

Y chromosome sequence and epigenomic reconstruction across human populations

Paula Esteller-Cucala^{*1}, Marc Palmada-Flores^{*1}, Lukas F. K. Kuderna¹, Claudia Fontserè¹, Aitor Serres-Armero¹, Marc Dabad², María Torralvo¹, Armida Faella¹, Luis Ferrández-Peral¹, Laia Llovera¹, Oscar Fornas^{3,4}, Eva Julià³, Erika Ramírez³, Irene González³, Jochen Hecht³, Esther Lizano^{1,5}, David Juan¹, Tomàs Marquès-Bonet^{1,2,4-6}

¹ Institut de Biologia Evolutiva, (CSIC-Universitat Pompeu Fabra), Doctor Aiguader 88, Barcelona, Spain

² CNAG-CRG, Centre for Genomic Regulation (CRG), Barcelona Institute of Science and Technology (BIST), Baldri i Reixac 4, Barcelona, Spain

³ Centre for Genomic Regulation (CRG), The Barcelona Institute for Science and Technology, Doctor Aiguader 88, Barcelona, Spain

⁴ Universitat Pompeu Fabra (UPF), Doctor Aiguader 88, Barcelona, Spain

⁵ Institut Català de Paleontologia Miquel Crusafont, Universitat Autònoma de Barcelona, Edifici ICTA-ICP, Cerdanyola del Vallès, Spain

⁶ Institució Catalana de Recerca i Estudis Avançats (ICREA), Passeig Lluís Companys 23, Barcelona, Spain

* These authors contributed equally

Correspondence to: tomas.marques@upf.edu and paula.esteller@upf.edu

Abstract

Recent advances in long-read sequencing technologies have allowed the generation and curation of more complete genome assemblies, enabling the analysis of traditionally neglected chromosomes, such as the human Y chromosome (chrY). Native DNA was sequenced on a MinION Oxford Nanopore Technologies sequencing device to generate genome assemblies for 7 major chrY human haplogroups. We analyzed and compared the chrY enrichment of sequencing data obtained using two different selective sequencing approaches: adaptive sampling and flow cytometry chromosome sorting. We show that adaptive sampling can produce data to create assemblies comparable to chromosome sorting while being a less expensive and time-consuming technique. We also assessed haplogroup-specific structural variants, which would be otherwise difficult to study using short-read sequencing data only. Finally, we took advantage of this technology to detect and profile epigenetic modifications amongst the considered haplogroups. Altogether, we provide a framework to study complex genomic regions with a simple, fast, and affordable methodology that could be applied to larger population genomics datasets.

Introduction

Human sex chromosomes have been traditionally excluded from genome-wide studies^{1,2}. This exclusion is particularly pronounced for the Y chromosome, the study of which could be key in understanding differences in disease susceptibility between men and women³⁻⁵. However, the Y chromosome is now considered important not only for male-specific traits but also for the study and characterization of common complex diseases⁴. Sex-limited chromosomes, defined as those unique to a heterogametic genome⁶, are usually harder to assemble since they are haploid and thus have half the sequencing depth when sequenced together with other autosomal and homogametic chromosomes. Moreover, their repetitive nature, filled with ampliconic regions and heterochromatin, poses an additional challenge for assemblers⁷.

The first Y chromosome assemblies were generated by means of bacterial artificial chromosomes (BACs) which are labor-intensive and time-consuming approaches⁸⁻¹⁰. Indeed, the Y chromosome sequences in the GRCh38 assembly¹¹⁻¹³ are a composite of BAC clones¹¹ from a male that belongs to the R1b haplogroup¹⁴ and pseudoautosomal (PAR) regions from the X-chromosome.

To facilitate the assembly process, and also to avoid the use of such costly techniques, one can decrease the potential interchromosomal assembly overlaps by specifically enriching the chromosome of interest. This can be done by physically isolating the chromosome using flow cytometry (chromosome sorting)^{6,15-17}. Alternatively, other selective sequencing methods such as adaptive sampling on Oxford Nanopore Technologies (ONT) devices¹⁸ can potentially be used.

Chromosome sorting allows the chromosome of interest to be sequenced on different platforms after its physical isolation by flow cytometry. This separation is possible because different chromosomes have specific fluorescence intensity¹⁹. On the other hand, adaptive sampling allows for the sequencing of specific DNA regions by locus-specific enrichment or depletion of off-target reads without the need for previous chromosome enrichment^{20,21}. To obtain a *de novo* assembly, it is also important to avoid whole genome amplification (WGA), as this process can introduce chimeras, bias the assembly process²² and prevent the detection of epigenetic modifications.

Long-read whole-genome sequencing enables the assessment of previously unsolved repeats, and thus allows generating more contiguous assemblies. Currently, ONT can achieve the longest read lengths compared to any other existing sequencing technology²³⁻²⁵. Moreover, ONT allows the detection of DNA (and RNA) modifications based on the different current signals of the nanopores^{26,27}. Taken together, this technology is able to resolve gaps, allowing for the true completion of chromosomes or even genomes²⁸⁻³⁰. Here we assess the performance of two enrichment methods to

sequence and assemble the Y chromosomes from 7 major human haplogroups. Moreover, we provide insights into their structural variation and epigenomic landscape showing that enrichment techniques coupled with ONT can be used to study variation between population datasets.

Results

Data production

Complete Y chromosomes from 6 different human haplogroups were isolated as previously described¹⁷ (Fig. 1A). In brief, chromosomes were obtained from lymphoblastoid cell lines (LCLs) used in the 1000 Genomes Project³¹ (1kGP) and sequenced on the ONT MinION. We also made use of the Y chromosome sorted ONT data generated by Kuderna et al.¹⁷, whose haplogroup (A0) represents one of the deepest-rooting known haplogroups. Additionally, we also generated Illumina short-read data for the same flow-sorted chromosomes (Supplementary Table 1).

The ONT data available for the cell lines ranged from 6.4 to 10.33 Gb, of which 7 to 33% mapped to the Y chromosome in the reference. Moreover, we also generated 6.4 to 35 Gb of Illumina data for all the chromosome sorting extractions. This is a notably high amount of data, especially considering that the Y chromosome sequence represents less than 1% of the known sequence in GRCh38.

The Y chromosome enrichment specificity was assessed by aligning the basecalled data to the human reference genome assembly GRCh38 and calculating the normalized coverage on each chromosome accounting for the gaps of the reference genome and the ploidy of each chromosome (Methods, Fig. 1B and Supplementary Figs. 1 and 2). The Y chromosome-specific enrichment factor of the six samples showed high variability, as it ranged from 15 to 50-fold, whereas the A0 haplogroup was over 100-fold enriched (Supplementary Table 2). As noted in Kuderna et al.¹⁷, we found that chromosome 22 partially co-sorts with chromosome Y, showing enrichment values slightly higher than 1.

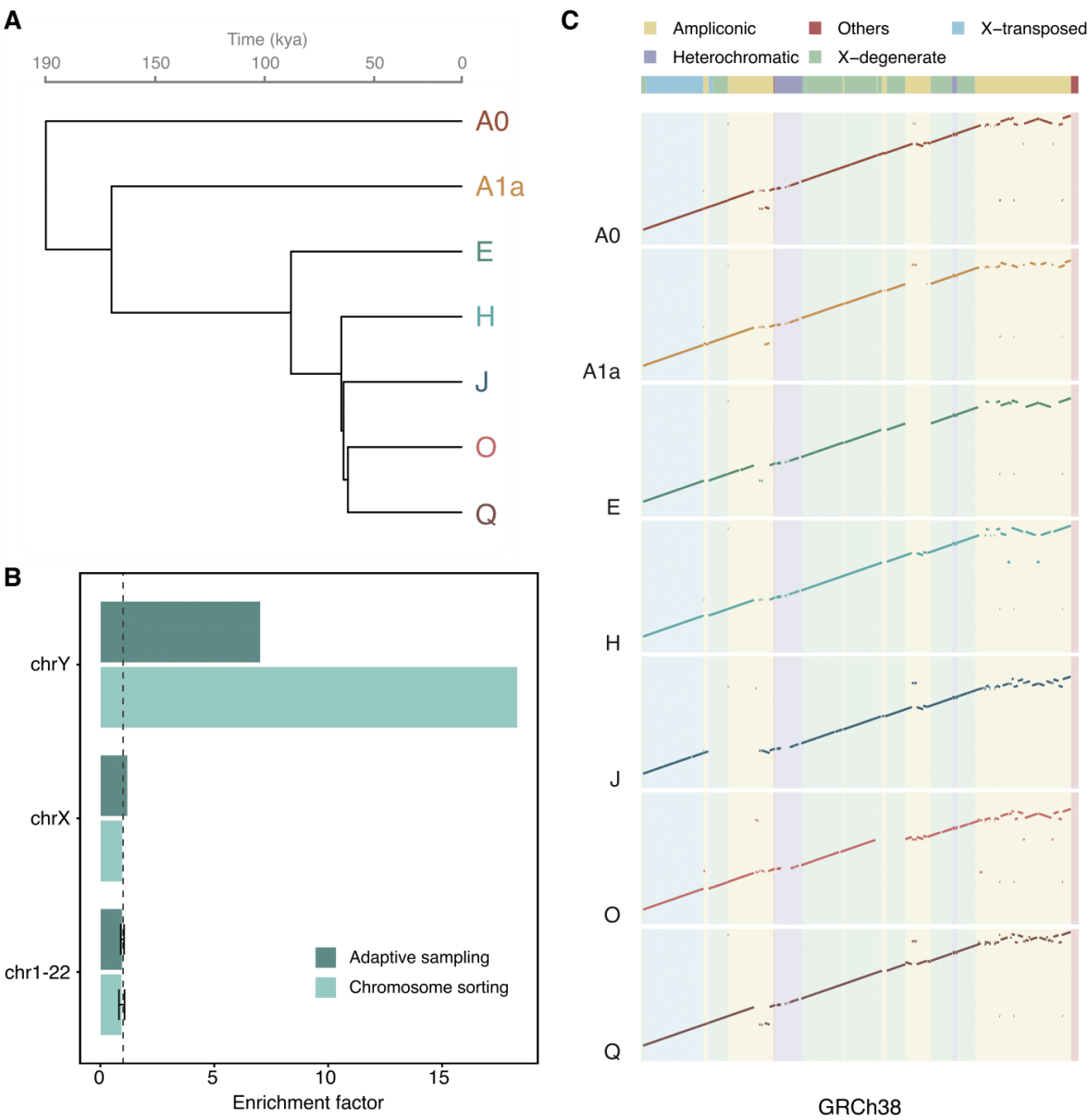


Figure 1. Study design, enrichments, and assemblies. (A) Phylogenetic tree of the human Y chromosomes used in the study. Split times taken from Jobling & Tyler-Smith, 2017³². kya, kilo years ago. (B) Enrichment factor values of the H haplogroup from data generated using chromosome sorting and adaptive sampling. The chrY shows higher enrichment with chromosome sorting than with adaptive sampling for the haplogroup compared. The dashed vertical line equal to 1 denotes no chromosomal enrichment. (C) Dot-plots of the manually scaffolded Y chromosomes compared to the resolved MSY region of GRCh38. The large-scale deletion in the J haplogroup is most likely due to its low coverage.

Adaptive sampling as a strategy to enrich specific chromosomes

A limiting factor in chromosome sorting is the need to culture hundreds of millions of cells in order to enrich the chromosome of interest effectively^{16,17}. To overcome this limitation, we explored the potential of adaptive sampling to specifically enrich the Y chromosome. This approach was done for one of the cell lines (haplogroup H) for which

chromosome sorting data had also been generated. We used the nucleotide sequences of the Y chromosome (chrY) and the contig *chrY_K1270740v1_random* (chrY_random) as provided in the GRCh38 assembly as the target sequences to enrich. To obtain comparable coverages using the two methodologies (~18x), we ran two ONT MinION flowcells with adaptive sampling. In both experiments, we showed that the Y chromosome was preferentially enriched to the other chromosomes (Supplementary Fig. 3). Although the Y chromosome enrichment factor value with chromosome sorting doubles the one in adaptive sampling for this cell line (Fig. 1B and Supplementary Table 2), adaptive sampling proves to be a cheaper and less time-consuming strategy.

Y chromosome assembly across haplogroups and enrichment techniques

We obtained Y chromosome assemblies corresponding to 7 different Y chromosome haplogroups using chromosome sorting on distinct cell lines (Methods, Supplementary Figs. 4 to 6 and Supplementary Table 3). The coverage used by the assembler to generate each assembly ranged from 13 to 50x, with a mean assembly coverage of over 28x. The resulting assemblies spanned from 18.95 - 22.23 Mb in length, being 16 to 28% shorter than the length of chromosome Y in GRCh38. We also observed that assemblies with higher continuity (contig N50) tend to have higher values of read length N50 and mean read lengths (Supplementary Table 3).

Our assemblies had a similar amount of contigs compared to a previously published African Y chromosome assembly (haplogroup A0)¹⁷, which is lower than the number of contigs of the GRCh38. The N50 across our assemblies ranged from 1.40 - 2.67 Mb, and are thus within the same order of magnitude as the Y chromosome in GRCh38 (6.9 Mb). These results suggest that creating *de novo* assemblies primarily based on long reads mapping to a reference chromosomal assembly might lead to shorter assemblies, less fragmented but with lower continuity values (such as lower contig N50). However, these results also show that by using ONT and Illumina platforms it is possible to generate assemblies almost as continuous as the current GRCh38 reference, for which much more effort and resources were devoted^{11-13,33}.

Furthermore, we manually scaffolded each haplogroup to create a single scaffold based on genome-to-genome alignments to the GRCh38 Y chromosome (Fig. 1C). Comparing our assemblies to the GRCh38 in the male-specific region of the Y chromosome (MSY) shows how most of the MSY sequence classes were assembled for the most part. As previously reported¹⁷, the ampliconic region was the most fragmented and the least complete, most likely due to collapsed repeats. Of note, for the J haplogroup, we observed a big gap in the region comprising 6.7 - 9.3 Mb of the Y chromosome. This region includes an X-degenerate region and most of its adjacent ampliconic region. When inspecting the reads mapping to this region, we found few reads present, thus possibly explaining why we could not accurately assemble this region.

Compared to the previously assembled African chrY that also made use of chromosome sorting data, we were able to generate a longer and more contiguous assembly starting from the same raw fast5 reads (Supplementary Table 3). This demonstrates the value of combining up-to-date basecalling and assembly tools, which are constantly evolving for long-read data^{34–36}.

Apart from chromosome sorting (CS), data from the H haplogroup (GM21113 cell line) was also generated using adaptive sampling (AS). In order to generate and compare the assemblies between the two enrichment methods, and given the unequal amount of data generated between them (83.5Mb difference), we restricted the comparison to assemblies generated using the same number of bases (Methods). The resulting assemblies showed similar values in metrics such as genome span (CS: 21.8 Mb, AS: 22.0 Mb), contig N50 Mb (CS: 2.7 Mb, AS: 2.6 Mb), and L50 (both L50 = 3 scaffolds). Moreover, the AS-based assembly led to a slightly more fragmented assembly (44 sequences) compared to the CS-based one (31 sequences) (Table 1 and Supplementary Fig. 7).

Table 1. Assembly metrics of Y chromosome assemblies for the GM21113 cell line (haplogroup H) using adaptive sampling and chromosome sorting, and an assembly using all the data available for GM21113.

Selective sequencing method to generate the long-read data	Assembly span (bp)	Scaffold N50 (bp)	Scaffold L50	Number of sequences
Adaptive sampling (AS)	21,955,745	2,612,207	3	44
Chromosome sorting (CS)	21,794,102	2,666,112	3	31
AS + CS	22,007,578	2,640,901	3	42

Altogether, we have generated assemblies for 7 Y chromosome haplogroups with similar contiguity to previously published assemblies. Moreover, we also show that adaptive sampling can be used for generating assemblies that are comparable to those generated by chromosome sorting.

The landscape of structural variants across the human Y chromosome phylogeny

As expected by the nature of these data, methods to detect structural variants which make use of long reads show an overall better performance than methods based on short-read data³⁷. Taking advantage of our data, we assessed the landscape of structural variants in the Y chromosome in the 7 haplogroups. For that, we used two approaches: one based on long-read mapping (*Sniffles*^{38,39}) and another based on assembly comparison (*Assemblytics*⁴⁰).

First, we identified different structural variants based on how the reads align to a reference genome using *Sniffles*. We used chrY and the chrY_random sequence from the GRCh38 as the reference. After merging the indel calls (Methods), we identified 803 unique variants (801 indels), including 194 structural variants (at least 50 bp in size, Supplementary Table 4). The number of structural variants ranged from 103 to 536 events per haplogroup. Moreover, *Sniffles* detected 1 translocation in the H haplogroup and 1 duplication event shared between 5 haplogroups (all but A0 and A1a, which are basal relative to the others). The detected duplication is located in the position *chrY:56,673,215*, at the end of a gap. This indicates that the reference is missing a region of around 98,295 bp, similar to the sequence after the gap. Most of the events were indels of 10 to 50 bp in size (Fig. 2A). Out of the 803 variants found, there were 320 insertions and 481 deletions (including 3 insertions and 9 deletions from chrY_random). We manually investigated the longest events detected using *Sniffles* and confirmed the longest insertion of 6,023 bp and the longest deletion of 6,314 bp. Both were found in haplogroup A0 and belonged to different X-degenerate regions (Supplementary Figs. 8 and 9). After merging and regenerating the panel of indels, we observed that the cell line belonging to haplogroup J was the one having more undetermined genotypes (ie. positions with no genotype information). This correlates with a lower sequencing depth for this sample. We also observed that 14% of the variants genotyped in all haplogroups shared the same genotype which was different from the reference (Fig. 2B). Haplogroups A0 and A1a harbour the most haplogroup-specific variants, concordant to their genetic distance to the reference. We also manually assessed previously reported events for the A0 haplogroup¹⁷ and confirmed that they were restricted to this cell line (Supplementary Figs 10 and 11). This indicates that structural variants found in only one haplogroup might not be representative of widespread structural variants of a chromosome, but rather delimited to one specific population or group of individuals.

Second, we identified structural variants based on the comparison of the obtained chrY assemblies to the reference chrY GRCh38 using *Assemblytics* v1.2.¹⁴⁰ (Supplementary Table 5). This allowed for the detection of 557 to 1,019 putative variants, for which 202 to 258 were at least 50 bp in size (Fig. 2A). We also found between 1 to 4 structural variants bigger than 50,000 bp for the cell lines studied, a type of variant that the mapping-based method may not detect because it would require constant coverage along a long region of the reference.

We observe 194 to 406 insertions per haplogroup with *Assemblytics*⁴⁰ compared to the 45 to 255 insertions detected with *Sniffles*^{38,39}, and 206 to 424 deletions against 57 to 288, respectively. However, similar amounts of structural variant indels are detected by *Sniffles* in all haplogroups (between 52 and 110) but for the J haplogroup (n = 21 variants) compared to *Assemblytics* (between 47 and 80, J haplogroup having 55 variants). These results, together with the fact that J haplogroup is the one with less data generated,

suggest that mapping-based structural variation detection methods may not be able to detect as many structural variants compared to assembly-to-assembly comparison-based methods when having limited sequencing depth. In that situation, generating a *de novo* assembly and using *Assemblytics* can lead to the identification of larger indels. Moreover, *Assemblytics* provides many other structural variant events such as tandem and repeat expansions or contractions, while *Sniffles* was able only to capture one duplication event and one translocation.

We further assess if the indels found using long reads could be similarly genotyped using short-read data. For that, we genotyped the variants confidently called by the ONT mapping-based approach in Illumina data generated for the same cell line extractions. With the Illumina data, we were able to replicate over one-quarter of the indels found in the nanopore data (214 out of the 801 indels). A significant positive association was seen between the predicted genotypes using ONT to those observed using Illumina data for each cell line (Fig. 2C). The observed phi coefficients (correlation values for binary variables) range between 0.35 and 0.54, not close to the highest correlation value of one^{41,42}. This is expected, given that most of the variants in our panel cannot be called by the Illumina data. The intra-haplogroup correlation is generally higher than that inter-haplogroup. However, when testing inter-haplogroup associations for genotypes derived from the same methodology, strong associations are also detected (Supplementary Fig. 12). To explore which of the indels could be observed in a panel of human variation we genotyped the same indels in the male samples present in the 1000 Genomes Project (Methods). As expected, intra-haplogroup associations were typically positive and significant, generally having stronger associations than inter-haplogroups comparisons (Supplementary Fig. 12).

259

260 CpG methylation across the phylogeny

261 ONT sequencing relies on the identification of different current signals when the DNA
262 passes through the pore, so it is possible to go beyond the identification of the four
263 canonical nucleotide bases and detect other modifications in the DNA. We used
264 *nanopolish v1.12*⁴³ to call the methylation status of 5-methylcytosines (5mC) at CpG
265 positions from the nanopore current signal. Assessing the Y chromosome methylome
266 using long reads is beneficial for exploring regions that are traditionally inaccessible using
267 short-read techniques, such as the PAR, X-transposed regions, and even the ampliconic
268 regions.

269 For that, we performed quantile normalization on the methylation values across samples
270 with a minimum coverage of 4x (Supplementary Fig. 13). We observed consistent
271 methylation patterns along Y chromosomes across samples, indicating a strong overall
272 correlation on the methylation status (Fig. 3A and Supplementary Figs. 14 and 15).
273 However, 5mC frequency values could not recapitulate the expected phylogeny, either
274 chromosome-wise or segregating by sequence class or epigenetic annotation
275 (Supplementary Fig. 16). Given that methylation levels might vary within the population,
276 age, environmental exposures, and cell culture conditions^{44–46}, and the absence of
277 replicates for each of the haplogroups considered, this observation could be due to
278 differences in any of these variables. However, given the uncertainties about the cell lines
279 and age of the individuals from which were generated, we are unable to discern the 5mC
280 variation which accounts for the different haplogroups from that which could be caused
281 by other factors. As expected by the nature of the sequence classes¹¹, the X-degenerate
282 region, which harbours single-copy genes and mostly ubiquitous expression, showed
283 5mC frequency values which resembled the most those normally seen in mammalian
284 autosomal chromosomes⁴⁷ (Supplementary Fig. 17). X-degenerate regions showed the
285 characteristic bimodal distribution of frequency values with a median close to 0.7,
286 whereas all other regions showed less defined distributions. We also inspected the
287 behavior of methylation according to the epigenetic annotation of the CpG of each of the
288 sequence classes. For that, we divided the CpGs into four mutually exclusive categories
289 (Fig. 3B and Supplementary Fig. 18): those in CpG islands (CGI), CpG shores, CpG shelves,
290 and other inter-CGI regions (open sea). CGI in the X-degenerate and X-transposed regions
291 were predominantly unmethylated, while all the other regions were mostly methylated.
292 Open sea regions showed intermediate methylation levels for all sequence classes but
293 the X-degenerate, whose median 5mC frequency reached 0.75. As expected by the
294 dynamic nature of the human methylome, CpG shores and shelves showed intermediate
295 values transitioning from CGI and open sea regions (Supplementary Fig. 19)^{48–50}.

DNA methylation is associated with gene expression⁵¹, and so we also inspected the 5mC frequency patterns across different gene annotations (Supplementary Fig. 20). Most annotated genes are present in the X-degenerate and ampliconic sequence classes (Supplementary Fig. 21), and consistent with the different expression profiles of the genes in LCLs (retrieved from GTEx⁵²) in each of these two sequence classes, we observed clear distinct methylation patterns in their TSS, UTRs, and intragenic CpGs (Fig. 3C). Not surprisingly, we found 5'UTRs to be the most constrained gene feature across samples, which directly link its methylation status to gene expression (Supplementary Fig. 22). Moreover, we found a direct relationship between upstream CGI methylation status with gene expression (Supplementary Fig. 21). Finally, we explored those cases in which differential methylation could have an effect on gene expression. We encountered a region with high methylation dispersion fully spanning a protein-coding gene (Fig. 3D and Supplementary Fig. 23). In that location, haplogroup A1a was found to be undermethyated compared to the other haplogroups, and though this difference was only modest, it could have an effect on the expression level of the gene located in this region. This gene is *NLGN4Y*, which is a long gene that spans over 300 kb and is expressed in brain and other tissues, including LCLs ($\tau_{NLGN4Y} = 0.714$). Interestingly, this gene has been proposed as a candidate for autism spectrum disorder^{53,54}. As expected, we found CGIs located upstream of this gene to be unmethylated (CGI_1 and CGI_2), whereas a CGI potentially regulating an overlapping non-coding gene in the opposite strand and which has no expression in LCLs was shown to be fully methylated in all cell lines (CGI_3).

Altogether, we show that ONT can be used to study 5mC across different cell lines, and it can prove to be helpful for the study of traditionally challenging genomic regions, particularly those present in the Y chromosome.

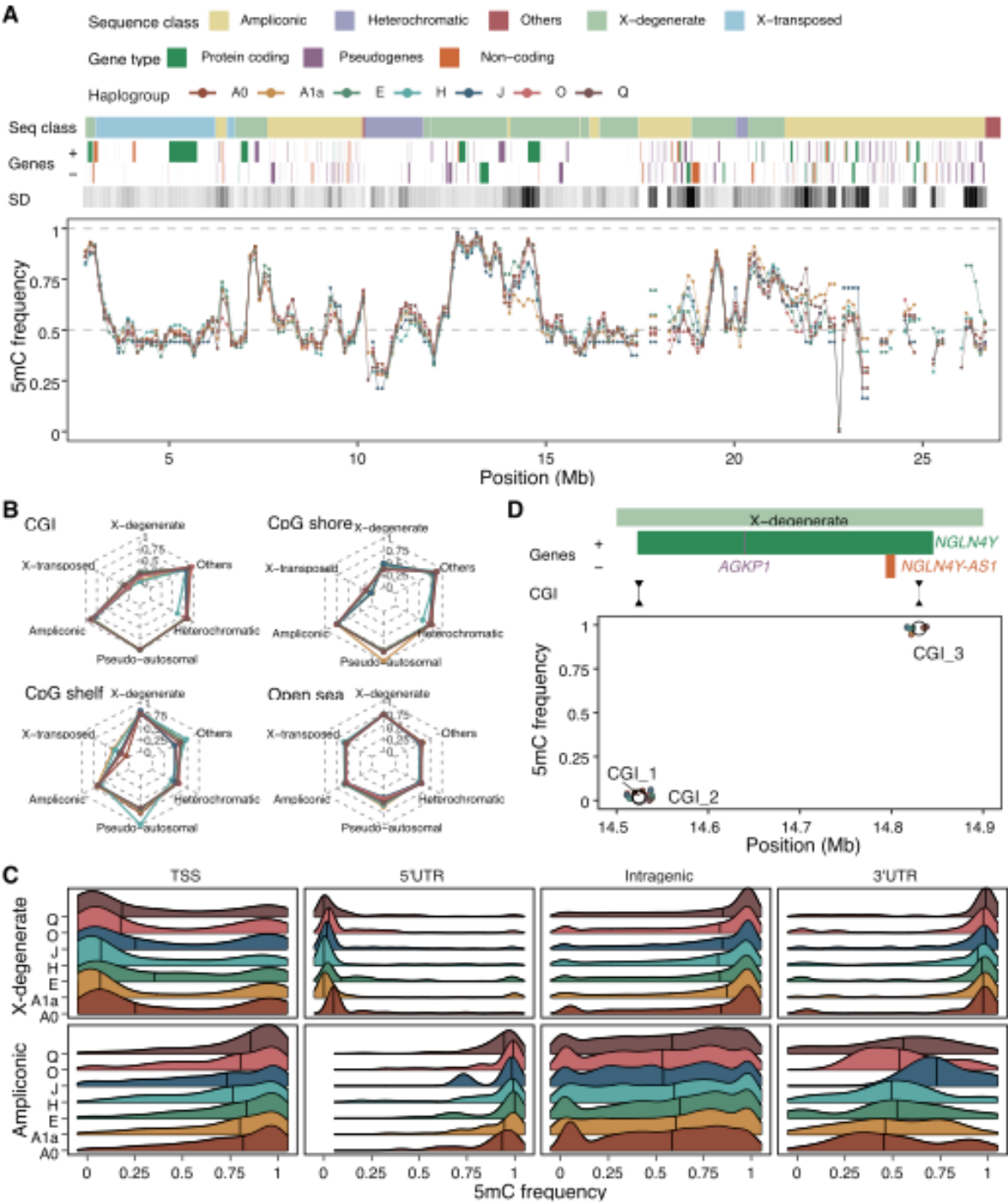


Figure 3. Methylation landscape across the Y chromosome phylogeny. (A) Frequency of 5mC in the seven cell lines along the resolved MSY of the GRCh38. The methylation levels are calculated as the median 5mC frequency value in 250kb sliding windows for each cell line. The sequence classes, the genes annotated and the standard deviation of the methylation levels across cell lines are also shown. The standard deviation of the 5mC frequency is represented in a white-to-black scale, in which a darker color denotes a higher standard deviation value. (B) Median methylation value per cell line segregated by CpG annotation and sequence classes. CpG annotations are mutually exclusive regions that comprise: CpG islands (CGI), CpG shores (up to 2kb away from the end of the CGI), CpG shelves (up to 2kb away from the end of the CpG shores), and inter-CGI or open sea regions (where all remaining CpG are allocated). (C) 5mC frequencies on different gene

features in X-degenerate and ampliconic sequence classes. Gene annotation features shown are TSS (region of 200 bp upstream of the transcription start site), both UTRs, and intragenic regions (which combine all exonic and intronic regions without considering the first gene exon). (D) Methylation frequencies in 3 CpG islands (CGI) surrounding the *NLGN4Y* and *NLGN4Y-AS1* genes. Empty circles show the mean 5mC frequency per CGI, whereas smaller colored points indicate the individual value in each cell line.

Discussion

Here, we present a panel of ONT data for 7 cell lines that represent the major human Y chromosome haplogroups. We have generated assemblies for each of them and studied their diversity, focusing on structural variation and methylation. To generate this resource, chromosome sorting data was employed and compared to adaptive sampling data, an enrichment technique that is compatible with ONT sequencing data. After generating and comparing the assemblies of the two enrichment techniques, we showed that both methods can lead to comparable assemblies, while they require different time, cost, and expertise. In terms of enrichment factor values, chromosome sorting shows co-enrichment with chromosome 22. This is mainly due to the fact that both chromosomes have similar sizes. However, this is not seen in adaptive sampling. In fact, samples enriched with adaptive sampling show the lowest standard deviation of the enrichment factor on autosomal chromosomes. Nevertheless, given the homology of the sex chromosomes, and the fact that adaptive sampling is performed by providing the genomic sequence of chromosome Y, chromosome X shows higher enrichment factor compared to the other chromosomes and the samples enriched by chromosome sorting. Altogether we show that adaptive sampling is a viable alternative strategy for the enrichment of specific genomic regions. We also emphasize the importance of using high molecular weight DNA or long DNA fragments, which are especially convenient for the enrichment of small chromosomes with adaptive sampling. As such, at longer DNA fragment sizes, the time the sequencer will be scanning for on-target regions (i.e., those that belong to the Y chromosome) will be reduced. Therefore, we realize that having started from higher DNA fragment sizes for the haplogroup H sample would have led to higher enrichment efficiencies in the adaptive sampling enrichment method.

One major limitation of our work is the conservative filtering we have used to generate the assemblies. Our approach uses the Y chromosome of the current genome of reference GRCh38 as a backbone. All data obtained using adaptive sampling relies heavily on the GRCh38 reference and may include a few reads of other chromosomes that start with a similar sequence. On the other side, chromosome sorting produces data on unresolved chromosomal regions but includes some undesired full chromosomes data. As such, restricting our assembly to only those reads that map to the reference leads to the loss of a fraction of Y chromosome potentially informative reads during the filtering process. Conversely, this approach minimizes the retention of non-chromosome Y data

and limits the resulting assembly to the Y chromosome only. Compared to a previous assembly created with the same data for the cell line that belongs to haplogroup A0, our approach yielded a more contiguous assembly. As such, it shows the potential that re-processing the same raw data with novel approaches might have in the future, especially in the context of the big data era^{55,56}.

Due to its large fraction of heterochromatin, around half of the sequence in the current Y chromosome assembly is unresolved. This limitation, together with the fact that we are using a partial reference genome to generate assemblies of a specific chromosome, hampers the possibility of reconstructing the totality of this chromosome. In the future, telomere-to-telomere Y chromosome sequencing would undoubtedly avoid reference-biases we encountered in this study²⁸.

We also took advantage of the long-read data generated to explore the landscape of structural variants in each cell line. For that, we used two different methods for structural variant calling: one based on long-read mapping and another based on assembly comparison. The former allows for two rounds of genotyping and so the final candidates are potentially more curated. The latter is based on genome-to-genome comparisons, so it is able to detect longer genomics variants. We consider that for low data samples the creation of a *de novo* Y chromosome assembly may allow the detection of structural variants that cannot be recognized with a mapping-based method, considering the low coverage of reads mapping in those regions.

Besides the potential to generate high-accuracy assemblies and resolve complex genomic regions like structural variants, ONT also allows for studying the epigenome. We have assessed the methylation status of cytosines in a CpG context for our panel of cell lines. Despite the fact that the epigenome of the Y chromosome has not been deeply studied, we were able to consistently replicate the methylation patterns that have been described in other human autosomal chromosomes^{57,58}. Not surprisingly, with the methylation values obtained, we were unable to recapitulate our samples' expected phylogeny. Two main factors can be attributed to this: the lack of replicates for each haplogroup and also within-population variability^{45,46,59,60} which, in our case, could also be confounded by epigenetic drift⁶¹. Still, methylation differences at the population level are to be expected to be small in magnitude⁴⁵. In this line, we were able to detect small differences in methylation in regions that could have an effect on the regulation of specific genes. This is the case of gene *NYGN4Y*, which we found to fully overlap with a region with consistently lower methylation in the cell line belonging to haplogroup A1a.

Nevertheless, we are using lymphoblastoid cell lines (LCLs), which are artificially transformed cells, so caution must be taken when extrapolating these findings. But the extent to which the generalization of our results could be biased is even more consequential when reporting those findings that are sample-specific. As such, an

increase in the number of replicates would help to discern which of our findings are artifacts from those which have a true biological meaning.

Taken together, here we provide a framework to study complex genomic regions. We applied this simple, fast and affordable technology to study diverse human population groups. Moreover, this approach can be applied to the generation of long-read data of other regions or chromosomes of interest. As such, it could be used for the characterization of virtually any species, although it would be especially advantageous for those rich in complex genomic features.

Methods

Flow chromosome sorting followed by ONT or Illumina sequencing

Chromosome preparation was performed as previously described^{16,17}. The libraries to obtain the Illumina paired-end data were constructed using a SureSelect V6-Post Library Kit. Raw data generated for the haplogroup A0 (HG02982 cell line) was retrieved from Kuderna et al.¹⁷ The data generated in each MinION run was basecalled using *Guppy v5.0.15*³⁵ with the super accuracy model *dna_r94.1_450bps_sup*.

Adaptive sampling for the enrichment of a specific chromosome

We extracted DNA from cultured cells of haplogroup H (GM21113 cell line) using the Qiagen MagAttract HMW Kit. DNA libraries for ONT sequencing were obtained using the Ligation Sequencing Kit (SQK-LSK110) and sequenced in two ONT MinION flowcells (FLO-MIN106 R9.4.1) using a MinION Mk1C with MinKNOW v21.02-beta4~xenial. We aimed for the specific enrichment of the chrY and the chrY_random, by adding their nucleotide sequence as provided in the GRCh38 assembly. This method bioinformatically labels the reads that are being sequenced for enrichment or depletion. After a DNA strand enters the pore, the sequencer only needs one second (around 420 bases) to decide whether to continue sequencing the DNA if it matches the region of interest or to eject it if it does not. Each of the strands that enter a pore will be labeled as *unblock* and *no decision* when they are rejected by the pore or they are so short that their status remains inconclusive, respectively. They will be labeled as *stop receiving* when they are on target, thus further sequenced. Only reads labeled as *stop receiving* were used in this project.

The enrichment obtained with adaptive sampling highly depends on the fragment length of the library. Longer DNA library lengths are preferred, as the adaptive sampling enrichment algorithm takes a fixed amount of time to recognize whether to enrich a DNA strand. Because of this, in order to target a specific region of the genome that is particularly small (the Y chromosome represents ~1% of the genome), it will always be better to have few long DNA fragments, rather than many short DNA fragments, as the time spent by the sequencer scanning for on-target regions will be reduced.

Assessing the performance of two different enrichment methodologies

The coverage and enrichment factor for each chromosome were calculated as follows:

$$\text{Coverage of chrN} = \frac{\text{Mapped bp in chrN}}{\text{Size chrN (bp, without N)}} \quad (\text{Eq. 1})$$

$$\text{Enrichment factor of chrN} = \frac{\frac{\text{Mapped bp in chrN}}{\text{Total mapped bp}}}{\frac{\text{Size chrN} \times n(\text{bp, without N})}{\text{Diploid genome size (bp, without N)}}} \quad (\text{Eq. 2})$$

Since more than 50% of the Y chromosome in GRCh38 is composed of long stretches of unknown sequence (that in the assembly is seen as N), it is important to exclude these regions from the coverage and enrichment calculations. Because of that, Eq. 1 and Eq. 2 only consider chromosome sizes without Ns. Moreover, for calculating the enrichment, and in order to account for the real target space of each chromosome, the size of each of them is multiplied by its ploidy.

Assembly generation

Basecalled passed reads ($Q > 10$) were mapped to the human GRCh38 genome assembly using *minimap2* v2.17-r941⁶² with the option *-x map-ont*. The resulting bam was indexed using *SAMTOOLS* v1.12^{62,63} and the reads mapping either to chrY or chrY_random (chrY specific reads) were retrieved.

We ran *Flye* v2.9³⁴ using the chrY specific reads with the option *--nano-hq* as suggested by the developers while using data basecalled using *Guppy* v5³⁵ onwards with the super accuracy model. We added the option *--scaffold* to enable scaffolding based on the assembly graph, and included two internal rounds of polishing with the argument *-i 2*.

As we used uncorrected long reads to obtain the draft assemblies, we polished the initial assemblies by first using ONT reads. We started with 2 rounds of *Racon* v1.3.1⁶⁴, using *minimap2* v2.9-r720⁶² with the option *-x ont* to obtain the mapping file, and adding to *Racon* the argument *-u* to keep any unpolished sequences. To further improve the assembly, we then ran *medaka* v1.4.1⁶⁵ using the *medaka_consensus* program with default settings and the model *-m r941_prom_sup_g507*.

Additionally, to polish the assemblies with Illumina data we used *HyPo* v1.0.3⁶⁶, mapping the Illumina reads to the polished assembly. For mapping short reads to the existing assembly we used *minimap2* v2.9-r720⁶² with the option *-x sr*.

Once we polished the assemblies we purged them using *purge_dups* v1.2.5^{62,67}, with default parameters and the *-2* option. This was done to remove any haplotig present in the assemblies. We obtained the mapping files using *minimap2* v2.14-r883⁶² with the option *-x ont* to map the ONT reads to the polished assembly, and with the options *-xasm5 -DP* to map the split polished assembly to itself.

For the comparison of the assemblies generated from the two selective sequencing methods (chromosome sorting and adaptive sampling), we downsampled the data of the adaptive sampling experiment. For that, we used *Filtlong* v0.2.0⁶⁸ with the option `--keep_percent 87.8` so as to retrieve 87.8% of the AS sequencing data. From that point on, the assembly process was the same as the one explained (Supplementary Fig. 5).

Genome-to-genome comparisons

To obtain genome-to-genome alignments we used MuMmer v3.23⁶⁹ *nucmer* tool with options `--maxmatch -l 100 -c 100`. To manually scaffold the Y chromosome assemblies, we used the dot-plot viewer *dot*⁷⁰.

Structural variant detection with long reads

Structural variation was called using *Sniffles* v2.0.2^{38,39} with a minimum number of reads that support an SV of `-s 10`, fed with the bam files for which we calculated and added MD tags using *SAMTOOLS* v1.9⁶³, with the program *samtools calmd* adding options `-uAr -Q`. We summarized the amount of SVs per type and filtered out the SVs considered 'IMPRECISE' by *Sniffles*.

We merged the insertions and deletion separately with a maximum permitted distance of 100 bp (so that indels located 100 bp upwards or downwards will be considered a single event) found independently in all the cell lines using *SURVIVOR* v1.0.7⁷¹. We removed any genotype with quality under 25 (MQ) and the events that were homozygous for the reference genotype in all samples.

We used *Assemblytics* v1.2.1⁴⁰ to find structural variants in the different assemblies generated by comparing them to the reference GRCh38. We looked for structural variants with sizes between 10 to 100,000 and the unique sequence length required to call a variant of 1,000.

Structural variant genotyping with short reads

We genotyped, using the indels obtained using *Sniffles* and *SURVIVOR* as reference, the structural variants based on Illumina data with the program *graph typer* v2.7.5 with the option `"genotype_sv"` and only kept the indels with a quality > 0 . We genotyped them with the Illumina data generated in this study and with the Illumina data of the *1kGP*³¹ available for the Y chromosome.

Correlation between structural variant detection using long or short reads

To assess the reproducibility of the structural variant calls obtained with ONT data in short-read data. We took the indels genotyped in the Illumina data and compared them between platforms. For each structural variant, the ONT genotypes were assumed to be

true positives and all genotype calls were binarised into presence (1) or absence (0). For each structural variant, and given that Y chromosomes are hemizygous, homozygous and heterozygous alternative calls were considered present, and homozygous reference genotypes absent.

As we wanted to study the correlation of binary variables, we made use of phi coefficients, also known as Matthews correlation coefficient or MCC. Phi coefficients should be interpreted similarly to a Pearson correlation coefficient. For the *1kgp* data, we considered a structural variant to be present if it was at a frequency higher than 0.2.

Studying methylation using ONT

The methylation status was called using *nanopolish v0.13.2*⁴³, which assigns a log-likelihood ratio to each individual CpG site. To avoid adding noise to the methylation results, we only used reads with the highest mapping quality as provided by minimap2 (mapQ = 60) and filtered out all others. We used the default log-likelihood threshold of 2 as implemented in *nanopolish v0.12* onwards. As suggested by the developers, we called methylation with the option `--min-separation 5` to help calling CpG dense regions. The methylation frequency was calculated for each site as the number of mapped reads predicted as methylated divided by the number of total mapped reads.

We filtered out the few instances in which alternative alleles were present in a genomic position with a cytosine in the reference sequence. We performed quantile normalization on the methylation values across samples with a minimum coverage of 4x using the R package *preprocessCore v1.56.0*⁷². CpG and genic annotations were obtained using the R package *annotatr v1.24.0*⁷³. Minor modifications were made to these annotations for different analyses. All these modifications have been specifically described when used in the text. For the overlapping regions in the genic annotations, the priority set was the following: promoters, UTRs (5', 3'), first exon, non-first exons, all introns, and upstream region.

Contributions

T.M.B and L.F.K.K. conceived the study. O.F., E.J., and E.V. cultured cells and performed the flow cytometry. A.F. cultured cells. P.E.-C., C.F., and L.L. performed adaptive sampling. I.G. and J.H. performed sequencing. L.F.K.K., A.S., L.F.-P., and E.L. provided analytical support. M.D. and M.T. helped with data curation and analyses. M.P.-F. generated assemblies and structural variant calls. P.E.-C. performed methylation analyses and supervised analyses. D.J. designed and supervised analyses. P.E.-C. and M.P.-F. wrote the manuscript with input from all co-authors.

Acknowledgments

MPF has support of an INPhINIT Retaining Fellowship from "La Caixa" Foundation (ID 100010434) with code LCF/BQ/DR20/11790032. TMB is supported by funding from the European Research Council (ERC) under the European Union's Horizon 2020 research and innovation programme (grant agreement No. 864203), "Unidad de Excelencia María de Maeztu", funded by the AEI (CEX2018-000792-M) and NIH 1R01HG010898-01A1.

References

1. Accounting for sex in the genome. *Nature Medicine* vol. 23 1243–1243 Preprint at <https://doi.org/10.1038/nm.4445> (2017).
2. Wise, A. L., Gyi, L. & Manolio, T. A. eXclusion: toward integrating the X chromosome in genome-wide association analyses. *Am. J. Hum. Genet.* **92**, 643–647 (2013).
3. Wilson, M. A. The Y chromosome and its impact on health and disease. *Human Molecular Genetics* vol. 30 R296–R300 Preprint at <https://doi.org/10.1093/hmg/ddab215> (2021).
4. Anderson, K., Cañadas-Garre, M., Chambers, R., Maxwell, A. P. & McKnight, A. J. The Challenges of Chromosome Y Analysis and the Implications for Chronic Kidney Disease. *Frontiers in Genetics* vol. 10 Preprint at <https://doi.org/10.3389/fgene.2019.00781> (2019).
5. Molina, E., Clarence, E. M., Ahmady, F., Chew, G. S. & Charchar, F. J. Coronary Artery Disease: Why We should Consider the Y Chromosome. *Heart, Lung and Circulation* vol. 25 791–801 Preprint at <https://doi.org/10.1016/j.hlc.2015.12.100> (2016).
6. Mank, J. E. The W, X, Y and Z of sex-chromosome dosage compensation. *Trends in Genetics* vol. 25 226–233 Preprint at <https://doi.org/10.1016/j.tig.2009.03.005> (2009).
7. Tomaszewicz, M., Medvedev, P. & Makova, K. D. Y and W Chromosome Assemblies: Approaches and Discoveries. *Trends in Genetics* vol. 33 266–282 Preprint at <https://doi.org/10.1016/j.tig.2017.01.008> (2017).
8. Hughes, J. F. *et al.* Chimpanzee and human Y chromosomes are remarkably divergent in structure and gene content. *Nature* **463**, 536–539 (2010).
9. Hughes, J. F. *et al.* Strict evolutionary conservation followed rapid gene loss on human and rhesus Y chromosomes. *Nature* vol. 483 82–86 Preprint at <https://doi.org/10.1038/nature10843> (2012).
10. Soh, Y. Q. S. *et al.* Sequencing the mouse Y chromosome reveals convergent gene acquisition and amplification on both sex chromosomes. *Cell* **159**, 800–813 (2014).
11. Skaletsky, H. *et al.* The male-specific region of the human Y chromosome is a mosaic of discrete sequence classes. *Nature* **423**, 825–837 (2003).
12. Lander, E. S. *et al.* Initial sequencing and analysis of the human genome. *Nature* **409**, 860–921 (2001).
13. International Human Genome Sequencing Consortium. Finishing the euchromatic

- sequence of the human genome. *Nature* **431**, 931–945 (2004).
14. Mendez, F. L., David Poznik, G., Castellano, S. & Bustamante, C. D. The Divergence of Neandertal and Modern Human Y Chromosomes. *Am. J. Hum. Genet.* **98**, 728–734 (2016).
15. Tomaszewicz, M. *et al.* A time- and cost-effective strategy to sequence mammalian Y Chromosomes: an application to the de novo assembly of gorilla Y. *Genome Res.* **26**, 530–540 (2016).
16. Kuderna, L. F. K. *et al.* Flow Sorting Enrichment and Nanopore Sequencing of Chromosome 1 From a Chinese Individual. *Front. Genet.* **10**, 1315 (2019).
17. Kuderna, L. F. K. *et al.* Selective single molecule sequencing and assembly of a human Y chromosome of African origin. *Nat. Commun.* **10**, 4 (2019).
18. Martin, S. *et al.* Nanopore adaptive sampling: a tool for enrichment of low abundance species in metagenomic samples. *Genome Biol.* **23**, 11 (2022).
19. Doležal, J. *et al.* Chromosomes in the flow to simplify genome analysis. *Funct. Integr. Genomics* **12**, 397–416 (2012).
20. Payne, A. *et al.* Readfish enables targeted nanopore sequencing of gigabase-sized genomes. *Nat. Biotechnol.* **39**, 442–450 (2021).
21. Kovaka, S., Fan, Y., Ni, B., Timp, W. & Schatz, M. C. Targeted nanopore sequencing by real-time mapping of raw electrical signal with UNCALLED. *Nat. Biotechnol.* **39**, 431–441 (2021).
22. Pinard, R. *et al.* Assessment of whole genome amplification-induced bias through high-throughput, massively parallel whole genome sequencing. *BMC Genomics* **7**, 216 (2006).
23. Udaondo, Z. *et al.* Comparative Analysis of PacBio and Oxford Nanopore Sequencing Technologies for Transcriptomic Landscape Identification of *Penaeus monodon*. *Life* vol. 11 862 Preprint at <https://doi.org/10.3390/life11080862> (2021).
24. Lang, D. *et al.* Comparison of the two up-to-date sequencing technologies for genome assembly: HiFi reads of Pacbio Sequel II system and ultralong reads of Oxford Nanopore. Preprint at <https://doi.org/10.1101/2020.02.13.948489>.
25. Tvedte, E. S. *et al.* Comparison of long-read sequencing technologies in interrogating bacteria and fly genomes. *G3 Genes|Genomes|Genetics* vol. 11 Preprint at <https://doi.org/10.1093/g3journal/jkab083> (2021).
26. Rand, A. C. *et al.* Mapping DNA methylation with high-throughput nanopore sequencing. *Nat. Methods* **14**, 411–413 (2017).
27. Simpson, J. T. *et al.* Detecting DNA cytosine methylation using nanopore sequencing. *Nat. Methods* **14**, 407–410 (2017).
28. Nurk, S. *et al.* The complete sequence of a human genome. *Science* **376**, 44–53 (2022).
29. Logsdon, G. A. *et al.* The structure, function and evolution of a complete human chromosome 8. *Nature* **593**, 101–107 (2021).
30. Miga, K. H. *et al.* Telomere-to-telomere assembly of a complete human X

- chromosome. *Nature* **585**, 79–84 (2020).
31. A global reference for human genetic variation. *Nature* **526**, 68–74 (2015).
32. Jobling, M. A. & Tyler-Smith, C. Human Y-chromosome variation in the genome-sequencing era. *Nat. Rev. Genet.* **18**, 485–497 (2017).
33. Tilford, C. A. *et al.* A physical map of the human Y chromosome. *Nature* **409**, 943–945 (2001).
34. Kolmogorov, M., Yuan, J., Lin, Y. & Pevzner, P. A. Assembly of long, error-prone reads using repeat graphs. *Nat. Biotechnol.* **37**, 540–546 (2019).
35. Wick, R. R., Judd, L. M. & Holt, K. E. Performance of neural network basecalling tools for Oxford Nanopore sequencing. *Genome Biol.* **20**, 129 (2019).
36. Dida, F. & Yi, G. Empirical evaluation of methods for *de novo* genome assembly. *PeerJ Computer Science* vol. 7 e636 Preprint at <https://doi.org/10.7717/peerj-cs.636> (2021).
37. Mahmoud, M. *et al.* Structural variant calling: the long and the short of it. *Genome Biol.* **20**, 246 (2019).
38. Sedlazeck, F. J. *et al.* Accurate detection of complex structural variations using single-molecule sequencing. *Nat. Methods* **15**, 461–468 (2018).
39. Smolka, M. *et al.* Comprehensive Structural Variant Detection: From Mosaic to Population-Level. *bioRxiv* 2022.04.04.487055 (2022) doi:10.1101/2022.04.04.487055.
40. Nattestad, M. & Schatz, M. C. Assemblytics: a web analytics tool for the detection of variants from an assembly. *Bioinformatics* **32**, 3021–3023 (2016).
41. Chicco, D. & Jurman, G. The advantages of the Matthews correlation coefficient (MCC) over F1 score and accuracy in binary classification evaluation. *BMC Genomics* **21**, 6 (2020).
42. Yule, G. U. On the methods of measuring association between two attributes. *J. R. Stat. Soc.* **75**, 579 (1912).
43. Loman, N. J., Quick, J. & Simpson, J. T. A complete bacterial genome assembled *de novo* using only nanopore sequencing data. *Nat. Methods* **12**, 733–735 (2015).
44. Varley, K. E. *et al.* Dynamic DNA methylation across diverse human cell lines and tissues. *Genome Res.* **23**, 555–567 (2013).
45. Fraser, H. B., Lam, L. L., Neumann, S. M. & Kobor, M. S. Population-specificity of human DNA methylation. *Genome Biol.* **13**, R8 (2012).
46. Husquin, L. T. *et al.* Exploring the genetic basis of human population differences in DNA methylation and their causal impact on immune gene regulation. *Genome Biol.* **19**, 222 (2018).
47. Bell, J. T. *et al.* DNA methylation patterns associate with genetic and gene expression variation in HapMap cell lines. *Genome Biol.* **12**, R10 (2011).
48. Illingworth, R. S. & Bird, A. P. CpG islands--'a rough guide'. *FEBS Lett.* **583**, 1713–1720 (2009).
49. Irizarry, R. A. *et al.* The human colon cancer methylome shows similar hypo- and hypermethylation at conserved tissue-specific CpG island shores. *Nat. Genet.* **41**, 178–

- 186 (2009).
50. Ziller, M. J. *et al.* Charting a dynamic DNA methylation landscape of the human genome. *Nature* **500**, 477–481 (2013).
51. Lowdon, R. F., Jang, H. S. & Wang, T. Evolution of Epigenetic Regulation in Vertebrate Genomes. *Trends Genet.* **32**, 269–283 (2016).
52. GTEx Consortium. The GTEx Consortium atlas of genetic regulatory effects across human tissues. *Science* **369**, 1318–1330 (2020).
53. Ross, J. L., Tartaglia, N., Merry, D. E., Dalva, M. & Zinn, A. R. Behavioral phenotypes in males with XYY and possible role of increased NLGN4Y expression in autism features. *Genes Brain Behav.* **14**, 137–144 (2015).
54. Chen, J., Yu, S., Fu, Y. & Li, X. Synaptic proteins and receptors defects in autism spectrum disorders. *Front. Cell. Neurosci.* **8**, 276 (2014).
55. Dall’Alba, G., Casa, P. L., Abreu, F. P. de, Notari, D. L. & de Avila E Silva, S. A Survey of Biological Data in a Big Data Perspective. *Big Data* **10**, 279–297 (2022).
56. Kamble, S. S., Gunasekaran, A., Goswami, M. & Manda, J. A systematic perspective on the applications of big data analytics in healthcare management. *Int. J. Healthc. Manag.* **12**, 226–240 (2019).
57. McCarthy, N. S. *et al.* Meta-analysis of human methylation data for evidence of sex-specific autosomal patterns. *BMC Genomics* **15**, 981 (2014).
58. Johansson, A., Enroth, S. & Gyllenstein, U. Continuous Aging of the Human DNA Methylome Throughout the Human Lifespan. *PLoS One* **8**, e67378 (2013).
59. Palumbo, D., Affinito, O., Monticelli, A. & Coccozza, S. DNA Methylation variability among individuals is related to CpGs cluster density and evolutionary signatures. *BMC Genomics* **19**, 229 (2018).
60. Galanter, J. M. *et al.* Differential methylation between ethnic sub-groups reflects the effect of genetic ancestry and environmental exposures. *Elife* **6**, (2017).
61. Horvath, S. & Raj, K. DNA methylation-based biomarkers and the epigenetic clock theory of ageing. *Nat. Rev. Genet.* **19**, 371–384 (2018).
62. Li, H. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics* **34**, 3094–3100 (2018).
63. Li, H. *et al.* The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**, 2078–2079 (2009).
64. Vaser, R., Sović, I., Nagarajan, N. & Šikić, M. Fast and accurate de novo genome assembly from long uncorrected reads. *Genome Res.* **27**, 737–746 (2017).
65. GitHub - nanoporetech/medaka: Sequence correction provided by ONT Research. *GitHub* <https://github.com/nanoporetech/medaka>.
66. Kundu, R., Casey, J. & Sung, W.-K. HyPo: Super Fast & Accurate Polisher for Long Read Genome Assemblies. *bioRxiv* 2019.12.19.882506 (2019) doi:10.1101/2019.12.19.882506.
67. Guan, D. *et al.* Identifying and removing haplotypic duplication in primary genome

703 assemblies. *Bioinformatics* **36**, 2896–2898 (2020).
704 68. rrwick/Filtlong. *GitHub* <https://github.com/rrwick/Filtlong>.
705 69. Kurtz, S. *et al.* Genome Biology. vol. 5 R12 Preprint at [https://doi.org/10.1186/gb-](https://doi.org/10.1186/gb-2004-5-2-r12)
706 2004-5-2-r12 (2004).
707 70. Website. GitHub - MariaNattestad/dot: Dot: An interactive dot plot viewer for
708 comparative genomics. GitHub <https://github.com/marianattestad/dot>.
709 71. Jeffares, D. C. *et al.* Transient structural variations have strong effects on quantitative
710 traits and reproductive isolation in fission yeast. *Nat. Commun.* **8**, 14061 (2017).
711 72. Ben Bolstad <bmb@bmbolstad.com>. *preprocessCore*. (Bioconductor, 2017).
712 doi:10.18129/B9.BIOC.PREPROCESSCORE.
713 73. Cavalcante, R. G. & Sartor, M. A. annotatr: genomic regions in context. *Bioinformatics*
714 **33**, 2381–2383 (2017).