

1     **A Novel Methodology to Recalibrate Pathogenic Range of SCA36 Repeat Expansions**  
2   **for PGT-M**

3                             Fulin Liu<sup>1</sup>, Wen Huang<sup>1</sup>, Ling Liao<sup>\*</sup>, Jiyun Yang<sup>\*</sup>

4     Sichuan Provincial Key Laboratory for Human Disease Gene Study, Center for Medical Gene  
5     tics, Department of Laboratory Medicine, Sichuan Academy of Medical Sciences & Sichuan  
6     Provincial People's Hospital, University of Electronic Science and Technology, Chengdu  
7     610072, China

8

9     <sup>\*</sup> Corresponding authors:

10    Email address: [yangjiyun@yeah.net](mailto:yangjiyun@yeah.net) (Jiyun Yang), [450766957@qq.com](mailto:450766957@qq.com) (Ling Liao)

11    <sup>1</sup> These authors contribute equally to this work.

12

13    **Abstract**

14    **Background:** Spinocerebellar ataxia-36 (SCA36) is an inherited neurodegenerative disorder  
15    caused by the heterozygous expansion of an intronic GGCCTG hexanucleotide repeat in the  
16    NOP56 gene on chromosome 20p13. Unaffected individuals typically carry 3 to 14 repeats,  
17    whereas affected individuals carry 650 to 2,500. However, based on a single study, this  
18    pathogenic range was conservatively established, limiting its extended clinical applicability  
19    such as preimplantation genetic testing (PGT). In this study, we propose a novel methodology  
20    to recalibrate the pathogenic range of SCA36 repeat expansion.

21    **Methods:** We conducted a comprehensive literature review and collected examination data  
22    from 2012 onward. We used the gamma distribution to describe the data distribution and  
23    applied Bayesian methods to update the prior distribution with data from recent publications.  
24    Based on the recalibrated distribution, the 95% confidence interval (CI) was used to  
25    determine the new lower boundary of the pathogenic range. A pedigree was collected to

26 validate the proposal with long-read sequencing (LRS) applied to detect the high GC content  
27 and long length of repeat expansions.

28 **Results:** Our results, based on 2 studies, indicate that the data distribution is well-described  
29 by gamma distribution. The prior, likelihood and posterior distributions within the 95% CI  
30 for the integrated research of SCA36 pathogenic repeat expansions were [446,  $+\infty$ ), [124,  
31  $+\infty$ ), and [484,  $+\infty$ ), respectively. These recalibrated pathogenic ranges were validated by an  
32 authentic case: a proband diagnosed with SCA36 carrying 418 repeats and her daughter with  
33 499 repeats, under the detection of LRS.

34 **Conclusions:** Therefore, we proposed a novel methodology that integrates updated data, 95%  
35 CI using Bayesian methods and LRS for accurate detection of repeat expansions of dynamic  
36 mutations to present an up-to-date pathogenic range of SCA36, as well as other similar  
37 diseases.

38

39 **Keywords:** Spinocerebellar ataxia-36; Repeat expansions; Bayesian methods; Confidence  
40 interval; Diagnosis; Long-read sequencing; Preimplantation genetic testing

41

## 42 **Background**

43

44 The human genome contains a substantial number of tandem repeats (TRs). TRs exhibit a  
45 high mutation rate within the genome due to their repetitive nature, resulting in enormous  
46 polymorphism and multiallelism. The variability of TRs in the length of the repeat sequences  
47 occurs in both non-coding and coding regions of the genome [1]. A certain range of TR  
48 expansions beyond a critical threshold can lead to disease onset and may increase or decrease  
49 across generations, which is the foundation of dynamic mutations [2]. The intergenerational

50 expansion of dynamic mutations is a well-documented phenomenon that underlies the  
51 pathogenesis of at least 50 known disorders [3].

52

53 Clinical diagnosis of dynamic mutations necessitates not only the observation of  
54 characteristic clinical manifestations but also the application of auxiliary diagnostic tests,  
55 with genetic testing being paramount. Diagnosis depends on identifying the number of  
56 dynamic mutations within a defined pathogenic range. When the mutation count falls within  
57 the defined range, a clinical diagnosis can be made; otherwise, a negative result is inferred.

58 However, there exists an uncertain range of clinical significance beyond the confirmatory and  
59 exclusion ranges. Currently, the evaluation of the uncertain range primarily relies on clinical  
60 diagnosis. If clinical diagnosis is definitive, the uncertain range is considered positive; if  
61 clinical phenotypes cannot be correlated, it is deemed negative. Nonetheless, clinical practice  
62 often encounters cases where phenotypes are ambiguous, and genetic testing results fall  
63 within the uncertain range. This poses ethical challenges, particularly for preimplantation  
64 genetic testing (PGT) candidates, where ethics committees may not approve the subsequent  
65 procedure according to the guidelines [4-6].

66

67 PGT, especially PGT for monogenic disorders (PGT-M), aims to enable pregnancies  
68 unaffected by specific genetic traits carried by one or both parents. PGT-M offers couples  
69 with genetic disorders the opportunity to have their offspring and avoid the birth of defective  
70 fetuses, yet its ethical standards are stringently ruled. This stringency stems from the current  
71 limitations in basic medical research to fully elucidate the relationship between genetic  
72 variants and clinical phenotypes, leading to poor gene-disease mapping. Thus, strict inclusion  
73 criteria are employed to minimize off-target resulting in defective fetuses. Namely, PGT-M  
74 can be only administered in a certain pathogenic range but not the uncertain range of repeat

75 expansions, which poses a significant challenge to the uncertainty. Such stringent criteria  
76 may result in too many considerations to perform the PGT-M. Therefore, this study proposes  
77 a novel but reasonable methodology to recalibrate the defined pathogenic range of repeat  
78 expansions so that more patients within the uncertain range win steady support for the PGT-  
79 M.

80

81 Taking spinocerebellar ataxia type 36 (SCA36) as an example, SCA36 was initially identified  
82 in Japanese and Spanish families as a new type of SCA, characterized by late-onset, slowly  
83 progressive cerebellar syndrome often involving motor neurons, or accompanied by sensory  
84 neural hearing loss, and cognitive, and emotional disturbances [7, 8]. It is caused by a  
85 heterozygous expansion of an intronic GGCCTG hexanucleotide repeat in the NOP56 gene  
86 on chromosome 20p13. Unaffected individuals carry 3 to 14 repeats, whereas affected  
87 individuals have 650 to 2,500 repeats (OMIM #3614153). However, the range was defined  
88 based on only one single research [7] under an obsolete technology and unsuitable statistics,  
89 making individuals under the uncertain range for SCA36 diagnosis between 14 to 650 repeats  
90 inaccessible to PGT-M.

91

92 We recommend an update of knowledge. This study, taking CA36 as an example, aims to  
93 integrate updated data, 95% CI under Bayesian methods and long-read sequencing (LRS) to  
94 present an up-to-date pathogenic range of SCA36, along with the validation of a specific  
95 clinical case, to propose an efficient novel methodology to recalibrate pathogenic repeat  
96 expansions.

97

98 **Methods**

## 99 **Participants**

100 The study subjects were members of one generation of a Han Chinese ataxia family pedigree  
101 in Sichuan province. A detailed medical history and physical examination record of the  
102 proband (II-2) and other family members were evaluated by our multidisciplinary teams  
103 including a geneticist, neurologist, otorhinolaryngologist, obstetrician, and gynecologist,  
104 reproductive endocrinologist, members of the ethical committee. The blood samples were  
105 obtained from participants who asked for carrier screening for SCA36 at the Center for  
106 Medical Genetics, Sichuan Academy of Medical Sciences & Sichuan Provincial People's  
107 Hospital. This study was in agreement with the Guidance of the Ministry of Science and  
108 Technology for the Review and Approval of Human Genetic Resources. All study procedures  
109 were approved by the Ethical Committee of Sichuan Academy of Medical Sciences &  
110 Sichuan Provincial People's Hospital in accordance with the Helsinki Declaration of 1975, as  
111 revised in 2000, and patients signed an informed consent form.

## 112 **Repeat-primed PCR**

113 Given the initial suspicion of SCA36, repeat-primed polymerase chain reaction (RP-PCR)  
114 screening combined with Southern blotting with capillary electrophoresis was employed to  
115 identify the presence of SCA36 in the proband and other family members. RP-PCR utilizes  
116 repeat primers that amplify alleles, producing PCR products that anneal and form a “ladder”  
117 pattern upon separation by capillary electrophoresis. For DNA fragment analysis, the RP-  
118 PCR products were processed using the ABI-Prism 3730XL Genetic Analyzer, and the  
119 resulting data were analyzed with GeneMarker software.

## 120 **Targeted gene capture and PacBio long-read Sequencing**

121 Using a g-tube tube (Covaris, Australia) and centrifugation at  $6000 \times g$ , 3  $\mu\text{g}$  of genomic  
122 DNA was physically sheared into 8-9 kb fragments after the purity, followed by the  
123 evaluation of the integrity of the DNA. The DNA fragments were then connected with  
124 barcodes after being repaired by the Damage Repair End Preparation Mix (Vazyme, China).  
125 Subsequently, Rapid DNA ligase (Vazyme, China) was used to attach adapters to the DNA  
126 fragments. After purifying the ligation products with Agencourt AMPure XP beads  
127 (Beckman Coulter, USA), amplification was carried out. The amplified products were  
128 captured by hybridization with biotinylated probes (Boke Bioscience, China) that targeted  
129 genes relevant to repeat expansion diseases. The SMRTbell Express Template Kit 2.0  
130 (PacBio, USA) was utilized to generate an SMRTbell library. The constructed DNA library  
131 was then subjected to long-read sequencing on the PacBio Sequel IIe platform.

132 After the sequencing data was evaluated and qualified by SMRT Link, bioinformatics  
133 analysis of the raw sequencing data was carried out. The particular processes are as follows.  
134 The raw sequencing data was analyzed using the official PacBio software package SMRT  
135 Link (version 12.0) to generate HiFi reads, and CCS (Circular Consensus Sequence) reads  
136 were also automatically generated by the PacBio SMRT Analysis Module, which is an apart  
137 of SMRT Link. The data of each sample was divided according to the barcode using the Lima  
138 software (version 2.7.1). The CCS reads were aligned to the reference genome GRCh37/hg19  
139 using the Minimap2 software (version 2.24), and the number of repeat units of the gene was  
140 calculated using GrandSTR, while the number of repeat units on each read was manually  
141 checked using the IGV software.

## 142 **Statistical analysis**

143

144 Considering unaffected individuals typically carry 3 to 14 repeats, whereas affected  
145 individuals carry 650 to 2,500 repeats, the pathogenic repeats exhibit a pronounced skewed  
146 distribution with one tail. Gamma distribution was hence adopted to describe the pathogenic  
147 repeat range. Based on this assumption, we first estimated the prior distribution based on the  
148 dataset from [7].

149

$$150 \quad \text{Prior: } \text{gamma}(\kappa|s, r) = \left(\frac{r^s}{\Gamma(s)}\right) \kappa^{s-1} e^{-r\kappa},$$

151

152 where  $\kappa$  is a shape parameter,  $s$  is a scale parameter and  $r$ , called a rate parameter, is an  
153 inverse scale parameter  $r = 1/s$ .  $\kappa$  and  $s$  are calculated based on the mean and variance of  
154 the old data.

155 Bayesian approaches were applied to update this posterior probability, based on the prior  
156 probability and new data. Given two events  $A$  and  $B$ , the conditional probability of  $A$  given  
157 that  $B$  is true is expressed as follows:

$$P(A|B) = P(B|A)P(A)/P(B)$$

158 where  $P(A)$  is the prior probability of A to present the experience,  $P(B|A)$  is the likelihood  
159 function to present the new data,  $P(A|B)$  is posterior probability to present the current  
160 knowledge, and  $P(B)$  is the total probability of the probability of B with  $P(B) \neq 0$ . As  $P(B)$   
161 does not change in the same analysis, the formula can be also interpreted as:

$$P(A|B) \propto P(B|A)P(A)$$

162 where the current knowledge is proportional to the product of the likelihood function and the  
163 new data under the same data distribution. Assume that the new observed dataset consists of  
164  $n$  observations, denoted as  $x_1, x_2, \dots, x_n$ . Assuming a conjugate prior of the same distribution,  
165 the posterior distribution also follows a gamma distribution with updated parameters:

166

$$s_{new} = s + ns_0,$$

167

168

$$k_{new} = k + \sum_{i=1}^n x_i,$$

169 where  $s_0$  is estimated by the likelihood from the new data [9]. The posterior distribution is

170 then given by:

171

$$Post: gamma(k_{new}|s_{new}, r_{new}) = \left( \frac{r_{new}^{s_{new}}}{\Gamma(s_{new})} \right) k_{new}^{s_{new}-1} e^{-r_{new}k_{new}},$$

172

173 where  $r_{new} = 1/s_{new}$ . Using this distribution, we calculated a 95% confidence interval (CI)

174 to estimate the range within which the true population parameter is expected to fall with high

175 confidence. This range also defines the new cutoff for normal observations, classifying

176 patients with repeats falling within this range as having the presence of SCA36. Data were

177 analyzed and visualized with R statistical software (v4.2.1; R Core Team 2023)[10].

178

## 179 **Results**

180

181 1 Literature review and summary

182

183 We conducted a comprehensive search with the keywords “SCA36” or “Spinocerebellar

184 ataxia-36” across several databases, including PubMed, Embase, Web of Science, and

185 additional platforms such as Dimensions and Semantic Scholar in May 2024, yielding a total

186 of 67 articles. The subsequent filter process is illustrated in Figure 1. To elaborate, we first

187 refined our search publication date from 2012 onwards to include literature published after

188 the research from which the pathogenic range of SCA36 was defined, retaining 55 records.

189 We then meticulously reviewed the abstracts of these articles, excluding 9 duplicates, 11

190 reviews, 18 basic medical research, 2 case reports, 7 comments, and 2 non-biomedical



191 articles. The remaining 9 articles encompassed detailed data and processing methodologies.  
192 Accordingly, Table 1 summarizes these 9 articles, detailing the authors, publication year,  
193 ethnic group, sample size, repeat range, mean age at onset, detection method, and references.  
194 Upon a thorough examination of the original data from these 9 articles [7, 9, 11-17], we  
195 noticed issues such as insufficient data [11-13], questionable data reliability [14], suboptimal  
196 data distribution [11], and poor data accessibility [15, 16], and non-Caucasian ethical  
197 population [11-14]. Consequently, we determined that only two articles provided data  
198 suitable for subsequent analysis [7, 9].

199

200 2 Recalibrating pathogenic range of SCA36 repeat expansions

201

202 According to the data published in 2012 by Garcia-Murias et al [7], which Online Mendelian  
203 Inheritance in Man (OMIM, <https://www.omim.org/>) used as a guide, the range of pathogenic  
204 repeat expansions is 650 to 2,500. By examining the original data, we found it skewed and  
205 inconsistent with the normal distribution. Taking the one-tailed and skewed distribution of  
206 the data into consideration, we employed the gamma distribution to simulate the prior  
207 distribution, obtaining a 95% CI of [446,  $+\infty$ ]. Since the prior distribution data shares a  
208 common original disease with the data published in 2014 by Obayashi et al.[9] we defined  
209 them as conjugate distributions. Therefore, we simulated the likelihood function, obtaining a  
210 95% CI of [124,  $+\infty$ ]. Using Bayesian methods, we inferred the posterior distribution and its  
211 95% CI of [484,  $+\infty$ ]. The fitted distribution curves and the integrated curve are depicted in  
212 Figure 2. The results indicate that the lower bound of the pathogenic repeat expansions in the  
213 posterior distribution is 484, which is less than the 650 as defined in the OMIM guidelines.

214

215 3 Case validation

216

217 In this study, the proband (II-2) initially presented with progressive symptoms in her 50s,  
218 including gait instability, frequent falls, difficulty in lifting the feet, swaying while walking,  
219 inability to walk in a straight line, slurred speech, choking on water, dizziness, and slight  
220 difficulty in writing. Magnetic resonance imaging indicated cerebellar atrophy (Figure 3B).  
221 She did not exhibit bradykinesia, hand tremors, or other Parkinsonian symptoms, and her  
222 pain, light touch, and proprioception were intact. Additionally, electromyography did not  
223 show significant reductions or the presence of large motor unit potentials. Genetic panel  
224 results were negative for SCA1, SCA2, SCA3, SCA6, SCA7, SCA8, SCA12, SCA17, and  
225 DRPLA, but revealed that the SCA36 repeats exceeded 14 (Figure 3B). Her daughter (III-1),  
226 seeking consultation for assisted reproductive technology, also underwent genetic panel  
227 testing, which similarly indicated that the SCA36 repeats exceeded 14 (Supplementary  
228 Table\_S1). Testing of I-2 was negative, while I-1 was absent due to passing away. Given III-  
229 1's request for PGT-M, we performed long-read sequencing (LRS) to ascertain the precise  
230 repeats. Figure 3B illustrates the span of long reads in LRS for II-2 and III-1. Results  
231 revealed that II-2 had 418 repeats, while III-1 had 499 repeats, with other potential disease-  
232 related gene mutations excluded (Supplementary Table\_S2&S3). These findings demonstrate  
233 that III-1's repeats fall within the previously mentioned 95% CI of posterior distribution [484,  
234  $+\infty$ ], thus confirming that III-1 meets the criteria for a positive SCA36 diagnosis within the  
235 updated pathogenic range of repeat expansion.

236

## 237 **Discussion**

238

239 The ethical standards for PGT-M in SCA36 patients are stringently ruled and confined to a  
240 certain pathogenic range. The overcautious criterion has prevented those SCA36 patients

241 within an uncertain range from performing the PGT-M. In this study, we propose a more  
242 inclusive and pragmatic diagnostic methodology to encompass more SCA36 patients in the  
243 uncertain range. Initially, we conducted a comprehensive literature review, identifying three  
244 original studies with accessible data. Subsequently, we employed Bayesian methods to infer  
245 the posterior distribution, obtaining a 95% CI for the distribution. Our results indicated that  
246 the lower bound fell below the defined pathogenic range. Lastly, we utilized the LRS  
247 technology in the clinical practice to identify the pathogenic repeat expansions in a case,  
248 whose repeat expansions also fell within the aforementioned 95% CI, thereby validating the  
249 necessity and efficacy of recalibration for repeat expansions.

250

251 SCA36 exemplifies cases within an uncertain range that encounter the ambiguity of  
252 mutations and phenotypes, displaying an uncertain diagnostic range of 14-650 repeat  
253 expansions, within which patients may not exhibit pronounced symptoms, particularly those  
254 of childbearing age who have not reached the average onset age of  $51.23 \pm 7.33$  years (male  
255 = 51.23 years, female = 51 years) [11]. In our case, III-1 is in her 30s, suggesting that she is  
256 likely far from the average age of disease onset. However, SCA36, as a dominant repeat  
257 expansion disorder, is characterized by genetic anticipation, where clinical manifestations  
258 appear earlier and/or become more severe in successive generations. Dominant repeat  
259 expansion disorders have a well-established correlation between expansion size and both age  
260 of onset and disease severity, with larger expansions associated with earlier onset and more  
261 severe phenotypes [18]. Compared to II-2's 418 repeats, that III-1 has 499 repeats suggests  
262 III-1 might experience earlier onset and more severe symptoms than II-2. To prevent the  
263 transmission and exacerbation of the serious condition to her children, the application of  
264 PGT-M in assisted reproductive technology is highly warranted.

265

266 While the stringent rules for PGT-M in assisted reproductive technology are  
267 understandable—primarily aimed at avoiding potential off-targets, especially for monogenic  
268 diseases of unknown etiology—such stringent criteria may not be appropriate for the  
269 uncertain definition of dynamic mutations. Taking SCA36 into consideration, PGT-M is only  
270 permissible when repeat expansions exceed 650. Nevertheless, in our case, both II-2 (with  
271 typical manifestations) and III-1 fall below the range. Therefore, it is reasonable to suggest  
272 that the defined pathogenic range should be extended more rationally to enhance the  
273 application of PGT-M in clinical practice. In this study, the Bayesian methods, along with 95%  
274 CI were applied to recalibrate the pathogenic range. Bayesian theorem allows for the  
275 integration of prior experience with updated data, produced by new research or novel  
276 technologies, to refresh the current knowledge [19]. Thus, an effective method to integrate  
277 past and future data is highly beneficial for updating clinical diagnostic guidelines. In our  
278 case, the defined pathogenic range of 650 to 2,500 was based on the single research from [7],  
279 but with our recalibration based on Bayesian methods the updated range is [484, +∞]. III-1's  
280 repeat expansion therefore falls in, meeting the ethical requirements for PGT-M in assisted  
281 reproductive technology. Indeed, we have also conducted a direct integration of the datasets  
282 from the two aforementioned studies, yielding a 95% CI of [337, +∞] (Supplementary  
283 Figure\_S1), which agreed with the conclusion but might undertake the potential  
284 overshadowing of information and weight bias from smaller datasets by those with larger  
285 sample sizes. Therefore, we advocate for the application of Bayesian methods, which  
286 amalgamates prior and likelihood function data, to refine and enhance the integral  
287 understanding.

288

289 It is important to note that SCA36 is characterized by a heterozygous expansion of an intronic  
290 GGCCTG hexanucleotide repeat in the NOP56 gene, resulting in high GC content and

291 extremely long repeat motifs in the repeat expansion [7, 8]. Conventional detection methods  
292 are limited. In our case, the patient initially underwent RP-PCR combined with Southern  
293 blotting following capillary electrophoresis to measure the GGCCTG repeat expansions, but  
294 specific repeat sequences were not obtained. As a matter of fact, over the past decades, the  
295 repetitive nature and abundance of short TRs in the human genome have posed significant  
296 challenges for genome-wide studies [18]. Short-read sequencing alone fails to accurately map  
297 long repeat sequences, leading to covered reads often mapping to multiple genomic regions,  
298 thereby being truncated or discarded. Recently, the LRS has emerged in clinical practice to  
299 perform a more precise detection of TR expansions [20]. The LRS can obtain reads  
300 exceeding 15 kb and sometimes up to 2 Mb, enabling the detection of thousands of structural  
301 variants, including repeat regions, in individuals [21]. In this study, compared to RP-PCR  
302 combined with Southern blotting, LRS was used for the accurate detection of dynamic  
303 mutations to determine whether the repeat expansions exceeded 14 by considering the high  
304 GC content and long length. The LRS technology, with its long-read length, high accuracy,  
305 and lack of GC bias, provided more specific repeats, thereby offering crucial assistance in  
306 determining whether the repeat expansion falls within the recalibrated pathogenic range.

307

308 This study has limitations. Firstly, the sample size is limited. We obtained merely 44 patients  
309 (27 from the prior and 17 from the likelihood) to recalibrate the pathogenic range. The  
310 limited sample size might lead to bias in estimating the true distribution of the data, despite of  
311 our attempt to request data from other corresponding authors. Notably, in these limited  
312 samples, low repeats between 20 and 30 occurred only in three individuals, casting doubt on  
313 whether low repeats were outliers. Fortunately, subsequent cases in Chinese and British  
314 cohorts [11, 15] with low repeats reaching 30 or 60 provided proofs to dispute the doubt of  
315 the previous data that repeats between 20 to 30 were probably outliers. The second limitation

316 pertains to ethnic differences. Caucasian participants were included in our work to describe  
317 the prior, likelihood, and posterior distributions. However, a Han Chinese pedigree as the  
318 clinical practice was adopted to validate the recalibration, which might raise the  
319 overestimation or underestimation of the pathogenic range given the diversity of gene  
320 variations in different ethnic populations [22]. This underscores the necessity for further  
321 follow-up studies in Caucasian populations, or alternatively, a retrospective analysis with  
322 East Asian population to reassess the conclusions.

323

## 324 **Conclusions**

325

326 To sum up, we proposed a novel methodology that integrates updated data, 95% CI using  
327 Bayesian methods and LRS for accurate detection of repeat expansions of dynamic mutations  
328 to present an up-to-date pathogenic range of SCA36. We sincerely hope this work would well  
329 serve as a strong evidence for those suffering from an uncertain range of repeat expansion.

330

## 331 **Ethics approval and consent to participate**

332 This study was in agreement with the Guidance of the Ministry of Science and Technology  
333 for the Review and Approval of Human Genetic Resources. All study procedures were  
334 approved by the Ethical Committee of Sichuan Academy of Medical Sciences & Sichuan  
335 Provincial People's Hospital in accordance with the Helsinki Declaration of 1975, as revised  
336 in 2000, and patients signed an informed consent form.

337

## 338 **Consent for publication**

339 Patients signed a consent form for publication.

340

341 **Availability of data and materials**

342 All data are incorporated into the article and its online supplementary material.

343

344 **Competing interests**

345 The authors declare that they have no conflict of interest.

346

347 **Funding**

348 This study was supported by the Youth Talent Foundation of Sichuan Academy of Medical  
349 Sciences & Sichuan Provincial People's Hospital (No. 2022QN37), Sichuan Science and  
350 Technology Program (No. 2023NSFSC1605).

351

352 **Authors' contributions**

353 Fulin Liu.: designed the study and wrote the first draft of the manuscript. Wen Huang:  
354 participated in the manuscript organization. Ling Liao: supervised the study and edited the  
355 first draft of the manuscript. Jiyun Yang: conceived the original idea, and edited the final  
356 version of the manuscript.

357

358 **Acknowledgements**

359 We thank Beijing GrandOmics Biosciences Co., Ltd. for the technical support on PacBio  
360 SMRT target sequencing.

361

362 **References:**

- 363 1. Fazal S, Danzi MC, Cintra VP, Bis-Brewer DM, Dolzhenko E, Eberle MA, Zuchner S:  
364 **Large scale in silico characterization of repeat expansion variation in human**  
365 **genomes. *Sci Data* 2020, 7:294.**

- 366 2. Balzano E, Pelliccia F, Giunta S: **Genome (in)stability at tandem repeats.** *Semin*  
367 *Cell Dev Biol* 2021, **113**:97-112.
- 368 3. Depienne C, Mandel JL: **30 years of repeat expansion disorders: What have we**  
369 **learned and what are the remaining challenges?** *Am J Hum Genet* 2021, **108**:764-  
370 785.
- 371 4. Yan L, Cao Y, Chen ZJ, Du J, Wang S, Huang H, Huang J, Li R, Liu P, Zhang Z, et al:  
372 **Chinese experts' consensus guideline on preimplantation genetic testing of**  
373 **monogenic disorders.** *Hum Reprod* 2023, **38**:ii3-ii13.
- 374 5. Committee EPCS, Carvalho F, Coonen E, Goossens V, Kokkali G, Rubio C, Meijer-  
375 Hoogeveen M, Moutou C, Vermeulen N, De Rycke M: **ESHRE PGT Consortium**  
376 **good practice recommendations for the organisation of PGT.** *Hum Reprod Open*  
377 2020, **2020**:hoaa021.
- 378 6. Ginoza MEC, Isasi R: **Regulating Preimplantation Genetic Testing across the**  
379 **World: A Comparison of International Policy and Ethical Perspectives.** *Cold*  
380 *Spring Harb Perspect Med* 2020, **10**.
- 381 7. Garcia-Murias M, Quintans B, Arias M, Seixas AI, Cacheiro P, Tarrío R, Pardo J,  
382 Millan MJ, Arias-Rivas S, Blanco-Arias P, et al: **'Costa da Morte' ataxia is**  
383 **spinocerebellar ataxia 36: clinical and genetic characterization.** *Brain* 2012,  
384 **135**:1423-1435.
- 385 8. Kobayashi H, Abe K, Matsuura T, Ikeda Y, Hitomi T, Akechi Y, Habu T, Liu W,  
386 Okuda H, Koizumi A: **Expansion of intronic GGCCTG hexanucleotide repeat in**  
387 **NOP56 causes SCA36, a type of spinocerebellar ataxia accompanied by motor**  
388 **neuron involvement.** *Am J Hum Genet* 2011, **89**:121-130.
- 389 9. Obayashi M, Stevanin G, Synofzik M, Monin ML, Duyckaerts C, Sato N,  
390 Streichenberger N, Vighetto A, Desestret V, Tesson C, et al: **Spinocerebellar ataxia**



- 391            **type 36 exists in diverse populations and can be caused by a short hexanucleotide**  
392            **GGCCTG repeat expansion.** *J Neurol Neurosurg Psychiatry* 2015, **86**:986-995.
- 393    10.    Dessau RB, Pipper CB: **["R"--project for statistical computing].** *Ugeskr Laeger*  
394            2008, **170**:328-330.
- 395    11.    Zou J, Wang F, Gong Z, Wang R, Chen S, Zhang H, Sun R, Gao C, Li W, Shang J,  
396            Zhang J: **A Chinese SCA36 pedigree analysis of NOP56 expansion region based**  
397            **on long-read sequencing.** *Front Genet* 2023, **14**:1110307.
- 398    12.    Wang Q, Zhang C, Liu S, Liu T, Ni R, Liu X, Zhong P, Wu Q, Xu T, Ke H, et al:  
399            **Long-read sequencing identified intronic (GGCCTG)<sub>n</sub> expansion in NOP56 in**  
400            **one SCA36 family and literature review.** *Clin Neurol Neurosurg* 2022, **223**:107503.
- 401    13.    Lee YC, Tsai PC, Guo YC, Hsiao CT, Liu GT, Liao YC, Soong BW:  
402            **Spinocerebellar ataxia type 36 in the Han Chinese.** *Neurol Genet* 2016, **2**:e68.
- 403    14.    Ikeda Y, Ohta Y, Kobayashi H, Okamoto M, Takamatsu K, Ota T, Manabe Y,  
404            Okamoto K, Koizumi A, Abe K: **Clinical features of SCA36: a novel**  
405            **spinocerebellar ataxia with motor neuron involvement (Asidan).** *Neurology* 2012,  
406            **79**:333-341.
- 407    15.    Lopez S, He F: **Spinocerebellar Ataxia 36: From Mutations Toward Therapies.**  
408            *Front Genet* 2022, **13**:837690.
- 409    16.    Valera JM, Diaz T, Petty LE, Quintans B, Yanez Z, Boerwinkle E, Muzny D,  
410            Akhmedov D, Berdeaux R, Sobrido MJ, et al: **Prevalence of spinocerebellar ataxia**  
411            **36 in a US population.** *Neurol Genet* 2017, **3**:e174.
- 412    17.    Zeng S, Zeng J, He M, Zeng X, Zhou Y, Liu Z, Xia K, Pan Q, Jiang H, Shen L, et al:  
413            **Genetic and clinical analysis of spinocerebellar ataxia type 36 in Mainland China.**  
414            *Clin Genet* 2016, **90**:141-148.

- 415 18. Chintalaphani SR, Pineda SS, Deveson IW, Kumar KR: **An update on the**  
416 **neurological short tandem repeat expansion disorders and the emergence of**  
417 **long-read sequencing diagnostics.** *Acta Neuropathol Commun* 2021, **9**:98.
- 418 19. Xu H, Wang S, Ma LL, Huang S, Liang L, Liu Q, Liu YY, Liu KD, Tan ZM, Ban H,  
419 et al: **Informative priors on fetal fraction increase power of the noninvasive**  
420 **prenatal screen.** *Genet Med* 2018, **20**:817-824.
- 421 20. **Tandem repeats in the long-read sequencing era.** *Nat Rev Genet* 2024, **25**:449.
- 422 21. van Dijk EL, Naquin D, Gorrichon K, Jaszczyszyn Y, Ouazahrou R, Thermes C,  
423 Hernandez C: **Genomics in the long-read sequencing era.** *Trends Genet* 2023,  
424 **39**:649-671.
- 425 22. Taylor DJ, Chhetri SB, Tassia MG, Biddanda A, Battle A, McCoy RC: **Sources of**  
426 **gene expression variation in a globally diverse human cohort.** *bioRxiv* 2023.

427

428 **Figure legends:**

429

430 **Figure 1.** Schematic process of literature review and screening.

431

432 **Figure 2.** Prior, likelihood, posterior distributions and corresponding 95% confidence  
433 interval.

434

435 **Figure 3.** Identification of Expanded GGCCTG Repeat within NOP56 in the SCA family  
436 pedigree. (A) The pedigree of a family with spinocerebellar ataxia 36. (B) Magnetic  
437 resonance imaging indicates cerebellar atrophy (white arrows) of II-2 and her corresponding  
438 conventional PCR shows one peak in an allele and another a characteristic ladder pattern with

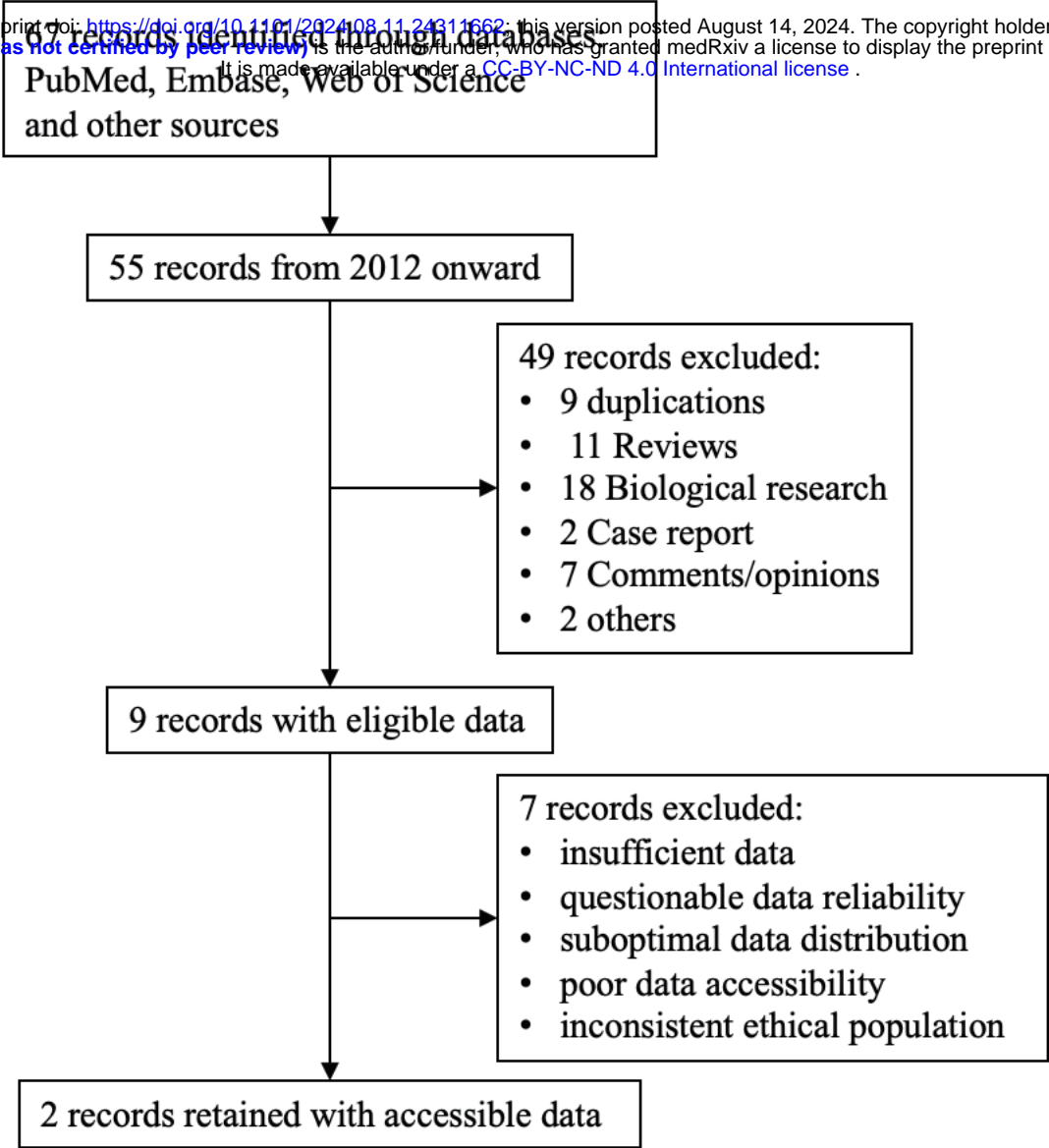
439 a typical 6-bp periodicity repeat expansion in NOP56. (C) Long-read sequencing to ascertain

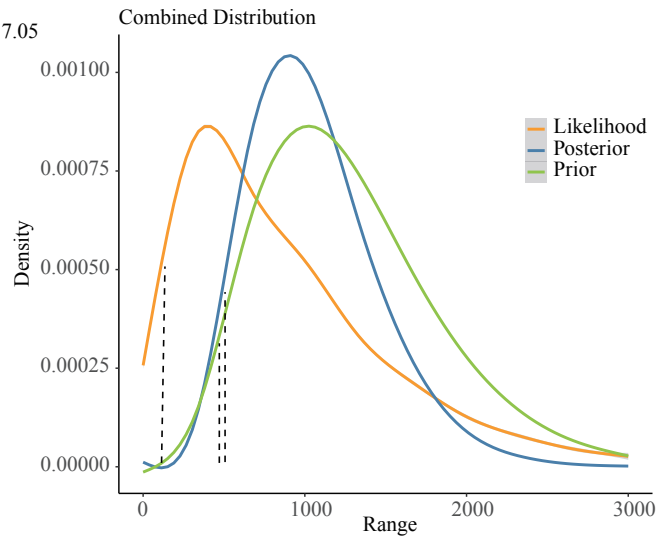
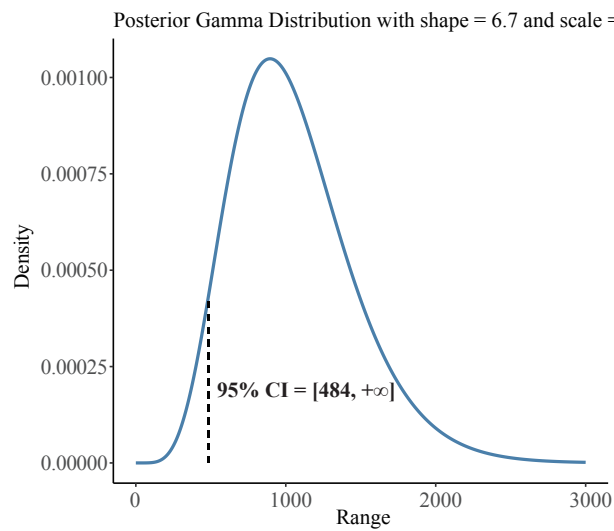
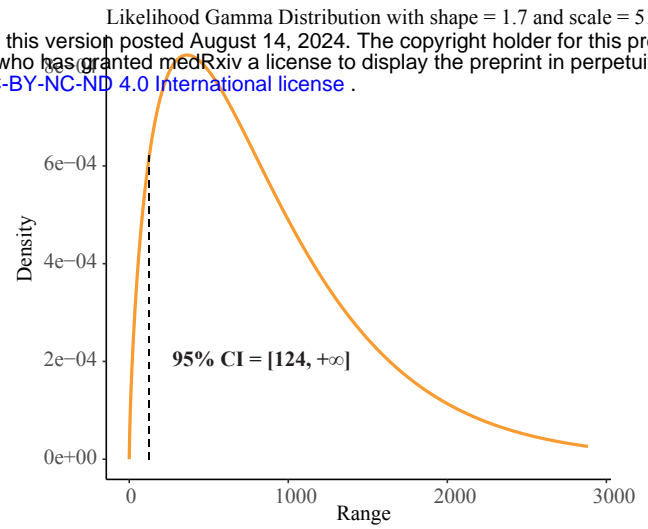
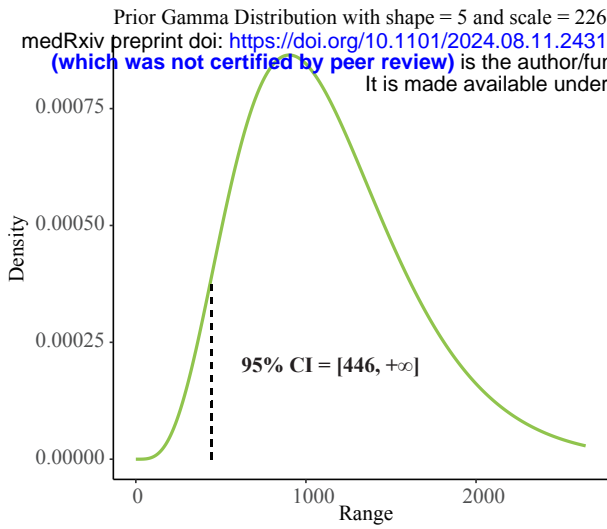
440 the precise repeat range of II-2 and III-1.

441

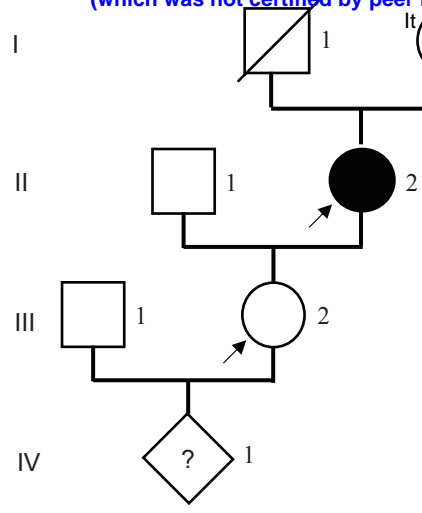
Table 1. Summary of 9 records with eligible data.

tim	Year	Ethnic group	Count	Pathogenic repeats	Mean age at onset (years)	Detection method	Reference
Garcia-Murias et al.	2012	Spain	27	650-2500	$52.5 \pm 7.2$	Southern blot	[7]
Obayashi et al.	2014	French/ German/Japanese	17	21-2000	$50.4 \pm 7.2$	PCR-fragment analyses	[9]
Zou et al.	2023	Chinese	3	60-2023	$44.5 \pm 3.8$	Long-read sequencing	[11]
Wang et al.	2022	Chinese	4	782 ~	$45 \pm /$	Long-read sequencing	[12]
Lee et al.	2016	Chinese	11	Not detected	$44.8 \pm 3.8$	PCR-fragment analysis	[13]
Ikeda et al.	2012	Japanese	18	1700-2300	$52.8 \pm 4.3$	Southern blot	[14]
Lam et al.	2022	British	7	30~	$48.4 \pm /$	Whole Genome Sequence	[15]
Valera et al.	2017	USA	6	1500~	$44.5 \pm 5.8$	PCR-fragment analyses	[16]
Zeng et al.	2015	Chinese	14	1000~	$50.82 \pm 5.02$	Southern blot	[17]





A

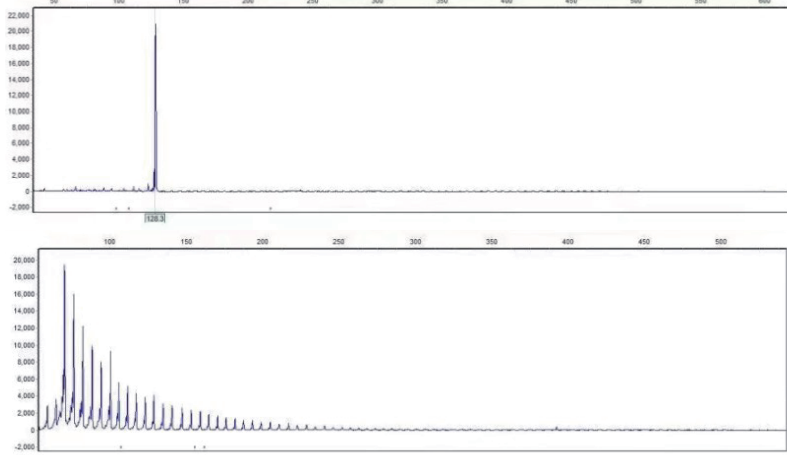
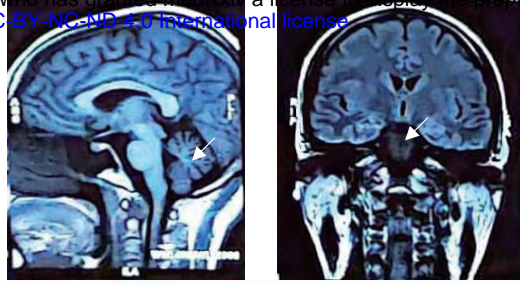


● Proband ○ Consultant  
 □ ○ Unaffected  
 ◻ ◊ Deceased ◊ Unborn

B

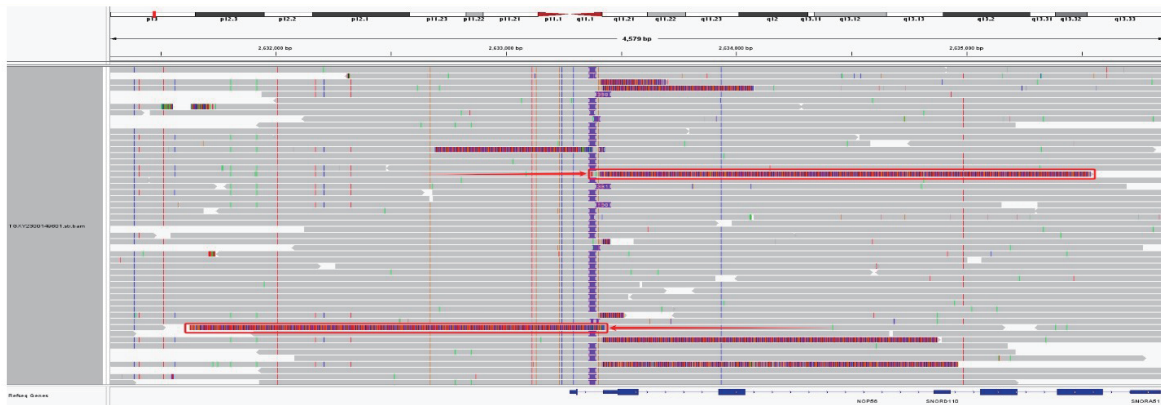
medRxiv preprint doi: <https://doi.org/10.1101/2024.08.11.24311662>; this version posted August 14, 2024. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted medRxiv a license to display the preprint in perpetuity. It is made available under a [CC BY-NC-ND 4.0 International license](https://creativecommons.org/licenses/by-nc-nd/4.0/).

II-2



C

II-2



III-2

