# Detection of *GBA* missense mutations and other variants using the Oxford Nanopore MinION

Melissa Leija-Salazar MSc<sup>1</sup>, Fritz J. Sedlazeck PhD<sup>2</sup>, Katya Mokretar PhD<sup>1</sup>, Stephen Mullin PhD MRCP<sup>1,3</sup>, Marco Toffoli MD<sup>1</sup>, Maria Athanasopoulou MSc<sup>4</sup>, Aimee Donald MRCP<sup>5</sup>, Reena Sharma MRCP<sup>6</sup>, Derralynn Hughes FRCP<sup>7</sup>, Anthony H Schapira DSc FRCP<sup>1</sup>, Christos Proukakis PhD FRCP<sup>1\*</sup>

<sup>1</sup>Department of Clinical Neuroscience, Institute of Neurology, Royal Free Hospital,

University College London, London, United Kingdom

<sup>2</sup>Human Genome Sequencing Center, Baylor College of Medicine, Houston, USA

<sup>3</sup>Institute of Translational and Stratified Medicine, Plymouth University Peninsula School of Medicine.

<sup>4</sup>Department of Molecular Neuroscience, Institute of Neurology - Queen Square,

University College London, London, United Kingdom

<sup>5</sup> Department of Paediatrics, Royal Manchester Children's Hospital

<sup>6</sup> The Mark Holland Metabolic Unit, Salford Royal Foundation NHS Trust, Salford, UK <sup>7</sup> Institute of Immunity and Transplantation, Lysosomal Storage Disorders Unit, Royal Free Hospital, London, UK

\*Corresponding author

c.proukakis@ucl.ac.uk

Tel +44 2077940500 x36701

## Abstract

## Purpose

Mutations in *GBA* cause Gaucher disease when biallelic, and are strong risk factors for Parkinson's disease when heterozygous. *GBA* analysis is complicated by the nearby pseudogene. We aimed to design and validate a method for sequencing *GBA* on the Oxford Nanopore MinION.

### Methods

We sequenced an 8.9 kb amplicon from DNA samples of 17 individuals, including patients with Parkinson's and Gaucher disease, on older and current (R9.4) flow cells. These included samples with known mutations, assessed in a blinded fashion on the R9.4 data. We used NanoOK for quality metrics, two different aligners (Graphmap and NGMLR), Nanopolish and Sniffles to call variants, and Whatshap for phasing.

#### Results

We detected all known mutations, including the common p.N409S (N370S) and p.L483P (L444P), and three rarer ones, at the correct zygosity, as well as intronic SNPs. In a sample with the complex RecNcil allele, we detected an additional coding mutation, and a 55-base pair deletion. We confirmed compound heterozygosity where relevant. False positives were easily identified.

## Conclusion

The Oxford Nanopore MinION can detect missense mutations and an exonic deletion in this difficult gene, with the added advantage of phasing and intronic analysis. It can be used as an efficient diagnostic tool.

# Introduction

The *GBA* gene encodes the lysosomal enzyme Glucocerobrosidase, deficiency of which leads to accumulation of glucosylceramide. Biallelic (homozygous or compound heterozygous) mutations in *GBA* cause Gaucher disease (GD), the most common lysosomal storage disorder.<sup>1</sup> Heterozygous *GBA* mutations are a significant risk factor for Parkinson's disease (PD),<sup>2,3</sup> with evidence of longitudinal changes in many carriers suggestive of prodromal PD.<sup>4</sup> *GBA* mutations are also associated with Dementia with Lewy bodies <sup>5</sup> and Multiple System Atrophy, <sup>6,7</sup> related conditions which also demonstrate aggregation of the alpha-synuclein protein. At present, more than 300 mutations have been linked to Gaucher disease,<sup>8</sup> and the number of studies analysing the prevalence and phenotype of *GBA* mutations in PD is rapidly increasing.<sup>9</sup> <sup>10,11</sup> <sup>12,13</sup>

*GBA* comprises eleven exons and ten introns over ~8 kb on chromosome 1q21. A nearby pseudogene *GBAP* has 96% exonic sequence homology to the *GBA* coding region.<sup>14</sup> The region also contains the Metaxin gene (*MTX1*), and its pseudogene. The existence of these two pseudogenes confers an increased risk for recombination between homologous regions, which can generate complex alleles. The homology between *GBA* and *GBAP* is highest between exons 8-11, where most of the pathogenic mutations have been reported, usually resulting from recombination events.<sup>8</sup>

The complex regional genomic structure complicates PCR and DNA sequencing, and some exons are also problematic in exome sequencing<sup>15</sup> and whole genome sequencing.<sup>16</sup> Established analysis protocols usually involve PCR of up to three

fragments, carefully designed to not amplify *GBAP*,<sup>17</sup> followed by Sanger sequencing of coding exons. Illumina targeted sequencing protocols have also recently been developed.<sup>9,12</sup> In recent years, long reads produced by Single Molecule Real Time (SMRT) sequencing have become commercially available, and have several advantages over short reads.<sup>18</sup> Oxford Nanopore sequencing technology analyses a single DNA molecule while it passes through a pore, producing characteristic changes in current depending on the sequence.<sup>19</sup> The Oxford Nanopore MinION is currently the most portable long-read sequencer. It can be plugged into a computer through a USB connection and provides sequencing data and run metrics data in real time. It has been used for applications ranging from pathogen sequencing in the field,<sup>20</sup> to sequencing a whole human genome.<sup>21</sup> It is still not routinely used in human disease diagnostics, but has been successfully used for SNV detection *in CYP2D6*, *HLA-A* and *HLA-B*,<sup>22</sup> *TP53* in cancer <sup>23</sup> and *BCR-ABL1* in leukemia,<sup>24</sup> and for chromosome 20 in a recent whole genome sequencing study.<sup>21</sup>

In the present study, we present an efficient laboratory and bioinformatic protocol for *GBA* analysis using the MinION. In addition to disease-causing variants, it can detect intronic ones, and provide phasing information. The MinION protocol can thus provide further insights into *GBA* than other sequencing technologies, and is ready to be considered for diagnostic use.

# Materials and methods

#### **Overview, DNA extraction and PCR**

Samples used in this study were derived from 17 individuals (Table S1). We used brain DNA from eight PD patients, one MSA, and one control, including two samples from brains of PD patients known to carry *GBA* mutations RecNcil and p.L483P (L444P).<sup>17</sup> Brain samples were provided by Queen Square and Parkinson's UK brain banks. We also used samples from saliva of seven living individuals, six of whom had previously been found to have at least one mutation in GBA, although these results were not known to MLS and CP, who performed the SNV analyses, until it had been completed. All individuals had given informed consent, and ethics approval was provided by the local research ethics committee. DNA was isolated from brain using Phenol-Chloroform,<sup>25</sup> and from saliva using Oragene-DNA kit.

We amplified an 8.9-kb *GBA* sequence, which covered all coding exons, the introns between them, and part of the 3' UTR region (chr1: 155202296-155211206; Figure S1). We customised previously reported primers<sup>26</sup> to carry Oxford Nanopore adapters and barcodes for multiplexing. Primer sequences were npGBA-F: 5'-TTTCTGTTGGTGCTGATATTGCTCCTAAAGTTGTCACCCATACATG-3' and npcMTX1: 5'-ACTTGCCTGTCGCTCTATCTTCCCAACCTTTCTTCCTTCTTCAA-3'.

Two DNA polymerases with appropriate optimised PCR conditions were used to amplify the *GBA* target region (Table S2): Expand Long Template PCR (Roche) and Kapa Hi-Fi polymerase (Kapa Biosystems). Amplicons were purified by Qiaquick PCR purification kit (Qiagen) and DNA concentration was measured by Qubit.

#### **Barcoding, Library Preparation and Sequencing**

For sample multiplexing, a barcoding step was carried out after generating the *GBA* amplicons with PCR Barcoding Kit I (Oxford Nanopore). Up to 12 samples were pooled in each case for library preparation according to the manufacturer amplicon sequencing protocol, starting with 1  $\mu$ g of DNA and 1%  $\lambda$ DNA CS spike-in for the dA-tailing step, followed by purification using AMPure beads. Nanopore adapters were ligated to the end-prepped DNA, using the NEB blunt/TA ligase master mix recommended by the manufacturer. Flow cell priming was performed according to the requirements of each flow cell version. We first used R7.3 and R9 flow cells with 2D reads, where a molecule passes through the pore in both directions. After recent technical advances, we used 1D reads from a R9.4 flow cell.

#### **Bioinformatic analysis**

MinKNOW versions 0.51.1.62 and later were used for data acquisition and run monitoring. Metrichor versions v2.38.1033 - v2.40.17 were used for basecalling, demultiplexing and fast5 file generation. The software divides reads into "pass" and "fail", and only "pass" reads were analysed. We used NanoOK (version 1.25)<sup>27</sup> to obtain a wide range of quality control metrics, with Graphmap (version 0.3.0)<sup>22</sup> alignment, using the precise region targeted as reference. We first converted fast5 files to fastq using NanoOK, or Poretools (version 0.6.0)<sup>28</sup> with a 2-kb size cut-off. NanoOK output included the N50 (the size at which reads of the same or greater length contain 50% of the bases sequenced), the commonest erroneous substitutions, and overall error estimates, notably the aligned base identity excluding indels (ABID), and identical bases per 100 aligned bases including indels (IBAB). We

aligned reads to the human genome (hg19) for detailed study and variant calling using GraphMap or NGMLR (version 0.2.6)<sup>29</sup>, both specifically developed for long reads. Samtools (version 1.3.1) were used where required to merge, sort and index bam files. Coverage was calculated using bedtools (version 2.25.0) coverageBed. Data were viewed on IGV.

We used Nanopolish (versions 0.6-dev and 0.8.4)<sup>20</sup> to call variants over our target region. Nanopolish was specially developed to improve accuracy by reanalysis of raw signals after alignment, and used in a recent whole genome study.<sup>21</sup> It relies on a hidden Markov model which calculates the probability of the MinION data at the signal-level for a given proposed sequence.<sup>30</sup> We called variants setting ploidy to 2, and invoked the "fix homopolymers" option. When using Nanopolish 0.8.4, we had to use Albacore (version 2.1.3, Oxford Nanopore) to re-generate fastq files for analysis. We filtered any indel calls smaller than 5 bases, due to the known problem of Nanopore in calling these, especially in homopolymer regions.<sup>21,29</sup> We reviewed the variant quality of all calls and visualised them on IGV. We used WhatsHap, designed for long reads,<sup>31</sup> to phase all true variants, and tag bam files for visualisation. We used Sniffles, another tool designed specifically for such data, to call structural variants.<sup>29</sup>

All nomenclature is based on the Human Genome Variation Society guidelines,<sup>32</sup> using reference sequence NM\_000157.3. The traditional numbering for *GBA* missense mutations, which omits the first 39 amino acids, is given in brackets to ensure easy comparability with previous literature. SNVs were annotated using ANNOVAR,<sup>33</sup> and viewed on <u>www.varsome.com</u>, which provides data from dbSNP, gnomAD <sup>34</sup> genomes and exomes where available, and other useful metrics.

# Statistical analysis

This was performed using Graphpad Prism v.6.0 (Graphpad, CA, USA) using paired

t-test and Spearman correlation analysis as indicated.

# Results

#### Validation of correct alignment and preliminary mutation detection

To test our method, we first performed sequencing with 2D reads on brain DNA, acquired over five runs, using two earlier chemistry versions (R7.3 and R9), and the Graphmap aligner. The basic metrics of all are shown in Tables S3- S4. Nine samples had GBA coverage >60 over these runs, and were analysed further. We confirmed that reads mapped to the GBA target region, with only  $\sim 1.1\%$  of reads aligning to pseudogene (Table S4). Two of these samples were known to have heterozygous GBA coding variants (three SNVs comprising the "RecNcil" allele in S5, and p.L483P, or L444P in S8). These were called by Nanopolish, and were clearly visible on IGV (Figure 1). One additional coding SNV was detected in S5 (g.155205518C>G; rs1064651; p.D448H, or D409H). This can occur in *cis* with the RecNcil allele.<sup>35</sup> There were non-coding variants in all samples. One apparently novel intronic variant in sample S1 (chr1:g.155207565C>T; intron 6: c.762-196G>A) was confirmed by Sanger sequencing (Figure S2). It has now been reported as a very rare SNP in dbSNP build 150 (rs979955939; minor allele frequency 0.0001) and gnomAD genomes (1 of 30,762 alleles). Several other non-coding SNPs were detected (Table S4). One additional candidate, which was not a known SNP, was detected in S3. Review on IGV revealed that the same base change was present in all eight other samples, with AF 10-20% in five of them (Figure S3). NanoOK showed that A>G is the second commonest erroneous substitution in this sample (15.94%). Sanger sequencing did not confirm this variant, demonstrating that a common base error, with frequent reads supporting the same variant in most samples, is a false positive.

#### Use of newer chemistry to test detection in Gaucher patients

For the next part of the study, given the rapid improvements in Nanopore chemistry, and availability of the newer R9.4 cells, we decided to test samples known to carry pathogenic mutations, to determine the potential for diagnostic use. We used the Kapa PCR protocol, because of a possible minimal error reduction (Table S4). We included DNA from the two previously tested PD brain samples carrying RecNcil and p.L483P (S5 and S8), two additional untested PD cases (one brain and one saliva), and six saliva samples which had previously been determined to carry heterozygous or biallelic mutations. These comprised four GD patients and two carriers, although their status and previously established genotypes were not revealed until after the analysis was performed. We multiplexed these 10 samples on a R9.4 flow cell. NanoOK analysis showed high base accuracy for all samples (mean 93.2%) (Tables S3 and S5). We aligned data using both Graphmap, and the newly developed NGMLR, with a mean *GBA* coverage >300, and minimal number of reads aligning to the pseudogene (average 0.78% and 1.97% of the reads aligning to gene with Graphmap and NGMLR respectively; Table S5).

#### Coding mutations are detected

We called variants using Nanopolish (version 0.8.4) on data aligned both with GraphMap and NGMLR. We first focused on coding SNVs, which were detected in eight of the ten samples, regardless of the aligner used (Table 1; Figure 2). The two untested PD patients S16 and S18 were negative. We detected all previously known coding missense mutations, at the correct zygosity. These included p.N409S (N370S) in three GD cases, in the homozygous state in two (S12, S14), and

heterozygous in one (S17) (Figure 2A), and the second mutation in S17 (p.L105P; Figure 2B). In another GD patient we detected two other heterozygous pathogenic mutations (p.R502C, p.R535C; Figure 2C,D). In the "RecNcil" carrier (S5), the additional p.D448H mutation was confirmed (Figure 2F), which would lead to this allele being designated as "RecTL". We also detected heterozygosity in three samples from individuals without GD for p.L483P (L444P) (Figure 2E), including the one tested earlier. The mean quality score for coding heterozygous SNVs was 638 (standard deviation 229), and the lowest 337.8. The lowest scores were in S17, which had the second lowest coverage (205.1). There was a non-significant trend for coding heterozygote SNVs in a sample to have higher mean quality scores with higher coverage (Spearman correlation r=0.77, p=0.10, for the six samples which had at least one). The cut-off for a true positive may therefore partly depend on coverage.

#### Non-coding SNVs are also detected, and false positives can be identified

We reviewed all other SNV calls, and noted several known SNPs present in the heterozygous or homozygous state, with quality scores also >500 (Table S6). We also noted seven SNVs that were reported in one or (usually) several samples with low quality scores (all but one <200), all but one intronic (Table S6). These were always transitions (G>A, A>G, or C>T). These base changes were identified as common errors by NanoOK (occurring in 13.31%,12.66%, and 11.95% on average of the relevant base respectively). Furthermore, review of these positions on IGV in all samples revealed a high percentage of reads with the aberrant base, including those where the SNV was not called (11-31%; Figure S4). We concluded that these were false positives. Some were shared by Graphmap and NGMLR alignments from the

same sample. Overall, however, the NGMLR alignments had significantly fewer false positives, mostly due to one SNP that was always called in Graphmap samples, but never in NGMLR (mean per sample 2.2 with Graphmap, and 1.2 with NGMLR; paired t-test p=0.0038). For the mean quality of NGMLR false positives, there was no correlation with coverage (r=0.07, p=0.91, for the seven samples which had at least one).

# Structural variant detection and mutation phasing provides additional relevant information

Sniffles and Nanopolish both reported a 55-bp exonic deletion in S5 in the NGMLRalignment only, clearly visible on IGV in this alignment (Figure S5). This sample had been previously designated "RecNcil" based on the presence of three pseudogene derived missense changes which comprise this genotype. Our detection of the additional missense change p.D448H, and the 55-bp deletion, both of which may coexist with the "RecNcil" mutations, would change the classification to a "c.1263del+RecTL allele", indicating a different site of recombination with the pseudogene than RecNcil.<sup>8</sup> Detecting this deletion can be difficult with Illumina targeted sequencing.<sup>9</sup> No other structural variants were reported.

We next phased all variants using Whatshap (Table S7). We verified that the four coding SNVs and the deletion in S5 were *in cis*, as well as five rare intronic SNPs already detected in the original analysis (Figure 3). We confirmed compound heterozygosity in two GD cases, S7, heterozygous for p.N409S and p.L105P, and S15, heterozygous for p.R502C and p.R535C. We noted a haplotype comprising 8

SNPs over 6.7 kb. This corresponds to the previously reported Pv1.1<sup>+/-</sup> haplotype,<sup>36</sup> later extended to a 70-kb haplotype designated 111.<sup>37</sup> One sample was homozygous and two heterozygous for Pv1.1<sup>+</sup> (Table S7). p.N409S (N370S) was always on the Pv1.1<sup>-</sup> background, as expected.<sup>8</sup> The p.L483P (L444P) mutation was on the Pv1.1<sup>-</sup> haplotype in two individuals and the Pv1.1<sup>+</sup> in one, consistent with the reported lack of founder effect.<sup>8</sup> p.L105P and the complex recombinant allele were on a Pv1.1<sup>-</sup> haplotype, and p.R502C and p.R535C on Pv1.1<sup>+</sup>.

# Discussion

We have sequenced a long-range *GBA* amplicon, covering all coding exons and introns, using the Oxford Nanopore Technologies MinION. We first validated our approach on brain DNA samples, using the early R7.3 and R9 chemistry. We then used the newer R9.4 chemistry, as used in human whole genome sequencing,<sup>21</sup> to analyse samples mostly known to carry biallelic or heterozygous mutations, in a blinded fashion. We confirmed common mutations in six samples (p.N409S, p.L483P), differentiating p.N409S homozygosity and heterozygosity. We also detected other mutations in two GD patients, and confirmed compound heterozygosity by phasing mutations where relevant. We further characterised the complex allele previously reported as RecNcil in one PD patient, finding another missense change and a 55-base pair deletion *in cis*, both reported with it before.<sup>8</sup>

Recent years have seen the introduction of single-molecule real time (SMRT) sequencing technologies by Oxford Nanopore and PacBio which can easily generate long reads of several kb,<sup>18</sup> and in the case of the Nanopore up to hundreds of kb.<sup>21</sup> Using long reads has several advantages, despite the lower accuracy at the base level,<sup>18</sup> some of which were evident here. The challenge of aligning short reads to regions with high homology is often not fully appreciated,<sup>15</sup> with false negatives in *GBA* targeted Illumina sequencing when the whole genome was used as a reference.<sup>9</sup> We observed minimal alignment to the pseudogene. We also detected a coding 55-bp deletion, which can be missed by Illumina data, <sup>9</sup> and intronic SNPs, an understudied area in *GBA* and other lysosomal disorders.<sup>9</sup> Finally, the long reads allowed the phasing of mutations, enabling a haplotype-resolved personalized

assessment. This helps overcome the frequent problem of phasing, which may require analysis of relatives.<sup>38</sup>

The nanopore chemistry, and bioinformatic tools available, have evolved considerably during the time in which this work was performed. We compared two aligners (Graphmap and the recently developed NGMLR), both of which gave negligible alignment to the pseudogene. NGMLR allowed detection of the 55-bp deletion, and halved the number of false positives. We thus recommend using NanoOK for guality control, NGMLR for alignment, Nanopolish for SNV calling, and Sniffles for structural variant calling. Nanopolish has been designed for SNV calling by correcting accuracy problems arising in nanopore default basecalling by reanalysing the raw signal data.<sup>21</sup> Nanopolish variant calling option uses a likelihoodbased method to generate haplotypes that serve as the reference sequence for the target region.<sup>20</sup> It has been instrumental in projects ranging from Ebola virus<sup>20</sup> to human genome sequencing.<sup>21</sup> A cut-off quality score of 320 in our work differentiated all true and false positives, although the quality score of true positives may partly depend on coverage, and calls with scores ~200-400 would require careful review, and possible Sanger analysis. Although even 120x coverage allowed detection of a p.N409S homozygote, we recommend 300x or more to ensure accurate determination of zygosity. In a human genome sequencing study with SNP analysis of chromosome 20, coverage of only 30x remarkably allowed SNP calling with accuracy ~95% against annotated variants, but zygosity was not always correctly determined.<sup>21</sup>

We were able to identify and filter false positives, based on (1) the low quality score on Nanopolish, (2) the high % of these changes occurring as errors based on

NanoOK, and (3) the significant percentage of aberrant bases at the same positions in all samples, even where not called as mutations. Notably, they were always transitions, which were also the main errors in whole genome sequencing using the MinION.<sup>21</sup> Current limitations include the inability to accurately resolve homopolymers and detect small insertions and deletions (indels),<sup>21,29</sup> and we did not attempt to do this, filtering any indel calls <5 bases. Sniffles can detect insertions and deletions, as demonstrated here, as well as complex structural variants.<sup>29</sup> Based on the rapid developments in the chemistry and bioinformatics, we expect calling of small indels and further reduction of false positive SNV calls in the very near future.

As treatments are now available, neonatal screening for lysosomal storage diseases is becoming commoner,<sup>39</sup> including in some cases Gaucher.<sup>40,41</sup> This relies on biochemical activity, often by blood-spot screening,<sup>42</sup> with several false positives in Gaucher, possibly due to carrier status.<sup>41</sup> Genetic confirmation is ultimately required, so a rapid and cost-effective method would be useful in this setting. The advantages of the MinION include the very low capital cost, space requirements, and turnaround time of the analysis. The cost per sample is likely to compare favourably with Sanger and Illumina sequencing in all settings, especially taking into account the ability to phase variants and detect structural variants in this complex region. Current R9.4 flow cells yields are at least 5 Gb of sequence, and often much more. For our 8.9 kb amplicon, 96 samples, which can be multiplexed on a single flow cell, would therefore achieve a mean coverage >1,000x, even if less than a fifth of the reads aligned successfully.

# Conclusion

Oxford Nanopore is a versatile single-molecule real-time sequencing technology that has been used in several innovative applications, from detection of Ebola to proof-of-principle human whole genome sequencing. Here we demonstrate that the MinION can detect and phase pathogenic variants in *GBA*, and intronic SNPs that would not be detected by Sanger sequencing of exons. The rapid evolution of specific bioinformatic methods, and the improvements in accuracy and data yield, combined with the minimal footprint and capital investment, make the MinION a suitable platform for long-read sequencing of difficult genes such as *GBA*, both in the diagnostic and research environments.

## Acknowledgements

Melissa Leija-Salazar is funded by CONACYT. FJS was supported by NHGRI grant UM1 HG008898. Additional funding was received by the Michael J Fox Foundation for Parkinson's research. We are grateful to the Queen Square and Parkinson's UK Brain Banks, and to all individual who donated their brains or DNA samples to research. The Queen Square Brain Bank is supported by the Reta Lila Weston Institute for Neurological Studies and the Medical Research Council UK. The Parkinson's UK Tissue Bank is funded by Parkinson's UK, a charity registered in England and Wales (258197) and in Scotland (SC037554). We are grateful to Atul Mehta and Sarah Cable for help recruiting participants. We thank Richard Leggett for support of NanoOK, and Jared Simpson for support of Nanopolish.

## **Competing interests**

F.J.S. has participated in PacBio sponsored meetings over the past few years and received travel reimbursement. CP is a participant of the Oxford Nanopore Technology early access MinION scheme.

# Table 1: Coding mutations detected.

The samples are separated into GD, where two mutations were expected, and others, where up to one was expected. \* indicates sample also included in early 2D read work. The old aminoacid notation is included. Zygosity is shown for each mutation in each sample (het=heterozygous, hom=homozygous). The Nanopolish quality score is shown for the NGMLR aligned data.

Sample	Genomic	Base change	Amino acid	Old notation	Zygosity	Nanopolish
	position		change			score
GD						
S12	155205634	c.1226A>G	p.N409S	N370S	hom	1618.9
S14	155205634	c.1226A>G	p.N409S	N370S	hom	702.1
S17	155205634	c.1226A>G	p.N409S	N370S	het	425.5
	155209547	c.314T>C	p.L105P	L66P	het	337.8
S15	155204794	c.1603C>T	p.R535C	R496C	het	975.5
	155204987	c.1504C>T	p.R502C	R463C	het	572.7
Other						
S5*	155205542	c.[1263_1317del55;	p.L422Pfs*4	c.1263del+RecTL	het	6326
	155205518	1342G>C;			het	557.1
	155205043	1448T>C;			het	382.1

	155205008	1483G>C;			het	545
	155204994	1497G>C]			het	991.9
S8*	155205043	c.1448T>C	p.L483P	L444P	het	630.3
S13	155205043	c.1448T>C	p.L483P	L444P	het	875.5
S19	155205043	c.1448T>C	p.L483P	L444P	het	728.5

Figure legends

## Figure 1. Detection of known variants in S5 and S8 in the validation phase.

The IGV trace over exon 10 is shown for all samples sequenced with 2D reads.

### Figure 2. Missense mutations detected with R9.4 chemistry.

The IGV trace is shown for each sample with a mutation. The mutated base is shown, with 20 bases on either side. The three SNPs which comprise RecNcil were shown in figure 1.

## Figure 3. Detection and phasing of a 55-base pair exonic deletion in S5.

The coverage track, with eight SNVs highlighted, and a selection of reads are shown, over exons 9 and 10 (chr1:155,204,981-155,205,661; NGMLR alignment). The deletion is clearly visible as a drop in coverage (red bracket). Reads are grouped and coloured by haplotype for these variants, which are all on the blue-coloured reads. The arrows point to the SNVs (red= coding, blue= non-coding) and the red box to the deletion.

#### References

- Schapira AHV, Chiasserini D, Beccari T, Parnetti L. Glucocerebrosidase in Parkinson's disease: Insights into pathogenesis and prospects for treatment. *Mov. Disord.* 2016;31(6):830-835. doi:10.1002/MDS.26616.
- Sidransky E, Nalls MA, Aasly JO, et al. Multicenter analysis of glucocerebrosidase mutations in Parkinson's disease. *N Engl J Med* 2009;361(17):1651-1661. doi:361/17/1651 [pii]10.1056/NEJMoa0901281.
- Mullin S, Schapira A. The genetics of Parkinson's disease. *Br. Med. Bull.* 2015;114(1):39-52. doi:10.1093/bmb/ldv022.
- Beavan M, McNeill A, Proukakis C, Hughes DA, Mehta A, Schapira AH V. Evolution of Prodromal Clinical Markers of Parkinson Disease in a GBA Mutation-Positive Cohort. *JAMA Neurol.* 2015;72(2):201-208. doi:10.1001/jamaneurol.2014.2950.
- Geiger JT, Ding J, Crain B, et al. Next-generation sequencing reveals substantial genetic contribution to dementia with Lewy bodies. *Neurobiol. Dis.* 2016;94:55-62. doi:10.1016/j.nbd.2016.06.004.
- Mitsui J, Matsukawa T, Sasaki H, et al. Variants associated with Gaucher disease in multiple system atrophy. *Ann. Clin. Transl. Neurol.* 2015;2(4):417-426. doi:10.1002/acn3.185.
- Sklerov M, Kang UJ, Liong C, et al. Frequency of *GBA* Variants in Autopsyproven Multiple System Atrophy. *Mov. Disord. Clin. Pract.* 2017;4(4):574-581. doi:10.1002/mdc3.12481.

- Hruska KS, LaMarca ME, Scott CR, Sidransky E. Gaucher disease: mutation and polymorphism spectrum in the glucocerebrosidase gene (GBA). *Hum. Mutat.* 2008;29(5):567-83. doi:10.1002/humu.20676.
- Zampieri S, Cattarossi S, Bembi B, Dardis A. GBA Analysis in Next-Generation Era. *J. Mol. Diagnostics* 2017;19(5):733-741. doi:10.1016/j.jmoldx.2017.05.005.
- Adler CH, Beach TG, Shill HA, et al. GBA mutations in Parkinson disease: earlier death but similar neuropathological features. *Eur. J. Neurol.* 2017;24(11):1363-1368. doi:10.1111/ene.13395.
- 11. Berge-Seidl V, Pihlstrøm L, Maple-Grødem J, et al. The GBA variant E326K is associated with Parkinson's disease and explains a genome-wide association signal. *Neurosci. Lett.* 2017;658:48-52. doi:10.1016/j.neulet.2017.08.040.
- Liu G, Boot B, Locascio JJ, et al. Specifically neuropathic Gaucher's mutations accelerate cognitive decline in Parkinson's. *Ann. Neurol.* 2016;80(5):674-685. doi:10.1002/ana.24781.
- Alcalay RN, Levy OA, Waters CC, et al. Glucocerebrosidase activity in Parkinson's disease with and without GBA mutations. *Brain* 2015;138(Pt 9):2648-58. doi:10.1093/brain/awv179.
- Bruce ME, McBride PA, Farquhar CF. Precise targeting of the pathology of the sialoglycoprotein, PrP, and vacuolar degeneration in mouse scrapie. *Neurosci. Lett.* 1989;102(1):1-6. doi:10.1016/0304-3940(89)90298-X.
- 15. Mandelker D, Schmidt RJ, Ankala A, et al. Navigating highly homologous

genes in a molecular diagnostic setting: a resource for clinical next-generation sequencing. *Genet. Med.* 2016;18(12):1282-1289. doi:10.1038/gim.2016.58.

- Bodian DL, Klein E, Iyer RK, et al. Utility of whole-genome sequencing for detection of newborn screening disorders in a population cohort of 1,696 neonates. *Genet. Med.* 2016;18(3):221-230. doi:10.1038/gim.2015.111.
- Neumann J, Bras J, Deas E, et al. Glucocerebrosidase mutations in clinical and pathologically proven Parkinson's disease. *Brain* 2009;132(7):1783-1794. doi:10.1093/brain/awp044.
- Goodwin S, McPherson JD, McCombie WR. Coming of age: ten years of nextgeneration sequencing technologies. *Nat. Rev. Genet.* 2016;17(6):333-51. doi:10.1038/nrg.2016.49.
- Ip CLC, Loose M, Tyson JR, et al. MinION Analysis and Reference Consortium: Phase 1 data release and analysis. *F1000Research* 2015;4(1075):1-35. doi:10.12688/f1000research.7201.1.
- Quick J, Loman NJ, Duraffour S, et al. Real-time, portable genome sequencing for Ebola surveillance. *Nature* 2016;530(7589):228-32. doi:10.1038/nature16996.
- Jain M, Koren S, Miga KH, et al. Nanopore sequencing and assembly of a human genome with ultra-long reads. *Nat. Biotechnol.* 2018. doi:10.1038/nbt.4060.
- 22. Sović I, Šikić M, Wilm A, Fenlon SN, Chen S, Nagarajan N. Fast and sensitive mapping of nanopore sequencing reads with GraphMap. *Nat. Commun.*

2016;7:11307. doi:10.1038/ncomms11307.

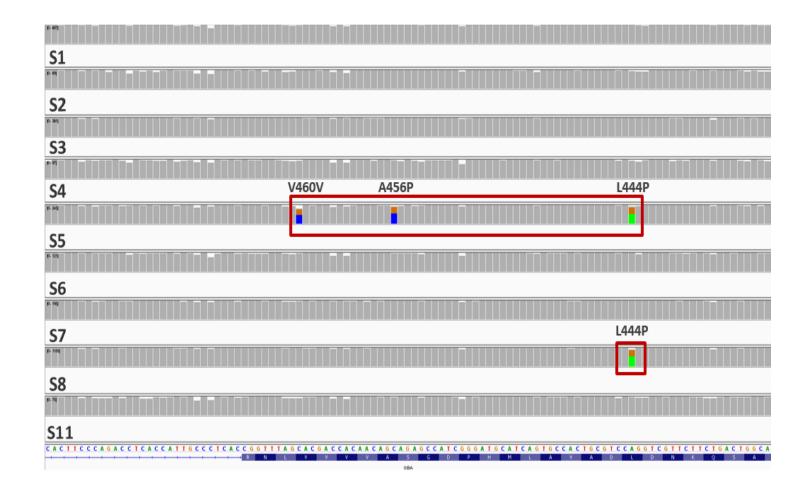
- Minervini CF, Cumbo C, Orsini P, et al. TP53 gene mutation analysis in chronic lymphocytic leukemia by nanopore MinION sequencing. *Diagn. Pathol.* 2016;11(1):96. doi:10.1186/s13000-016-0550-y.
- Minervini CF, Cumbo C, Orsini P, et al. Mutational analysis in BCR ABL1 positive leukemia by deep sequencing based on nanopore MinION technology. *Exp. Mol. Pathol.* 2017;103(1):33-37. doi:10.1016/j.yexmp.2017.06.007.
- 25. Nacheva E, Mokretar K, Soenmez A, et al. DNA isolation protocol effects on nuclear DNA analysis by microarrays, droplet digital PCR, and whole genome sequencing, and on mitochondrial DNA copy number estimation. *PLoS One* 2017;12(7):e0180467. doi:10.1371/journal.pone.0180467.
- Jeong S-Y, Kim S-J, Yang J-A, Hong J-H, Lee S-J, Kim HJ. Identification of a novel recombinant mutation in Korean patients with Gaucher disease using a long-range PCR approach. *J. Hum. Genet.* 2011;56(6):469-71. doi:10.1038/jhg.2011.37.
- Leggett RMM, Heavens D, Caccamo M, Clark MDD, Davey RPP. NanoOK: multi-reference alignment analysis of nanopore sequencing data, quality and error profiles. *Bioinformatics* 2015;32(1):142-4. doi:10.1093/bioinformatics/btv540.
- 28. Loman NJ, Quinlan AR. Poretools: a toolkit for analyzing nanopore sequence data. *Bioinformatics* 2014;30(23):3399-401. doi:10.1093/bioinformatics/btu555.
- 29. Sedlazeck FJ, Rescheneder P, Smolka M, et al. Accurate detection of complex

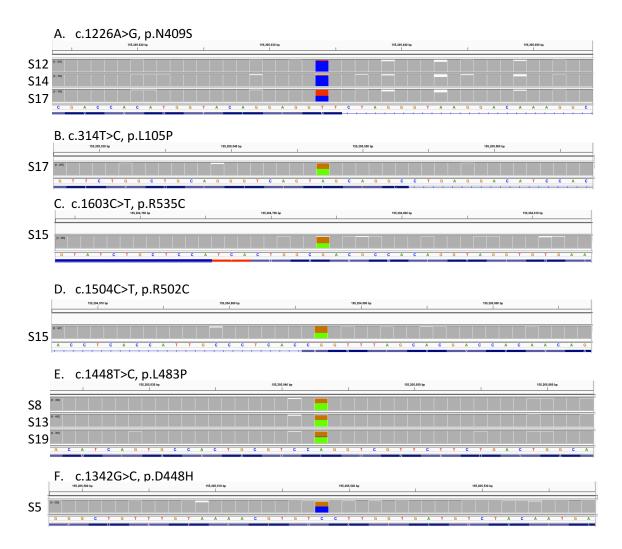
structural variations using single molecule sequencing. *bioRxiv* 2017. Available at: http://www.biorxiv.org/content/early/2017/07/28/169557. Accessed September 5, 2017.

- Loman NJ, Quick J, Simpson JT. A complete bacterial genome assembled de novo using only nanopore sequencing data. *Nat. Methods* 12(8). doi:10.1038/nMeth.3444.
- 31. Martin M, Patterson M, Garg S, et al. WhatsHap: fast and accurate read-based phasing. *bioRxiv* 2016:85050. doi:10.1101/085050.
- den Dunnen JT, Dalgleish R, Maglott DR, et al. HGVS Recommendations for the Description of Sequence Variants: 2016 Update. *Hum. Mutat.* 2016;37(6):564-569. doi:10.1002/humu.22981.
- Wang K, Li M, Hakonarson H. ANNOVAR: functional annotation of genetic variants from high-throughput sequencing data. *Nucleic Acids Res.* 2010;38(16):e164-e164. doi:10.1093/nar/gkq603.
- Lek M, Karczewski KJ, Minikel E V., et al. Analysis of protein-coding genetic variation in 60,706 humans. *Nature* 2016;536(7616):285-291. doi:10.1038/nature19057.
- Tayebi N, Stubblefield BK, Park JK, et al. Reciprocal and Nonreciprocal Recombination at the Glucocerebrosidase Gene Region: Implications for Complexity in Gaucher Disease. *Am. J. Hum. Genet.* 2003;72(3):519-534. doi:10.1086/367850.
- 36. Beutler E, West C, Gelbart T. Polymorphisms in the human

glucocerebrosidase gene. *Genomics* 1992;12(4):795-800. Available at: http://www.ncbi.nlm.nih.gov/pubmed/1572652. Accessed January 23, 2018.

- Mateu E, Pérez-Lezaun A, Martínez-Arias R, et al. PKLR-GBA region shows almost complete linkage disequilibrium over 70 kb in a set of worldwide populations. *Hum. Genet.* 2002;110(6):532-544. doi:10.1007/s00439-002-0734-2.
- Tewhey R, Bansal V, Torkamani A, Topol EJ, Schork NJ. The importance of phase information for human genomics. *Nat. Rev. Genet.* 2011;12(3):215-23. doi:10.1038/nrg2950.
- 39. Minter Baerg MM, Stoway SD, Hart J, et al. Precision newborn screening for lysosomal disorders. *Genet. Med.* 2017. doi:10.1038/gim.2017.194.
- 40. Burton BK, Charrow J, Hoganson GE, et al. Newborn Screening for Lysosomal Storage Disorders in Illinois: The Initial 15-Month Experience. *J. Pediatr.* 2017;190:130-135. doi:10.1016/j.jpeds.2017.06.048.
- Hopkins P V., Campbell C, Klug T, Rogers S, Raburn-Miller J, Kiesling J. Lysosomal Storage Disorder Screening Implementation: Findings from the First Six Months of Full Population Pilot Testing in Missouri. *J. Pediatr.* 2015;166(1):172-177. doi:10.1016/j.jpeds.2014.09.023.
- 42. Johnson BA, Dajnoki A, Bodamer O. Diagnosis of Lysosomal Storage Disorders: Gaucher Disease. In: *Current Protocols in Human Genetics*.Vol 82. Hoboken, NJ, USA: John Wiley & Sons, Inc.; 2014:17.15.1-17.15.6. doi:10.1002/0471142905.hg1715s82.





155,205,000 bp	155,205,100 bp	155,205,200 bp	681 bp 155,205,300 bp	155,205,400 bp	155,205,500 bp	■ 155,205,600 bp
(0 = 22i)					[	
11 1		11	1		1	1