



UNIVERSITY OF LEEDS

This is a repository copy of *Characterization and Genomic Localization of a SMAD4 Processed Pseudogene*.

White Rose Research Online URL for this paper:  
<http://eprints.whiterose.ac.uk/121101/>

Version: Accepted Version

---

**Article:**

Watson, CM, Camm, N, Crinnion, LA et al. (6 more authors) (2017) Characterization and Genomic Localization of a SMAD4 Processed Pseudogene. *Journal of Molecular Diagnostics*, 19 (6). pp. 933-940. ISSN 1525-1578

<https://doi.org/10.1016/j.jmoldx.2017.08.002>

---

(c) 2017, Published by Elsevier Inc. on behalf of the American Society for Investigative Pathology and the Association for Molecular Pathology. This manuscript version is made available under the CC BY-NC-ND 4.0 license  
<https://creativecommons.org/licenses/by-nc-nd/4.0/>

**Reuse**

This article is distributed under the terms of the Creative Commons Attribution-NonCommercial-NoDerivs (CC BY-NC-ND) licence. This licence only allows you to download this work and share it with others as long as you credit the authors, but you can't change the article in any way or use it commercially. More information and the full terms of the licence here: <https://creativecommons.org/licenses/>

**Takedown**

If you consider content in White Rose Research Online to be in breach of UK law, please notify us by emailing [eprints@whiterose.ac.uk](mailto:eprints@whiterose.ac.uk) including the URL of the record and the reason for the withdrawal request.

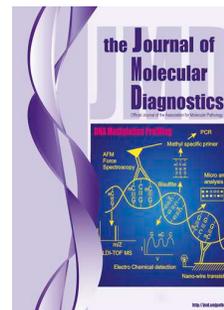


[eprints@whiterose.ac.uk](mailto:eprints@whiterose.ac.uk)  
<https://eprints.whiterose.ac.uk/>

# Accepted Manuscript

Characterization and Genomic Localization of a *SMAD4* Processed Pseudogene

Christopher M. Watson, Nick Camm, Laura A. Crinnion, Agne Antanaviciute, Julian Adlard, Alexander F. Markham, Ian M. Carr, Ruth Charlton, David T. Bonthron



PII: S1525-1578(17)30316-1

DOI: [10.1016/j.jmoldx.2017.08.002](https://doi.org/10.1016/j.jmoldx.2017.08.002)

Reference: JMDI 633

To appear in: *The Journal of Molecular Diagnostics*

Accepted Date: 16 August 2017

Please cite this article as: Watson CM, Camm N, Crinnion LA, Antanaviciute A, Adlard J, Markham AF, Carr IM, Charlton R, Bonthron DT, Characterization and Genomic Localization of a *SMAD4* Processed Pseudogene, *The Journal of Molecular Diagnostics* (2017), doi: 10.1016/j.jmoldx.2017.08.002.

This is a PDF file of an unedited manuscript that has been accepted for publication. As a service to our customers we are providing this early version of the manuscript. The manuscript will undergo copyediting, typesetting, and review of the resulting proof before it is published in its final form. Please note that during the production process errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

1 Characterisation and genomic localisation of a *SMAD4* processed pseudogene

2

3 Characterisation of a *SMAD4* pseudogene

4

5 Christopher M. Watson,\*†‡ Nick Camm,\* Laura A. Crinnion,\*†‡ Agne Antanaviciute,†

6 Julian Adlard,\* Alexander F. Markham,† Ian M. Carr,†‡ Ruth Charlton,\* and David T.

7 Bonthron\*†‡

8

9 From the Yorkshire Regional Genetics Service,\* St. James's University Hospital, Leeds;  
10 the MRC Medical Bioinformatics Centre,† Leeds Institute for Data Analytics, and the  
11 MRC Single Cell Functional Genomics Centre,‡ University of Leeds, St. James's University  
12 Hospital, Leeds, United Kingdom

13

14 Corresponding author:

15 Dr. Christopher M. Watson

16 6.2 Clinical Sciences Building

17 Yorkshire Regional Genetics Service

18 St James's University Hospital

19 Leeds, LS9 7TF

20 United Kingdom

21 Email: c.m.watson@leeds.ac.uk

22

23 This work was supported by grants MR/M009084/1 and MR/L01629X/1 awarded by  
24 the UK Medical Research Council.

25 **ABSTRACT**

26 Like many clinical diagnostic laboratories, we undertake routine investigation of  
27 cancer-predisposed individuals by high-throughput sequencing of patient DNA that has  
28 been target-enriched for genes associated with hereditary cancer. Accurate diagnosis  
29 using such reagents requires alertness against rare non-pathogenic variants that may  
30 interfere with variant calling. In a cohort of 2,042 such cases, we identified five that  
31 initially appeared to be carriers of a 95-bp deletion of *SMAD4* intron 6. More detailed  
32 analysis indicated that these individuals all carried one copy of a *SMAD4* processed  
33 gene. Because of its interference with diagnostic analysis, we characterized this  
34 processed gene in detail. Whole genome sequencing and confirmatory Sanger  
35 sequencing of junction PCR products were used to show that in each of the five cases,  
36 the *SMAD4* processed gene was integrated at the same position on chromosome 9,  
37 located within the last intron of the *SCAI* gene. This rare polymorphic processed gene  
38 therefore reflects the occurrence of a single ancestral retrotransposition event.  
39 Compared to the reference *SMAD4* mRNA sequence NM\_005359.5  
40 (<https://www.ncbi.nlm.nih.gov/nucleotide/>), the 5' and 3' UTR regions of the processed  
41 gene are both truncated, but its open reading frame is unaltered. Our experience leads  
42 us to advocate the use of an RNA-seq aligner, as part of diagnostic assay quality  
43 assurance, since this allows their recognition in a comparatively facile automated  
44 fashion.

45

46

47 **INTRODUCTION**

48 The availability of diagnostic molecular genetic assays has increased significantly in  
49 recent years. This has largely been due to the ubiquitous adoption of next generation  
50 sequencing (NGS) instruments, which have replaced comparatively low-throughput  
51 Sanger sequencing technology, as the standard technique for mutation detection. New  
52 laboratory assays combined with ever-increasing automation is resulting in increased  
53 patient throughput and more efficient workflows. That several genes can be analysed  
54 concurrently has enabled an expansion of testing for heterogeneous genetic disorders  
55 which may have previously been considered too rare for a *bona-fide* genetic test to have  
56 been established and offered in a routine clinical laboratory. To be able to request a  
57 comprehensive analysis of all genes that correspond to a patient's phenotype is  
58 transforming diagnostic referral pathways, by eliminating costly 'test and review'  
59 processes that are necessary when referrals are made in a consecutive manner.

60  
61 Operational requirements associated with test portfolios that can accommodate varying  
62 combinations of target genes have necessitated a fundamental transformation in assay  
63 design. Typically, a far larger range of targets are selected for sequencing than is  
64 suggested *a priori* from the patient's presenting phenotype. An *in silico* virtual gene  
65 panel is applied to these data thus masking inappropriate results from those requested  
66 by the referring clinician. Although this approach generates unnecessary sequence data,  
67 laboratories are able to reduce the complexity of wet-laboratory processes thereby  
68 streamlining their workflows. As the cost of DNA sequencing continues to fall the  
69 number of genes that can be feasibly targeted, while maintaining iteratively comparable  
70 test sensitivity, will continue to increase. Indeed, our originally reported 36-gene

71 reagent has been periodically revised and presently targets the coding exons of 155  
72 cancer-associated genes [1].

73

74 For many commentators, the long-held aspiration that custom-designed panels will be  
75 replaced by exome- and subsequently whole genome-sequencing, is being expedited by  
76 large-scale, population based, sequencing projects. Nevertheless, the prevailing  
77 approach for performing target enrichment, using probe-based hybridisation, has  
78 overcome the need to design and optimise long-range PCR amplicons [2]. This has  
79 improved the scalability of targeted loci, as previously only a finite number of long-  
80 range PCR primer pairs could be handled by a single laboratory. Despite this advance,  
81 hybridisation capture methods have a lower specificity for target enrichment due to the  
82 capture of 'off-target' sequences. A comparatively greater number of reads it therefore  
83 required to achieve the same depth of coverage (although this is typically off-set by no  
84 longer needing to sequence a gene's introns).

85

86 Off-target sequences are captured for reasons that may include hybridisation of probes  
87 to low-diversity nucleotide sequences, sequence homology between the targeted region  
88 and that of a related gene family member or an interfering pseudogene, or reaction  
89 kinetics. Although off-target reads are typically ignored, a number of studies have  
90 demonstrated their utility for the inadvertent identification of single nucleotide  
91 polymorphisms [3] and as a source of low-coverage whole genome sequencing reads for  
92 genomewide copy-number analysis [4]. Less useful is the capture and sequencing of  
93 DNA fragments that are highly homologous to target loci; it is usually not possible to  
94 determine the true genomic origin of these resulting data. As pseudogene sequences

95 may therefore affect the interpretation of clinical assays their identification and  
96 characterisation is of particular importance to the diagnostic community.

97

98 A *SMAD4* processed pseudogene was recently detected in a subset of patients referred  
99 for diagnostic analysis of hereditary cancer predisposition genes [5]. *SMAD4* is  
100 associated with both juvenile polyposis syndrome (OMIM: 174900) and combined  
101 juvenile polyposis/hereditary hemorrhagic telangiectasia syndrome (OMIM: 175050).  
102 Here we corroborate this observation and assess the frequency of the *SMAD4*  
103 pseudogene in our cohort of 2,042 diagnostically referred hereditary cancer cases. We  
104 further define the genomic integration site and report the transcript structure following  
105 end-to-end sequencing of the identified *SMAD4* processed pseudogene.

106

107

108

109

110

111

112

113

114

115

116

117

118

119

120 **MATERIALS AND METHODS**

121 Patients were referred to the Leeds Genetics Laboratory for diagnostic testing of one or  
122 more hereditary cancer predisposition genes using a custom-designed SureSelect  
123 hybridisation enrichment assay (Agilent Technologies, Wokingham, UK). The original  
124 36-gene reagent has been iteratively redesigned since the service was launched in 2013  
125 [1] and now targets the exons and immediate flanking sequence of 155 hereditary  
126 cancer genes.

127  
128 DNA was isolated from blood lymphocytes using either a standard salting out method or  
129 the Chemagic™ 360 automated extractor (PerkinElmer, Seer Green, UK). For each  
130 sample, an Illumina-compatible sequencing library was generated. Initially, 3 µg of  
131 genomic DNA was sheared using a Covaris S2 or E220 (Covaris Inc., Woburn, MA, USA)  
132 before whole genome library preparation was undertaken using SureSelect XT reagents  
133 (Agilent Technologies, Wokingham, UK). This consisted of end-repair, (A)-addition,  
134 adaptor ligation and PCR enrichment. A custom RNA probeset was used to perform a  
135 targeted capture hybridisation on each of the whole genome libraries, following  
136 manufacturer's protocols throughout. Samples were initially prepared manually, but a  
137 fully automated solution has since been introduced using a Sciclone G3 liquid handling  
138 workstation (PerkinElmer, Seer Green, UK). The quality and concentration of final  
139 libraries were confirmed using either an Agilent Bioanalyser or Agilent TapeStation  
140 (Agilent Technologies, Wokingham, UK) before, typically, 16 samples were combined  
141 into a single batch for sequencing. Each batch was either sequenced on a single lane of  
142 an Illumina HiSeq2500 rapid-mode flow cell (2 × 101 bp sequencing reads) or pooled  
143 with two additional batches and sequenced on an Illumina NextSeq500 (2 × 151 bp  
144 sequencing reads) using a High Output flow cell using version 2 chemistry (Illumina

145 Inc., San Diego, CA, USA). Raw sequence data was converted to FASTQ.gz format using  
146 bcl2fastq v.2.17.1.14.  
147  
148 A common data processing pipeline, running on the Leeds high-performance computer  
149 MARC1 (<http://arc.leeds.ac.uk/systems/marc1/>), was applied to each of the per-sample  
150 directories from the SureSelect target enrichment assay. Initially, adaptor sequences  
151 and low-quality bases (Q score  $\leq 10$ ) were trimmed from reads using Cutadapt v.1.9.1  
152 (<https://github.com/marcelm/cutadapt>) [6]. The resulting analysis-ready reads were  
153 assessed using FastQC v.0.11.5  
154 (<http://www.bioinformatics.babraham.ac.uk/projects/fastqc/>). Reads were next  
155 aligned to an indexed human reference genome (hg19) using BWA MEM v.0.7.13  
156 (<https://sourceforge.net/projects/bio-bwa/files/>) [7] before being sorted by  
157 chromosome coordinate and having PCR duplicates marked by Picard v.2.1.1  
158 (<http://broadinstitute.github.io/picard/>) to create a processed.bam file. These data  
159 were realigned using ABRA v.0.97 (<https://github.com/mozack/abra>) [8] and the  
160 Genome Analysis Toolkit (GATK) v.3.6-0 was used to perform variant calling following  
161 best practice guidelines. This involved indel realignment, base quality score  
162 recalibration and variant calling using the Haplotypecaller to generate a per-sample VCF  
163 file [9]. These variant data were annotated using Alamut Batch Standalone v.1.4.4  
164 (database v.2016.03.04) (Interactive Biosoftware, Rouen, France). Coverage metrics  
165 were determined using the GATK walkers DepthOfCoverage, CallableLoci and  
166 CountReads. Visualisation of aligned sequence reads was performed with the  
167 Integrative Genome Viewer v.2.3.80 (<http://software.broadinstitute.org/software/igv/>)  
168 [10]. The analysis-ready reads for five samples with apparent *SMAD4* intron 6 deletions  
169 were aligned to an indexed hg19 reference genome annotated using GENCODE Release

170 26 using the RNA-seq aligner STAR v.2.5.3a with default settings  
171 (<https://github.com/alexdobin/STAR/>) [11].  
172  
173 Illumina-compatible whole genome sequencing libraries were subsequently prepared  
174 for the same five samples. Approximately 3 µg DNA was sheared using a Covaris S2  
175 prior to end-repair, (A)-addition and adaptor ligation steps being undertaken using  
176 NEBNext® Ultra™ reagents, following manufacturer's protocols (New England Biolabs,  
177 Ipswich, MA, USA). An ampure size selection ratio for a 300-bp to 400-bp insert and a 6-  
178 cycle enrichment PCR was performed. Following an assessment of library quality, the  
179 final libraries were pooled in equimolar concentrations and the pooled batch was  
180 sequenced using an Illumina NextSeq500 High Output flow cell generating 2 × 151 bp  
181 read lengths. For each sample, a processed.bam file was generated using the same  
182 bioinformatics pipeline described above. Sequence reads mapping to the *SMAD4* locus  
183 (chr18:48550000-48620000) were extracted from the coordinate-sorted duplicate-  
184 marked bam file using samtools v.0.1.18 with the options -q 1 and -F 14 [12]. These  
185 filters ensured that the mapped read quality score was greater than 0, that neither read  
186 in the pair was unmapped and that the pair was not considered to be a "proper pair".  
187 Read pairs with one read mapping outside the *SMAD4* locus and whose non-*SMAD4* read  
188 clustered within 500 bp of the nearest of non-*SMAD4* read were reviewed and  
189 compared between patients.  
190  
191 Three PCR amplicons were generated to amplify across the breakpoints identified by  
192 medium coverage whole genome sequencing. The specificity of the amplicons was  
193 evaluated for each reaction; one primer was located within *SCAI* intron 18, and the  
194 second primer within the *SMAD4* pseudogene.

195

196 Two amplicons spanning the 5' end of the *SMAD4* pseudogene were designed and  
197 amplified. The first comprised a common *SCAI*-bound forward primer 5'-  
198 CTGAGCTTGTGATCTGCCTG-3' and the *SMAD4* exon 2-located reverse primer 5'-  
199 TGAAGCCTCCCATCCAATGT-3'. Each PCR reaction consisted of 7.46 µl nuclease-free  
200 H<sub>2</sub>O, 1.2 µl, 10× Buffer + Mg, 0.12 µl dNTPs, 2.4 µl GC-rich buffer, 0.12 µl Faststart Taq  
201 polymerase, 0.1 µl 10 µM forward primer, 0.1 µl 10 µM reverse primer and 0.5 µl of  
202 approximately 100 ng/µl DNA (Roche Diagnostics Ltd., Burgess Hill, U.K.).  
203 Thermocycling conditions were 96°C for 5 minutes, followed by 35 cycles of 96°C for 30  
204 seconds, 55°C for 30 seconds, 72°C for 2 minutes and a final 72°C extension step for 10  
205 minutes. The second amplicon was amplified using the same common *SCAI*-bound  
206 forward primer and a reverse primer specific to *SMAD4* exon 8 5'-  
207 TGGAAATGGGAGGCTGGAAT-3'. PCR reagents and volumes were equivalent to the first  
208 reaction. Thermocycling conditions were the same, but an additional 5 cycles were  
209 performed. PCR products from the second reaction were gel-extracted and purified  
210 using a QIAquick column following manufacturer's protocols (Qiagen GmbH, Hilden,  
211 Germany). Sanger sequencing was performed on PCR products from both reactions  
212 using amplification primers and, for the second reaction, a further two internally sited  
213 *SMAD4* exon 2 primers (5'-TTCCTTGCAACGTTAGCTGT-3' and 5'-  
214 ACATTGGATGGGAGGCTTCA-3') with an ABI3730 following manufacturer's instructions  
215 (Applied Biosystems, Paisley, UK).

216

217 The 3' end of the *SMAD4* processed pseudogene was amplified using a forward primer  
218 bridging the *SMAD4* exon 5/6 junction 5'-ACAAGTCAGCCTGCCAGTAT-3' and an *SCAI*  
219 reverse primer 5'-CAGGAAACAGCTATGACCTGCAATGACTCGATCTCAGC-3'. The reverse

220 primer contained a universal tag (underlined) for Sanger sequencing using our routine  
221 diagnostic workflow. Each reaction consisted of 12.74 µl nuclease-free H<sub>2</sub>O, 2 µl  
222 SequaPrep™ 10× Reaction Buffer, 0.36 µl SequaPrep™ 5U/µl Long Polymerase, 0.4 µl  
223 dimethyl sulfoxide (DMSO), 2 µl SequaPrep™ 10× Enhancer A, 1 µl 10 µM forward  
224 primer, 1 µl 10 µM reverse primer and 0.5 µl of approximately 100 ng/µl DNA  
225 (Invitrogen, Paisley, UK). Thermocycling conditions comprised a denaturation step of  
226 94°C for 2 minutes, followed by 10 cycles of 94°C for 10 seconds, 60°C for 30 seconds  
227 and 68°C for 3 minutes then 25 cycles of 94°C for 10 seconds 60°C for 30 seconds, 68°C  
228 for 3 minutes with an additional 20 seconds added per cycle, before a final extension  
229 step at 72°C for 5 minutes. PCR products of approximately 2 kb were gel extracted and  
230 purified using the QIAquick column following the manufacturer's protocol. Sanger  
231 sequencing was performed using the amplification forward primer, universal reverse  
232 primer and the following internally sited primers: 5'-AGCCATTGAGAGAGCAAGGT -3'  
233 (SMAD4 exon 9/10 forward), 5'-CCTCCAGCTCCTAGACGAAG-3' (SMAD4 exon 12  
234 forward), 5'-CCATGTGGGTGAGTTAATTTTACC-3' (SMAD4 exon 12 forward), 5'-  
235 TGGAAATGGGAGGCTGGAAT-3' (SMAD4 exon 8 reverse), 5'-  
236 AAAGCAGCGTCACTCTACCT-3' (SMAD4 exon 12 forward) and 5'-  
237 TCAGTTTTTGTATCTTGGGGCA-3' (SMAD4 exon 12 forward).

238

239 Sequence chromatograms for all Sanger sequencing reactions were analysed using  
240 4Peaks v.1.8 (<http://nucleobytes.com/4peaks/index.html>).

241

## 242 **RESULTS**

243 Since 2013, we have used a custom-hybridisation enrichment assay and NGS pipeline  
244 for the diagnostic analysis of hereditary cancer genes [1]. In the present study, we

245 retrospectively examined 2,042 patient libraries that had been sequenced in 131  
246 batches. We noticed five cases in which our standard variant-calling pipeline identified  
247 an apparent 95-bp deletion, corresponding to the entire *SMAD4* intron 6 nucleotide  
248 sequence (c.787+1\_788-1del, NM\_005359.5,  
249 <https://www.ncbi.nlm.nih.gov/nucleotide/>). Assay performance metrics for each of  
250 these five libraries are displayed in Supplemental Table S1.

251

252 Visualisation of *SMAD4* read coverage charts for the five cases with an apparent intron 6  
253 deletion revealed plots with prominent ‘cliff-edge’ shaped profiles, the discontinuities in  
254 which aligned with the *SMAD4* exon-intron boundaries (Supplemental Figure S1). This  
255 was particularly conspicuous for *SMAD4* exon 8. Close inspection of these data  
256 established that reads at the exon-intron boundaries had been “soft-clipped”. To further  
257 investigate whether these soft-clipped reads spanned *SMAD4* exon-to-exon splice  
258 junctions, sequence reads were mapped to a transcript-annotated human genome, using  
259 the RNA-seq aligner STAR. Resulting Sashimi plots displaying splice junction read  
260 counts were consistent with the presence of a spliced *SMAD4* sequence whose exon  
261 structure matched that of the reference mRNA sequence NM\_005359.5  
262 (<https://www.ncbi.nlm.nih.gov/nucleotide/>) (Figure 1). These data thus suggested the  
263 presence of a processed (intron-lacking) *SMAD4* pseudogene in these five individuals.

264

265 The existence of such a pseudogene was indeed recently reported [5], although its  
266 structure was not characterized in detail. The frequency ( $5/2042 = 0.24\%$ ) of cases we  
267 observed carrying the *SMAD4* pseudogene was in keeping with that reported by Millson  
268 et al. ( $12/4672 = 0.26\%$ ). The likely interference of the pseudogene with diagnostic

269 testing prompted us to define its exact structure, and address the question of whether  
270 its sequence and location are identical among carriers.

271

272 To assess the relative number of copies of the *SMAD4* pseudogene, we determined the  
273 ratio of gapped (pseudogene-derived) to non-gapped (non-pseudogene) read  
274 alignments spanning intron 6. (Although this region was not specifically targeted when  
275 designing the capture enrichment probes, its small size and proximity to *SMAD4* exons 5  
276 and 6 ensured that the intron was fortuitously sequenced.) The ratio of gapped:non-  
277 gapped reads was approximately 1:2, suggesting that only a single copy of the  
278 pseudogene was present in each case (Table 1). Further, by comparing the normalised  
279 read-depths of these cases to controls from the same sequencing batches, we  
280 determined relative dosage values for each *SMAD4* exon. These results indicated the  
281 presence of three copies of most of the *SMAD4* exons, again indicative of a single copy of  
282 the *SMAD4* pseudogene (Supplemental Table S2). Although data for exons 4 and 8  
283 deviate from this interpretation, this is probably due to the small genomic intervals  
284 represented by these exons (30 bp and 51 bp, respectively). Additionally, the greater  
285 variability displayed by sample 1 is probably attributable to the reduced number of  
286 available intra-batch controls (9 samples, vs. 15 samples for the other 4 cases).

287

288 Retrospective variant calling was undertaken using VarScan2 [13], to assess the allelic  
289 ratios of coding and non-coding variants. No non-reference coding variants were  
290 identified. However, for sample 4, two variants c.905-52A>G (rs948589) and  
291 c.955+58C>T (rs948588) were present, in introns 7 and 8 respectively. The non-  
292 reference read frequencies were 47% for c.905-52A>G (681 of 1455 reads) and 46% for  
293 c.955+58C>T (722 of 1562 reads). This diploid allelic ratio further supports the

294 inference that the *SMAD4* pseudogene is processed, allelic ratios of intronic SNPs being  
295 unaffected by the presence of the pseudogene.

296

297 To determine whether a common *SMAD4* pseudogene integration site was shared  
298 between the five cases, medium-coverage whole genome sequencing (approximately 9×  
299 per sample) was performed. Evidence for the integration site being located on  
300 chromosome 9, within intron 18 of the *SCAI* gene, was provided by the 16 read pairs  
301 detailed in Table 2. These data characterise DNA fragments whose opposite ends were  
302 each mapped to (a) *SCAI* intron 18 and (b) either the 5' (14 read pairs) or 3' (2 read  
303 pairs) end of *SMAD4*. Soft-clipped reads spanning the precise integration site indicated  
304 that this was identical, at least among samples 2-5. (For sample 1, no supporting read-  
305 pairs were identified, despite there being no obvious difference between the assay  
306 performance metrics, as displayed in Supplemental Table S3.) *SMAD4* mapped reads  
307 indicated that the pseudogene sequences for exons 1 and 12 were shorter than those  
308 reported in transcript record NM\_005359.5  
309 (<https://www.ncbi.nlm.nih.gov/nucleotide/>). However, the precise terminal nucleotide  
310 of the 3'-UTR could not be determined from this dataset. This was probably due to the  
311 presence of the poly-A tail, hindering DNA sequencing and mapping, and resulting in an  
312 underrepresentation of exon 12 mapped read pairs. Interestingly, the library insert for  
313 the sample 5 read pair 4:23601:11116:11521 was sufficiently large that the *SMAD4*-  
314 mapped read spanned the exon 1-2 splice junction.

315

316 To confirm the identified integration site, and establish the terminal nucleotide of the  
317 *SMAD4* 3'-UTR, three overlapping PCR amplicons, each anchored at one end by a primer  
318 bound to *SCAI* intron 18, were amplified and sequenced (Figure 2). All five cases were

319 confirmed to have the same genomic integration site, at which the inserted pseudogene  
320 is flanked by a 4-nt microduplication (TTTC). The exon-exon arrangement was identical  
321 to transcript record NM\_005359.5 (<https://www.ncbi.nlm.nih.gov/nucleotide/>), and no  
322 nucleotide sequence variants were identified in any of the pseudogene exons. Compared  
323 to the mRNA reference sequence, 41 nt are missing from the beginning of the *SMAD4* 5'-  
324 UTR and 5,265 nucleotides are absent from the end of the 3'-UTR. A schematic  
325 representation of the integration site and scale drawing of the gene structure are  
326 displayed in Figure 3.

327

## 328 **DISCUSSION**

329 In recent years, the significantly increased number of genes that are attributable to  
330 clinically recognisable phenotypes have resulted in far greater scope for genetic testing.  
331 Laboratories typically create target enrichment panels that sequence more loci than are  
332 requested by the referring clinician and the unwanted variant data is masked by  
333 creating virtual gene panels *in silico*. While this approach facilitates the creation of  
334 efficient wet-laboratory processes, it also generates sequence data that is not routinely  
335 analysed. For the purposes of this study we harnessed these data to determine the  
336 frequency of a reported *SMAD4* processed pseudogene in our cohort of patients that had  
337 been referred for hereditary cancer testing. We determined the pseudogene to be  
338 present at a frequency of 1 in 408, which is consistent with the previously reported  
339 frequency of 1 in 389 [5]. That the integration site was common to all five patients  
340 suggests that this reflects a single ancestral founder event. Given that the majority of  
341 our laboratory's referrals are of northern European ancestry it will be interesting to  
342 determine whether this variant is also detected in more diverse ethnic populations.  
343 Unsurprisingly, many other polymorphic processed genes have been found to be

344 restricted to certain ethnic groups [14]. Polymorphic processed genes of the present  
345 type have been revealed by large-scale sequencing surveys to be a frequent feature of  
346 the human (and mouse) genome. Although the insertion site of the *SMAD4* processed  
347 gene was not determined in the large-scale studies of Ewing et al., (2013) and Shriver et  
348 al., (2013) [14, 15].

349

350 Most processed genes in the reference human genome are known to be non-functional  
351 (*i.e.* they are processed pseudogenes), either because they lack promoter sequences  
352 (“dead on arrival”) or have acquired inactivating mutations subsequent to  
353 retrotransposition. However, processed genes whose existence is polymorphic within  
354 the normal population are likely to have been recently transposed, and therefore (as in  
355 the present case) not to have acquired many inactivating mutations. There is  
356 population-level evidence that new processed genes are frequently subject to positive  
357 or negative evolutionary selection [15] as well as anecdotal examples of individual  
358 functional effects of processed genes (discussed in Richardson et al., (2014)) [16].

359

360 Since the coding region of the *SMAD4* processed gene is unaltered in comparison to its  
361 parent gene, we cannot be completely certain that it is non-functional (*i.e.* that it really  
362 is a processed pseudogene). We have been unable to address this question, since RNA is  
363 not available from any of the five carrier individuals, to permit analysis of whether the  
364 processed gene is transcribed. For the same reason, we cannot address any possible  
365 effect of the retrotransposed gene on the splicing of the *SCAI* gene, within which it is  
366 integrated. A newly transposed processed gene can be disease-causing as a result of  
367 disruption of splicing of its target gene [17].

368

369 *SCAI* itself is a nuclear protein that was first characterized for its suppressive effects  
370 upon tumour cell invasiveness, through regulation of beta1-integrin expression [18]. It  
371 has also been shown to be a TP53BP1 interaction partner with an important role in  
372 double-strand break repair [19]. It has been reported that the *SCAI* 3'-UTR contains a  
373 binding site for miR-1228. When bound, this microRNA is capable of down-regulating  
374 endogenous *SCAI* protein [20]. Furthermore, *SCAI* levels have been observed to be  
375 down-regulated in human tumours leading to reports of its tumour suppressor  
376 characteristics. RNA interference experiments of *SCAI* have shown an upregulation of  
377  $\beta$ 1-integrin gene expression and a resulting increase in invasive cell migration. Despite  
378 these observations, we were unable to obtain relevant tissue specimens from our  
379 patients to determine whether *SCAI* expression is perturbed by the presence of the  
380 *SMAD4* pseudogene.

381

382 Pseudogenes commonly interfere with the diagnostic analysis of clinically important  
383 genes. In extreme cases, unambiguous analysis may be impossible without resort to  
384 highly specialized methodologies; such is the case for mutations in *PMS2*, which in the  
385 heterozygous or biallelic state cause low-penetrance colorectal cancer predisposition  
386 (Lynch syndrome; OMIM: 614337), and a young-onset mismatch repair cancer  
387 syndrome (OMIM: 276300), respectively [21, 22]. Typically, however, because  
388 pseudogenes are not polymorphic, assay designs can be tailored to avoid interference  
389 and allow robust and reliable clinical diagnosis.

390

391 The ad-hoc discovery of polymorphic processed pseudogenes is likely to become more  
392 frequent as an increasingly genomic approach is applied to molecular diagnostic  
393 investigations. It is perhaps therefore surprising that given the clinical importance of

394 *SMAD4* [23], no comprehensive analysis of the *SMAD4* pseudogene integration site had  
395 hitherto been undertaken.

396

397 While the initial identification of the *SMAD4* pseudogene stemmed from aberrant MLPA  
398 result, the clinical adoption of NGS-based hybridisation enrichment panels is outpacing  
399 the production of gene-specific MLPA kits. Consequently, the per-exon cost of  
400 performing MLPA to detect novel pseudogenes, on a large-scale, would likely be cost-  
401 prohibitive. Our study demonstrates a convenient approach of using an RNA-seq aligner  
402 to detect processed pseudogenes from hybridisation capture data. We also report how  
403 comparative read depth methods can effectively determine the allelic copy number of  
404 novel pseudogene sequences. Increased demand for genetic testing has meant  
405 laboratories are becoming ever-more reliant on automated variant calling pipelines that  
406 do not involve visualisation of the directly sequenced reads, and clinical scientists are  
407 required to interpret sequence variants for unfamiliar genes. To maintain quality  
408 assurance of these tests, we advocate the inclusion of an RNA-seq aligner into  
409 laboratory pipelines as a means of detecting as-yet unreported polymorphic processed  
410 pseudogenes which, if they remain undetected, could interfere with the interpretation  
411 of clinical results.

412

413 In summary, we report a common genomic integration site for the polymorphic *SMAD4*  
414 processed pseudogene. We demonstrate how alignment of these data using an RNA-seq  
415 aligner can confirm the presence of splice-junction containing reads. And advocate that  
416 as the number of genes analysed by clinical laboratories continues to expand this would  
417 provide a worthwhile quality assurance approach for target enrichment experiments.

418

419 **REFERENCES**

- 420 1. Watson CM, Crinnion LA, Morgan JE, Harrison SM, Diggle CP, Adlard J, Lindsay  
421 HA, Camm N, Charlton R, Sheridan E, Bonthron DT, Taylor GR, Carr IM: Robust  
422 diagnostic genetic testing using solution capture enrichment and a novel variant-  
423 filtering interface. *Hum Mutat* 2014, 35:434-441
- 424
- 425 2. Morgan JE, Carr IM, Sheridan E, Chu CE, Hayward B, Camm N, Lindsay HA,  
426 Mattocks CJ, Markham AF, Bonthron DT, Taylor GR: Genetic diagnosis of familial  
427 breast cancer using clonal sequencing. *Hum Mutat* 2010, 31:484-491
- 428
- 429 3. Guo Y, Long J, He J, Li CI, Cai Q, Shu XO, Zheng W, Li C: Exome sequencing  
430 generates high quality data in non-target regions. *BMC Genomics* 2012, 13:194
- 431
- 432 4. Bellos E, Coin LJ: cnvOffSeq: detecting intergenic copy number variation using  
433 off-target exome sequencing data. *Bioinformatics* 2014, 30:i639-645
- 434
- 435 5. Millson A, Lewis T, Pesaran T, Salvador D, Gillespie K, Gau CL, Pont-Kingdon G,  
436 Lyon E, Bayrak-Toydemir P: Processed Pseudogene Confounding  
437 Deletion/Duplication Assays for SMAD4. *J Mol Diagn* 2015, 17:576-582
- 438
- 439 6. Martin M: Cutadapt removes adapter sequences from high-throughput  
440 sequencing reads. *EMBnet.journal* 2011, 17:10-12
- 441
- 442 7. Li H, Durbin R: Fast and accurate short read alignment with Burrows-Wheeler  
443 transform. *Bioinformatics* 2009, 25:1754-1760

444

445 8. Mose LE, Wilkerson MD, Hayes DN, Perou CM, Parker JS: ABRA: improved coding  
446 indel detection via assembly-based realignment. *Bioinformatics* 2014, 30:2813-  
447 2815

448

449 9. DePristo MA, Banks E, Poplin R, Garimella KV, Maguire JR, Hartl C, Philippakis  
450 AA, del Angel G, Rivas MA, Hanna M, McKenna A, Fennell TJ, Kernytzsky AM,  
451 Sivachenko AY, Cibulskis K, Gabriel SB, Altshuler D, Daly MJ: A framework for  
452 variation discovery and genotyping using next-generation DNA sequencing data.  
453 *Nat Genet* 2011, 43:491-498

454

455 10. Thorvaldsdóttir H, Robinson JT, Mesirov JP. Integrative Genomics Viewer (IGV):  
456 high-performance genomics data visualization and exploration. *Brief Bioinform*  
457 2013, 14:178-192

458

459 11. Dobin A, Davis CA, Schlesinger F, Drenkow J, Zaleski C, Jha S, Batut P, Chaisson M,  
460 Gingeras TR: STAR: ultrafast universal RNA-seq aligner. *Bioinformatics* 2013,  
461 29:15-21.

462

463 12. Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G,  
464 Durbin R; 1000 Genome Project Data Processing Subgroup: The Sequence  
465 Alignment/Map format and SAMtools. *Bioinformatics* 2009, 25:2078-2079

466

- 467 13. Koboldt DC, Zhang Q, Larson DE, Shen D, McLellan MD, Lin L, Miller CA, Mardis  
468 ER, Ding L, Wilson RK: VarScan 2: somatic mutation and copy number alteration  
469 discovery in cancer by exome sequencing. *Genome Res* 2012, 22:568-576  
470
- 471 14. Ewing AD, Ballinger TJ, Earl D; Broad Institute Genome Sequencing and Analysis  
472 Program and Platform, Harris CC, Ding L, Wilson RK, Haussler D:  
473 Retrotransposition of gene transcripts leads to structural variation in  
474 mammalian genomes. *Genome Biol* 2013, 14:R22  
475
- 476 15. Schrider DR, Navarro FC, Galante PA, Parmigiani RB, Camargo AA, Hahn MW, de  
477 Souza SJ: Gene copy-number polymorphism caused by retrotransposition in  
478 humans. *PLoS Genet* 2013, 9:e1003242  
479
- 480 16. Richardson SR, Salvador-Palomeque C, Faulkner GJ: Diversity through  
481 duplication: whole-genome sequencing reveals novel gene retrocopies in the  
482 human population. *Bioessays* 2014, 36:475-481  
483
- 484 17. de Boer M, van Leeuwen K, Geissler J, Weemaes CM, van den Berg TK, Kuijpers  
485 TW, Warris A, Roos D: Primary immunodeficiency caused by an exonized  
486 retroposed gene copy inserted in the CYBB gene. *Hum Mutat* 2014, 35:486-496  
487
- 488 18. Brandt DT, Baarlink C, Kitzing TM, Kremmer E, Ivaska J, Nollau P, Grosse R: SCAI  
489 acts as a suppressor of cancer cell invasion through the transcriptional control of  
490 beta1-integrin. *Nat Cell Biol* 2009, 11:557-568  
491

- 492 19. Hansen RK, Mund A, Poulsen SL, Sandoval M, Klement K, Tsouroula K, Tollenaere  
493 MA, Räschle M, Soria R, Offermanns S, Worzfeld T, Grosse R, Brandt DT, Rozell B,  
494 Mann M, Cole F, Soutoglou E, Goodarzi AA, Daniel JA, Mailand N, Bekker-Jensen S:  
495 SCAI promotes DNA double-strand break repair in distinct chromosomal  
496 contexts. *Nat Cell Biol* 2016, 18:1357-1366  
497
- 498 20. Lin L, Liu D, Liang H, Xue L, Su C, Liu M: MiR-1228 promotes breast cancer cell  
499 growth and metastasis through targeting SCAI protein. *Int J Clin Exp Pathol* 2015,  
500 8:6646-6655  
501
- 502 21. De Vos M, Hayward BE, Picton S, Sheridan E, Bonthron DT: Novel PMS2  
503 pseudogenes can conceal recessive mutations causing a distinctive childhood  
504 cancer syndrome. *Am J Hum Genet* 2004, 74:954-964  
505
- 506 22. Goodenberger ML, Thomas BC, Riegert-Johnson D, Boland CR, Plon SE,  
507 Clendenning M, Win AK, Senter L, Lipkin SM, Stadler ZK, Macrae FA, Lynch HT,  
508 Weitzel JN, de la Chapelle A, Syngal S, Lynch P, Parry S, Jenkins MA, Gallinger S,  
509 Holter S, Aronson M, Newcomb PA, Burnett T, Le Marchand L, Pichurin P, Hampel  
510 H, Terdiman JP, Lu KH, Thibodeau S, Lindor NM: PMS2 monoallelic mutation  
511 carriers: the known unknown. *Genet Med* 2016, 18:13-19  
512
- 513 23. Kalia SS, Adelman K, Bale SJ, Chung WK, Eng C, Evans JP, Herman GE, Hufnagel  
514 SB, Klein TE, Korf BR, McKelvey KD, Ormond KE, Richards CS, Vlangos CN,  
515 Watson M, Martin CL, Miller DT: Recommendations for reporting of secondary  
516 findings in clinical exome and genome sequencing, 2016 update (ACMG SF v2.0):

517 a policy statement of the American College of Medical Genetics and Genomics.

518 Genet Med 2017, 19:249-255

519

520

521

522

523

524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

540

541

542 **FIGURE LEGENDS**

543 **Figure 1:** *SMAD4* Sashimi plots generated following alignment of targeted capture data  
544 using the RNA-seq aligner STAR. Each arc's corresponding value records the number of  
545 reads crossing the reported splice junction. Alignment coverage data is displayed with  
546 y-axis values ranging from 0-20,000 for Sample 1 and 0-6,000 for all other samples.

547

548 **Figure 2:** DNA sequence at the common *SMAD4* processed gene integration site, located  
549 within *SCAI* intron 18 (using reference transcript NM\_173690.4,  
550 <https://www.ncbi.nlm.nih.gov/nucleotide/>). Genomic coordinates refer to human  
551 reference genome build hg19. **(A)** The dashed red line marks the breakpoint 5' to the  
552 processed gene. **(B)** The last nucleotide matching the *SMAD4* 3' untranslated region is  
553 identified, immediately to the left of the vertical dashed line. To the right of this line is a  
554 poly(A) sequence. **(C)** The *SCAI* intron 18 integration site beyond the poly(A) tail. This  
555 sequence was generated using a reverse strand primer. The four nucleotides located  
556 between the dashed red lines are duplicated from the proximal breakpoint.

557

558 **Figure 3:** A schematic representation of the *SCAI* locus, displaying the exon  
559 arrangement of the *SMAD4* processed pseudogene, which is consistent with that  
560 reported for NM\_005359.5 (<https://www.ncbi.nlm.nih.gov/nucleotide/>). Exons (green  
561 boxes) are drawn to scale using GeneDrawer  
562 ([www.insilicase.com/Desktop/GeneDrawer.aspx](http://www.insilicase.com/Desktop/GeneDrawer.aspx), last accessed August 18, 2017).

563

**Table 1: The ratio of gapped to non-gapped sequence alignments in cases with an apparent *SMAD4* intron 6 deletion.**

<b>Sample</b>	<b>Gapped read alignments spanning intron 6</b>	<b>Mean per-base read depth for intron 6 nucleotides</b>	<b>Ratio of gapped to non-gapped reads</b>
1	1,620	3,069	1:1.89
2	739	1,750	1:2.37
3	571	1,183	1:2.07
4	1,198	2,250	1:1.88
5	770	1,323	1:1.72

Intron numbering determined according to NM\_005359.5 (<https://www.ncbi.nlm.nih.gov/nucleotide/>).

**Table 2: Characteristics of whole genome sequencing reads supporting the intragenic *SCAI* integration site.**

Sample	Read pair ID	Read 1					Read 2				
		Locus	Chr.	Start	Str.	CIGAR	Locus	Chr.	Start	Str.	CIGAR
2	4:13608:15564:14605	5'-SMAD4	18	48,556,641	-	151M	SCAI	9	127,732,358	+	150M
2	2:13210:1908:2419	SCAI	9	127,732,501	+	151M	5'-SMAD4	18	48,556,701	-	6S143M
2	4:22402:24489:4305	5'-SMAD4	18	48,556,700	-	151M	SCAI	9	127,732,506	+	150M
2	1:11204:8894:18752	5'-SMAD4	18	48,556,624	-	60S91M	SCAI	9	127,732,633	+	81M70S
2	4:21606:19277:14299	SCAI	9	127,732,700	-	24S127M	3'-SMAD4	18	48,605,924	+	101M49S
3	3:22511:22720:13254	5'-SMAD4	18	48,556,622	-	151M	SCAI	9	127,732,422	+	151M
3	2:11311:9152:15694	5'-SMAD4	18	48,556,624	-	60S91M	SCAI	9	127,732,437	+	151M
3	2:21212:20833:4797	SCAI	9	127,732,556	+	150M	5'-SMAD4	18	48,556,624	-	64S87M
3	1:21211:11117:11719	3'-SMAD4	18	48,605,995	+	150M	SCAI	9	127,732,700	-	87S62M
4	3:13407:5904:6688	5'-SMAD4	18	48,556,711	-	151M	SCAI	9	127,732,459	+	150M
4	1:11210:11536:3015	5'-SMAD4	18	48,556,624	-	58S93M	SCAI	9	127,732,487	+	151M
4	2:21206:5607:15745	5'-SMAD4	18	48,556,624	-	48S103M	SCAI	9	127,732,604	+	110M41S
5	4:13501:18475:17089	5'-SMAD4	18	48,556,624	-	45S106M	SCAI	9	127,732,552	+	150M
5	4:11605:22916:12006	SCAI	9	127,732,569	+	145M6S	5'-SMAD4	18	48,556,635	-	151M
5	4:12410:4727:7249	5'-SMAD4	18	48,556,800	-	151M	SCAI	9	127,732,582	+	132M19S
5	4:23601:11116:11521	SCAI	9	127,732,607	+	107M44S	5'-SMAD4	18	48,556,882	-	2S114M35S

Str.: Strand. CIGAR: The mapping defined by the BWA alignment. All coordinates are provided according to human genome build hg19. Locus represents the read mapping to one of three possible loci, either the 5' end of the *SMAD4* pseudogene (5'-SMAD4), the 3' end of the *SMAD4* pseudogene (3'-SMAD4), or the *SCAI* integration site (SCAI).

