

Long-read sequencing identifies the first retrotransposon insertion and resolves structural variants causing antithrombin deficiency.

Journal:	<i>Thrombosis and Haemostasis</i>
Manuscript ID	TH-21-09-0539.R1
Manuscript Type:	Original Article: New Technologies, Diagnostic Tools and Drugs
Category:	Clinical Studies
Date Submitted by the Author:	n/a
Complete List of Authors:	<p>de la Morena-Barrio, Belén; Servicio de Hematología y Oncología Médica, Hospital Universitario Morales Meseguer, Centro Regional de Hemodonación, Universidad de Murcia, IMIB-Arrixaca, CIBERER</p> <p>Stephens, Jonathan ; Department of Haematology, University of Cambridge, NHS Blood and Transplant Centre, Cambridge, CB2 0PT, UK; NIHR BioResource, Cambridge University Hospitals NHS Foundation Trust, Cambridge Biomedical Campus, Cambridge, CB2 0QQ, UK.</p> <p>de la Morena-Barrio, Maria Eugenia; Servicio de Hematología y Oncología Médica, Hospital Universitario Morales Meseguer, Centro Regional de Hemodonación, Universidad de Murcia, IMIB-Arrixaca, CIBERER, Luca, Stephanucci; Department of Haematology, University of Cambridge, NHS Blood and Transplant Centre, Cambridge, CB2 0PT, UK; National Health Service Blood and Transplant (NHSBT), Cambridge Biomedical Campus, Cambridge, CB2 0PT, UK.; BHF Centre of Excellence, Division of Cardiovascular Medicine, Addenbrooke's Hospital, Cambridge Biomedical Campus, Cambridge, CB2 0QQ, UK</p> <p>Padilla, Jose; Centro Regional de Hemodonación, Miñano, Antonia; Servicio de Hematología y Oncología Médica, Hospital Universitario Morales Meseguer, Centro Regional de Hemodonación, Universidad de Murcia, IMIB-Arrixaca</p> <p>Gleadall, Nicholas ; Department of Haematology, University of Cambridge, NHS Blood and Transplant Centre, Cambridge, CB2 0PT, UK; NIHR BioResource, Cambridge University Hospitals NHS Foundation Trust, Cambridge Biomedical Campus, Cambridge, CB2 0QQ, UK.</p> <p>García, Juan Luis; Servicio de Hematología, Hospital Universitario de Salamanca, Salamanca, Spain.</p> <p>Lopez Fernandez, M Fernanda; C.Hospitalario Universitario A Coruña, Hematology</p> <p>Morange, Pierre; Inserm, Hematology</p> <p>Puurunen, Marja; Boston University School of Medicine, Framingham Heart Study</p> <p>Undas, Anetta; Jagiellonian University School of Medicine, Department of Medicine</p> <p>Vidal, Francisco; Blood and Tissue Bank (BST), Congenital Coagulopathies Laboratory; Vall d'Hebron University Hospital, Barcelona, Spain, Haemophilia Unit</p>

	<p>Raymond, F Lucy; NIHR BioResource, Cambridge University Hospitals NHS Foundation Trust, Cambridge Biomedical Campus, Cambridge, CB2 0QQ, UK.; Department of Medical Genetics, University of Cambridge, Cambridge Biomedical Campus, Cambridge, UK</p> <p>Vincente, Vicente; Servicio de Hematología y Oncología Médica, Hospital Universitario Morales Meseguer, Centro Regional de Hemodonación, Universidad de Murcia, IMIB-Arrixaca, CIBERER</p> <p>Ouwehand, Willem; Department of Haematology, University of Cambridge, NHS Blood and Transplant Centre, Cambridge, CB2 0PT, UK; NIHR BioResource, Cambridge University Hospitals NHS Foundation Trust, Cambridge Biomedical Campus, Cambridge, CB2 0QQ, UK.</p> <p>Corral, Javier; Servicio de Hematología y Oncología Médica, Hospital Universitario Morales Meseguer, Centro Regional de Hemodonación, Universidad de Murcia, IMIB-Arrixaca, CIBERER</p> <p>Sanchis-Juan, Alba; Department of Haematology, University of Cambridge, NHS Blood and Transplant Centre, Cambridge, CB2 0PT, UK; NIHR BioResource, Cambridge University Hospitals NHS Foundation Trust, Cambridge Biomedical Campus, Cambridge, CB2 0QQ, UK.</p>
Keywords:	Long-read sequencing, Antithrombin deficiency, Structural variants, SVA retrotransposon
Abstract:	<p>The identification of inherited antithrombin deficiency (ATD) is critical to prevent potentially life-threatening thrombotic events. Causal variants in SERPINC1 are identified for up to 70% of cases, the majority being single-nucleotide variants and indels. The detection and characterization of structural variants (SVs) in ATD remain challenging due to the high number of repetitive elements in SERPINC1. Here, we performed long-read whole-genome sequencing on 10 familial and 9 singleton cases with type I ATD proven by functional and antigen assays, that were selected from a cohort of 340 patients with this rare disorder because genetic analyses were either negative, ambiguous, or not fully characterized. We developed an analysis workflow to identify disease-associated SVs. This approach resolved, independently of its size or type, all 8 SVs detected by MLPA, and identified for the first time a complex rearrangement previously misclassified as a deletion. Remarkably, we identified the mechanism explaining ATD in 2 out of 11 cases with previous unknown defect: the insertion of a novel 2.4Kb SINE-VNTR-Alu retroelement, which was characterized by de novo assembly and verified by specific PCR amplification and sequencing in the probands and affected relatives. The nucleotide-level resolution achieved for all SVs allowed breakpoint analysis, which revealed repetitive elements and microhomologies supporting a common replication-based mechanism for all the SVs. Our study underscores the utility of long-read sequencing technology as a complementary method to identify, characterize and unveil the molecular mechanism of disease-causing SVs involved in ATD, and enlarges the catalogue of genetic disorders caused by retrotransposon insertions.</p>

Title

Long-read sequencing identifies the first retrotransposon insertion and resolves structural variants causing antithrombin deficiency.

Running title

Long-read sequencing to resolve structural variants in *SERPINC1*

Author list

Belén de la Morena-Barrio,¹ Jonathan Stephens,^{2,3} María Eugenia de la Morena-Barrio,¹ Luca Stefanucci,^{2,4,5} José Padilla,¹ Antonia Miñano,¹ Nicholas Gleadall,^{2,3} Juan Luis García,⁶ María Fernanda López-Fernández,⁷ Pierre-Emmanuel Morange,⁸ Marja K Puurunen,⁹ Anetta Undas,¹⁰ Francisco Vidal,¹¹ NIHR BioResource,³ F Lucy Raymond,^{3,12} Vicente Vicente García,¹ Willem H Ouweland,^{2,3} Javier Corral,¹ Alba Sanchis-Juan^{2,3}

Affiliations

1. Servicio de Hematología y Oncología Médica, Hospital Universitario Morales Meseguer, Centro Regional de Hemodonación, Universidad de Murcia, Instituto Murciano de Investigación Biosanitaria (IMIB-Arixaca), Centro de Investigación Biomédica en Red de Enfermedades Raras (CIBERER), Murcia, Spain.
2. Department of Haematology, University of Cambridge, NHS Blood and Transplant Centre, Cambridge, CB2 0PT, UK.
3. NIHR BioResource, Cambridge University Hospitals NHS Foundation Trust, Cambridge Biomedical Campus, Cambridge, CB2 0QQ, UK.

- 1
- 2
- 3
- 4 4. National Health Service Blood and Transplant (NHSBT), Cambridge Biomedical
- 5 Campus, Cambridge, CB2 0PT, UK.
- 6
- 7
- 8 5. BHF Centre of Excellence, Division of Cardiovascular Medicine, Addenbrooke's
- 9 Hospital, Cambridge Biomedical Campus, Cambridge, CB2 0QQ, UK.
- 10
- 11
- 12 6. Servicio de Hematología, Hospital Universitario de Salamanca, Salamanca, Spain.
- 13
- 14 7. Servicio de Hematología, Complejo Hospitalario Universitario de A Coruña, A
- 15 Coruña, Spain.
- 16
- 17
- 18 8. Laboratory of Haematology, La Timone Hospital, Marseille, France; C2VN,
- 19 INRAE, INSERM, Aix-Marseille Université, Marseille, France.
- 20
- 21
- 22
- 23 9. National Heart, Lung and Blood Institute's. The Framingham Heart Study,
- 24 Framingham, MA, US.
- 25
- 26
- 27
- 28 10. Institute of Cardiology, Jagiellonian University Medical College and John Paul II
- 29 Hospital, 80 Prądnicka St, Kraków, Poland.
- 30
- 31
- 32 11. Banc de Sang i Teixits, Barcelona, Spain; Vall d'Hebron Research Institute,
- 33 Universitat Autònoma de Barcelona (VHIR-UAB), Barcelona, Spain; CIBER de
- 34 Enfermedades Cardiovasculares, Madrid, Spain.
- 35
- 36
- 37
- 38 12. Department of Medical Genetics, University of Cambridge, Cambridge Biomedical
- 39 Campus, Cambridge, UK.
- 40
- 41
- 42
- 43
- 44
- 45
- 46

47 **Correspondence**

48
49 Alba Sanchis-Juan, University of Cambridge, Department of Haematology, NHS Blood
50 and Transplant Centre, Cambridge, CB2 0PT, UK. Email: as2635@cam.ac.uk. Phone
51 number: +44(0)1223588035.
52
53
54
55
56
57
58
59
60

1
2
3 Javier Corral, University of Murcia, Centro Regional de Hemodonación, Calle Ronda
4 de Garay s/n, Murcia 30003, Spain. Email: javiercorraldelacalle@gmail.com. Phone
5
6
7 number: +34968341990.
8
9

11 12 **Authorship Contributions**

13
14 BMB, WHO, JC and ASJ designed the study.

15
16 MMB, LS, JP, AM, NG, FLR and VV helped with study design.

17
18 BMB, MMB, JP, AM performed laboratory experiments and analyzed the experimental
19
20
21 data.

22
23 JS performed sample preparation and executed long-read sequencing.

24
25 ASJ developed the analysis workflow for long-read sequencing, applied this to data
26
27
28 processing and performed the computational and statistical analyses.

29
30 BMB performed computational analyses and variant validation.

31
32 JM, FV, provided valuable insight into microarray and NGS data analysis.

33
34 AU, MF, MP and PM recruited participants and collected the clinical data and samples.

35
36 BMB, WHO, JC and ASJ wrote the manuscript.

37
38 All authors read and approved the final manuscript.
39
40
41
42
43

44 **Data and Code Availability**

45
46 The workflow developed for the detection of structural variants is publicly available at
47
48 <http://github.com/who-blackbird/magpie>.

51 52 **Funding**

53
54 This work was supported by the National Institute for Health Research England (NIHR)
55
56
57 for the NIHR BioResource project (grant numbers RG65966 and RG94028), the
58
59
60

1
2
3
4 PI18/00598, PI21/00174 and PMP21/00052 projects (ISCIII & FEDER) and the
5 19873/GERM/15 project (Fundación Séneca).
6
7
8

9 **Declaration of Interests**

10
11 The authors declare that they have no conflicts of interest.
12
13
14
15

16 **Patient consent statement**

17
18 All included subjects gave their written informed consent to enter the study.
19
20
21
22

23 **Ethics**

24
25 This study was approved by the Ethics Committee of Morales Meseguer Hospital and
26 the East of England Cambridge South national institutional review board (13/EE/0325).
27
28

29 The research conforms with the principles of the Declaration of Helsinki and their later
30 amendments.
31
32
33
34
35

36 **Keywords**

37
38 Long-read sequencing; Antithrombin deficiency; Structural variants; SVA
39 retrotransposon.
40
41
42
43
44
45

46 **Abstract**

47
48 The identification of inherited antithrombin deficiency (ATD) is critical to prevent
49 potentially life-threatening thrombotic events. Causal variants in *SERPINC1* are
50 identified for up to 70% of cases, the majority being single-nucleotide variants and
51 indels. The detection and characterization of structural variants (SVs) in ATD remain
52 challenging due to the high number of repetitive elements in *SERPINC1*. Here, we
53 performed long-read whole-genome sequencing on 10 familial and 9 singleton cases
54
55
56
57
58
59
60

with type I ATD proven by functional and antigen assays, that were selected from a cohort of 340 patients with this rare disorder because genetic analyses were either negative, ambiguous, or not fully characterized. We developed an analysis workflow to identify disease-associated SVs. This approach resolved, independently of its size or type, all 8 SVs detected by MLPA, and identified for the first time a complex rearrangement previously misclassified as a deletion. Remarkably, we identified the mechanism explaining ATD in 2 out of 11 cases with previous unknown defect: the insertion of a novel 2.4Kb SINE-VNTR-Alu retroelement, which was characterized by *de novo* assembly and verified by specific PCR amplification and sequencing in the probands and affected relatives. The nucleotide-level resolution achieved for all SVs allowed breakpoint analysis, which revealed repetitive elements and microhomologies supporting a common replication-based mechanism for all the SVs. Our study underscores the utility of long-read sequencing technology as a complementary method to identify, characterize and unveil the molecular mechanism of disease-causing SVs involved in ATD, and enlarges the catalogue of genetic disorders caused by retrotransposon insertions.

WC: 250

Introduction

Antithrombin deficiency is the most severe congenital thrombophilia firstly identified in 1965 by O Egeberg¹. The key hemostatic role of this anticoagulant serpin explains the high risk of thrombosis associated to congenital antithrombin deficiency (OR: 20-30), which is mainly caused by haploinsufficiency of *SERPINC1*, the coding gene.² Accurate genetic diagnosis of antithrombin deficiency facilitates the management of both symptomatic and asymptomatic carriers^{3,4}, and increases the antithrombotic arsenal of

1
2
3 carriers with antithrombin concentrates.⁵ Routine investigation of antithrombin
4
5 deficiency combines functional assays, antigen quantification and genetic analyses to
6
7 determine the molecular base. However, most studies do not reach a molecular
8
9 characterization, despite it could contribute to a better definition of the thrombotic risk.²
10
11
12

13
14 In genetic diagnostic centers, causal **Single Nucleotide Variants (SNVs) and small**
15
16 **insertions or deletions (indels)** are routinely identified in *SERPINC1* by Sanger
17
18 sequencing, and copy number changes are investigated by multiple ligation-dependent
19
20 probe amplification (MLPA).² Only few cases with gross gene defects have been
21
22 analyzed by microarray to determine the extent of the variants. These methods identify
23
24 causal mutations in *SERPINC1* for 70% of cases, whilst 5% of patients harbor defects in
25
26 other genes and 25% remain without a genetic diagnosis.² To date, 441 causal variants
27
28 in *SERPINC1* have been reported,⁶ and these adhere to the typical spectrum observed in
29
30 disorders with a dominant inheritance, being 63% SNVs, 28% indels and 9% structural
31
32 variants (SVs).^{7,8}
33
34
35
36
37
38
39

40 However, there are important limitations to these techniques, including that neither
41
42 MLPA nor microarray consider the full spectrum of SVs and do not provide nucleotide-
43
44 level resolution, which is important for confirming causality and reveal insights into
45
46 SVs formation.^{7,9,10} These limitations may now be addressed by long-reads, that can
47
48 span repetitive or other problematic regions, allowing identification and characterization
49
50 of SVs.¹⁰⁻¹⁴ This is particularly advantageous for antithrombin deficiency due to the
51
52 high number of repetitive elements in and around *SERPINC1* (where 35% of sequence
53
54 are interspersed repeats)¹⁵, that hinders SVs identification by other methods.
55
56
57
58
59
60

1
2
3 Here, we report on the results of long-read whole-genome sequencing (LR-WGS) on 19
4 unrelated cases with antithrombin deficiency, selected from one of the largest cohort of
5 patients with this disorder based on negative or ambiguous results, as well as not fully
6 characterized SVs provided by routine molecular tests. Our aim was to identify new
7 causal variants, resolve ambiguous ones and investigate the most likely mechanism of
8 formation of SVs involved in this severe thrombophilia.
9
10
11
12
13
14
15
16
17
18

19 **Methods**

20 *Cohort*

21
22
23 Nineteen unrelated individuals with antithrombin deficiency were selected from our
24 cohort of 340 cases, recruited between 1994 and 2019 and largely characterized by
25 functional, biochemical, and molecular analyses. Selection was done based on negative
26 results from multiple genetic studies evaluating *SERPINC1* gene, including Sanger
27 sequencing followed by Next-Generation Sequencing (NGS) and MLPA, as well as
28 negative glycosylation analysis (N= 11). Additionally, individuals with SVs that could
29 not be characterized or that were identified by MLPA but had ambiguous results from
30 other approaches (such as microarray and/or long-range PCR) were also selected (N= 8)
31 (Table 1). Detailed information of these procedures is shown in Supplemental Methods.
32
33 Measurement of antithrombin levels and function were performed for all participants as
34 previously described.^{16,17}
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50

51 *Long-read whole-genome sequencing*

52
53 Long-read whole-genome sequencing (LR-WGS) of DNAs purified from peripheral
54 blood leukocytes using Gentra Puregene Qiagen kit, used to reduce the fragmentation of
55 DNA, was done using the PromethION platform (Oxford Nanopore Technologies).
56
57
58
59
60

1
2
3 Samples were prepared using the 1D ligation library prep kit (SQK-LSK109) and
4 genomic libraries were sequenced on R9 flow cell. Read sequences were extracted from
5
6
7 base-called FAST5 files by Guppy (versions 3.0.4 to 3.2.8; 3.0.4+e7dbc23 to
8
9
10 3.2.8+bd67289) to generate FASTQ files, that were then merged per sample.

11 12 13 14 15 *Data processing and SV identification*

16
17 We used the Snakemake library to develop an *in-house* multi-modal analysis workflow
18
19 for the sensitive detection of SVs,¹⁸ which is publicly available at
20
21 <https://github.com/who-blackbird/magpie>. An overview of the workflow is shown in
22
23 Figure 1A. Detailed information is provided in Supplemental Methods.

24 25 26 27 28 *De novo assembly of the SINE-VNTR-Alu retroelement*

29
30 Local *de novo* assembly was performed to characterize the SINE-VNTR-Alu
31
32 retroelement insertion in P9. Reads within the region [GRCh38/hg38]
33
34 Chr1:173,840,000-174,820,000 were extracted from the alignment of this individual and
35
36 converted to a FASTQ file using Samtools.¹⁹ *De novo* assembly was performed with
37
38 wtdbg2 v2.5, using the parameters ``-x ont -g 980k -X 10 -e 3``.²⁰ The *de*
39
40
41
42 *nov*o contig was then aligned to the reference genome using minimap2²¹ with default
43
44 parameters for nanopore reads. The genomic sequence containing the SINE-VNTR-Alu
45
46 retroelement was then extracted from the alignment and analyzed with RepeatMasker
47
48 (<http://www.repeatmasker.org>) to characterize the type of SINE-VNTR-Alu and its
49
50
51 subelements.
52
53
54
55

56 57 58 59 60 *Validations and breakpoint flanking sequence analysis*

1
2
3 All candidate SV junctions were confirmed by PCR amplification and Sanger
4 sequencing to verify all variant configurations at nucleotide level resolution. We then
5 manually identified the presence of microhomology, insertions and deletions at the
6 breakpoints as previously described.²² The percentage of repetitive sequence was also
7 calculated for each junction (+/- 150 bps) by intersecting these regions with the human
8 genomic repeat library (hg38) from RepeatMasker version open-4.0.5 using bedtools.²³
9
10
11
12
13
14
15
16
17
18

19 **Results**

20 *Long-read sequencing identifies SVs involving SERPINC1*

21
22
23
24
25
26 Nanopore sequencing in 21 runs produced reads with an average length of 4,499bp and
27 median genome coverage of 16x (Figure 1B). After a detailed quality control analysis
28 (Figure S1), 83,486 SVs were identified, consistent with previous reports using LR-
29 WGS (Figure S2).¹¹ Focusing on rare variants (allele count <= 10 in gnomAD v3, NIH
30 BioResource and NGC project)^{11,24,25} in *SERPINC1* and flanking regions, 10 candidate
31 heterozygous SVs were observed in 9 individuals (Figure 1C). Visual inspection of read
32 alignments identified an additional heterozygous SV in a region of low coverage
33 involving *SERPINC1* in an additional patient (Table 1).
34
35
36
37
38
39
40
41
42
43
44
45
46

47 *Resolution of causal SVs: identification of the first complex SV*

48
49
50
51 Nanopore sequencing resolved the precise configuration of all SVs previously identified
52 by MLPA in 8 individuals (P1-P8). Structural variants were identified independently of
53 their size (from 7Kb to 968Kb, restricted to *SERPINC1* or involving neighboring genes)
54 and their type (six deletions, one tandem duplication and one complex SV) (Figure 2,
55
56
57
58
59
60

1
2
3 **Table 1**). In all the cases the extension of the variants was determined, and nucleotide
4
5 level resolution **of breakpoints** was achieved by the long reads **(Table 1)**. Importantly,
6
7 nanopore sequencing facilitated the resolution of the SVs identified in two patients (P2
8
9 and P6) that presented inconsistent or ambiguous results from MLPA and long-range
10
11 PCR and NGS results **(Table 1)**.
12
13
14
15
16

17 For the first case (P2) MLPA detected a deletion of exon 1, but long-range PCR
18
19 followed by NGS suggested a deletion of exons 1 and 2. The discordant results were
20
21 explained by nanopore sequencing, as this method revealed a complex SV in
22
23 *SERPINC1* resulting in a dispersed duplication of exons 2 and 3 and a deletion spanning
24
25 exons 1 and 2, both in the same allele (Figure 3). Specific PCR amplification and
26
27 Sanger sequencing validated this complex structural variant in the proband and his
28
29 affected daughter also with antithrombin deficiency.
30
31
32
33
34

35 For the second case (P6) MLPA detected a duplication of exons **2, 3 and 5** and a
36
37 deletion of exon 6. Here, our sequencing approach identified a tandem duplication of
38
39 exons 1 to 5, which was confirmed by long-range PCR (Figure 4). The tandem
40
41 duplication of exons 1-5 was observed to be present in the affected son of P6, also with
42
43 antithrombin deficiency.
44
45
46
47
48

49 *A **SINE-VNTR-Alu** retroelement insertion is identified in two previously unresolved*
50
51 *cases and characterized by de novo assembly*
52
53
54
55

56 We aimed to identify new disease-causing variants in the remaining 11 participants with
57
58 negative results using current molecular methods. Remarkably, two cases (P9 and P10)
59
60

1
2
3 presented an insertion of 2,440bp in intron 6. Blast analysis of the inserted sequence
4 revealed a new **SINE-VNTR-Alu** retroelement (Figure 2, **Table 1**). Local *de novo*
5 assembly using the data from P9 revealed an antisense-oriented **SINE-VNTR-Alu**
6 element flanked by a target site duplication (TSD) of 14bp (Figure 2C), consistent with
7 a target-primed reverse transcription mechanism of insertion into the genome.^{26,27}
8 Interestingly, the TSD in both individuals was also the same. The inserted sequence was
9 aligned to the canonical **SINE-VNTR-Alu** A-F sequences (Figure S3A) and it was
10 observed to be closest to the **SINE-VNTR-Alu** E in the phylogenetic tree (Figure S3B).
11 Moreover, the VNTR sub-element harbored 1,449bp, which was longer than the typical
12 ~520bp-long VNTR in the canonical sequences. Multiple PCRs covering the
13 retroelement were attempted to validate this insertion, but all PCRs using flanking
14 primers failed **due to the highly repetitive sequence of this element, specially the VNTR**
15 **sub-element, which is longer in this new SINE-VNTR-Alu**. Only one specific PCR
16 using an internal **SINE-VNTR-Alu** primer, which designed was facilitated by the
17 Nanopore data, was able to amplify the breakpoint (Figure S4). This method was used
18 to confirm the insertion in P9, P10 and to confirm the Mendelian inheritance of this
19 **SINE-VNTR-Alu**, as it was also present in two affected relatives, both with
20 antithrombin deficiency (Figure S4).
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46

Breakpoint analysis supports a replication-based mechanism for the majority of SVs

47
48
49
50
51 Breakpoint analysis was performed to investigate the mechanism underlying the
52 formation of these SVs involving *SERPINC1*. Nanopore sequencing facilitated primer
53 design to perform Sanger sequencing confirmations for all the new formed junctions,
54 demonstrating a 100% accuracy in 7/10 (70%) SVs called. Repetitive elements (RE)
55
56
57
58
59
60

1
2
3 were detected in all the SVs, with Alu elements being the most frequent (16/24, 67%)
4 (Table S1). Additionally, breakpoint analysis identified microhomologies (7/11, 64%)
5 and insertions, deletions or duplications (7/11, 64%) (Figure S5 and Table S2).
6
7 Importantly, we observed a non-random formation driven by the presence of REs in
8 some of the SVs. We point out an Alu element in intron 5, involved in SVs of P6, P7
9 and P8 (Figure 2B, Table S1).
10
11
12
13
14
15
16
17
18

19 Discussion

20
21 In this study we aimed to resolve the precise configuration of SVs involved in
22 antithrombin deficiency using nanopore, to identify new candidate variants in
23 previously unresolved cases and to investigate the possible mechanisms of formation of
24 these SVs by breakpoint analysis. We have characterized disease-causing SVs in eight
25 individuals with previous positive findings from MLPA and other methods but with
26 unresolved variants, in two cases with previous contradictory results. Additionally, we
27 reported a new causal SINE-VNTR-Alu retroelement insertion in two unrelated
28 individuals that we characterized by local *de novo* assembly. Finally, we presented
29 evidence for a replication-based mechanism of formation for most of the SVs causing
30 this severe thrombophilia.
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46

47 We show new evidence of how LR-WGS can be used to identify SVs causal of a
48 genetic disease, in this case antithrombin deficiency, independently of its length or type.
49 LR-WGS also gives information for the exact extension of the event involved and
50 resolves conflictive data obtained by other methods. Additionally, we show how this
51 approach is particularly powerful to investigate complex SVs, which are genomic
52 rearrangements typically composed of three or more breakpoint junctions. Since these
53
54
55
56
57
58
59
60

1
2
3 are particularly challenging to detect and interpret by other methods, complex SVs are
4 typically missed or misclassified in research and clinical diagnostic pipelines, although
5 they have been reported as associated with multiple Mendelian diseases.¹⁰ Here we
6 show for the first time a complex SV in a patient with antithrombin deficiency,
7 expanding the landscape of SV types involved in this disorder. Further investigations
8 will be required to elucidate the exact mechanism of formation, since it remains unclear
9 if this event occurred by one or multiple mutational events.

10
11
12
13
14
15
16
17
18
19
20
21 Additionally, we identified an intronic **SINE-VNTR-Alu retroelement** insertion in 2/11
22 (18%) previously unresolved individuals (P9 and P10). **SINE-VNTR-Alu** retroelements,
23 along with other retrotransposons, are a source of regulatory variation in the human
24 genome, but can also cause disease.²⁸ Although the number of pathogenic retroelements
25 has increased during the last years with the use of WGS technologies,^{25,29–31} these are
26 usually missed by routine diagnostic methods. With LR-WGS we have not only
27 identified the causal mutation in two previously unresolved families, but also performed
28 local *de novo* assembly to characterize the exact sequence and length of its sub-
29 elements, which might be relevant for future studies to investigate their possible role in
30 severity and age of disease onset as other studies have shown.³²

31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47 Furthermore, the genomic heterogeneity observed between the causal **SINE-VNTR-Alu**
48 **retroelement** and the canonical sequences highlights the diverse genomic landscape of
49 these retroelements and underscores the importance of their characterization in order to
50 obtain a reliable catalogue of novel mobile elements to identify and interpret this type of
51 causal variants in other patients and other disorders where retrotransposon insertions
52 might also be involved.^{27,33,34} This characterization has been historically challenging by
53
54
55
56
57
58
59
60

1
2
3 the application of classic technologies, but here we show that it can be achieved by *de*
4
5 *novo* assembly of long-reads.
6
7

8
9
10 The decreased levels of antithrombin in plasma of P9 and P10 might be consistent with
11
12 transcriptional interference of *SERPINC1* induced by the SINE-VNTR-Alu
13
14 retroelement, as reported for other cases with pathogenic SINE-VNTR-Alu insertions.²⁸
15
16 Besides, the 2.4Kb insertion of a retroelement in intron 6 could introduce splicing
17
18 signals affecting the normal splicing of *SERPINC1* RNA. However, the specific hepatic
19
20 expression of *SERPINC1* hinders investigation of the exact mechanism, but the co-
21
22 segregation of this variant with antithrombin deficiency observed in family studies of
23
24 both probands supports the pathogenic consequences of this insertion. The identification
25
26 of the same retrotransposon in two unrelated families from different regions of Spain
27
28 (570 km far from each other) with the same TSD, does not only support the germline
29
30 transmission of this SV, but also suggests a shared mechanism of formation or a founder
31
32 effect, which must be confirmed by further studies.
33
34
35
36
37
38
39

40 In antithrombin deficiency, the detection and characterization of SVs remain
41
42 particularly challenging due to the high number of repetitive elements in and around
43
44 *SERPINC1* (35% of sequence in these gene are interspersed repeats). Specific
45
46 mutational signatures can yield insights into the mechanisms by which the SVs are
47
48 formed. Our breakpoint analysis suggested for most of the cases (P1-P8) a replication-
49
50 based mechanism (such as BIR/MMBIR/FoSTeS),³⁵ consistent with previous studies
51
52 done in antithrombin deficiency,^{36,37} but importantly, we observed a non-random
53
54 formation in some instances given the recurrent involvement of specific REs such as
55
56 *Alu* elements in intron 5 of *SERPINC1*. It has been suggested that RE may provide
57
58
59
60

1
2
3 larger tracks of microhomologies, also termed ‘microhomology islands’, that could
4 assist strand transfer or stimulate template switching during repair by a replication-
5 based mechanism.³⁵ These microhomology islands were present in the SVs of 3 cases (
6 P6, P7, P8), highlighting the important role that RE play in the formation of non-
7 recurrent, but non-random, SVs. These results highlight that *SERPINC1* might be a hot-
8 spot for SVs given the high number of REs in this gene and shows how LR-WGS can
9 be used to investigate and resolve events occurring in repetitive genes and regions.
10
11
12
13
14
15
16
17
18
19
20

21 In total, 9 cases in this cohort remain yet unresolved, three of whom reported to have
22 familial disease. An explanation may be that the causal variant was missed due to low
23 coverage, or alternatively the variant is located in an unidentified transacting gene or in
24 a regulatory element for *SERPINC1*, as we have recently reported for other genes.¹³ The
25 observation that the antithrombin deficiency in patients without causal SVs have
26 significantly higher anti-FXa activity than those with SVs (Figure 1D) is supportive of
27 the notion that causal variants may regulate gene expression, which must be analyzed in
28 future studies.
29
30
31
32
33
34
35
36
37
38
39
40
41

42 Altogether this study provides insight into the molecular mechanism of SVs causing
43 antithrombin deficiency and highlights the importance of identifying a new class of
44 causal variants to improve diagnostic rates, lead to new therapeutic opportunities, and
45 provide accurate family counselling, as decisions about long term anticoagulant
46 prophylaxis are complex and carry significant morbidity and mortality risks. Moreover,
47 our study suggests that SVs, which are often overlooked or misclassified by
48 conventional methods, may be more common than anticipated as a genetic mechanism
49 of antithrombin deficiency.
50
51
52
53
54
55
56
57
58
59
60

Description of Supplemental Data

Supplemental Material includes additional methods information, five figures and two tables.

Acknowledgments

We thank the participants involved in this study and their families. We thank NIHR BioResource volunteers for their participation, and gratefully acknowledge NIHR BioResource centers, NHS Trusts and staff for their contribution. We thank the National Institute for Health Research, NHS Blood and Transplant, and Health Data Research UK as part of the Digital Innovation Hub Programme. The views expressed are those of the author(s) and not necessarily those of the NHS, the NIHR or the Department of Health and Social Care.

References

1. EGEBERG O. THROMBOPHILIA CAUSED BY INHERITABLE DEFICIENCY OF BLOOD ANTITHROMBIN. *Scand J Clin Lab Invest*. 1965;17(1):92. ed. & tran. Accessed January 21, 2021.
2. Corral J, de la Morena-Barrio ME, Vicente V. The genetics of antithrombin. *Thromb Res*. 2018;169:23-29. ed. & tran. Accessed March 29, 2021.
3. Lijfering WM, Brouwer JLP, Veeger NJGM, et al. Selective testing for thrombophilia in patients with first venous thrombosis: Results from a retrospective family cohort study on absolute thrombotic risk for currently known thrombophilic defects in 2479 relatives. *Blood*. 2009;113(21):5314-5322.

- 1
2
3 ed. & tran. Accessed February 4, 2021.
4
5
6 4. Mahmoodi BK, Brouwer J-LP, Ten Kate MK, et al. A prospective cohort study on
7
8 the absolute risks of venous thromboembolism and predictive value of
9
10 screening asymptomatic relatives of patients with hereditary deficiencies of
11
12 protein S, protein C or antithrombin. *J Thromb Haemost.* 2010;8(6):1193-1200.
13
14 ed. & tran. Accessed August 2, 2018.
15
16
17 5. Bravo-Pérez C, Vicente V, Corral J. Management of antithrombin deficiency: an
18
19 update for clinicians. *Expert Rev Hematol.* 2019;12(6):397-405. ed. & tran.
20
21 Accessed January 22, 2021.
22
23
24 6. Stenson PD, Ball E V., Howells K, Phillips AD, Mort M, Cooper DN. The Human
25
26 Gene Mutation Database: providing a comprehensive central mutation database
27
28 for molecular diagnostics and personalized genomics. *Hum Genomics.*
29
30 2009;4(2):69-72.
31
32
33
34 7. Ordulu Z, Kammin T, Brand H, et al. Structural Chromosomal Rearrangements
35
36 Require Nucleotide-Level Resolution: Lessons from Next-Generation Sequencing
37
38 in Prenatal Diagnosis. *Am J Hum Genet.* 2016;99(5):1015-1033. ed. & tran.
39
40 Accessed March 29, 2021.
41
42
43
44 8. Beauchamp NJ, Makris M, Preston FE, Peake IR, Daly ME. Major structural
45
46 defects in the antithrombin gene in four families with type I antithrombin
47
48 deficiency. Partial/complete deletions and rearrangement of the antithrombin
49
50 gene. *Thromb Haemost.* 2000;83(5):715-721. ed. & tran. Accessed February 25,
51
52 2021.
53
54
55
56 9. Lam HYK, Mu XJ, Stütz AM, et al. Nucleotide-resolution analysis of structural
57
58 variants using BreakSeq and a breakpoint library. *Nat Biotechnol.* 2010;28(1):47-
59
60

- 1
2
3 55. ed. & tran. Accessed March 29, 2021.
4
5
6 10. Sanchis-Juan A, Stephens J, French CE, et al. Complex structural variants in
7
8 Mendelian disorders: identification and breakpoint resolution using short- and
9
10 long-read genome sequencing. *Genome Med.* 2018;10(1):95. ed. & tran.
11
12 Accessed March 18, 2020.
13
14
15 11. Beyter D, Ingimundardottir H, Eggertsson HP, et al. Long read sequencing of
16
17 1,817 Icelanders provides insight into the role of structural variants in human
18
19 disease. *bioRxiv.* November 2019:848366. ed. & tran. Accessed January 30,
20
21 2021.
22
23
24
25 12. Sedlazeck FJ, Rescheneder P, Smolka M, et al. Accurate detection of complex
26
27 structural variations using single-molecule sequencing. *Nat Methods.*
28
29 2018;15(6):461-468.
30
31
32
33 13. Cretu Stancu M, Van Roosmalen MJ, Renkens I, et al. Mapping and phasing of
34
35 structural variation in patient genomes using nanopore sequencing. *Nat*
36
37 *Commun.* 2017;8(1):1-13. ed. & tran. Accessed January 27, 2021.
38
39
40 14. French CE, Delon I, Dolling H, et al. Whole genome sequencing reveals that
41
42 genetic conditions are frequent in intensively ill children. *Intensive Care Med.*
43
44 2019;45(5):627-636. ed. & tran. Accessed March 1, 2021.
45
46
47 15. de Koning APJ, Gu W, Castoe TA, Batzer MA, Pollock DD. Repetitive elements
48
49 may comprise over Two-Thirds of the human genome. *PLoS Genet.* 2011;7(12).
50
51
52 16. de la Morena-Barrio M, Sandoval E, Llamas P, et al. High levels of latent
53
54 antithrombin in plasma from patients with antithrombin deficiency. *Thromb*
55
56 *Haemost.* 2017;117(5):880-888.
57
58
59 17. de la Morena-Barrio ME, Martínez-Martínez I, de Cos C, et al. Hypoglycosylation
60

- 1
2
3 is a common finding in antithrombin deficiency in the absence of a SERPINC1
4
5
6 gene defect. *J Thromb Haemost.* 2016;14(8):1549-1560.
7
- 8 18. De Coster W, De Rijk P, De Roeck A, et al. Structural variants identified by Oxford
9
10 Nanopore PromethION sequencing of the human genome. *Genome Res.*
11
12 2019;29(7):1178-1187. ed. & tran. Accessed June 2, 2021.
13
14
- 15 19. Li H, Handsaker B, Wysoker A, et al. The Sequence Alignment/Map format and
16
17 SAMtools. *Bioinformatics.* 2009;25(16):2078-2079. ed. & tran. Accessed
18
19 February 9, 2021.
20
21
- 22 20. Ruan J, Li H. Fast and accurate long-read assembly with wtdbg2. *Nat Methods.*
23
24 2020;17(2):155-158. ed. & tran. Accessed March 29, 2021.
25
26
- 27 21. Li H. Minimap2: Pairwise alignment for nucleotide sequences. *Bioinformatics.*
28
29 2018;34(18):3094-3100. ed. & tran. Accessed January 27, 2021.
30
31
- 32 22. Stankiewicz P, Lupski JR. Structural variation in the human genome and its role
33
34 in disease. *Annu Rev Med.* 2010;61(1):437-455. ed. & tran. Accessed November
35
36 28, 2018.
37
38
- 39 23. Quinlan AR, Hall IM. BEDTools: A flexible suite of utilities for comparing genomic
40
41 features. *Bioinformatics.* 2010;26(6):841-842. ed. & tran. Accessed February 9,
42
43 2021.
44
45
- 46 24. Köster J, Rahmann S. Snakemake - A scalable bioinformatics workflow engine.
47
48 *Bioinformatics.* 2018;34(20):3600. ed. & tran. Accessed February 9, 2021.
49
50
- 51 25. Turro E, Astle WJ, Megy K, et al. Whole-genome sequencing of patients with
52
53 rare diseases in a national health system. *Nature.* 2020;583(7814):96-102. ed. &
54
55 tran. Accessed February 14, 2021.
56
57
- 58 26. Vogt J, Bengesser K, Claes KBM, et al. SVA retrotransposon insertion-associated
59
60

- 1
2
3 deletion represents a novel mutational mechanism underlying large genomic
4
5 copy number changes with non-recurrent breakpoints. *Genome Biol.*
6
7 2014;15(6):15:R80. ed. & tran. Accessed March 29, 2021.
8
9
10 27. Payer LM, Burns KH. Transposable elements in human genetic disease. *Nat Rev*
11
12 *Genet.* 2019;20(12):760-772. ed. & tran. Accessed February 22, 2021.
13
14
15 28. Huang CRL, Burns KH, Boeke JD. Active transposition in genomes. *Annu Rev*
16
17 *Genet.* 2012;46:651-675. ed. & tran. Accessed February 22, 2021.
18
19
20 29. Nakamura Y, Murata M, Takagi Y, et al. SVA retrotransposition in exon 6 of the
21
22 coagulation factor IX gene causing severe hemophilia B. *Int J Hematol.*
23
24 2015;102(1):134-139. ed. & tran. Accessed May 31, 2021.
25
26
27 30. Van der Klift HM, Tops CM, Hes FJ, Devilee P, Wijnen JT. Insertion of an SVA
28
29 element, a nonautonomous retrotransposon, in PMS2 intron 7 as a novel cause
30
31 of lynch syndrome. *Hum Mutat.* 2012;33(7):1051-1055. ed. & tran. Accessed
32
33 May 31, 2021.
34
35
36
37 31. Aneichyk T, Hendriks WT, Yadav R, et al. Dissecting the Causal Mechanism of X-
38
39 Linked Dystonia-Parkinsonism by Integrating Genome and Transcriptome
40
41 Assembly. *Cell.* 2018;172(5):897-909.e21. ed. & tran. Accessed May 31, 2021.
42
43
44
45 32. Bragg DC, Mangkalaphiban K, Vaine CA, et al. Disease onset in X-linked dystonia-
46
47 parkinsonism correlates with expansion of a hexameric repeat within an SVA
48
49 retrotransposon in TAF1. *Proc Natl Acad Sci U S A.* 2017;114(51):E11020-
50
51 E11028. ed. & tran. Accessed May 31, 2021.
52
53
54
55 33. Hancks DC, Kazazian HH. Roles for retrotransposon insertions in human disease.
56
57 *Mob DNA.* 2016;7(1).
58
59
60 34. Kazazian HH, Moran J V. Mobile DNA in Health and Disease. *N Engl J Med.*

- 1
2
3 2017;377(4):361-370. ed. & tran. Accessed January 26, 2021.
4
5
6 35. Carvalho CMB, Lupski JR. Mechanisms underlying structural variant formation in
7
8 genomic disorders. *Nat Rev Genet.* 2016;17(4):224-238.
9
10
11 36. Kato I, Takagi Y, Ando Y, et al. A complex genomic abnormality found in a patient
12
13 with antithrombin deficiency and autoimmune disease-like symptoms. *Int J*
14
15 *Hematol.* 2014;100(2):200-205. ed. & tran. Accessed May 31, 2021.
16
17
18 37. Picard V, Chen J-M, Tardy B, et al. Detection and characterisation of large
19
20 SERPINC1 deletions in type I inherited antithrombin deficiency. *Hum Genet.*
21
22 2010;127(1):45-53. ed. & tran. Accessed August 2, 2018.
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Figure Titles and Legends

Figure 1. Long-read sequencing workflow and results. A) Overview of the general stages of the SVs discovery workflow. Algorithms used are depicted in yellow boxes. B) Nanopore sequencing results. i) Sequence length template distribution. Average read length was 4,499 bp (sd \pm 4,268); the maximum read length observed was 2.5Mb. ii) Genome median coverage per participant. The average across all samples was 16x (sd \pm 7.7). C) Filtering approach and number of SVs obtained per step. *SERPINC1* + promoter region corresponds to [GRCh38/hg38] Chr1:173,903,500-173,931,500. D) Anti-FXa percentage levels for the participants with a variant identified (P1-P10), cases without a candidate variant (P11-P19) and 300 controls from our internal database. The statistical significance is denoted by asterisks (*), where *** $P < 0.001$, **** $P \leq 0.0001$. p-values calculated by one-way ANOVA with Tukey's post-hoc test for repeated measures. ATD=Antithrombin Deficiency; ONT=Oxford Nanopore Technologies; SV=Structural Variant.

Figure 2. Candidate SVs identified by long-read sequencing. A) Schematic of chromosome 1 followed by protein coding genes falling in the zoomed region (1q25.1). SVs for each participant (P) are colored in red (deletions) and blue (duplications). The insertion identified in P9 and P10 is shown with a black line. B) Schematic of *SERPINC1* gene (NM_000488) followed by repetitive elements (RE) in the region. SINEs and LINEs are colored in light and dark grey respectively. Asterisks are present where the corresponding breakpoint falls within a RE. C) Characteristics of the antisense-oriented **SINE-VNTR-Alu (SVA)** retroelement (respect to the canonical sequence) observed in P9. Lengths of the fragments are subject to errors from Nanopore sequencing. TSD=Target site duplication.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

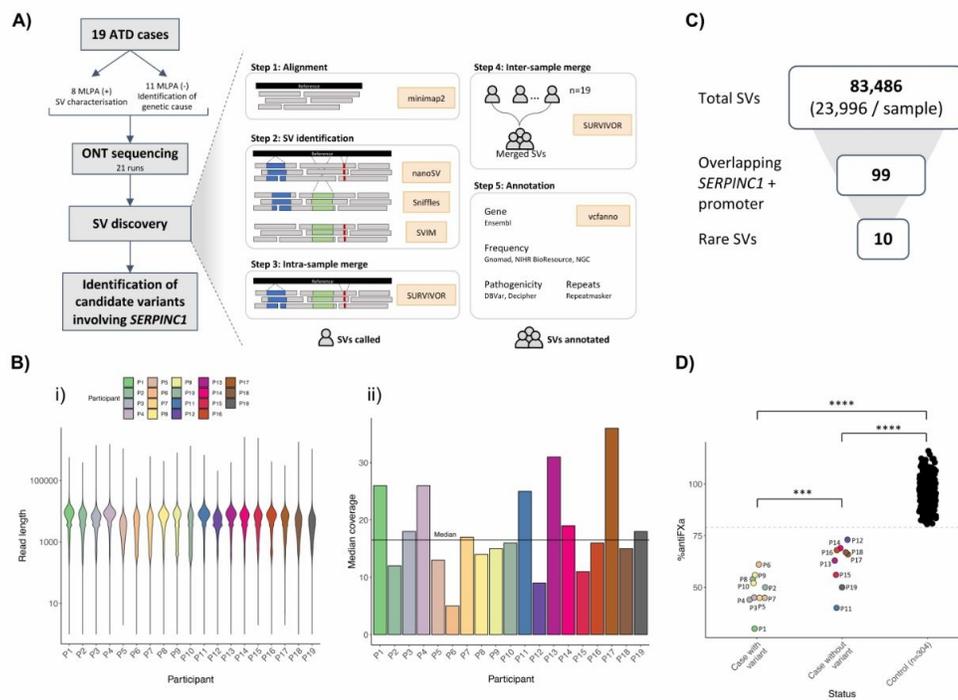
Figure 3. Resolution of a complex SV. Schematic representation of genetic diagnostic methods used to characterize the SVs in participant P2. Results from MLPA, LR-PCR and Nanopore are shown in white boxes. Primers used for both LR-PCR and Sanger validation experiments are shown representing the genetic location of each one with orange and green arrows respectively. *SERPINC1* gene in the IGV screenshot is represented in blue and exons are indicated. J1 and J2 correspond to the new formed junctions described in Figure S5. J=New junction; M1k=1 kb Molecular weight marker; M=100bp molecular weight marker; P=patient; C=control; B=Blank.

Figure 4. Schematic representation of genetic diagnostic methods used to characterize the SVs in participant P6. Results from MLPA, LR-PCR and Nanopore are shown in white boxes. Primers used for both LR-PCR and Sanger validation experiments are shown representing the genetic location of each one with orange and green arrows respectively. *SERPINC1* gene in the IGV screenshot is represented in blue and exons are indicated. J1 corresponds to the new formed junctions described in Figure S5. J=New junction; M=Molecular weight marker 1Kb or 100b; P=patient; C=control; B=Blank. For the LR-PCR results, C1 and P1 correspond to PCR 1 (done with Primer F + Primer R), and C2 and P2 correspond to PCR2 (done with Primer F + Primer R2).

1
2
3 **Table 1. Cohort of individuals included in this study, demographic, antithrombin**
4 **values and genetic results.** *SERPINC1* gene driven tests include MLPA, PGM
5 sequencing (Ion Torrent) and long range PCR (LR-PCR) amplification and Myseq
6 sequencing (Illumina). Genome wide tests are CGHa and whole genome sequencing
7 (WGS) using nanopore technology (ONT). Coordinates have been confirmed by Sanger
8 sequencing. Length refers to the extension of the structural variants. Het=
9 Heterozygous; Ag= Antigen; bp= Base pair.
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

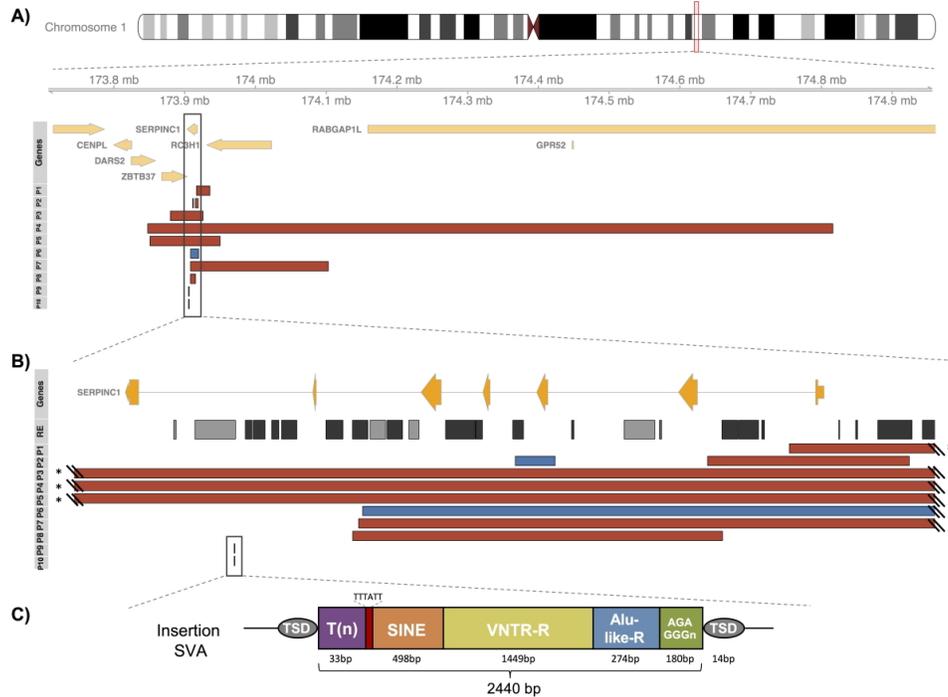
For Peer Review

Participant	Antithrombin		Family history	Gender	MLPA <i>SERPINC1</i>	PGM	CGHa	LR-PCR & Illumina sequencing	WGS ONT	Algorithm	Genotype	Coordinates	Length (bp)
	antiFXa%	Ag (%)											
P1	30	30	Yes	M	Deletion exon 1	-	Negative	Deletion exon 1	Deletion exon 1	nanosv;sniffles;svim	Het	1:173916704-173935703	18999
P2	54	41	Yes	M	Deletion exon 1	-	Negative	Deletion exons 1,2	CxSV(Deletion exon 1;Duplication exon 3)	nanosv;sniffles	Het;Het	1:173911379-173915115;1:173912151-173919034	3737;6884
P3	44	41	Yes	F	Complete deletion	-	Deletion 2 genes	-	Deletion 2 genes	nanosv;sniffles	Het	1:173879820-173925989	46169
P4	45	38	No	M	Complete deletion	-	Deletion 20 genes	-	Deletion 20 gene	nanosv;sniffles	Het	1:173847847-174816147	968005
P5	36	50	Yes	F	Complete deletion	-	-	-	Deletion 5 genes	nanosv	Het	1:173850996-173950174	99178
P6	61	46	Yes	M	Duplication exons 1,2 and 4; Deletion exon 6	-	Negative	Tandem duplication exons 1-5	Tandem Duplication exons 1-5	nanosv	Het	1:173908412-173919816	11404
P7	45	38	No	M	Deletion exons 1-5	-	Deletion exons 1-5 + 1 gene	-	Deletion 2 genes	nanosv;sniffles	Het	1:173908334-174103015	194389
P8	52	37	Yes	F	Deletion exons 2-5	-	-	-	Deletion exons 2-5	nanosv;sniffles	Het	1:173908218-173915405	7187
P9	56	61	Yes	F	Negative	Negative	-	Negative	Insertion SVA	nanosv	Het	1:173905922	2440
P10	50	46	Yes	F	Negative	Negative	-	Negative	Insertion SVA	visual inspection	Het	1:173905922	2440
P11	40	41	Yes	F	Negative	Negative	-	Negative	Negative				
P12	73	62	No	F	Negative	Negative	-	Negative	Negative				
P13	63	58	No	M	Negative	-	-	Negative	Negative				
P14	69	NA	No	F	Negative	Negative	-	Negative	Negative				
P15	56	45	Yes	F	Negative	-	-	-	Negative				
P16	68	54	No	M	Negative	Negative	-	Negative	Negative				
P17	66	67	No	M	Negative	Negative	-	Negative	Negative				
P18	67	70	No	F	Negative	-	-	-	Negative				
P19	50	70	Yes	M	Negative	-	-	-	Negative				



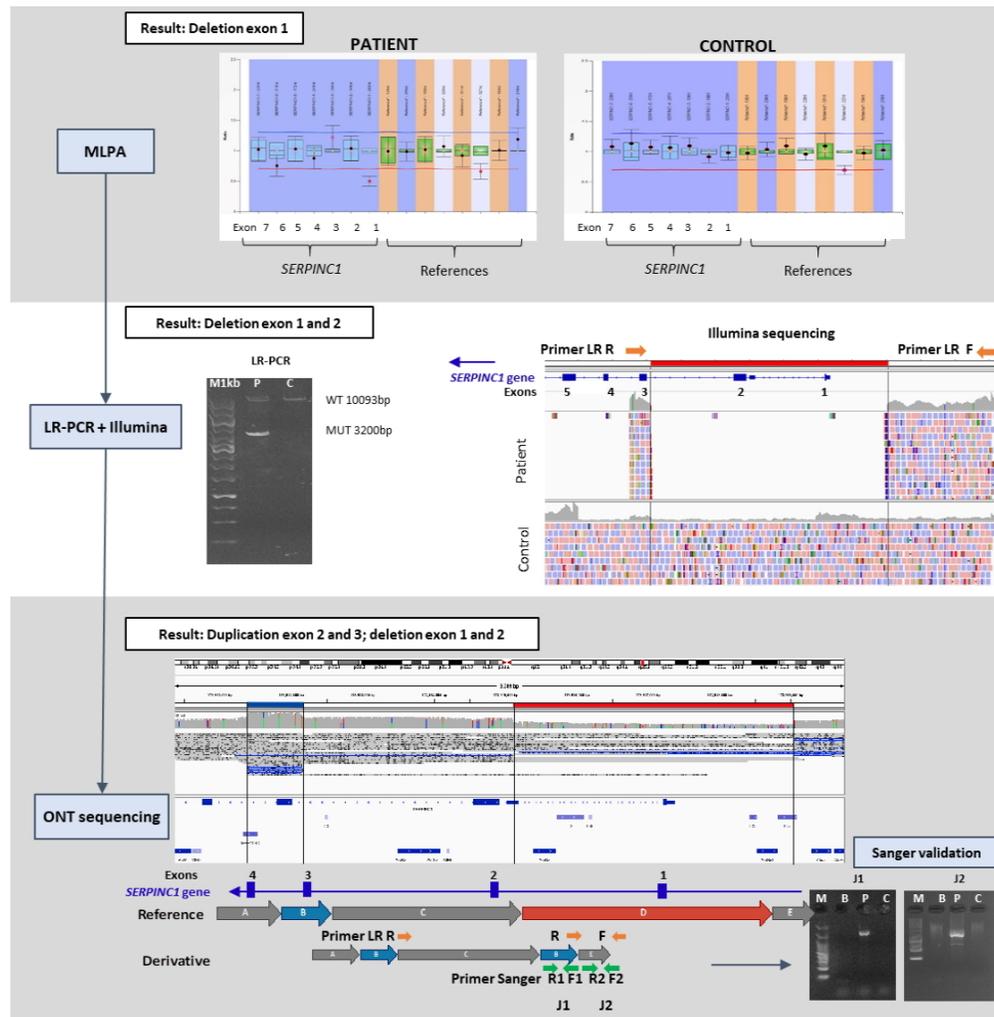
Long-read sequencing workflow and results. A) Overview of the general stages of the SVs discovery workflow. Algorithms used are depicted in yellow boxes. B) Nanopore sequencing results. i) Sequence length template distribution. Average read length was 4,499 bp (sd \pm 4,268); the maximum read length observed was 2.5Mb. ii) Genome median coverage per participant. The average across all samples was 16x (sd \pm 7.7). C) Filtering approach and number of SVs obtained per step. *SERPINC1* + promoter region corresponds to [GRCh38/hg38] Chr1:173,903,500-173,931,500. D) Anti-FXA percentage levels for the participants with a variant identified (P1-P10), cases without a candidate variant (P11-P19) and 300 controls from our internal database. The statistical significance is denoted by asterisks (*), where *** $P < 0.001$, **** $P \leq 0.0001$. p-values calculated by one-way ANOVA with Tukey's post-hoc test for repeated measures. ATD=Antithrombin Deficiency; ONT=Oxford Nanopore Technologies; SV=Structural Variant.

331x245mm (96 x 96 DPI)



Candidate SVs identified by long-read sequencing. A) Schematic of chromosome 1 followed by protein coding genes falling in the zoomed region (1q25.1). SVs for each participant (P) are colored in red (deletions) and blue (duplications). The insertion identified in P9 and P10 is shown with a black line. B) Schematic of SERPINC1 gene (NM_000488) followed by repetitive elements (RE) in the region. SINEs and LINEs are colored in light and dark grey respectively. Asterisks are present where the corresponding breakpoint falls within a RE. C) Characteristics of the antisense-oriented SINE-VNTR-Alu (SVA) retroelement (respect to the canonical sequence) observed in P9. Lengths of the fragments are subject to errors from Nanopore sequencing. TSD=Target site duplication.

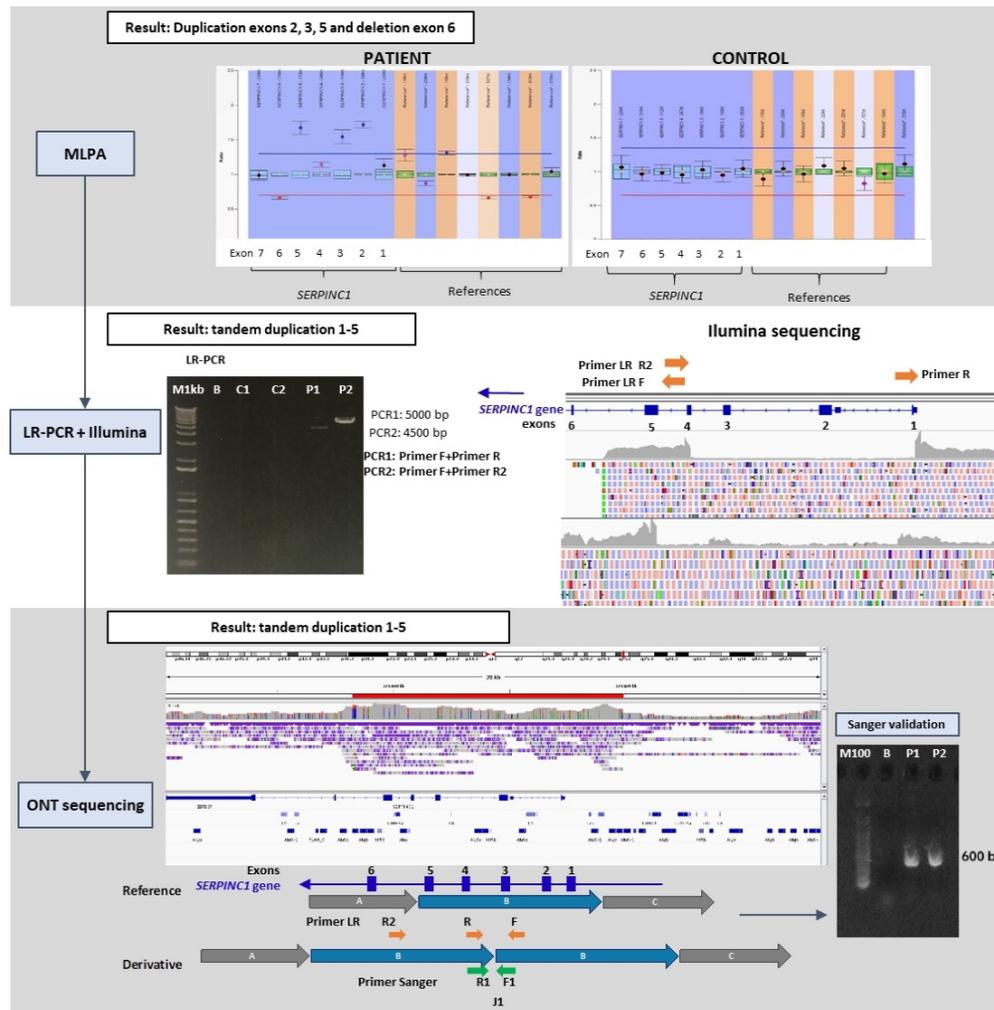
331x243mm (96 x 96 DPI)



39 Resolution of a complex SV. Schematic representation of genetic diagnostic methods used to characterize the SVs in participant P2. Results from MLPA, LR-PCR and Nanopore are shown in white boxes. Primers used for both LR-PCR and Sanger validation experiments are shown representing the genetic location of each one with orange and green arrows respectively. *SERPINC1* gene in the IGV screenshot is represented in blue and exons are indicated. J1 and J2 correspond to the new formed junctions described in Figure S5. J=New junction; M1k=1 kb Molecular weight marker; M=100bp molecular weight marker; P=patient; C=control; B=Blank.

46 295x301mm (96 x 96 DPI)

47
48
49
50
51
52
53
54
55
56
57
58
59
60



Schematic representation of genetic diagnostic methods used to characterize the SVs in participant P6. Results from MLPA, LR-PCR and Nanopore are shown in white boxes. Primers used for both LR-PCR and Sanger validation experiments are shown representing the genetic location of each one with orange and green arrows respectively. SERPINC1 gene in the IGV screenshot is represented in blue and exons are indicated. J1 corresponds to the new formed junctions described in Figure S5. J=New junction; M=Molecular weight marker 1Kb or 100b; P=patient; C=control; B=Blank. For the LR-PCR results, C1 and P1 correspond to PCR 1 (done with Primer F + Primer R), and C2 and P2 correspond to PCR2 (done with Primer F + Primer R2).

292x297mm (96 x 96 DPI)

SUPPLEMENTAL MATERIAL

Supplemental methods

Genetic diagnostic methods

Genetic diagnostic methods used to evaluate *SERPINC1* gene included: i) Sanger sequencing of exons and flanking regions, ii) Multiplex Ligation-dependent Probe Amplification (MLPA) covering the 7 exons of this gene (SALSA MLPA Kit P227 SerpinC1; MRC Holland, Amsterdam, The Netherlands), iii) whole gene sequencing by Ion Torrent technology (PGM; Thermo Fisher Scientific, Waltham, MA, USA), iv) long-range PCR (LR-PCR) amplification of the whole gene followed by Next Generation Sequencing (NGS) MiSeq platform (Illumina, San Diego, CA, USA) and / or v) Comparative Genomic Hybridization array (CGHa; CytoScan® HD Array; Thermo Fisher Scientific). These methods were performed as previously described.^{1,2}

SV identification

Reads were aligned against the GRCh38/hg38 human reference genome using minimap2 (2.17-r941)³ with default parameters for nanopore data ('-ax map-ont' parameter). SV discovery was done using a combination of three different algorithms:

- Sniffles v1.0.11⁴ was executed with a supporting read evidence of 4 ('-s 4' parameter) due to coverage variability.
- NanoSV v1.2.4⁵ was executed with default parameters. It was run on each independent chromosome in parallel to optimize compute time, with the limitation that inter-chromosomal variants were not detected.
- SVIM v1.2.0⁶ was executed with default parameters. Resulting SV calls with a quality score of less than 10 were filtered out of the dataset and not used for analysis. SV calls were merged at two different levels: intra-sample merge, to merge calls within individuals that had been identified by all three algorithms, and inter-sample merge, to merge SV calls across individuals.

Intra-sample merge. For each of the 19 samples, SV calls from all three different algorithms were merged using SURVIVOR v.1.0.7.⁷ VCF files were concatenated into one using bcftools⁸ and SURVIVOR was run using the command 'SURVIVOR merge in.fofn 500 1 -1 -1 -1 -1 out.vcf', requiring a maximum distance of 500bp between breakpoints. Additionally, intra-sample merge was done independently of the SV type, since different SV detection algorithms determine the type in different ways. For example, Sniffles determines canonical (DEL, DUP, INS, INV, TRA) and some complex SV types, while NanoSV calls only breakends (BND). The following options were turned-off: take the strands of SVs into account, estimate distance based on the size of SV, minimum size of SVs to be taken into account. After running SURVIVOR,

1
2
3 an *in-house* script was used to select the most common SV type. If there was no common type,
4 the order of selection was NanoSV (if the SV type was not a BND) > Sniffles > SVIM type.
5

6 Inter-sample merge. Then, all the 19 samples were merged using SURVIVOR, taking the SV type
7 into account. The command run was: 'SURVIVOR merge in.fofn 500 1 1 -1 -1 -1
8 out.vcf'.
9
10

11 **Identification of candidate SVs in SERPINC1**

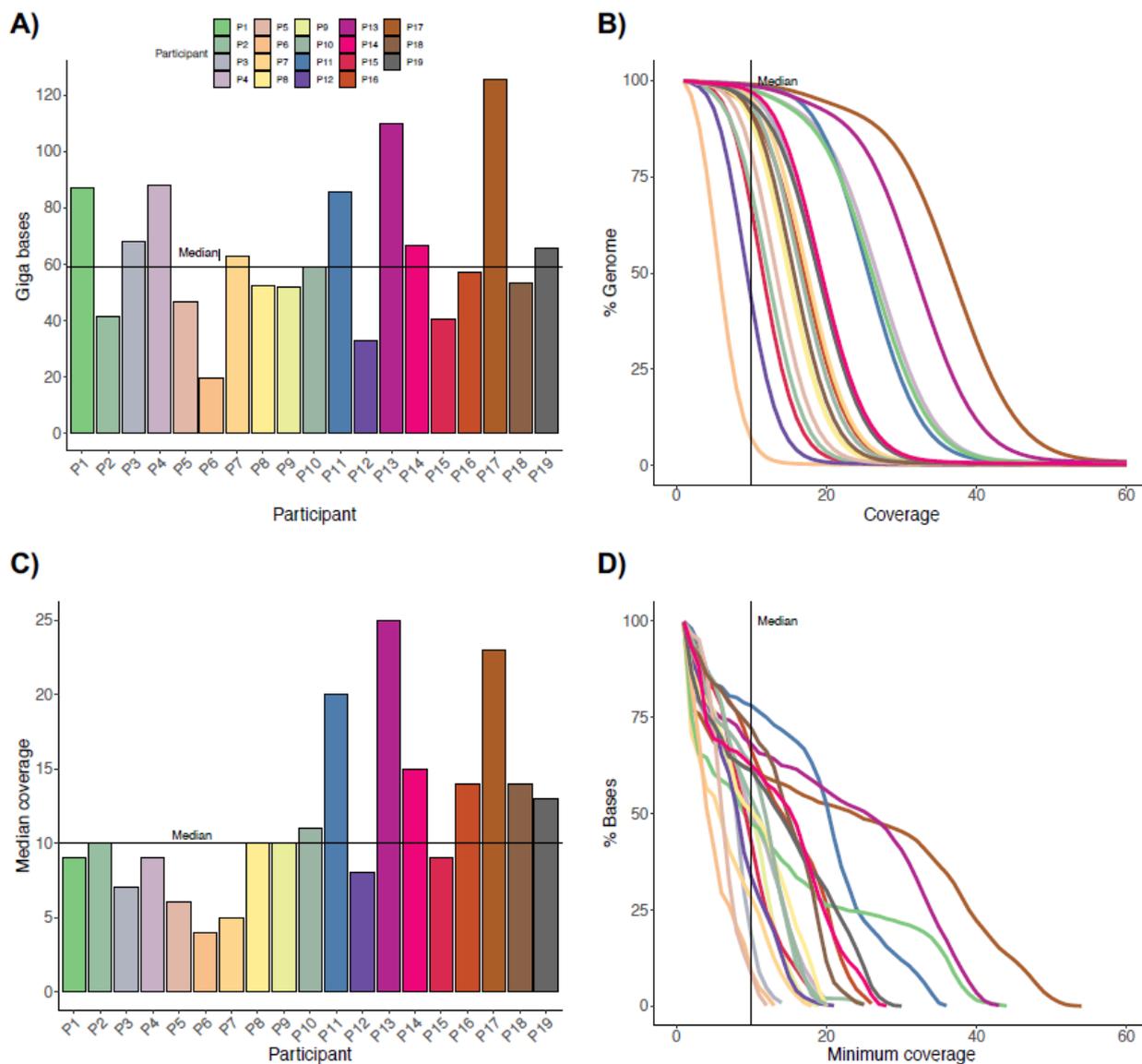
12
13
14 After filtering out variants that overlapped bad mapping quality regions (obtained from ⁹), a total
15 number of 83,486 SVs were identified across all participants, with a median of 23,996 (sd ±
16 3,431) SVs per sample. In order to identify disease-causing SVs associated with antithrombin
17 deficiency, we filtered for SVs overlapping the region [GRCh38/hg38] chr1:173,903,500-
18 173,931,500, which includes *SERPINC1* gene and its promoter region. A total number of 99 SVs
19 were observed, of which 10 were absent in gnomAD, NGC and the NIHR BioResource (5-7). These
20 10 SVs were observed in 9 samples: 6 were deletions, 1 was a tandem duplication, 1 was a
21 complex SV formed by 1 duplication and 1 deletion, and 1 was a SINE-VNTR-Alu (SVA)-type
22 retrotransposon insertion.
23
24
25
26
27

28 **Manual inspection of SVs**

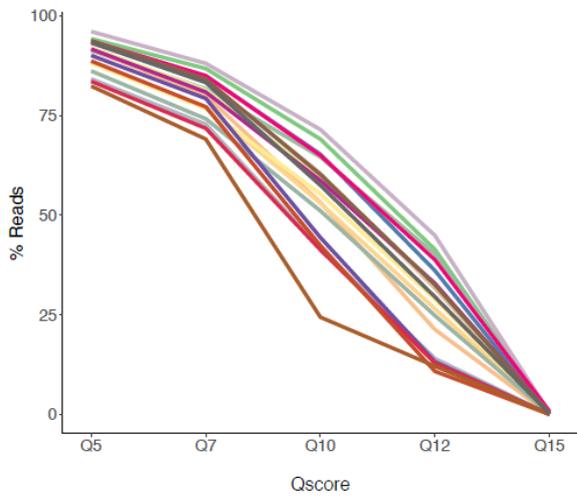
29
30 The alignments for all the cases without a candidate variant were manually inspected at the
31 above locus of interest using IGV. ¹⁰ An additional SVA insertion was observed in P10, at the exact
32 same position than P9. There were only two reads supporting the alternate, hence explaining
33 why it had not been called by any variant caller. Running Sniffles with a read evidence of 2 ('-s 2'
34 parameter) on the P10 data resulted in the SVA insertion being called.
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Supplemental figures

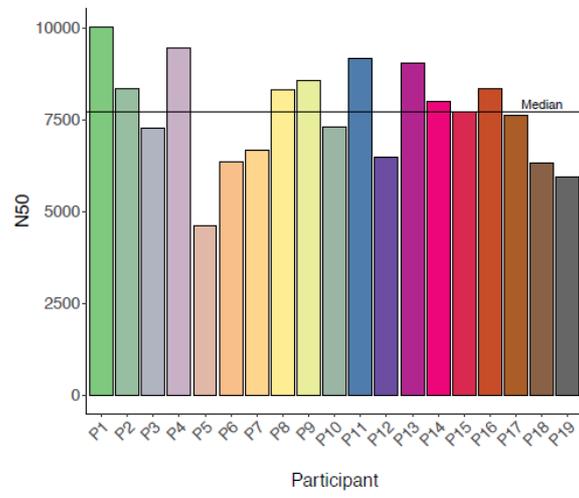
Figure S1. Sequencing results colored by participant. (A) Giga bases sequenced **(B)** Percentage of bases of the genome covered at a specific minimum coverage **(C)** Median coverage in *SERPINC1* + promoter region **(D)** Coverage distribution in *SERPINC1* + promoter region **(E)** Percentage of reads with a minimum Q score **(F)** Read N50, which refers to a value where half of the data is contained within reads with alignable lengths greater than this.



E)



F)



Peer Review

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Figure S2. Structural variants metrics. (A) Number of SVs identified by type and participant **(B)** Number of SVs by SV size **(C)** Fraction of SVs per allele count in our internal cohort of 62 individuals with long-read sequencing data **(D)** Number of SVs by median coverage and participant.

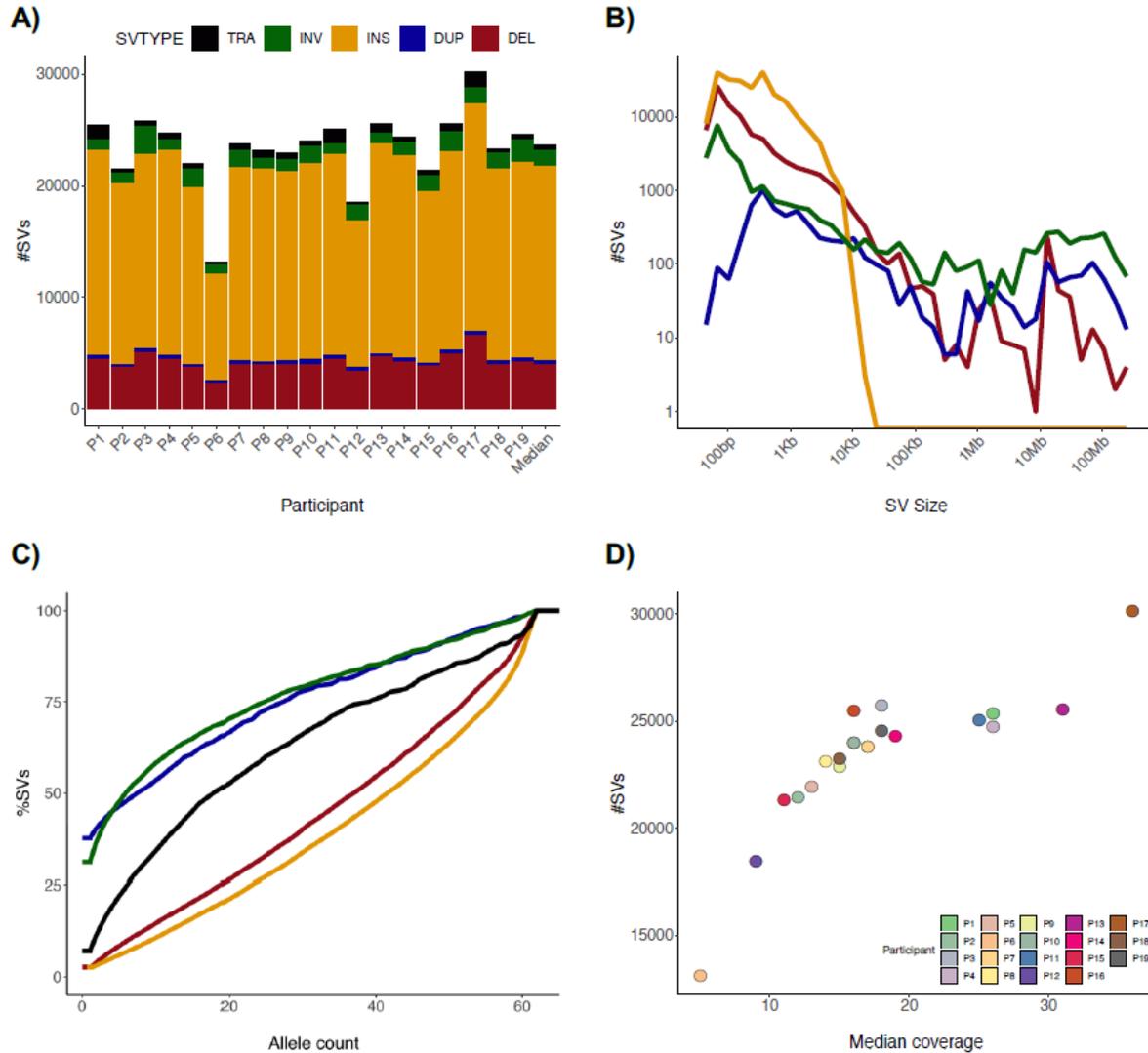
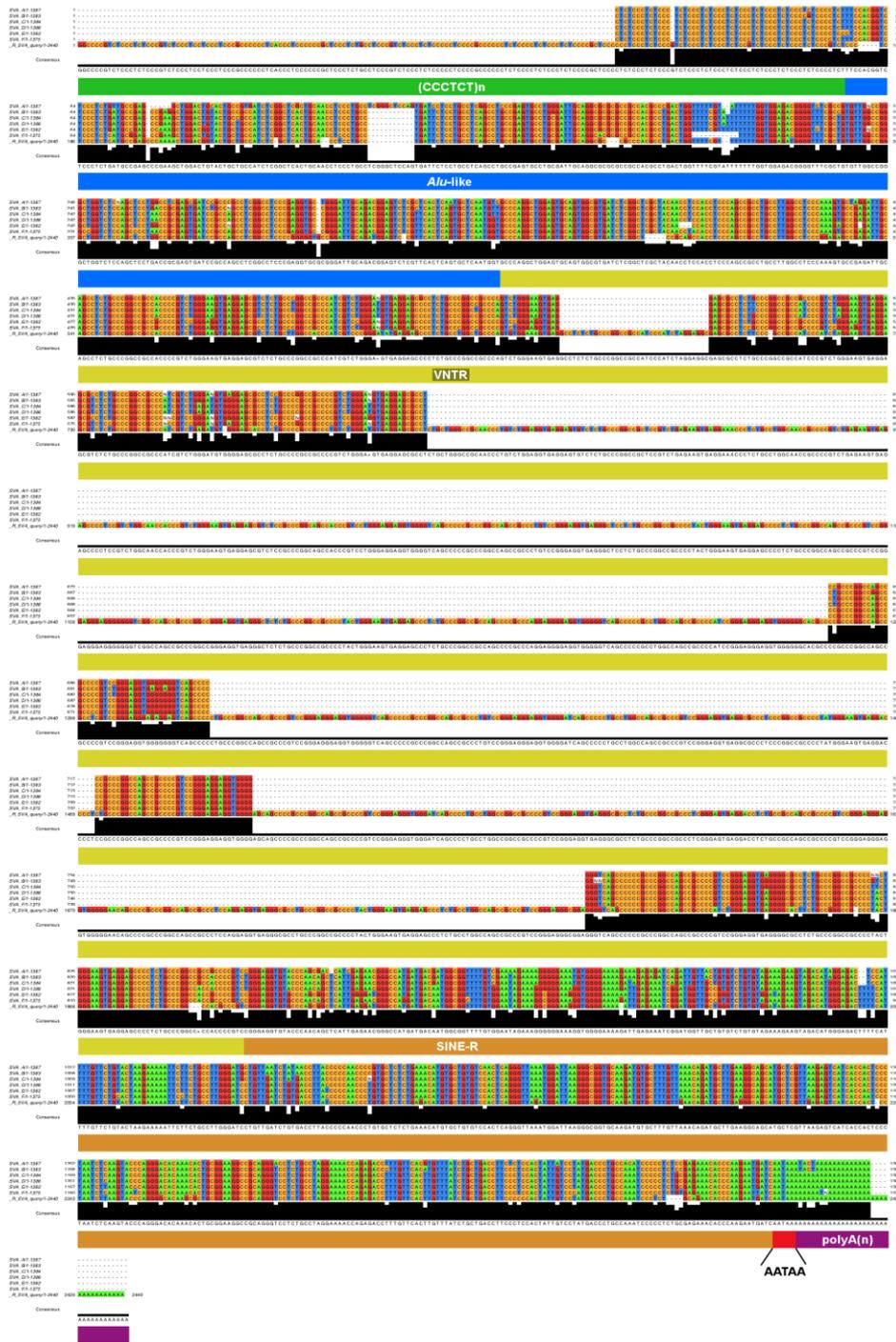


Figure S3. SVA sequence alignments. (A) The consensus sequences of SVA-A, -B, -C, -D, -E, and -F were taken from RepeatMasker (<http://www.repeatmasker.org>), then aligned using MAFFT¹¹ with default parameters. SVA_query corresponds to the SVA insertion in P9. Alignments were visually inspected and coloured by nucleotide using JalView.¹² Sub-elements of the SVAs are indicated underneath the consensus sequence matching colours in Figure 2C. (B) Phylogenetic tree was constructed with the Neighbour-Joining (NJ) algorithm using the Jukes-Cantor substitution model and visualized with iTol¹³. The SVA insertion in P9 was observed to be closest to the SVA E in the phylogenetic tree.

(A)



1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

(B)

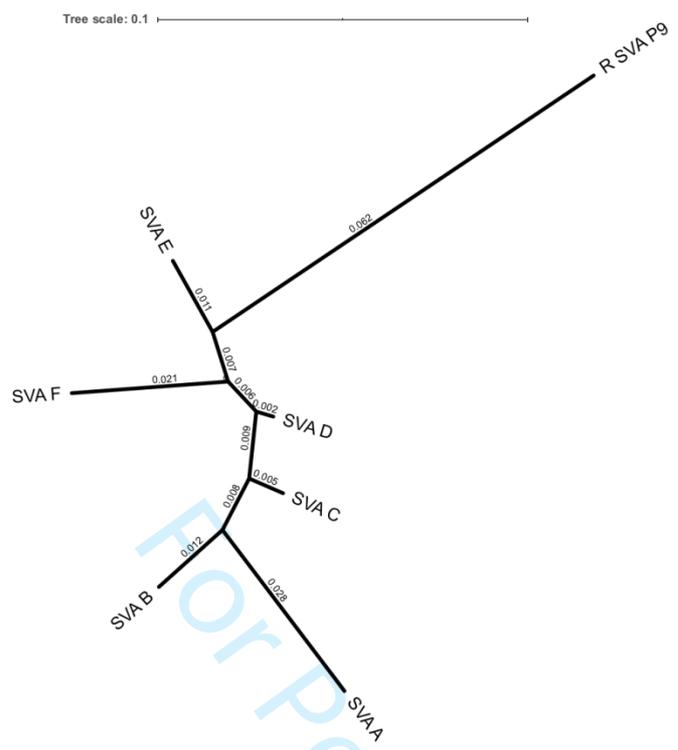


Figure S4. PCR validation of the SVA insertion in P9 and P10. (A) Schematic of *SERPINC1* gene (NM_000488) with zoom to intron 6 showing the SVA structure. Primers used in the long- and short-range PCRs are shown in orange and green respectively. Primer 8* was specifically designed within the inserted SVA sequence. Briefly, four reads of the retrotransposon present in the nanopore data for P9 and P10 were aligned to identify regions without any mismatch in order to select a 20 nucleotide sequence to be used as primer. That sequence was also checked to be present in *de novo* assembly alignment. **(B)** Primer combinations for PCR amplifications and expected sizes for wild type and mutated alleles are shown in the table. PCRs 1-4 were tested under different experimental conditions, and although in all cases the wild type allele was always amplified, no amplification of the mutated allele containing the SVA was obtained in P9 or P10. **(C)** The amplification of PCR 4 in agarose gel is shown. Only the 800 base pairs (bp) of the wild type allele was amplified in P9, P10 but also a healthy control. Only PCR 5, using the primer specific of the SVA rendered positive results and a specific 550bp band was obtained in P9, P10 and two relatives. **(D)** Family pedigrees of P9 and P10, including clinical information, the diagnosis of antithrombin deficiency (semi-filled symbols) and the anti-FX activity (as % of a reference plasma). B=Blank; M=Molecular Weight Marker; DVT=Deep vein thrombosis; PE=Pulmonary embolism.

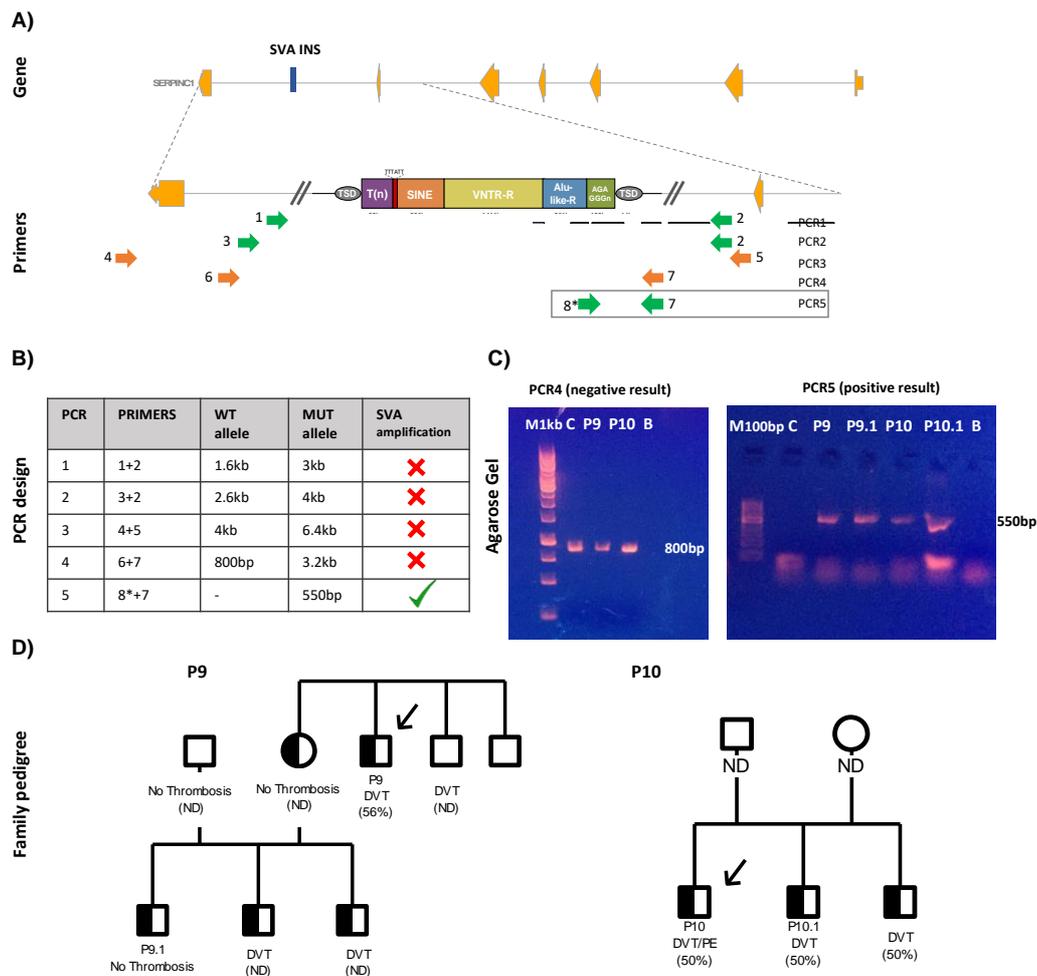


Figure S5. Nucleotide level characterization of the candidate SVs. Breakpoint junction sequence is aligned to the proximal and distal genomic reference sequence or sequence of breakpoint junction in hypothetical intermediate, as shown. Alignment is only shown for novel breakpoint junctions in the derivative chromosome. Microhomology at the breakpoint is indicated in red. Sequence in blue indicates inserted sequences at the breakpoint junction. Underline indicates repetitive elements in the reference, specified in *Italic*. J=Junction. **(A-H)** P1-8, **(I)** P9 and P10.

(A) P1

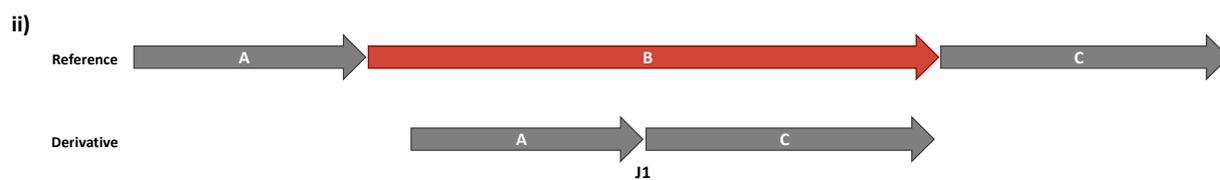
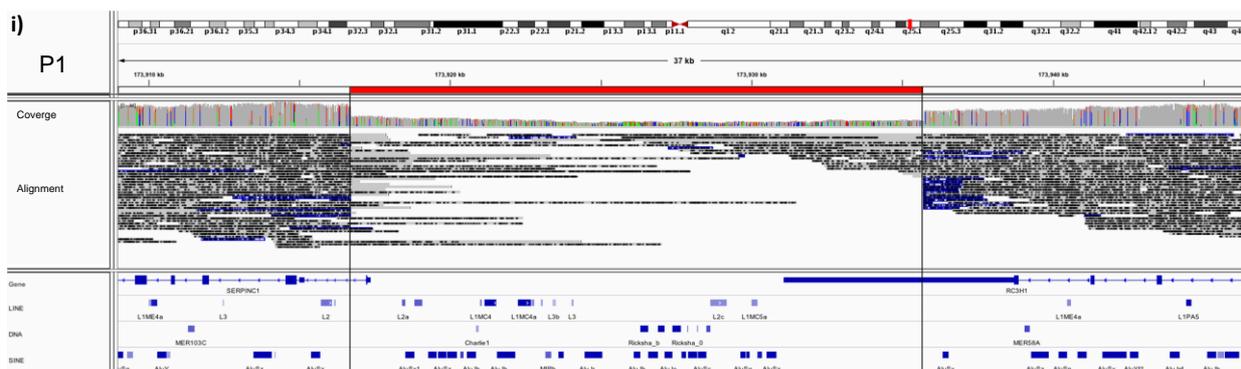


Figure S5(iii) shows a nucleotide alignment:

```

1-173916654  CCAAGACAGTTGTGCCACCACCCCGTTTCTAGCCTCCTGTGACAACTCA*GATGAAGAAGGGAGTGTGTGTGGCATTGAGGGAGGACAACCTCTCATTGTG 1739167 54
1-173916654  |||
J1-1         CCAAGACAGTTGTGCCACCACCCCGTTTCTAGCCTCCTGTGACAACTCA* GGAGCAAGGAACAAGGAAAAAAAAAGCAAAAAACAAAAACAATTTGGGGG 100
1-173916654  |||
1-173935659  AAATATGCCATTTGTCCTTACAGGCTCATGATATCACAGAGAAAAGGAA* -----AAGGAAAAAAAAAGCAAAAAACAAAAACAATTTGGGGG 173935759
    
```

A-rich

EMM

(B) P2

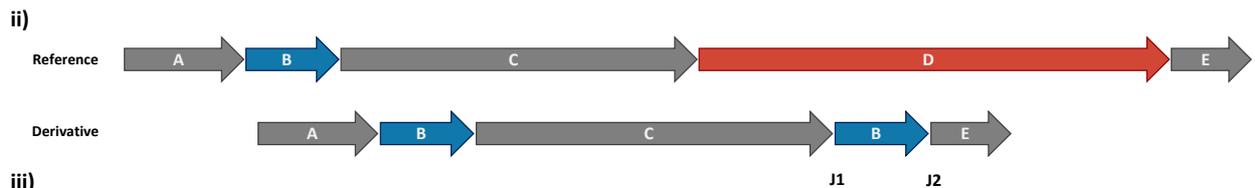
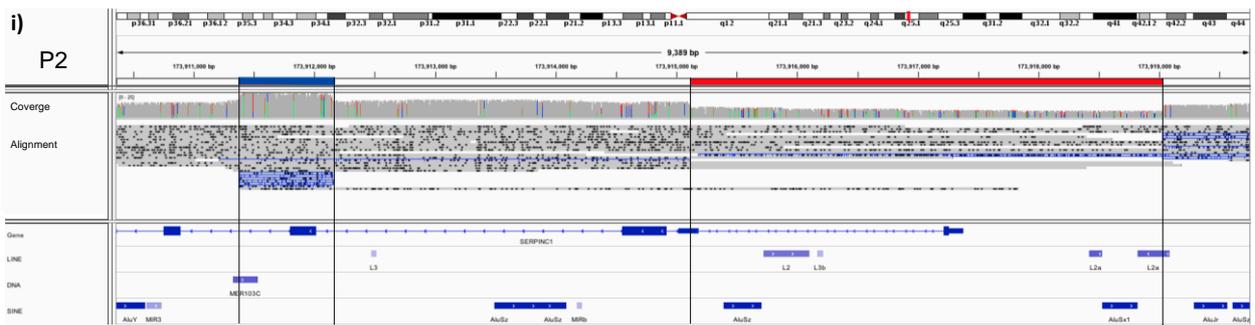


Figure B(iii) shows sequence alignments for reference and derivative sequences. The reference sequence is 1-173915066 to 173915165. The derivative sequence J1-1 is 1-173911329 to 173911428. The reference sequence 1-173912101 to 173912201 and derivative J2-1 (1-173918984 to 173919084) are also shown. The alignments highlight the MER103C and L2a regions. The derivative sequences show a deletion of segment D and a duplication of segment B.

(C) P3

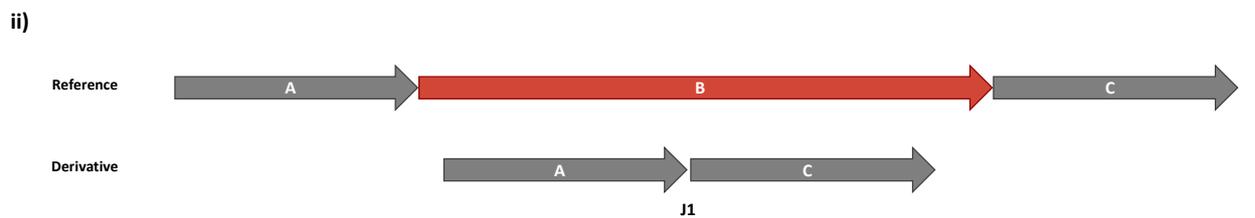
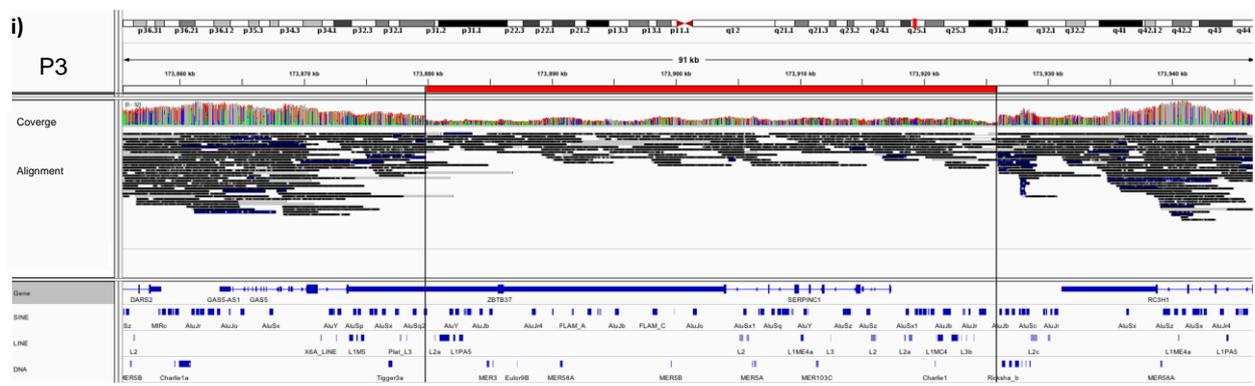
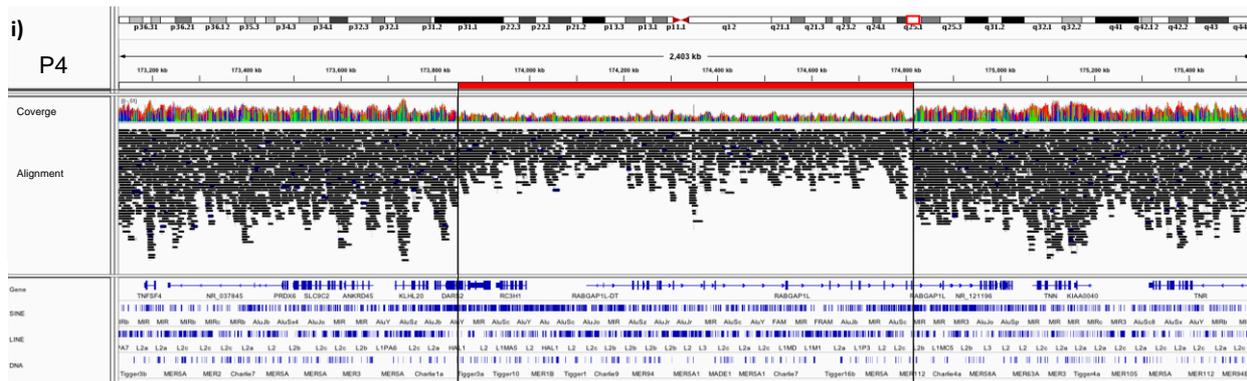
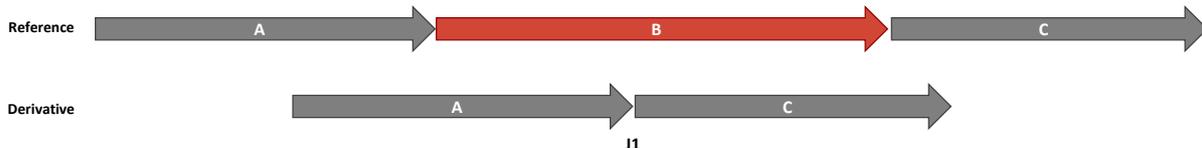


Figure C(iii) shows sequence alignments for reference and derivative sequences. The reference sequence is 1-173879770 to 173879870. The derivative sequence J1-1 is 1-173925939 to 173926039. The alignments highlight the AluSx1 regions. The derivative sequence shows a deletion of segment B.

(D) P4



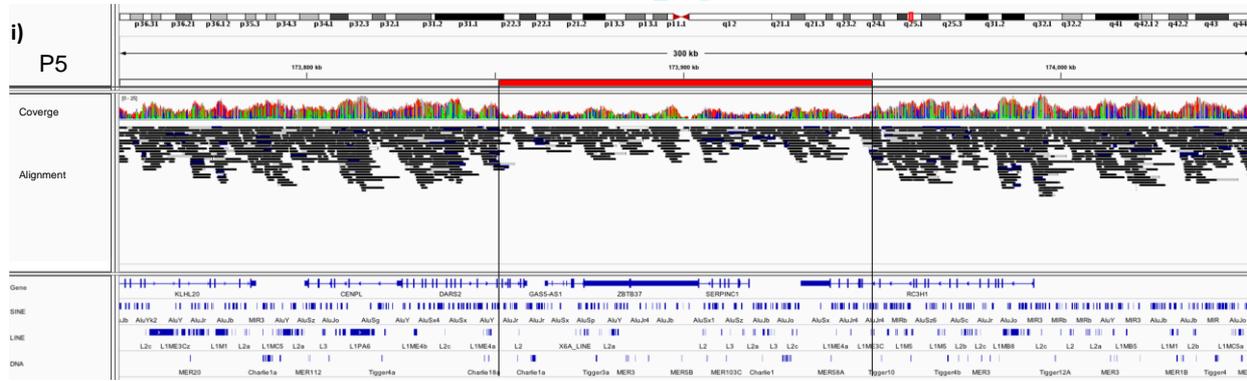
ii)



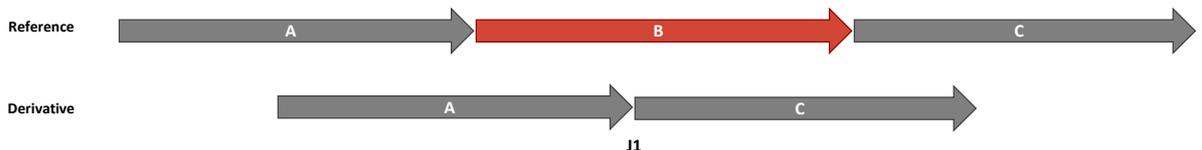
iii)

	<i>HAL1</i>	<i>AluY</i>	
1-173847797	CCTCACCTAGTTACTCTGTACCTTCTGCACAGCTGCTTTCTGAACCTT* TTTTTTTTTTTTTTTTTTTTGGAGCGGGAGTCTCACTCTGTGCCCCAG	CGGGAGTCTCACTCTGTGCCCCAG	173847897
J1-1	CCTCACCTAGTTACTCTGTACCTTCTGCACAGCTGCTTTCTGAAC ---* TTTTTTTTTTTTTTTTTTTTGGAGTAGAGTTTGCTCTGTTGCCCCAG	TAGAGTTTGCTCTGTTGCCCCAG	97
1-174816097	TTGTATGAACATATGAGCATTTTGTTTTCTTAGTAACATAATACATAG* TTTTTTTTTTTTTTTTTTTTGGAGTAGAGTTTGCTCTGTTGCCCCAG	TAGAGTTTGCTCTGTTGCCCCAG	174816197
		<i>AluSp</i>	

(E) P5



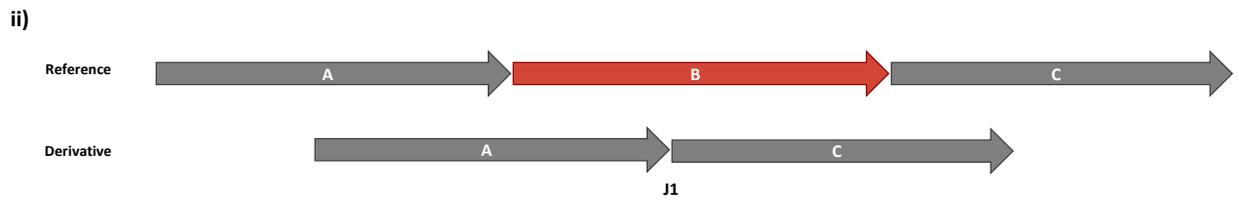
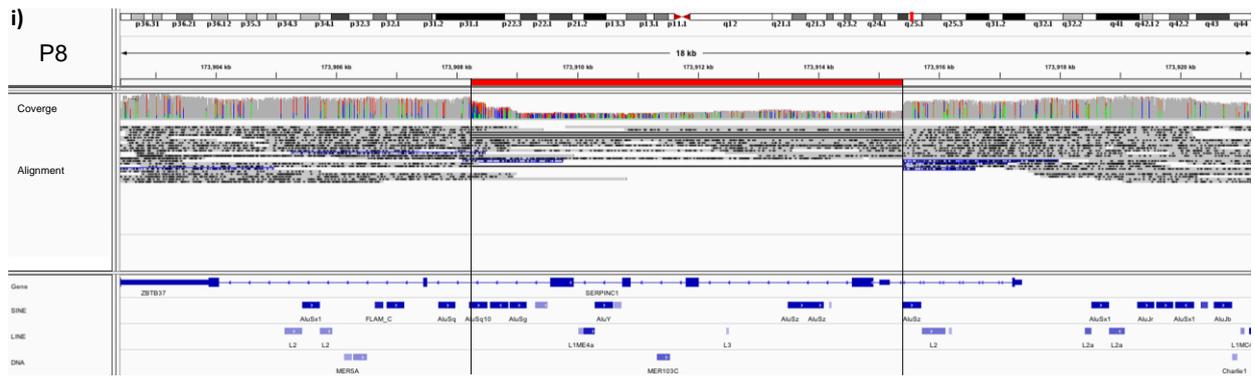
ii)



iii)

	<i>AluSx1</i>	<i>AluSx1</i>	
1-173850946	TTAAGGTTGGCCGGGCATGGTGGCTAACGCCTGTAATCCAGCAGCTTT*GGGAGCCGAGGTTGGTGGATCGCCTGAGGTGAGGTTGAGACCAGCCT	GGGAGCCGAGGTTGGTGGATCGCCTGAGGTGAGGTTGAGACCAGCCT	173851046
J1-1	TTAAGGTTGGCCGGGCATGGTGGCTAACGCCTGTAATCCAGCAGCTTT* -AAGTACACATCTTAAAGCAATGCTAAGTGAAAAATGAAAATGCAAGCT	TAAAGCAATGCTAAGTGAAAAATGAAAATGCAAGCT	99
1-173950124	TGCCACTGCACCTCCAGCTGGGTACAGAGTGGAACTCCGCTCAAAAA*AAA GTACACATCTTAAAGCAATGCTAAGTGAAAAATGAAAATGCAAGCT	TAAAGCAATGCTAAGTGAAAAATGAAAATGCAAGCT	173950224
	<i>AluSx</i>	<i>L1ME3C</i>	

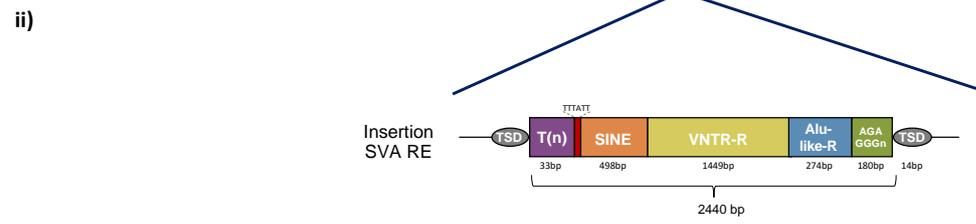
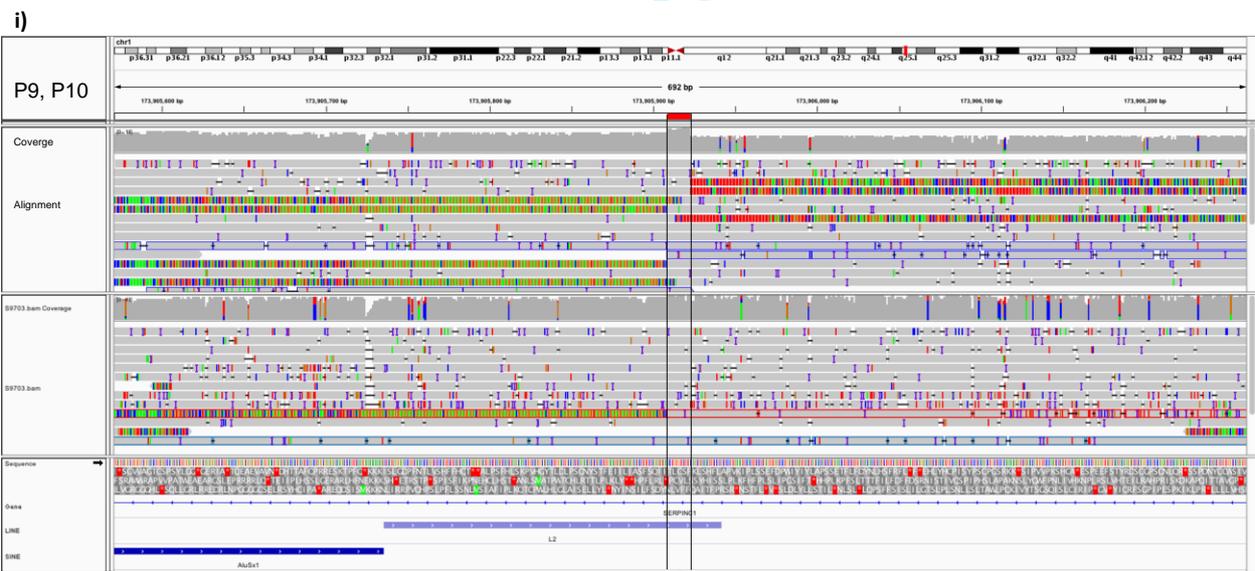
(H) P8



iii)

		<i>AluSx1</i>	
1-173908168	GGCAACATTTTAGTTTTCAAACCTAGCTTAAAGACAGCCCAAT AGCAC*AGTGGCTCAGCCTGTAATCCAGCACTTTGGGAGGC CAAGGCAGGTAGAT		173908268
J1-1	GGCAACATTTTAGTTTTCAA -CTAGCTTAAAGACAGCCCAATAGCAC* AGTGGCTCAGCCTGTAATCCAGCACTTTGGGAGGC TGAGGCATGCAGAT		99
1-173915355	AGAAGTAAGAAGAAATGGATGGAGCAAAGCAATATC TGGCCGGGTGCG* -GTGGCTCAGCCTGTAATCCAGCACTTTGGGAGGC TGAGGCATGCAGAT	<i>AluSp</i>	173915454

(I) P9 and P10



iii)

Participant 9

- P9 39-52: TSD
- P9 53-84: (T)n
- P9 85-90: TTTATT sequence
- P9 91-589: SINE
- P9 590-2038: VNTR-R
- P9 2039-2312: Alu-like-R
- P9 2313-2492: (AGAGGG)n
- P9 2493-2506: TSD

1-173905870	TTGAAATTATATTAATAGCATCCTTTTCTCAGATTATAACCTTGTGTTCTTT-----	173905926
P9-1	TTTATTGATCATTCT	100
1-173905926	-----	173905926
P9-101	TGGGTGTTCTGCAGAGGGATTGGCAGGGTCATAGGACAATAGTGGGGGAAGGTCAGCAGATAAACAAGTGAACAAAGTCTCTGGTTTTCTAGGCAG	200
1-173905926	-----	173905926
P9-201	AGGACCCCTGCGGCCTCCGCAGCGCTTGTGCCCTGGGTACTTAAGATTAGGAGTGGTGATGACTCTCAACGAGCATGCTGCCCTCAAGCATCTGTTCAA	300
1-173905926	-----	173905926
P9-301	CAAAGCACATCTTGACCCGCCCTTAATCCATCTAACCCCTGAGTGGACACAGCACATGTCTCAGAGAGCACAGGGTTGGGGATAAGGTCACAGATCAACAG	400
1-173905926	-----	173905926
P9-401	GATCCCAAGGCAGAAGAATTTTCTTAGTACAGAAACAAATGAAAAGTCTCCCATGTCTACTTCTATCCACACAGACCCAGCAACCATCCGATTTCTCAA	500
1-173905926	-----	173905926
P9-501	TTTTTCCCACCCCTTCCCGCCTTTCTATTCCACAAAACCGCCATTGTTCATCATGCCCCATCTCAATGAGCCGCTGGGCACACCTCCCAGCGGGCGTGGC	600
1-173905926	-----	173905926
P9-601	CGGGCAGAGGGGCTCCTCACTTCCCAGTAGGGCGGGCCGAGAAAGTGCCCTCACCTCCCAGATGGGGCGGCTGGCCGGCGGGGGCTGACCCCTCCGC	700
1-173905926	-----	173905926
P9-701	CCTCCCGGACGGGGCGGCTGGCCAGGCAGAGGGCTCCTCACTTCCCAGTAGGGCGGGCCGGCAGGCGCCCTCACCTCCTGGAGGGCGGCTGGCCGGGGCGG	800
1-173905926	-----	173905926
P9-801	GGCTGTTCCCCACCTCCCTCCCGACGGGGCGGCTGGCGGCAGAGGTCTCACTCCCGAGGCGGGCCGAGAGGCGCCCTCACCTCCCGACGGGGC	900
1-173905926	-----	173905926
P9-901	GGCCGGCCAGGCAGGGGCTGATCCACCTCCCGGACGGGGCGGCTGGCCGGCGGGGCTGCTCCACCTCCTCCCGGACGGGGCGGCTGGCCGGGGCAG	1000
1-173905926	-----	173905926
P9-1001	AGGGGTCTCACTTCCCATAGGGGCGGGCGGGAGGGCGGCTCACCTCCCGGACGGGGCGGCTGGCCAGGCGGGGGCTGATCCCCACCTCCCTCCCGGACA	1100
1-173905926	-----	173905926
P9-1101	GGGCGGCTGGCCGGGCGGGGGCTGACCCCACTCCCTCCCGACGGGGCGGCTGGCCGGCAGGGGGCTGACTCCTCCTCCCTCCCGACGAGGCGGCTGG	1200
1-173905926	-----	173905926
P9-1201	CCGGGCGGGGCGTGCCTCCACCTCCCTCCCGATGGGGCGGCTGGCCAGGCGGGGGCTGACCCCACTCCCTCCTGGCGGGGCTGGCGGGCGGGCA	1300
1-173905926	-----	173905926
P9-1301	GAGGGCTCCTCACTTCCCAGTAGGGGCGGGCGGGCAGAGAGCCCTCACCTCCCGGCGGGCGGCTGGCCAGCCCCCTCCCTCCCGGACGGGGCGGCTGG	1400
1-173905926	-----	173905926
P9-1401	CCGGGCGAGGGGGCTCCTCACTTCCCAGTAGGGGCGGGCGGGCAGAGGAGCCCTCACCTCCCGGACGAGGCGGCTGGCCGGGCGGGGGCTGACCCCACT	1500
1-173905926	-----	173905926
P9-1501	CCTCCAGGACGGGTGGCTGCCGGGCGGAGACGCTCCTCACTTCCAGACGGGTGGTTGCCAGACGAGGGGCTCCTCACTTCTCAGACGGGGCGGTTGCC	1600
1-173905926	-----	173905926
P9-1601	AGGCAGAGGGTTTCTCACTTCTCAGACGGAGCGGGCGGGCAGAGACTCCTCACCTCCAGACAGGTTGCGGCCAGCAGAGGGGCTCCTCACATCCC	1700
1-173905926	-----	173905926
P9-1701	AGACAGGGCGGGCGGGCAGAGGTGCTCCACATCTCAGACGATGGGGCGGGCGGGCAGAGACGCTCCTCACTTCTTAGATGGGATGGCGGGGGAAGAGGC	1800

SUPPLEMENTAL BIBLIOGRAPHY

1. de la Morena-Barrio, M.E., Martínez-Martínez, I., de Cos, C., Wypasek, E., Roldán, V., Undas, A., van Scherpenzeel, M., Lefeber, D.J., Toderici, M., Sevivas, T., et al. (2016). Hypoglycosylation is a common finding in antithrombin deficiency in the absence of a SERPINC1 gene defect. *J. Thromb. Haemost.* *14*, 1549–1560.
2. De la Morena-Barrio, B., Borràs, N., Rodríguez-Alén, A., Morena-Barrio, M.E., García-Hernández, J.L., Padilla, J., Bravo-Pérez, C., Miñano, A., Rollón, N., Corral, J., et al. (2019). Identification of the first large intronic deletion responsible of type I antithrombin deficiency not detected by routine molecular diagnostic methods. *Br. J. Haematol.* *186*, e82–e86.
3. Li, H. (2018). Minimap2: Pairwise alignment for nucleotide sequences. *Bioinformatics* *34*, 3094–3100.
4. Sedlazeck, F.J., Rescheneder, P., Smolka, M., Fang, H., Nattestad, M., Von Haeseler, A., and Schatz, M.C. (2018). Accurate detection of complex structural variations using single-molecule sequencing. *Nat. Methods* *15*, 461–468.
5. Cretu Stancu, M., Van Roosmalen, M.J., Renkens, I., Nieboer, M.M., Middelkamp, S., De Ligt, J., Pregno, G., Giachino, D., Mandrile, G., Espejo Valle-Inclan, J., et al. (2017). Mapping and phasing of structural variation in patient genomes using nanopore sequencing. *Nat. Commun.* *8*, 1–13.
6. Heller, D., and Vingron, M. (2019). SVIM: Structural variant identification using mapped long reads. *Bioinformatics* *35*, 2907–2915.
7. Jeffares, D.C., Jolly, C., Hoti, M., Speed, D., Shaw, L., Rallis, C., Balloux, F., Dessimoz, C., Bähler, J., and Sedlazeck, F.J. (2017). Transient structural variations have strong effects on quantitative traits and reproductive isolation in fission yeast. *Nat. Commun.* *8*, 1–11.
8. Li, H., Handsaker, B., Wysoker, A., Fennell, T., Ruan, J., Homer, N., Marth, G., Abecasis, G., and Durbin, R. (2009). The Sequence Alignment/Map format and SAMtools. *Bioinformatics* *25*, 2078–2079.
9. Sanchis-Juan, A., Stephens, J., French, C.E., Gleadall, N., Mégy, K., Penkett, C., Shamardina, O., Stirrups, K., Delon, I., Dewhurst, E., et al. (2018). Complex structural variants in Mendelian disorders: identification and breakpoint resolution using short- and long-read genome sequencing. *Genome Med.* *10*, 95.
10. Robinson, J.T., Thorvaldsdóttir, H., Wenger, A.M., Zehir, A., and Mesirov, J.P. (2017). Variant review with the integrative genomics viewer. *Cancer Res.* *77*, e31–e34.
11. Katoh, K., Misawa, K., Kuma, K.I., and Miyata, T. (2002). MAFFT: A novel method for rapid multiple sequence alignment based on fast Fourier transform. *Nucleic Acids Res.* *30*, 3059–3066.
12. Waterhouse, A.M., Procter, J.B., Martin, D.M.A., Clamp, M., and Barton, G.J. (2009). Jalview Version 2-A multiple sequence alignment editor and analysis workbench. *Bioinformatics* *25*, 1189–1191.
13. Letunic, I., and Bork, P. (2019). Interactive Tree of Life (iTOL) v4: Recent updates and new developments. *Nucleic Acids Res.* *47*, W256–W259.

Participant	Bkp in SV	Segments	Window	Repeat coordinates	Repeat	Family	Start in query	End in query	Repeat length
P1	Deletion start	A3'-B5'	1:173916554-173916854	-					
	Deletion end	B3'-C5'	1:173935553-173935853	1:173935699-173935738	A-rich	Low_complexity	147	186	40
P2	Duplication start	A3'-B5'	1:173911229-173911529	1:173911329-173911536	MER103C	DNA/hAT-Charlie	101	300	201
	Duplication end	B3'-C5'	1:173912001-173912301	-					
	Deletion start	C3'-D5'	1:173914965-173915265	-					
	Deletion end	D3'-E5'	1:173918884-173919184	1:173918822-173919087	L2a	LINE/L2	1	204	204
P3	Deletion start	A3'-B5'	1:173879670-173879970	1:173879670-173879957	AluSx1	SINE/Alu	1	288	288
	Deletion end	B3'-C5'	1:173925839-173926139	1:173926103-173926323	AluJb	SINE/Alu	265	300	37
P4	Deletion start	A3'-B5'	1:173847697-173847997	1:173847804-173847844	HAL1	LINE/L1	108	148	41
				1:173847845-173848152	AluY	SINE/Alu	149	300	153
	Deletion end	B3'-C5'	1:174815997-174816297	1:174816148-174816445	AluSp	SINE/Alu	152	300	150
P5	Deletion start	A3'-B5'	1:173850846-173851146	1:173850956-173851264	AluSx1	SINE/Alu	111	300	191
				1:173850606-173850904	AluSx	SINE/Alu	1	59	59
	Deletion end	B3'-C5'	1:173950024-173950324	1:173949885-173950176	AluSx	SINE/Alu	1	153	153
				1:173950177-173950297	L1ME3C	LINE/L1	154	121	121
				1:173950298-173950593	AluJr	SINE/Alu	275	300	27
P6	Duplication start	A3'-B5'	1:173908262-173908562	1:173908214-173908512	AluSx1	SINE/Alu	1	251	251
				1:173908555-173908859	AluJb	SINE/Alu	294	300	8
	Duplication end	B3'-C5'	1:173919640-173919940	1:173919609-173919890	AluSz	SINE/Alu	1	251	251
				1:173919924-173920236	AluSx1	SINE/Alu	285	300	17
P7	Deletion start	A3'-B5'	1:173908184-173908484	1:173908214-173908512	AluSx1	SINE/Alu	31	300	271
	Deletion end	B3'-C5'	1:174102865-174103165	1:174102793-174102889	Tigger1	DNA/TcMar-Tigger	1	25	25
				1:174102890-174103191	AluSp	SINE/Alu	26	300	276
P8	Deletion start	A3'-B5'	1:173908068-173908368	1:173908214-173908512	AluSx1	SINE/Alu	147	300	155
	Deletion end	B3'-C5'	1:173915255-173915555	1:173915394-173915705	AluSz	SINE/Alu	140	300	162
P9	Insertion site		1:173905771-173906071	1:173905736-173905941	L2	LINE/L2	1	171	171
P10	Insertion site		1:173905771-173906071	1:173905736-173905941	L2	LINE/L2	1	171	171

Table S1. Repetitive elements at the SV breakpoints. Segments column match to those in Figure S5, and refers to the relative location in the segments where the SV breakpoint is, where 5' and 3' correspond to the upstream and downstream parts respectively. Window was done for +/- 150 bp from the breakpoint. Start and End in query are the relative positions to the query sequence where the repeat starts and ends respectively. Bkp=Breakpoint.

Participant	New junction	Microhomology at breakpoint	Microhomology at breakpoint length	Deletion at breakpoint	Deletion at breakpoint length	Insertion at breakpoint	Insertion at breakpoint length	Duplication at breakpoint	Duplication at breakpoint length	Expected reference sequence	Observed sequence (Sanger)
P1	J1	-	-	-	-	GGAGCAAGGAAC	12	-	-	CCAAGACAGTTGTGCCACCACCCCGTTTCTAGCCTCTGTGACAACCTCA* AAGGAAAAGGAAAAAGCAAAAAACAAAAACAA AAATTTTGGGGGTGTAG	CCAAGACAGTTGTGCCACCACCCCGTTTCTAGCCTCTGTGACAACCTCA*(GGAGCAAGGAAC) AAGGAAAAGGAAAAACAAAAACAA ATTTTGGGGGTGTAG
P2	J1	GTTGC	5	GTTGC	5	-	-	-	-	CCCGGGACAGGTTCAAGTCTAGACTTCTGCCAGGGACAGTTCAG GTGC*GTGC AGAGTTCTCAGAACCTTCATAGGCC CATGTGCTCATGAATCTCC	CCCGGGACAGGTTCAAGTCTAGACTTCTGCCAGGGACAGTTCAGTGTGC*----- AGAGTTCTCAGAACCTTCATAGGCCCATGTGCTCATGAATCTCC
	J2	AGC	3	AGC	3	-	-	-	-	AATAGCCACCTCAGTTGTTAACTCCTTCAAGCAGTACTGGGAC AGC*AGC GTGTCCAAAGAGAAGCTCAGAAATCAGAAAG CTTGAGGTTTCAGGAACA	AATAGCCACCTCAGTTGTTAACTCCTTCAAGCAGTACTGGGACAGC*--- TGTCCAAAGAGAAGCTCAGAAATCAGAAAGCTTGAGGTTTCAGGAACA
P3	J1	TGGN	3	-	-	GCC	3	-	-	GTGAAACCCCATCTCTACTAAATATAAAACATTAGCCAGGCATG GTGG*TTGG ACAAGGCTATTCTTTGGAGGCTGAGCAAGA GAGAGTTGTGTGGGTCA	GTGAAACCCCATCTCTACTAAATATAAAACATTAGCCAGGCATG GTGG(GCC) TTGGACAAGGCTATTCTTTGGAGGCTGAGCAAGAGA GTTGTGTGGGTCA
P4	J1	TTTTTTTTTTTTTTTTTTGAGA	26	TTT	3	-	-	-	-	CCTCACCTAGTTACTCTGTACCTCTGCACAGCTGCTTTCTGAACCTTT* TTTTTTTTTTTTTTTTTTTGGAGATAGAGTTTTG CTCTTGTGCCCAG	CCTCACCTAGTTACTCTGTACCTCTGCACAGCTGCTTTCTGAACCTTT* TTTTTTTTTTTTTTTTTTTGGAGATAGAGTTTTGCTCTTGTGCCCAG
P5	J1	-	-	-	-	-	-	-	-	CTTAAGGGTTGGCCGGCATGGTGGCTAACGCCTGTAATCCAGCAGCTTT* AAAGTACACATCTTAAAGCAATGCTAAGTGA AAACATGAAATGCAAGCT	CTTAAGGGTTGGCCGGCATGGTGGCTAACGCCTGTAATCCAGCAGCTTT* AAAGTACACATCTTAAAGCAATGCTAAGTGA AAACATGAAATGCAAGCT
P6	J1	GGAGGCTGAGGCAGGAGAACTCACTTGA	28	-	-	-	-	-	-	ACCTGTAATCTCAGCTACTCA GGAGGCTGAGGCAGGAGAACTCACTTGA *TCTGGGGGTGGAGGTTGCAGTGAGCCAGAT TGCCACCCTGCAC	GCCTGTAATCTCAGCTACTCAGGAGGCTGAGGCAGGAGAACTCACTTGA* TCTGGGGGTGGAGGTTGCAGTGAGCCAGAT TGCCACCCTGCAC
P7	J1	CAAAAATTAGCC	12	-	-	-	-	-	-	TTCAAGACCAGCCGGCCCAACATAGTGAACCCCGTTTCTACTAAAAAT* ACAAAAATAGCC AGGCATGGTGGTGAAGCCCT GTAATCCAGCTACTCCA	TTCAAGACCAGCCGGCCCAACATAGTGAACCCCGTTTCTACTAAAAAT*(A)CAAAAATTAGCCAGGCATGGTGGTGAAGCCAGC AGCTACTCCA
P8	J1	GTGGCTCAGCCTGTAATCCAGCAGCTTTGGGAGGC	36	-	-	-	-	-	-	GGCAACATTTTAGTTTTCAAACTAGCTTAAAGACAGCCCAATAGCAC* GTGGCTCAGCCTGTAATCCAGCAGCTTTGGGA GGCTGAGGCATGCAGA	GGCAACATTTTAGTTTTCAAACTAGCTTAAAGACAGCCCAATAGCAC* GTGGCTCAGCCTGTAATCCAGCAGCTTTGGGAGGC ATGCAGA
P9	J1	-	-	-	-	-	-	ACCTTGTGTTCTTT	14	AAATATATTAATAGCATCCTTTCTCAGATTATAACCTGTGTTCTTT* CAAGCTATCACATTTCCGCTCCCGTTAAAAITCCA CGTTAAA-TTCCACTTCTCTGT	AAATATATTAATAGCATCCTTTCTCAGATTATAACCTGTGTTCTTT*(INS_SVA) ACCTTGTGTTCTTT CAAGCTATCACATTTCCCTGCTCC CGTTAAA-TTCCACTTCTCTGT
P10	J1	-	-	-	-	-	-	ACCTTGTGTTCTTT	14	AAATATATTAATAGCATCCTTTCTCAGATTATAACCTGTGTTCTTT* CAAGCTATCACATTTCCGCTCCCGTTAAAAITCCA CGTTAAA-TTCCACTTCTCTGT	AAATATATTAATAGCATCCTTTCTCAGATTATAACCTGTGTTCTTT*(INS_SVA) ACCTTGTGTTCTTT CAAGCTATCACATTTCCCTGCTCC CGTTAAA-TTCCACTTCTCTGT

Table S2. Structural variants breakpoint details. New junction numbers correspond to Figure S5. Expected reference sequence refers to the putative derivative sequence as shown in Figure S5. Microhomology sequences (from 3bp to 36bp) at breakpoints are highlighted in red, inserted sequences are in blue and duplicated fragments are in bold. Breakpoint is marked with *.



What is known about this topic.

Antithrombin deficiency is mainly caused by SNV, small indels, and structural variants in *SERPINC1*, usually identified by sequencing and MLPA. Up to 25% of cases had an unknown molecular base. Nanopore sequencing is an emerging 4th generation sequencing method that obtains long reads, which are ideal for identification and characterization of gross gene defects.

What does this paper add.

Long-read whole-genome nanopore sequencing resolved all types and sizes of structural variants causing antithrombin deficiency, and identified the first causal complex structural variant. This method also found a new causing mechanism: the insertion of a new SVA retrotransposon in 2 out of 11 unknown cases. This result enlarges the catalogue of genetic disorders caused by retrotransposon insertions.

Peer Review