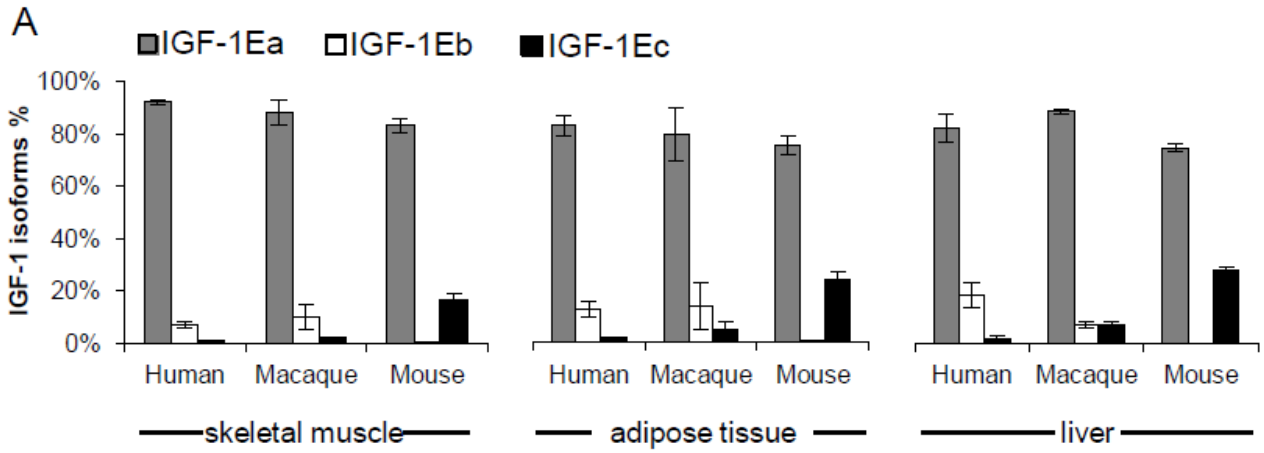
Conserved regions are highlighted in grey. Asterisks (\*) indicate conserved positions. The 3'ss site is `yagG` and the c5'ss site is `AGgtragt`.



10  
11  
12  
13

1  
2  
3

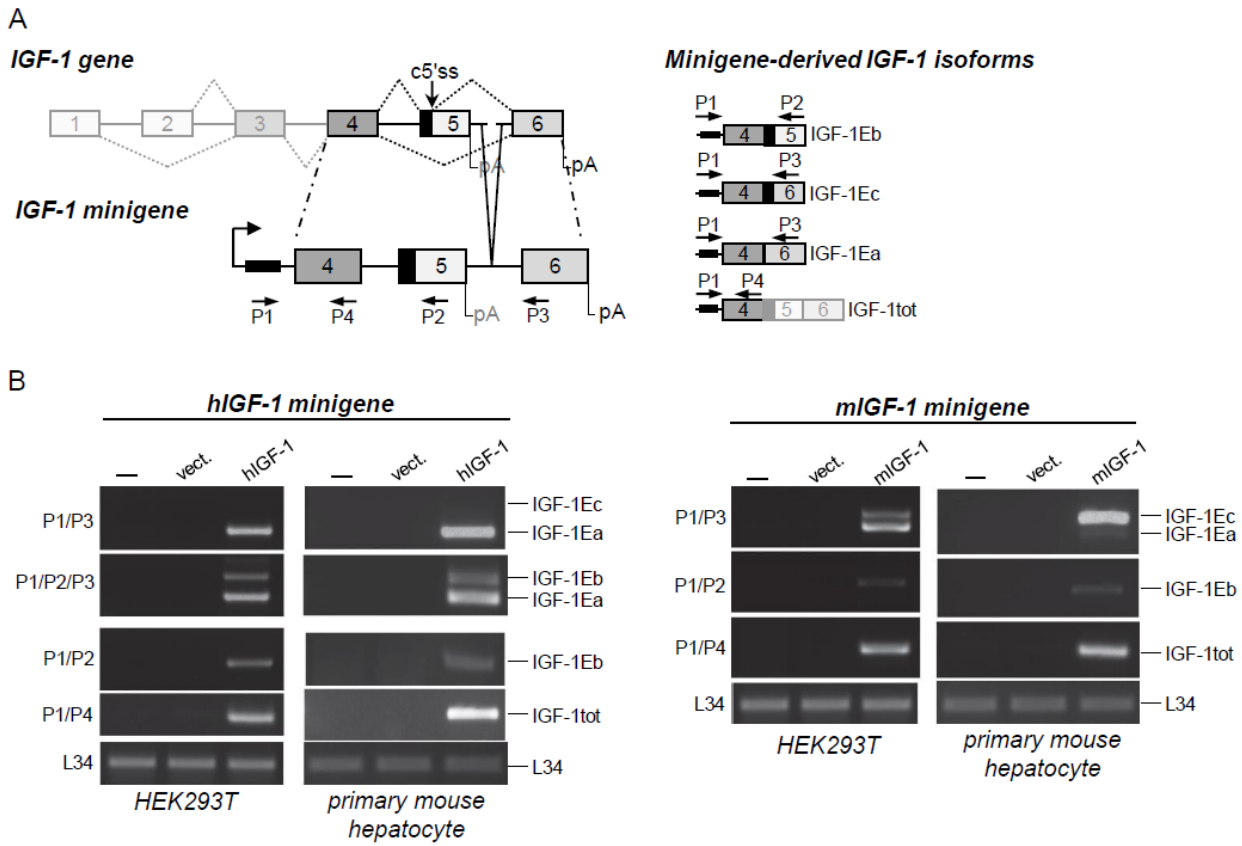
FIGURE 3



4  
5  
6  
7

1  
2  
3

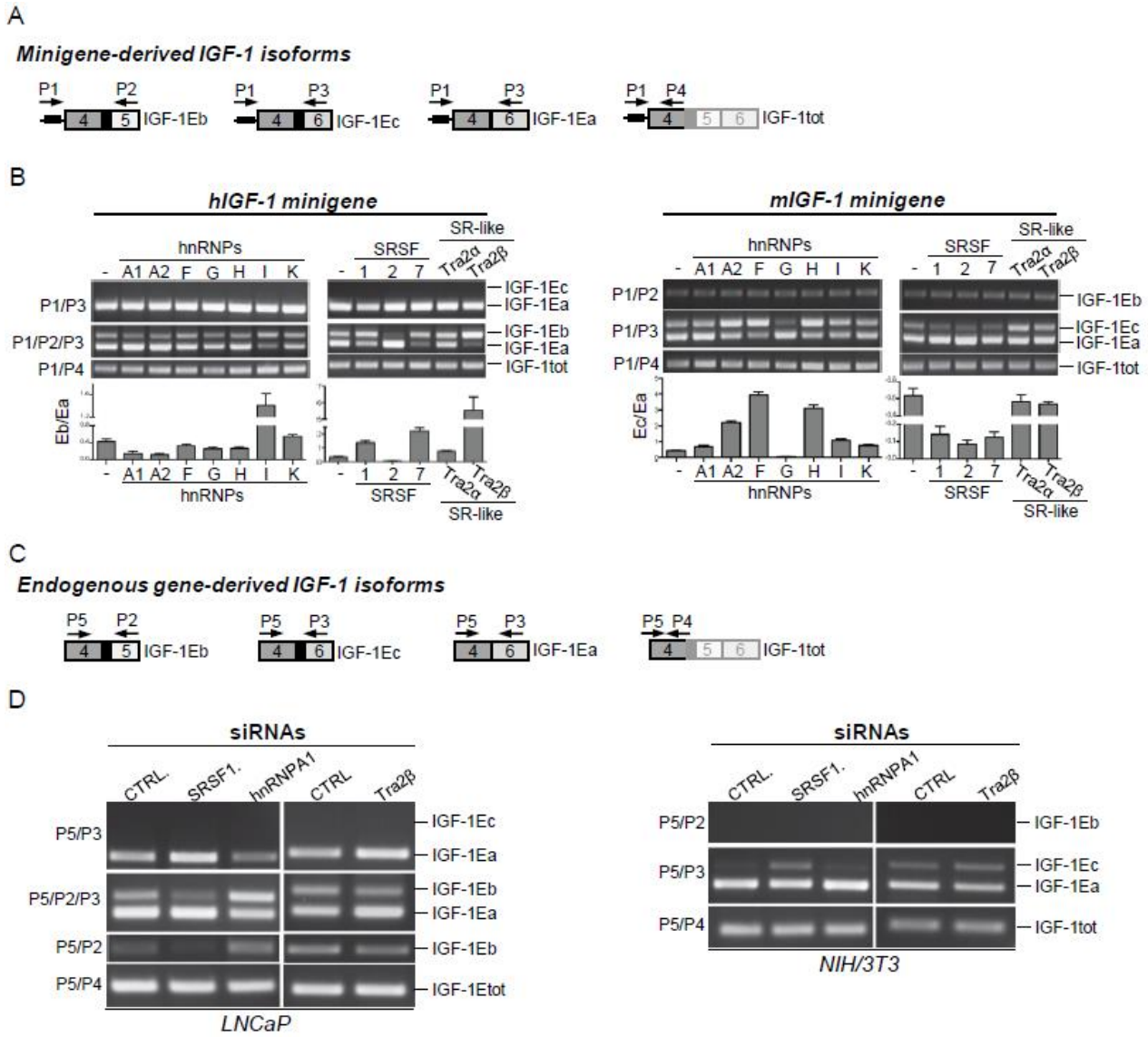
FIGURE 4



4  
5  
6

1  
2

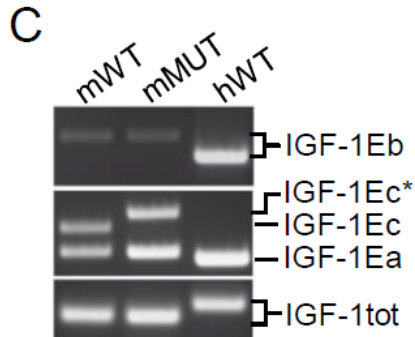
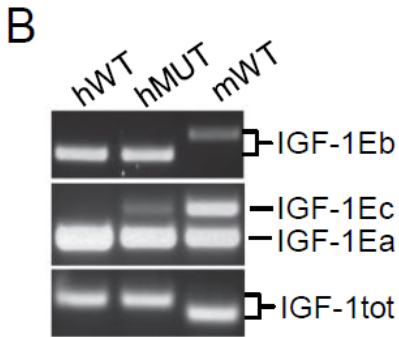
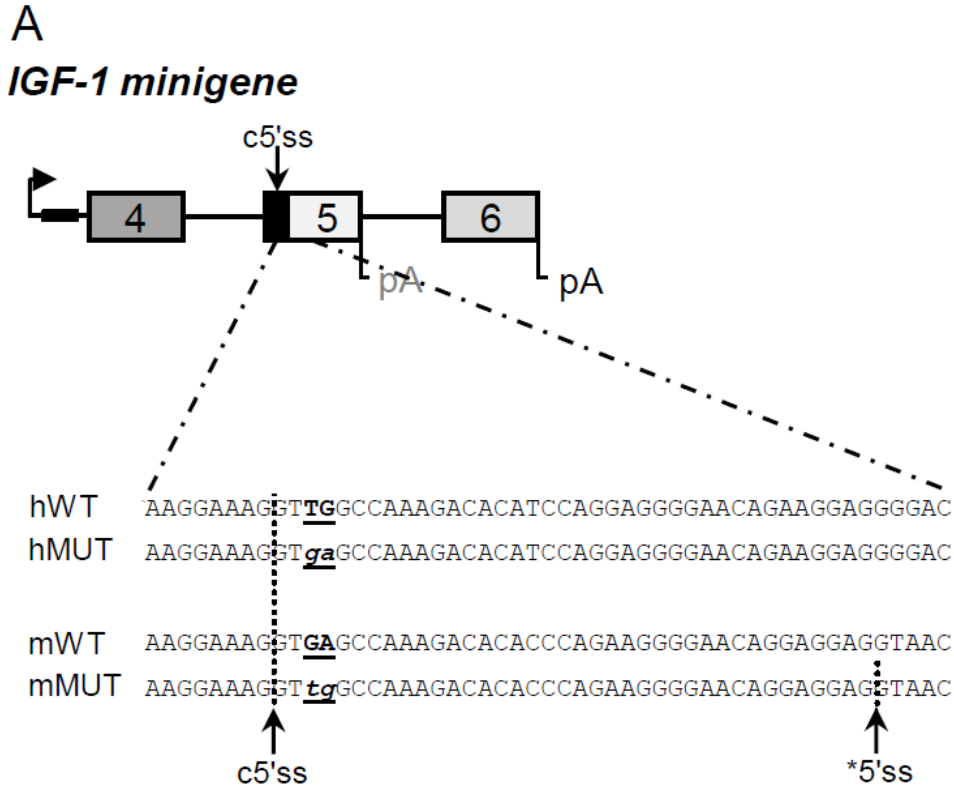
FIGURE 5



3  
4  
5

FIGURE 6

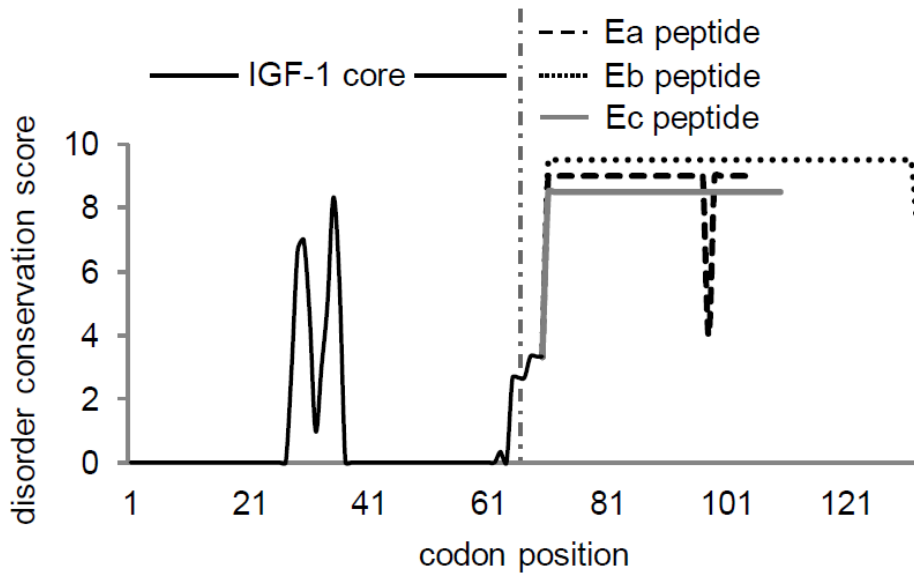
1  
2  
3  
4  
5



6  
7  
8

1  
2  
3  
4

FIGURE 7



5