

1 T2T-CHM13 improves read mapping
2 and detection of clinically relevant
3 genetic variation in the Swedish
4 population

5 Daniel Schmitz¹, Adam Ameer¹ and Åsa Johansson¹

6 ¹Department of Immunology, Genetics and Pathology, Science for Life Laboratory, Uppsala
7 University, Box 815, 751 08 Uppsala, Sweden

8 Corresponding Authors: Daniel Schmitz <Daniel.schmitz@igp.uu.se> and

9 Åsa Johansson <asa.johansson@igp.uu.se>. Department of Immunology, Genetics and
10 Pathology, Science for Life Laboratory, Uppsala University, Box 815, 751 08 Uppsala,
11 Sweden

12 Running title: T2T-CHM13 improves mapping and variant calling

13

1 Abstract

2 The T2T-CHM13 reference genome, released in March 2022, fills in the 8% of the human
3 genome that were not resolved in GRCh38 and reconstructs large parts of the known genome.
4 The more accurate and complete reference genome is expected to improve the quality of read
5 mapping and variant calling. Even though whole genome sequencing (WGS)-based
6 approaches have become the golden standard in medical genetics, the extent of these benefits
7 remains unclear. In this study, we aim to evaluate alignment and variant call performance
8 with T2T-CHM13 as a reference using a cross-sectional Swedish cohort (SweGen)
9 comprising 1000 individuals with short-read Illumina WGS data available. Remapping and
10 variant calling was performed using the nf-core/sarek pipeline. T2T-CHM13 improved a
11 wide range of mapping and variant calling related metrics, including a higher fraction of
12 properly paired reads, lower mismatch rate, and more uniform coverage of coding regions.
13 Moreover, the fraction of ambiguous alignments was higher, reflecting segmental
14 duplications that were incorrectly collapsed in GRCh37 and GRCh38. In comparison to
15 GRCh38, we identified 10 million additional variants in the cohort, including 5.5 million
16 singletons, and observed an increased sensitivity for rare variants. SnpEff assigned impact
17 ratings of moderate or high to 13% more variants in T2T-CHM13 than GRCh38. In
18 summary, we conclude that T2T-CHM13 improves alignment metrics with higher alignment
19 quality, better variant calling performance and confidence, including for rare and deleterious
20 variants. The T2T-CHM13 genome reference thus facilitates enhanced discovery of new
21 disease-causing variation, benefiting, for example, rare-disease diagnostics.

22 Introduction

23 The first draft of the human reference genome was published by the Human Genome
24 Sequencing Consortium in 2001 (International Human Genome Sequencing Consortium

1 2001). Since then, the reference genome has been continuously improved with patch releases,
2 which add missing sequences, correct assembly errors, and augment the genome with
3 alternative haplotypes. In 2022, the Telomere-to-Telomere (T2T) Consortium published the
4 first gapless T2T assembly of a human genome (Nurk et al. 2022), and in 2023, an update
5 which includes Chromosome Y (Rhie et al. 2023). This assembly fills in more than 200 Mbp
6 of missing sequences from GRCh38, which corresponds to ca. 8% of the human genome.
7 Among the regions that are now resolved are centromeres, the short arms of the acrocentric
8 Chromosomes 13, 14, 15, 21 and 22 as well as the majority of Chromosome Y. In addition to
9 resolving these unknown regions, they reported novel findings about the structure of the
10 human genomes. The T2T Consortium discovered 3604 new genes and identified around
11 50 Mbp of segmental duplications (SDs), which make up half of all resolved gaps. Overall,
12 SDs are a common genomic feature, representing 7% of the whole human genome, which are
13 especially common on the short arms of acrocentric chromosomes (Vollger et al. 2022).
14 Moreover, the T2T consortium benchmarked their assembly's performance for variant calling
15 and found it to allow for the detection of many novel variants, especially in the previously
16 unresolved regions. They also found that their novel reference reduces the number of false-
17 positive calls mainly in protein-coding genes and they highlight the increased sensitivity for
18 detection of rare variants and singletons (Aganezov et al. 2022). The T2T-CHM13 genome
19 has already been seen use in the evaluation of novel tools (Formenti et al. 2022; Mc Cartney
20 et al. 2022), characterization of SVs (Prodanov and Bansal 2022) and epigenetics (Gershman
21 et al. 2022).

22 During recent years, whole-genome sequencing (WGS) has been introduced as a clinical
23 diagnostic tool, especially for rare disease, which affect between 3.5% and 5.9% of the
24 population (Nguengang Wakap et al. 2019) and are often caused by a single nucleotide
25 variant (SNV) or a structural variant (SV). However, a minority of the rare disease patients

1 receive a definite diagnosis and often must go through a so-called “diagnostic odyssey” of
2 repeated visits to specialists before a conclusive result can be obtained (Molster et al. 2016;
3 Yan et al. 2020). A recent meta-analysis found WGS to significantly improve diagnostic
4 yield over other established methods, including whole-exome sequencing (WES), but still
5 estimated it to only be around 38.6% (Nurchis et al. 2023). This contributes to the increasing
6 pool of evidence that variants in non-coding regions contribute to rare-disease development
7 (Whiffin et al. 2020; Turro et al. 2020). Consequently, even with WGS as a first-line test, the
8 fraction of patients receiving a diagnosis remains low. Possible explanations for this low
9 yield could be insufficient annotation of non-coding variants, errors in the analysis pipeline,
10 including mapping and variant calling, and misassemblies of the reference genome. These
11 issues could be addressed by applying an improved reference genome but to which extent
12 T2T-CHM13 improves detection of rare and deleterious variants over the established
13 reference genomes GRCh37 and GRCh38 remains to be investigated.

14 SweGen is a cross-sectional cohort from Sweden providing a resource of genetic variability
15 in the local population (Ameur et al. 2017). It consists of 1000 individuals who have
16 undergone whole-genome sequencing (WGS). Since its inception in 2017, SweGen has
17 facilitated clinical genomics research and diagnostics by providing a representative set of
18 genetic variation in healthy individuals and, for example, by serving as a control cohort in
19 SWEDEGENE, a study that has identified genetic variants causing serious adverse drug
20 reactions (Hallberg et al. 2020). In another study, transposable elements (TEs) in protein-
21 coding genes have been characterized in SweGen and used to diagnose two unsolved cases of
22 rare diseases (Bilgrav Saether et al. 2023). SV frequency data from SweGen has been used to
23 map complex chromosomal rearrangements (Eisfeldt et al. 2019). Several studies have used
24 the frequency data from SweGen for identifying potentially pathogenic rare mutations usually
25 not found in healthy individuals e.g., to identify germline mutations associated with breast

1 cancer risk (Helgadóttir et al. 2021) and somatic variation in relapsed pediatric acute
2 lymphoblastic leukemia (Sayyab et al. 2021).

3 In this study, we aim to evaluate the improvement of mapping short read WGS data to the
4 new T2T-CHM13 reference compared to the reference genomes that are currently used in
5 clinical decision tools. We remapped the WGS data from 1000 individuals from SweGen and
6 compared the mapping, variant calling, and annotation performance as well as extrapolated
7 from what impact a transition from the current reference to the T2T-CHM13 will have for
8 clinical diagnostics. We also provide an open resource with updated summary data for
9 variants in the SweGen cohort based on the T2T-CHM13 reference to be used for research
10 and clinical applications.

11 Results

12 We obtained the original WGS alignments from SweGen, which were BAM files mapped to
13 GRCh37 and contained all generated reads, including those that were not mapped. We
14 reanalyzed the SweGen dataset with the freely available nf-core/sarek pipeline using BWA-
15 MEM to T2T-CHM13 v2.0 (Garcia et al. 2020a; Md et al. 2019). We performed variant
16 calling using GATK HaplotypeCaller in joint germline calling mode (Van der Auwera and
17 O'Connor 2020). The same data had also been analyzed using GRCh38.p13 so we were able
18 to leverage the results from that analysis for comparison with our results.

19 T2T-CHM13 improves quality of alignments

20 In comparison to GRCh37, the number of mapped reads increased by, on average, 558,955
21 reads per individual, which corresponded to 0.0692% (Figure 1A,B, Supplemental Table S1).
22 On the other hand, fewer reads were mapped when comparing to GRCh38, with a mean
23 difference of 128,094 per individual, corresponding to 0.0158% of all reads (Figure 1A,B,
24 Supplemental Table S1). However, we observed that a larger portion of read pairs were

1 properly paired according to BWA-MEM when using T2T-CHM13 compared to the other
2 two references (Figure 1C, Supplemental Table S1). BWA-MEM considers a read pair
3 properly paired if it consists of two reads where one mate is mapped in forward and the other
4 in reverse direction, both mates' read directions point towards each other and the distance
5 between mates does not deviate too much from the mean insert size. In this case, this
6 threshold is approximately six to seven standard deviations from the mean insert size. On
7 average, 99.4% of reads per individual were paired properly (IQR: 99.4% – 99.5%) when
8 using T2T-CHM13 (Figure 1C). In comparison, 97.9% of read pairs had proper orientation
9 when using GRCh37 and 98.2% when using GRCh38. Furthermore, the alignments to T2T-
10 CHM13 showed a lower per-sample mismatch rate (Figure 1D, Supplemental Table S1). The
11 mean mismatch rate was 0.54% (IQR: 0.48% – 0.58%) for T2T-CHM13, compared to 0.78%
12 (IQR: 0.73% – 0.82%) for GRCh37 and 0.65% (IQR: 0.60% – 0.69%) for GRCh38. This
13 means that reads were mapped to T2T-CHM13 with higher accuracy.

14 [Alignments previously considered unique are ambiguous](#)

15 We noticed an increase in the number of alignments that received mapping quality 0 (MQ0)
16 i.e., they mapped equally well to more than one position (Figure 1E) as well as a higher
17 percentage of mapped reads with MQ0 (Figure 1F) with the T2T-CHM13 reference
18 (Supplemental Table S1). MQ0 reads are also significantly overrepresented in SDs (χ^2 p-
19 value $< 2.2 \times 10^{-16}$). Overall, 7.98% of all reads were assigned MQ0 but among those that
20 map to SDs, 31.9% were assigned MQ0. This is likely a result of T2T-CHM13 resolving
21 many SDs that appeared as unique regions in GRCh37 and GRCh38 as well as low-
22 complexity regions not included in GRCh37 and GRCh38.

1 Coverage of genes is more uniform

2 To investigate the coverage of genes in our sequencing data and assess potential
3 improvements to detection of functionally relevant variants, we calculated the read depth for
4 all non-overlapping 500-bp bins in gene regions annotated by the T2T Consortium. While we
5 did not observe a change in the mean coverage in genes, there was a reduction in the per-
6 individual standard deviation of the read counts compared to GRCh37 and GRCh38,
7 indicating a more uniform coverage of genes when aligning to T2T-CHM13 (Figure 2).

8 Reads from unplaced contigs map mostly to acrocentric short arms

9 We investigated the mapping positions of reads which were mapped to unplaced contigs in
10 GRCh38. Virtually all reads (99.99974%) that uniquely mapped to unplaced contigs in
11 GRCh38 could be mapped to canonical chromosomes in T2T-CHM13. The short arms of
12 Chromosomes 13, 14, 15, 21, 22 and Y were overrepresented among the regions where these
13 reads mapped, accounting for 84% of reads (Figure 3, Supplemental Table S2). However, we
14 could only assign a unique canonical chromosome to five unplaced contigs out of 126. Reads
15 from the remaining contigs mapped to positions on multiple chromosomes.

16 More rare variants with T2T-CHM13

17 To assess the general characteristics of the T2T-CHM13 variant calls, we collected statistics
18 about number, frequencies, and types on the cohort level. Overall, we detected 47,744,487
19 unique high-quality variants in the cohort, of which 19,077,986 (40%) were singletons; an
20 increase from GRCh37 and GRCh38 (Table 1, Figure 4A). This difference is not explained
21 by the higher number of resolved bases, as the T2T-CHM13 call set shows a higher density
22 of variants overall (15.3 variants/kbp) and singletons in particular (6.12 singletons/kbp) than
23 the older call sets (Table 1). Moreover, we observed higher numbers of rare (allele frequency
24 (AF) < 1%) and low-frequency variants (AF < 5%) (Table 1, Figure 4C, Supplemental Table

1 S3). On the other hand, the number of common variants ($AF \geq 5\%$) remained roughly the
2 same, and the number of non-reference alleles with an $AF > 50\%$ was generally lower with
3 T2T-CHM13 than the other references (Figure 4C). SweGen was established as a resource to
4 capture genetic variability within the Swedish cohort, and it has been shown that these
5 samples cluster closely with other European populations, although they are distinct from
6 central European populations (Ameur et al. 2017). Since T2T-CHM13 matches European
7 ancestries more closely than GRCh37 and GRCh38 (Aganezov et al. 2022), we would expect
8 to observe fewer variants where the non-reference allele is the major allele in our cohort
9 when T2T-CHM13 is used as the reference. Noticeable is also the peak of variants with $AF \approx$
10 0.5. with GRCh37 and GRCh38 (Figure 4C) which was not seen with T2T-CHM13. When
11 removing variants that deviated strongly from Hardy-Weinberg equilibrium ($P < 10^{-20}$), this
12 peak disappeared (Supplemental Figure S1). This points to it being caused by previously
13 unresolved segmental duplications, that give rise to stretches of heterozygous variant calls in
14 GRCh37 and GRCh38.

15 Proportions of SNVs and indels remain unchanged

16 To assess whether T2T-CHM13 affected our power to call SNVs or indels, we compared the
17 numbers and proportions of these variant types between references and in previously
18 unresolved regions. Around 87.8% of all called variants were SNVs, which is similar to what
19 was found for GRCh37 (86.7%) and GRCh38 (87.3%) (Figure 4A). The overall fraction of
20 indels remained roughly the same as well, constituting 13.9% of variants in T2T-CHM13,
21 14.3% in GRCh37 and 13.7% in GRCh38. In general, low-frequency variants skewed more
22 towards indels and mixed variation (i.e. indels and SNVs at the same site) than common
23 variants in T2T-CHM13 (Figure 4D). Indels were generally skewed towards deletions in all
24 assemblies. Overall, more bases were deleted than inserted, resulting in an excess of deleted
25 bases of 11,028,492 bases in T2T-CHM13, which was an increase from GRCh37 (1,065,087)

1 and GRCh38 (816,639) (Supplemental Figure S2). This excess was also visible in the number
2 of 1-bp deletions, which accounted for 1,425,614 variants. 729,804 of these deletions were
3 situated in genes, independent of predicted impact. 1-bp deletions were more likely to be
4 singletons than variants in general; 43.4% (619,334) were private to one individual, opposed
5 to 40% overall; but less likely to be fixed in the cohort (0.004% (170) vs 0.036% (16967)
6 overall).

7 Higher numbers of loss-of-function variants and variants with predicted 8 functional impact with T2T-CHM13

9 To confirm our previous assessment of improved detection of variants with functional impact
10 thanks to better gene coverage, we predicted variant effects using SnpEff, including loss of
11 function (LoF). We observed an increased number of variants with an impact rating, i.e., a
12 prediction of having impact on the protein function by, for example, causing truncation, loss
13 of function or triggering nonsense mediated decay (Figure 4B). Overall, using T2T-CHM13
14 we identified 474,050 variants with an impact rating of “low”, “moderate” or “high”,
15 corresponding to approximately 0.993% of all called variants. For comparison, 354,087
16 variants called with GRCh37 (0.996%) and 359,977 called with GRCh38 (0.949%) received
17 an impact rating. Among the variants with an impact rating with T2T-CHM13, 48.4% and
18 45.2% of variants with were predicted to have moderate and low impact respectively and the
19 remaining 6.4% of were rated high impact (Figure 4B). Our results thus suggest that
20 improved gene coverage facilitates detection of functionally relevant variation.

21 We identified 17,446 LoF variants using T2T-CHM13, an increase from GRCh37 and
22 GRCh38 (Table 1). Of these, 17,297 variants were annotated to 8,287 genes with existing
23 ClinVar annotations. Most affected genes were shared among all three references (Figure 5).
24 Across all three references, *MUC4*, which encodes mucin 4, was the gene with the highest

1 number of LoF variants (Table 2, Supplemental Table S4). However, the majority of LoF
2 variants in *MUC4* fall into a 48-bp variable tandem repeat in exon 2, which is problematic for
3 short reads and might have led to an inflated number of variants. Tandem repeats are a
4 general characteristic of the mucin gene family, which made up five of the ten top affected
5 genes across all assemblies (Ferez-Vilar and Hill 1999). Additionally, the *HLA-DRB1*,
6 belonging to the human leukocyte antigen (HLA) gene family, was present across all
7 references. Another *HLA* gene, *HLA-DRB5*, was among the top ten for GRCh37 and GRCh38
8 but not T2T-CHM13. However, as with the mucin gene family, *HLA* genes tend to be
9 difficult to analyze with short reads since they are highly polymorphic.

10 With these considerations in mind, we curated a random list of ten novel LoF variants which
11 were called in genes which had no such variants in GRCh37 and GRCh38. We could confirm
12 the existence of six but not the other four (Supplemental Figure S3). Previous studies have
13 suggested that an increased number of high-impact predictions are likely to represent artifacts
14 arising from lower quality or incomplete assemblies (Cooper and Shendure 2011). It is
15 therefore plausible that some of the novel LoF variants we did not curate are artifacts rather
16 than true LoF variants, particularly given the improved quality and completeness of the T2T-
17 CHM13 assembly. Moreover, we validated these LoF variants in two individuals
18 (SweGen0945 and SweGen0970) with available *de-novo* assemblies from 75× PacBio
19 continuous long-read (CLR) sequencing, where we could confirm 15 of 130 (11.5%) and 16
20 of 139 (11.5%) variants, respectively (Supplemental Table S5a,b). For a third individual
21 (SweGen0969), we could make use of available 35× WGS PacBio CLR data and confirm 21
22 of 143 variants (14.7%) (Supplemental Table S5c).

23 Among the top ten affected genes in T2T-CHM13, only four did not belong to the
24 aforementioned gene families. Melanoma antigen gene C1 (*MAGEC1*) has close links to

1 melanoma and myeloma development (Caballero et al. 2010; Jungbluth et al. 2005). AHNAK
2 Nucleoprotein 2 (AHNAK2) may play a role in calcium signalling and is up regulated in lung
3 cancer (Liu et al. 2020; Komuro et al. 2004). Nascent polypeptide associated complex
4 (NACA) binds to newly created polypeptides to prevent wrong translocation to the
5 endoplasmatic reticulum and plays a role in bone-formation, red blood cell differentiation and
6 has been linked to dermatitis (Lauring et al. 1995; Natter et al. 1998; Lopez et al. 2005).
7 Immunoglobulin-Like And Fibronectin Type III Domain-Containing Protein 1 (IGFN1) may
8 play a role in skeletal muscle development and has been associated with polypoidal choroidal
9 vasculopathy (Wen et al. 2018; Cracknell et al. 2020).

10 **Previously unresolved regions are not enriched for singletons**

11 To assess whether the observed increase in singleton calls was driven by regions that were
12 unresolved in GRCh37 and GRCh38, we collected general variant call statistics restricted to
13 these regions. There were 8,321,170 variants in regions unresolved in GRCh37 and in regions
14 which were unresolved in GRCh38, we identified 7,826,115 variants (Table 3). Within these
15 regions, we observed a higher density of variants (unresolved in GRCh37: 30.9 variants /
16 kbp, unresolved in GRCh38: 31.1 variants / kbp) than in the remaining genome (13.9 variants
17 / kbp). By that, we can conclude that these variants contribute to a majority, but not all, of the
18 excess in the number of variants called with T2T-CHM13. This was also apparent in unique
19 regions syntenic between GRCh38 and T2T-CHM13, where T2T-CHM13 showed a higher
20 variant density (13.9 variants / kbp) than GRCh38 (12.5 variants / kbp).

21 Variants in regions that were previously unresolved in GRCh37 and GRCh38 were more
22 biased towards SNVs (which made up 92.6% and 92.8% of variants there, respectively) than
23 other regions (Figure 6A). However, the density of deletions was higher in these regions
24 (unresolved in GRCh37: 2.45 deletions / kbp, unresolved in GRCh38: 2.46 deletions / kbp)

1 than in the rest of the genome (1.65 deletions / kbp). We detected an excess of 4,705,752
2 deleted bases in regions which were unresolved in GRCh38, (4,947,747 bp in regions not in
3 GRCh37) despite only contributing 13% of all indels. The excess in the regions syntenic with
4 GRCh38 is still higher, however, with 6,322,740 bp more bases deleted than inserted.

5 The fraction of singletons was lower than in the genome in general (χ^2 p value $< 2.2 \times$
6 10^{-16}) (Table 3), suggesting that the newly resolved regions do not contribute
7 disproportionately to the number of singletons despite a higher density of singletons (Table
8 3). However, they were significantly enriched for rare variants (χ^2 p value $< 2.2 \times 10^{-16}$)
9 and contributed the majority of the excess of rare alleles that were not observed in GRCh37
10 (61.1% of 12,594,438) or in GRCh38 (68.6% of 10,661,072).

11 Previously unresolved regions contain a low rate of variants with predicted 12 impact on a protein

13 Focusing on previously unresolved regions, we found that they had a much lower fraction of
14 variants with predicted impact on a protein according to SnpEff. The fractions of variants
15 with at least low impact were significantly lower than those of all variants called with T2T-
16 CHM13. In the regions unresolved in GRCh37, there were 3,823 variants with a rating of at
17 least low impact, corresponding to 0.045% (χ^2 P-value $< 2.2 \times 10^{-16}$). In the parts of the
18 genome, which were unresolved in GRCh38, there were 1,922 variants with assigned impact
19 rating, which corresponded to 0.025% (χ^2 P-value $< 2.2 \times 10^{-16}$). The distribution of impact
20 ratings among these variants also differed slightly between regions that were unresolved in
21 GRCh37 compared to GRCh38, with a larger fraction of high-impact variants in regions that
22 were unresolved in GRCh37 (Figure 6B). Together, this suggests that variations in the
23 previously unresolved regions generally have lower impact and might therefore be under
24 purifying selective pressure. However, this difference could also be due to a lower density of

1 genes or a lack of annotations compared to regions that have been available to functional
2 annotation before.

3 Proportion of SNVs of per-individual variants decreases with use of T2T- 4 CHM13

5 Considering the differences we observed on the cohort level between GRCh37, GRCh38 and
6 T2T-CHM13, we also investigated how variant calls differed at the individual level. Despite
7 the higher total number of called variants in the cohort with T2T-CHM13, we did not observe
8 this increase in the number of variants called per individual (Figure 7A, Table 4). After
9 limiting the comparison to syntenic regions with high uniqueness in GRCh38 and T2T-
10 CHM13, we observed a slight decrease in the number of variants per individual (mean
11 difference = -212,696.8, p value $< 2.2 \times 10^{-16}$) and their density (mean difference = -0.039
12 variants / kbp, p value $< 2.2 \times 10^{-16}$) (Supplemental Table S6).

13 We also assessed whether the proportions of SNVs and indels changed when calling variants
14 with T2T-CHM13. The number of SNVs per individual decreased and the number of indels
15 increased slightly. (Figure 7B, Table 4). These changes also became apparent in the fraction
16 of SNVs and indels in all called variants per individual (Figure 7C). The mean fraction of
17 SNVs decreased whereas we observed an increased proportion of indels. This discrepancy
18 can be explained by the previously described observation that the fraction of indels tends to
19 be higher among rare variants and that we detected fewer high frequency but a larger number
20 of rare variants. The called indels also tended to skew more towards deletions than the rest of
21 the genome.

22 Increased power to detect rare and loss-of-function variants per individual

23 The number of singletons called per individual was higher with T2T-CHM13 (Table 4,
24 Figure 7D). There were four individuals with singleton counts of more than three standard

1 deviations above the mean, which we considered outliers. These four individuals were
2 different from the four individuals with the largest number of singletons with GRCh38.
3 However, they made up four of the six individuals with the lowest fraction of properly paired
4 reads and showed anomalies in certain other QC metrics, such as the number of read pairs on
5 different chromosomes. Considering that this was consistent across references, we believe
6 these high numbers of singletons to be caused by technical issues with these four samples,
7 which had not been identified before (Supplemental Table S7). While these individuals
8 accounted for approximately 645,000 singletons, they did not drive the observed increase in
9 per-individual singletons, as every individual in the cohort had an increased number of
10 singleton variants (Supplemental Figure S4). Furthermore, we observed a slightly increased
11 number of variants with predicted LoF effects per individual over GRCh38 (mean difference
12 = 4.778, p value = 0.000722), which agreed with the higher number in the cohort as a whole
13 and showed the benefit of improved sequencing coverage in gene regions (Figure 7E, Table
14 4). However, in comparison to GRCh37, the per-individual number of LoF variants actually
15 decreased (mean difference = -174.429, p value < 2.2×10^{-16}). This excess can be explained
16 by variants whose reference alleles were the LoF alleles, which were corrected in GRCh38
17 and appeared as common or fixed variants in SweGen. For instance, there were 194 LoF
18 variants in the GRCh37 call set with AF > 90%, as compared to 83 in GRCh38 and 24 in
19 T2T-CHM13.

20 Discussion

21 In this study, we have shown that using the new T2T-CHM13 reference can improve
22 alignments, variant calls, and sensitivity for rare, singleton and LoF variants over GRCh37
23 and GRCh38, which are currently in widespread use, including in genes with clinical
24 relevance. This is an important prospect for clinical genetics as rare, ultra-rare and *de-novo*

1 variants are often the cause of rare disease. Continuing efforts are aiming to bring
2 sequencing, and especially WGS, into the diagnostic routine because of its power to detect
3 these variants (Tesi et al. 2023). For instance, the UK 100,000 Genomes Project found that
4 genome sequencing significantly increases the diagnostic yield of rare diseases (Smedley et
5 al. 2021). The Genomics England project is finding an increasing number of rare variants
6 with clinical significance and enabling diagnoses of previously undiagnosed patients (Joyce
7 et al. 2022; Han et al. 2022; Kojic et al. 2023; Vadgama et al. 2022). By employing T2T-
8 CHM13 as a reference for first-line WGS in the clinic, the number of people receiving a
9 diagnosis for a rare disease is likely to improve further. Moreover, it is becoming increasingly
10 clear that reanalysis of sequencing data can lead to better diagnoses and improved patient
11 outcomes (Machini et al. 2019; Kingsmore et al. 2022). Through reanalysis, it is possible to
12 detect new causal variations using improved tools and annotation and provide treatment to
13 affected individuals. Considering this, *FixItFelix*, which realigns reads falling into medically
14 relevant genes using a modified GRCh38 reference with certain genes replaced by their T2T-
15 CHM13 assemblies, was developed (Behera et al. 2023). This tool only requires reanalysis of
16 smaller regions and, as such, provides a faster turnaround. However, it requires limiting the
17 scope of the analysis to manually curated regions and, therefore, is not suitable for a whole-
18 genome study. As such, a re-evaluation of inconclusive clinical samples using T2T-CHM13,
19 including reads that failed to map to canonical chromosomes, can provide valuable diagnoses
20 to patients.

21 Rare variants are also becoming more prominent in genetic epidemiology. GWAS have only
22 been able to explain part of the observed heritability. Recent studies showed that part of this
23 missing heritability can be explained by rare variants and successfully quantified their effect
24 on complex traits and disease risk (Höglund et al. 2019; Kierczak et al. 2022; Wang et al.
25 2021; The UK10K Consortium 2015). Consequently, the possibility of detecting more rare

1 variants, especially in regions which lacked an assembly but are functionally relevant, such as
2 centromeres, promises to close the gap in quantifying the missing heritability.

3 The resources provided by the T2T Consortium, such as lifted over versions of dbSNP, allow
4 for the adoption of existing workflows to the new reference as well as alignments and variant
5 calls generated using it. However, there is currently a lack of original annotation and
6 databases based on T2T-CHM13. Reliance on annotations generated using older references
7 and lifted over to T2T-CHM13 might negatively impact results, as not all variants can be
8 lifted over. We expect this issue to diminish as more resources, such as this one, become
9 available. Another limitation was that we were unable to unambiguously elucidate why four
10 individuals were outliers with regards to number of singletons. As these four individuals
11 showed some anomalies regarding some QC measures, a possible explanation might have
12 been a higher sensitivity to sample quality or technical issues when aligning to T2T-CHM13.
13 However, we could not conclusively confirm this hypothesis at this point. Furthermore, as
14 the version of T2T-CHM13 we used included unmasked pseudoautosomal regions of
15 Chromosome Y, variants called in these regions suffer from lower accuracy. However, there
16 were only 415 variants with a predicted functional impact, so these issues did not
17 significantly impact our results (Supplemental Table S8). In addition, we did not restrict our
18 analyses to short-read accessible regions. Even though variants outside these regions have a
19 higher risk of being false positives, general trends regarding variant impacts and types within
20 them, agree with our findings from the whole genome (Supplemental Table S9, Supplemental
21 Figure S6).

22 We observed, in general, improvements in mapping and variant calling performance similar
23 to those presented by the T2T Consortium in their investigation (Aganezov et al. 2022).
24 These included an increase in quality of alignments as both the fraction of properly paired

1 and oriented read pairs increased as well as the mismatch rate decreased. The newly
2 assembled regions also allowed us to identify reads that appeared to map uniquely in older
3 references, but that actually map to duplicated regions which were previously unknown. This
4 allowed for greater confidence in the variant calls since spurious variants, which were called
5 as common alleles but originated from separate copies of the same SD, could be removed.
6 Furthermore, common variants might have disappeared because of altered alleles in T2T-
7 CHM13.

8 On the other hand, we identified generally more variants overall, in terms of absolute number
9 and density, including more deletions than with previous references, which might be a sign of
10 reference bias associated with the increased number of duplicated sequences in T2T-CHM13.
11 This difference was also present in syntenic regions between GRCh38 and T2T-CHM13 with
12 high uniqueness, where we expected calls to be largely consistent between references.
13 However, at the individual level, we observed a slight decrease in variant number and density
14 in these regions. Therefore, we argue that the majority of variants are called consistently and
15 that the observed opposite trends of these metrics between cohort and individual level are due
16 to minor differences in the syntenic regions and the changes to variant composition
17 mentioned above.

18 The markedly elevated number of LoF variants in specific genes, many of which are
19 implicated in disease-related pathways, suggests that variation in these genes have been
20 systematically underrepresented in previous reference assemblies. This finding underscores
21 the importance of T2T-based mapping, particularly for diseases where accurate variant
22 detection in these pathways is critical for understanding pathogenesis and improving
23 diagnostic resolution. However, some additional LoF variants we identified, including in
24 syntenic regions, were technical artifacts. However, through our curation of LoF variants in

1 T2T-CHM13 in genes without any detected LoF variation in GRCh37 and GRCh38, and
2 additional validation using PacBio CLR data, which, despite its higher error rate, allowed us
3 to confirm up to 14.7% of novel LoF variants, we could show that there was, in fact, a higher
4 sensitivity for functionally relevant variation. As such, the balance between detection of
5 novel LoF variants and technical artifact is an important consideration for future projects
6 using T2T-CHM13, especially considering that, as before, LoF variants should be confirmed
7 experimentally before they can be used to inform clinical decisions.

8 In conclusion, the T2T-CHM13 reference provides a valuable resource that gives improved
9 alignments and variant calling. Based on our results we recommend, going forward, the
10 adoption of T2T-CHM13 for novel sequencing studies. Re-evaluation of sequenced samples,
11 especially inconclusive ones, using T2T-CHM13 can provide further insights into genetic
12 variation, including high-impact events, that might otherwise be missed. Considering the
13 computational complexity of remapping and recalling large cohorts, however, it can be
14 sensible to take an iterative approach. But as more resources based on T2T-CHM13 become
15 available, we can expect increased power and precision in statistical and clinical genetics
16 applications.

17 **Methods**

18 **The SweGen Cohort**

19 The SweGen cohort is a cross-sectional cohort of the Swedish population consisting of 1000
20 individuals that all passed quality control as described in the original SweGen paper (Ameur
21 et al. 2017). The majority of SweGen comprises 942 unrelated individuals from the Swedish
22 Twin Register (STR) which were selected to mirror the overall distribution of genetic
23 principal components (PCs) of all 10,000 STR individuals that underwent SNV genotyping
24 (Lichtenstein et al. 2002). The remaining 58 individuals were selected from the Northern

1 Swedish Population Health Study (NSPHS) (Igl et al. 2010). NSPHS individuals were
2 selected based on their genetic PCs, as well. Library preparation and WGS were
3 responsibility of Sweden’s National Genomics Infrastructure (NGI) in Stockholm (NGI-S)
4 and Uppsala (NGI-U). After fragmenting to 350 bp inserts and library preparation, DNA
5 samples underwent WGS on the Illumina HiSeq X platform using v2.5 chemistry (Ameur et
6 al. 2017). The 1000 individuals included in the dataset all passed QC. The original mappings
7 were obtained by aligning the raw reads to GRCh37.p13 using BWA-MEM 0.7.12. The
8 mappings underwent quality control, which was described previously, including base quality
9 recalibration and duplication marking (Ameur et al. 2017). In addition, alignments and
10 variant calls based on GRCh38.p14 including alternative haplotypes and decoys were
11 available. The final GRCh37 BAM files served as the input for our analysis pipeline as they
12 contained all reads, including unmapped ones, and allowed us to take advantage of the
13 already performed quality control.

14 **Remapping and Variant Calling**

15 We obtained the T2T-CHM13 reference genome version 2.0, a lifted over version of dbSNP
16 and a BED file of segmental duplications from the T2T Consortium’s GitHub page
17 (<https://github.com/marbl/CHM13>). This version of the reference genome contains an
18 assembly of Chromosome Y and softmasked repeats. We performed mapping and variant
19 calling using nf-core/sarek version 3.2.0, which uses GATK 4.3.0.0, with default parameters
20 corresponding to GATK best practices (Garcia et al. 2020b; DI Tommaso et al. 2017; Ewels
21 et al. 2020; McKenna et al. 2010). Shortly, the reads contained in the delivered BAM files,
22 i.e., the post-QC alignments to GRCh37p13, which contained all generated reads, even those
23 that failed original QC or were not aligned, underwent read QC using FastQC 0.11.9 and
24 fastp 0.23.2 (Chen et al. 2018; Anders 2010). Reads failed QC if more than 40% of bases had
25 a quality score below 15, or more than five base calls were indeterminate, i.e., they were

1 called as “N”. On average, 4.1% of reads (inter-quartile range (IQR): 3.6% – 4.5%) failed this
2 filter. The remaining reads were aligned using BWA-MEM version 0.7.17 (Md et al. 2019).
3 The resulting CRAM files were merged, sorted, and indexed using SAMtools 1.16.1 (Li et al.
4 2009).

5 Variants were called using GATK HaplotypeCaller with the joint germline calling feature
6 enabled. For Variant QC, GATK VariantRecalibrator was used with lifted over versions of
7 dbSNPv155, 1000 Genomes Omni 2.5 variants, as well as 1000 Genomes phase I and gold
8 standard indels. We included all variants passing the more lenient default threshold (truth
9 sensitivity 99.9% – 100%) for increased sensitivity and because those were available in the
10 call sets of the older assemblies. General statistics on the performance of the T2T-CHM13
11 reference were based on the reports automatically generated by SAMtools and BCFtools 1.17
12 as part of the sarek pipeline (Danecek et al. 2021; Li 2011). We added annotations using
13 BCFtools and as well as functional effect predictions, such as loss of function (LoF) using
14 SnpEff based on the lifted over version of dbSNP and curated transcripts from RefSeq 5.1
15 provided by the T2T Consortium (Cingolani et al. 2012a, 2012b). We created a subset of
16 putatively clinically relevant LoF variants by filtering for annotation to genes in ClinVar
17 (version 20240215) (Landrum et al. 2014). We identified variants in previously unresolved
18 regions, both in relation to GRCh37 and GRCh38, by subsetting the joint VCF with BCFtools
19 using the “CHM13 unique” track in UCSC Genome Browser (Kent et al. 2002). We obtained
20 already existing alignments to GRCh37 and GRCh38 from the SweGen dataset as well as
21 variant calls, both per-individual and joint calls for the whole cohort, generated from these
22 alignments. We generated mapping and variant reports using SAMtools and BCFtools. We
23 calculated read depth in gene regions by creating non-overlapping 500-bp bins covering these
24 regions. Additionally, we created a subset of variants expected to be consistent between
25 GRCh38 and T2T-CHM13 by subsetting for syntenic regions and the “UMAP S100” track

1 from UCSC Genome Browser, as these regions are consistent between assemblies and allow
2 for high-confidence short-read mapping (Karimzadeh et al. 2018).

3 Validation of putative LoF variants

4 We visually inspected ten randomly selected predicted LoF variants from genes that had no
5 reported LoF variants in the GRCh37 and GRCh38 call sets in IGV 2.8.13. We validated
6 putative LoF variants in the same set of genes in individuals SweGen0945 and SweGen0970
7 using available *de-novo* assemblies generated from 75× coverage PacBio continuous long
8 read (CLR) sequencing (Ameur et al. 2018). We mapped the assembled contigs to T2T-
9 CHM13 using minimap2 2.26 and called variants using the included version of paftools (Li
10 2018). We compared the variants found in the *de-novo* assemblies with those in the short
11 reads using RTG Tools vcfEval 3.12.1, restricting the analysis to regions deemed reliable for
12 variant calling by paftools i.e., covered by exactly one contig, excluding alignments shorter
13 than 10 kbp (Cleary et al. 2015).

14 We validated putative LoF variants in individual SweGen0969 using previously generated
15 PacBio CLR sequencing data (Schmitz et al. 2023). We aligned the reads to T2T-CHM13
16 using minimap2 version 2.26 and called variants using Longshot 1.0 (Edge and Bansal 2019).
17 We compared the variants found in the PacBio data with those in the short reads using RTG
18 Tools vcfEval 3.12.1, restricting the analysis to regions with coverage higher than 15× and
19 lower than 60× to exclude copy-number variable regions.

20 Ethics Declaration

21 All participants in the Swedish Twin Register gave informed consent and the study was
22 approved by the regional ethics committee in Stockholm (Regionala Etikprövningsnämnden,
23 Stockholm, dnr 2007-644-31 and 2014/521-32). Participants of the Northern Swedish
24 Population Health Study gave their informed consent and the study was approved by the

1 regional ethics committee in Uppsala (Regionala Etikprövningsnämnden, Uppsala, dnr
2 2005:325 and 2016-03-09).

3 Data Access

4 Cohort-level allele frequencies and variant annotation are publicly available on the website of
5 SweGen (<https://swefreq.nbis.se/dataset/SweGen>) and Zenodo
6 (<https://doi.org/10.5281/zenodo.13739348>). All data generated in this study have been
7 deposited in FEGA Sweden and made findable through the European Genome-Phenome
8 Archive web portal (EGA, <https://ega-archive.org>) under study number EGAS50000000906
9 (<https://ega-archive.org/studies/EGAS50000000906>).

10 Competing Interest Statement

11 The authors declare no competing interests.

12 Acknowledgements

13 AA and DS conceived the study. DS performed the analyses, interpreted the results, wrote the
14 initial manuscript and generated figures and tables under supervision of ÅJ and AA. All
15 authors discussed and interpreted the results and contributed to the final manuscript.

16 This project was funded by the Swedish Research Council (Dnr. 2019-01497). The
17 computations were enabled by resources in project sens2016003 and sens2022031 provided
18 by the National Academic Infrastructure for Supercomputing in Sweden (NAISS) at
19 UPPMAX, funded by the Swedish Research Council through grant agreement no. 2022-
20 06725.

1 References

- 2 Aganezov S, Yan SM, Soto DC, Kirsche M, Zarate S, Avdeyev P, Taylor DJ, Shafin K,
3 Shumate A, Xiao C, et al. 2022. A complete reference genome improves analysis of
4 human genetic variation. *Science (1979)* **376**. doi:10.1126/SCIENCE.ABL3533.
- 5 Ameur A, Che H, Martin M, Bunikis I, Dahlberg J, Höijer I, Häggqvist S, Vezzi F, Nordlund
6 J, Olason P, et al. 2018. De Novo Assembly of Two Swedish Genomes Reveals Missing
7 Segments from the Human GRCh38 Reference and Improves Variant Calling of
8 Population-Scale Sequencing Data. *Genes 2018, Vol 9, Page 486* **9**: 486.
9 doi:10.3390/GENES9100486.
- 10 Ameur A, Dahlberg J, Olason P, Vezzi F, Karlsson R, Martin M, Viklund J, Kähäri AK,
11 Lundin P, Che H, et al. 2017. SweGen: a whole-genome data resource of genetic
12 variability in a cross-section of the Swedish population. *European Journal of Human*
13 *Genetics 2017 25:11* **25**: 1253–1260. doi:10.1038/ejhg.2017.130.
- 14 Anders S. 2010. Babraham Bioinformatics - FastQC A Quality Control tool for High
15 Throughput Sequence Data. *Soil* **5**: <http://www.bioinformatics.babraham.ac.uk/projects/>.
- 16 Behera S, LeFaive J, Orchard P, Mahmoud M, Paulin LF, Farek J, Soto DC, Parker SCJ,
17 Smith A V, Dennis MY, et al. 2023. FixItFelix: improving genomic analysis by fixing
18 reference errors. *Genome Biol* **24**: 31. doi:10.1186/s13059-023-02863-7.
- 19 Bilgrav Saether K, Nilsson D, Thonberg H, Tham E, Ameur A, Eisfeldt J, Lindstrand A.
20 2023. Transposable element insertions in 1000 Swedish individuals. *PLoS One* **18**:
21 e0289346. doi:10.1371/JOURNAL.PONE.0289346.

- 1 Caballero OL, Zhao Q, Rimoldi D, Stevenson BJ, Svobodová S, Devalle S, Róhrig UF,
2 Pagotto A, Michielin O, Speiser D, et al. 2010. Frequent MAGE Mutations in Human
3 Melanoma. *PLoS One* **5**: e12773. doi:10.1371/JOURNAL.PONE.0012773.
- 4 Chen S, Zhou Y, Chen Y, Gu J. 2018. fastp: an ultra-fast all-in-one FASTQ preprocessor.
5 *Bioinformatics* **34**: i884–i890. doi:10.1093/BIOINFORMATICS/BTY560.
- 6 Cingolani P, Patel VM, Coon M, Nguyen T, Land SJ, Ruden DM, Lu X. 2012a. Using
7 *Drosophila melanogaster* as a model for genotoxic chemical mutational studies with a
8 new program, SnpSift. *Front Genet* **3**: 35. doi:10.3389/FGENE.2012.00035/BIBTEX.
- 9 Cingolani P, Platts A, Wang LL, Coon M, Nguyen T, Wang L, Land SJ, Lu X, Ruden DM.
10 2012b. A program for annotating and predicting the effects of single nucleotide
11 polymorphisms, SnpEff. *Fly (Austin)* **6**: 80–92. doi:10.4161/fly.19695.
- 12 Cleary JG, Braithwaite R, Gaastra K, Hilbush BS, Inglis S, Irvine SA, Jackson A, Littin R,
13 Rathod M, Ware D, et al. 2015. Comparing Variant Call Files for Performance
14 Benchmarking of Next-Generation Sequencing Variant Calling Pipelines. *bioRxiv*
15 023754. doi:10.1101/023754.
- 16 Cooper GM, Shendure J. 2011. Needles in stacks of needles: finding disease-causal variants
17 in a wealth of genomic data. *Nat Rev Genet* **12**: 628–640. doi:10.1038/nrg3046.
- 18 Cracknell T, Mannsverk S, Nichols A, Dowle A, Blanco G. 2020. Proteomic resolution of
19 IGFN1 complexes reveals a functional interaction with the actin nucleating protein
20 COBL. *Exp Cell Res* **395**: 112179. doi:10.1016/J.YEXCR.2020.112179.

- 1 Danecek P, Bonfield JK, Liddle J, Marshall J, Ohan V, Pollard MO, Whitwham A, Keane T,
2 McCarthy SA, Davies RM. 2021. Twelve years of SAMtools and BCFtools.
3 *Gigascience* **10**: 1–4. doi:10.1093/GIGASCIENCE/GIAB008.
- 4 DI Tommaso P, Chatzou M, Floden EW, Barja PP, Palumbo E, Notredame C. 2017.
5 Nextflow enables reproducible computational workflows. *Nat Biotechnol* **35**: 316–319.
6 doi:10.1038/NBT.3820.
- 7 Edge P, Bansal V. 2019. Longshot enables accurate variant calling in diploid genomes from
8 single-molecule long read sequencing. *Nature Communications* 2019 10:1 **10**: 1–10.
9 doi:10.1038/s41467-019-12493-y.
- 10 Eisfeldt J, Pettersson M, Vezzi F, Wincent J, Källér M, Gruselius J, Nilsson D, Syk Lundberg
11 E, Carvalho CMB, Lindstrand A. 2019. Comprehensive structural variation genome map
12 of individuals carrying complex chromosomal rearrangements. *PLoS Genet* **15**:
13 e1007858. doi:10.1371/JOURNAL.PGEN.1007858.
- 14 Ewels PA, Peltzer A, Fillinger S, Patel H, Alneberg J, Wilm A, Garcia MU, Di Tommaso P,
15 Nahnsen S. 2020. The nf-core framework for community-curated bioinformatics
16 pipelines. *Nature Biotechnology* 2020 38:3 **38**: 276–278. doi:10.1038/s41587-020-0439-
17 x.
- 18 Ferez-Vilar J, Hill RL. 1999. The structure and assembly of secreted mucins. *Journal of*
19 *Biological Chemistry* **274**: 31751–31754. doi:10.1074/jbc.274.45.31751.
- 20 Formenti G, Rhie A, Walenz BP, Thibaud-Nissen F, Shafin K, Koren S, Myers EW, Jarvis
21 ED, Phillippy AM. 2022. Merfin: improved variant filtering, assembly evaluation and
22 polishing via k-mer validation. *Nature Methods* 2022 19:6 **19**: 696–704.
23 doi:10.1038/s41592-022-01445-y.

- 1 Garcia M, Juhos S, Larsson M, Olason PI, Martin M, Eisfeldt J, DiLorenzo S, Sandgren J,
2 Díaz De Ståhl T, Ewels P, et al. 2020a. Sarek: A portable workflow for whole-genome
3 sequencing analysis of germline and somatic variants. *F1000Research* 2020 9:63 **9**: 63.
4 doi:10.12688/f1000research.16665.2.
- 5 Garcia M, Juhos S, Larsson M, Olason PI, Martin M, Eisfeldt J, DiLorenzo S, Sandgren J,
6 Díaz De Ståhl T, Ewels P, et al. 2020b. Sarek: A portable workflow for whole-genome
7 sequencing analysis of germline and somatic variants. *F1000Research* 2020 9:63 **9**: 63.
8 doi:10.12688/f1000research.16665.2.
- 9 Gershman A, Sauria MEG, Guitart X, Vollger MR, Hook PW, Hoyt SJ, Jain M, Shumate A,
10 Razaghi R, Koren S, et al. 2022. Epigenetic patterns in a complete human genome.
11 *Science (1979)* **376**.
12 doi:10.1126/SCIENCE.ABJ5089/SUPPL_FILE/SCIENCE.ABJ5089_MDAR_REPROD
13 UCIBILITY_CHECKLIST.PDF.
- 14 Hallberg P, Yue QY, Eliasson E, Melhus H, Ås J, Wadelius M. 2020. SWEDEGENE—a
15 Swedish nation-wide DNA sample collection for pharmacogenomic studies of serious
16 adverse drug reactions. *The Pharmacogenomics Journal* 2020 20:4 **20**: 579–585.
17 doi:10.1038/s41397-020-0148-3.
- 18 Han JH, Ryan G, Guy A, Liu L, Quinodoz M, Helbling I, Lai-Cheong JE, Barwell J, Folcher
19 M, McGrath JA, et al. 2022. Mutations in the ribosome biogenesis factor gene LTV1 are
20 linked to LIPHAK syndrome, a novel poikiloderma-like disorder. *Hum Mol Genet* **31**:
21 1970–1978. doi:10.1093/HMG/DDAB368.

- 1 Helgadóttir HT, Thutkawkorapin J, Lagerstedt-Robinson K, Lindblom A. 2021. Sequencing
2 for germline mutations in Swedish breast cancer families reveals novel breast cancer risk
3 genes. *Scientific Reports* 2021 11:1 **11**: 1–7. doi:10.1038/s41598-021-94316-z.
- 4 Höglund J, Rafati N, Rask-Andersen M, Enroth S, Karlsson T, Ek WE, Johansson Å. 2019.
5 Improved power and precision with whole genome sequencing data in genome-wide
6 association studies of inflammatory biomarkers. *Scientific Reports* 2019 9:1 **9**: 1–14.
7 doi:10.1038/s41598-019-53111-7.
- 8 Igl W, Johansson A, Gyllensten U. 2010. The Northern Swedish Population Health Study
9 (NSPHS)--a paradigmatic study in a rural population combining community health and
10 basic research. *Rural Remote Health* **10**: 1363. doi:10.22605/RRH1363.
- 11 International Human Genome Sequencing Consortium. 2001. Initial sequencing and analysis
12 of the human genome. *Nature* **409**: 860–921. doi:10.1038/35057062.
- 13 Joyce KE, Onabanjo E, Brownlow S, Nur F, Olupona K, Fakayode K, Sroya M, Thomas GA,
14 Ferguson T, Redhead J, et al. 2022. Whole genome sequences discriminate hereditary
15 hemorrhagic telangiectasia phenotypes by non-HHT deleterious DNA variation. *Blood*
16 *Adv* **6**: 3956–3969. doi:10.1182/BLOODADVANCES.2022007136.
- 17 Jungbluth AA, Ely S, DiLiberto M, Niesvizky R, Williamson B, Frosina D, Chen YT,
18 Bhardwaj N, Chen-Kiang S, Old LJ, et al. 2005. The cancer-testis antigens CT7
19 (MAGE-C1) and MAGE-A3/6 are commonly expressed in multiple myeloma and
20 correlate with plasma-cell proliferation. *Blood* **106**: 167–174. doi:10.1182/BLOOD-
21 2004-12-4931.

- 1 Karimzadeh M, Ernst C, Kundaje A, Hoffman MM. 2018. Umap and Bimap: quantifying
2 genome and methylome mappability. *Nucleic Acids Res* **46**: e120–e120.
3 doi:10.1093/NAR/GKY677.
- 4 Kent WJ, Sugnet CW, Furey TS, Roskin KM, Pringle TH, Zahler AM, Haussler and D. 2002.
5 The Human Genome Browser at UCSC. *Genome Res* **12**: 996–1006.
6 doi:10.1101/gr.229102.
- 7 Kierczak M, Rafati N, Höglund J, Gourelé H, Lo Faro V, Schmitz D, Ek WE, Gyllensten U,
8 Enroth S, Ekman D, et al. 2022. Contribution of rare whole-genome sequencing variants
9 to plasma protein levels and the missing heritability. *Nat Commun* **13**: 2532.
10 doi:10.1038/s41467-022-30208-8.
- 11 Kingsmore SF, Smith LD, Kunard CM, Bainbridge M, Batalov S, Benson W, Blincow E,
12 Caylor S, Chambers C, Del Angel G, et al. 2022. A genome sequencing system for
13 universal newborn screening, diagnosis, and precision medicine for severe genetic
14 diseases. *Am J Hum Genet* **109**: 1605–1619. doi:10.1016/j.ajhg.2022.08.003.
- 15 Kojic M, Abbassi NEH, Lin TY, Jones A, Wakeling EL, Clement E, Nakou V, Singleton M,
16 Dobosz D, Kaliakatsos M, et al. 2023. A novel ELP1 mutation impairs the function of
17 the Elongator complex and causes a severe neurodevelopmental phenotype. *Journal of*
18 *Human Genetics* 2023 68:7 **68**: 445–453. doi:10.1038/s10038-023-01135-3.
- 19 Komuro A, Masuda Y, Kobayashi K, Babbitt R, Gunel M, Flavell RA, Marchesi VT. 2004.
20 The AHNAKs are a class of giant propeller-like proteins that associate with calcium
21 channel proteins of cardiomyocytes and other cells. *Proceedings of the National*
22 *Academy of Sciences* **101**: 4053–4058. doi:10.1073/PNAS.0308619101.

- 1 Landrum MJ, Lee JM, Riley GR, Jang W, Rubinstein WS, Church DM, Maglott DR. 2014.
2 ClinVar: public archive of relationships among sequence variation and human
3 phenotype. *Nucleic Acids Res* **42**: D980–D985. doi:10.1093/nar/gkt1113.
- 4 Lauring B, Sakai H, Kreibich G, Wiedmann M. 1995. Nascent polypeptide-associated
5 complex protein prevents mistargeting of nascent chains to the endoplasmic reticulum.
6 *Proceedings of the National Academy of Sciences* **92**: 5411–5415.
7 doi:10.1073/PNAS.92.12.5411.
- 8 Li H. 2011. A statistical framework for SNP calling, mutation discovery, association mapping
9 and population genetical parameter estimation from sequencing data. *Bioinformatics* **27**:
10 2987–2993. doi:10.1093/BIOINFORMATICS/BTR509.
- 11 Li H. 2018. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics* **34**:
12 3094–3100. doi:10.1093/BIOINFORMATICS/BTY191.
- 13 Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G, Durbin
14 R. 2009. The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**:
15 2078–2079. doi:10.1093/BIOINFORMATICS/BTP352.
- 16 Lichtenstein P, De Faire U, Floderus B, Svartengren M, Svedberg P, Pedersen NL. 2002. The
17 Swedish Twin Registry: a unique resource for clinical, epidemiological and genetic
18 studies. *J Intern Med* **252**: 184–205. doi:10.1046/J.1365-2796.2002.01032.X.
- 19 Liu G, Guo Z, Zhang Q, Liu Z, Zhu D. 2020. Ahnak2 promotes migration, invasion, and
20 epithelial-mesenchymal transition in lung adenocarcinoma cells via the $\text{tgf-}\beta\text{/smad3}$
21 pathway. *Onco Targets Ther* **13**: 12893–12903. doi:10.2147/OTT.S281517.

- 1 Lopez S, Stuhl L, Fichelson S, Dubart-Kupperschmitt A, St Arnaud R, Galindo JR, Murati A,
2 Berda N, Dubreuil P, Gomez S. 2005. NACA is a positive regulator of human erythroid-
3 cell differentiation. *J Cell Sci* **118**: 1595–1605. doi:10.1242/JCS.02295.
- 4 Machini K, Ceyhan-Birsoy O, Azzariti DR, Sharma H, Rossetti P, Mahanta L, Hutchinson L,
5 McLaughlin H, Green RC, Lebo M, et al. 2019. Analyzing and Reanalyzing the
6 Genome: Findings from the MedSeq Project. *Am J Hum Genet* **105**: 177–188.
7 doi:10.1016/j.ajhg.2019.05.017.
- 8 Mc Cartney AM, Shafin K, Alonge M, Bzikadze A V., Formenti G, Functammasan A, Howe
9 K, Jain C, Koren S, Logsdon GA, et al. 2022. Chasing perfection: validation and
10 polishing strategies for telomere-to-telomere genome assemblies. *Nature Methods* 2022
11 *19*:6 **19**: 687–695. doi:10.1038/s41592-022-01440-3.
- 12 McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernytsky A, Garimella K,
13 Altshuler D, Gabriel S, Daly M, et al. 2010. The Genome Analysis Toolkit: a
14 MapReduce framework for analyzing next-generation DNA sequencing data. *Genome*
15 *Res* **20**: 1297–1303. doi:10.1101/GR.107524.110.
- 16 Md V, Misra S, Li H, Aluru S. 2019. Efficient architecture-aware acceleration of BWA-
17 MEM for multicore systems. *Proceedings - 2019 IEEE 33rd International Parallel and*
18 *Distributed Processing Symposium, IPDPS 2019* 314–324.
19 doi:10.1109/IPDPS.2019.00041.
- 20 Molster C, Urwin D, Di Pietro L, Fookes M, Petrie D, Van Der Laan S, Dawkins H. 2016.
21 Survey of healthcare experiences of Australian adults living with rare diseases.
22 *Orphanet J Rare Dis* **11**: 1–12. doi:10.1186/S13023-016-0409-Z/TABLES/9.

- 1 Natter S, Seiberler S, Hufnagl P, Binder BR, Hirschl AM, Ring J, Abeck D, Schmidt T,
2 Valent P, Valenta R. 1998. Isolation of cDNA clones coding for IgE autoantigens with
3 serum IgE from atopic dermatitis patients. *FASEB J* **12**: 1559–1569.
4 doi:10.1096/FASEBJ.12.14.1559.
- 5 Nguengang Wakap S, Lambert DM, Olry A, Rodwell C, Gueydan C, Lanneau V, Murphy D,
6 Le Cam Y, Rath A. 2019. Estimating cumulative point prevalence of rare diseases:
7 analysis of the Orphanet database. *European Journal of Human Genetics* 2019 28:2 **28**:
8 165–173. doi:10.1038/s41431-019-0508-0.
- 9 Nurchis MC, Altamura G, Riccardi MT, Radio FC, Chillemi G, Bertini ES, Garlasco J,
10 Tartaglia M, Dallapiccola B, Damiani G. 2023. Whole genome sequencing diagnostic
11 yield for paediatric patients with suspected genetic disorders: systematic review, meta-
12 analysis, and GRADE assessment. *Archives of Public Health* **81**: 93.
13 doi:10.1186/s13690-023-01112-4.
- 14 Nurk S, Koren S, Rhie A, Rautiainen M, Bzikadze A V., Mikheenko A, Vollger MR,
15 Altemose N, Uralsky L, Gershman A, et al. 2022. The complete sequence of a human
16 genome. *Science (1979)* **376**: 44–53. doi:10.1126/science.abj6987.
- 17 Prodanov T, Bansal V. 2022. Robust and accurate estimation of paralog-specific copy
18 number for duplicated genes using whole-genome sequencing. *Nature Communications*
19 *2022 13:1* **13**: 1–12. doi:10.1038/s41467-022-30930-3.
- 20 Rhie A, Nurk S, Cechova M, Hoyt SJ, Taylor DJ, Altemose N, Hook PW, Koren S,
21 Rautiainen M, Alexandrov IA, et al. 2023. The complete sequence of a human Y
22 chromosome. *Nature* **621**: 344–354. doi:10.1038/s41586-023-06457-y.

- 1 Sayyab S, Lundmark A, Larsson M, Ringnér M, Nystedt S, Marincevic-Zuniga Y, Tamm KP,
2 Abrahamsson J, Fogelstrand L, Heyman M, et al. 2021. Mutational patterns and clonal
3 evolution from diagnosis to relapse in pediatric acute lymphoblastic leukemia. *Scientific*
4 *Reports* 2021 11:1 **11**: 1–17. doi:10.1038/s41598-021-95109-0.
- 5 Schmitz D, Li Z, Lo Faro V, Rask-Andersen M, Ameer A, Rafati N, Johansson Å. 2023.
6 Copy number variations and their effect on the plasma proteome. *Genetics* **225**.
7 doi:10.1093/GENETICS/IYAD179.
- 8 Smedley D, Smith KR, Martin A, Thomas EA, McDonagh EM, Cipriani V, Ellingford JM,
9 Arno G, Tucci A, Vandrovcova J, et al. 2021. 100,000 Genomes Pilot on Rare-Disease
10 Diagnosis in Health Care — Preliminary Report. *New England Journal of Medicine* **385**:
11 1868–1880. doi:10.1056/NEJMoa2035790.
- 12 Tesi B, Boileau C, Boycott KM, Canaud G, Caulfield M, Choukair D, Hill S, Spielmann M,
13 Wedell A, Wirta V, et al. 2023. Precision medicine in rare diseases: What is next? *J*
14 *Intern Med* **294**: 397–412. doi:10.1111/joim.13655.
- 15 The UK10K Consortium. 2015. The UK10K project identifies rare variants in health and
16 disease. *Nature* **526**: 82–90. doi:10.1038/nature14962.
- 17 Turro E, Astle WJ, Megy K, Gräf S, Greene D, Shamardina O, Allen HL, Sanchis-Juan A,
18 Frontini M, Thys C, et al. 2020. Whole-genome sequencing of patients with rare
19 diseases in a national health system. *Nature* 2020 583:7814 **583**: 96–102.
20 doi:10.1038/s41586-020-2434-2.
- 21 Vadgama N, Ameen M, Sundaram L, Gaddam S, Gifford C, Nasir J, Karakikes I. 2022. De
22 novo and inherited variants in coding and regulatory regions in genetic

- 1 cardiomyopathies. *Hum Genomics* **16**: 1–20. doi:10.1186/S40246-022-00420-
2 0/FIGURES/4.
- 3 Van der Auwera GA, O'Connor BD. 2020. *Genomics in the cloud: using Docker, GATK, and*
4 *WDL in Terra*. O'Reilly Media.
- 5 Vollger MR, Guitart X, Dishuck PC, Mercuri L, Harvey WT, Gershman A, Diekhans M,
6 Sulovari A, Munson KM, Lewis AP, et al. 2022. Segmental duplications and their
7 variation in a complete human genome. *Science (1979)* **376**.
8 doi:10.1126/science.abj6965.
- 9 Wang Q, Dhindsa RS, Carss K, Harper AR, Nag A, Tachmazidou I, Vitsios D, Deevi SVV,
10 Mackay A, Muthas D, et al. 2021. Rare variant contribution to human disease in 281,104
11 UK Biobank exomes. *Nature* 2021 597:7877 **597**: 527–532. doi:10.1038/s41586-021-
12 03855-y.
- 13 Wen X, Liu Y, Yan Q, Liang M, Tang M, Liu R, Pan J, Liu Q, Chen T, Guo S, et al. 2018.
14 Association of IGFN1 variant with polypoidal choroidal vasculopathy. *J Gene Med* **20**:
15 e3007. doi:10.1002/JGM.3007.
- 16 Whiffin N, Karczewski KJ, Zhang X, Chothani S, Smith MJ, Evans DG, Roberts AM, Quaife
17 NM, Schafer S, Rackham O, et al. 2020. Characterising the loss-of-function impact of 5'
18 untranslated region variants in 15,708 individuals. *Nature Communications* 2020 11:1
19 **11**: 1–12. doi:10.1038/s41467-019-10717-9.
- 20 Yan X, He S, Dong D. 2020. Determining How Far an Adult Rare Disease Patient Needs to
21 Travel for a Definitive Diagnosis: A Cross-Sectional Examination of the 2018 National
22 Rare Disease Survey in China. *International Journal of Environmental Research and*
23 *Public Health* 2020, Vol 17, Page 1757 **17**: 1757. doi:10.3390/IJERPH17051757.

1

Fig. 1

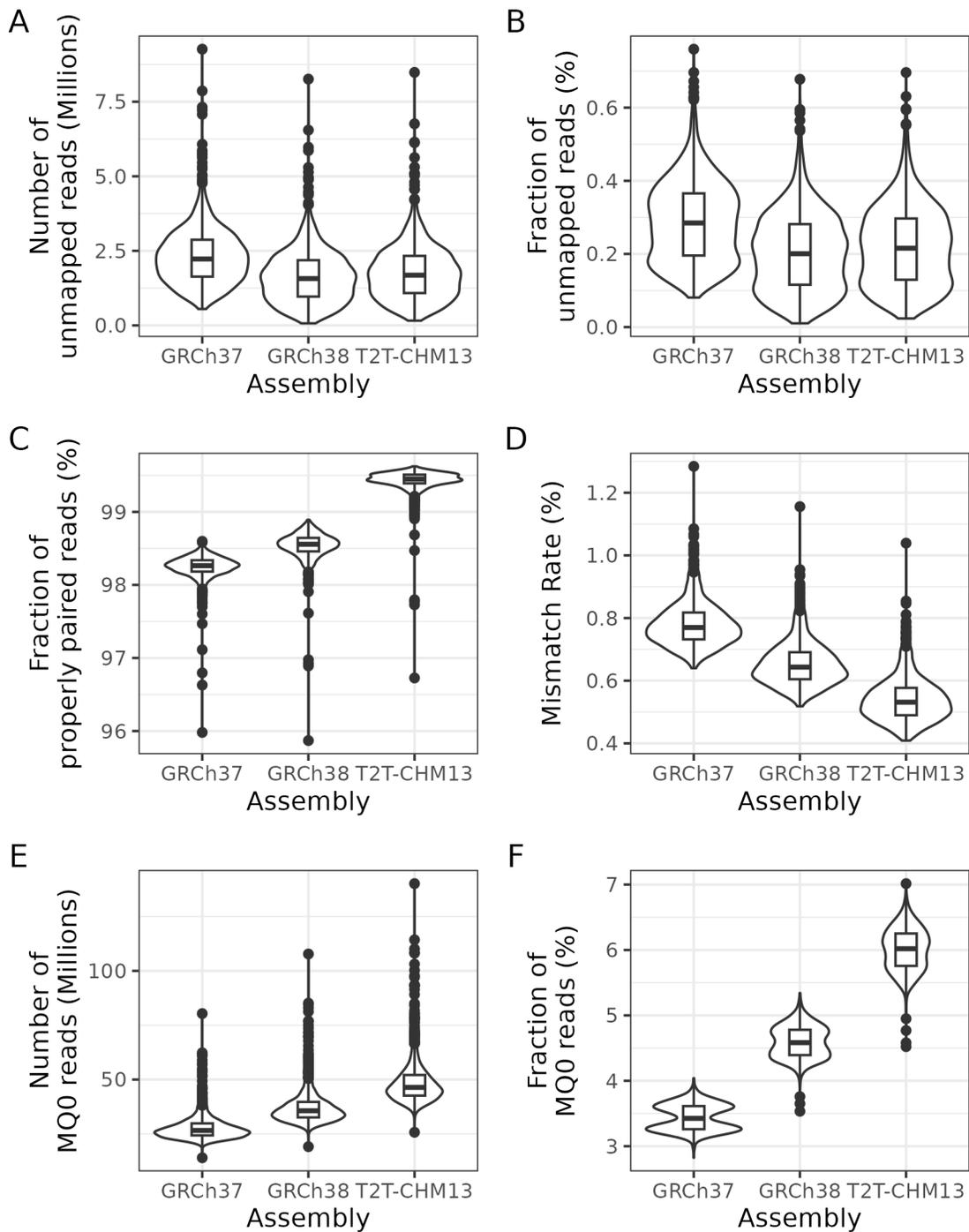


Fig. 2

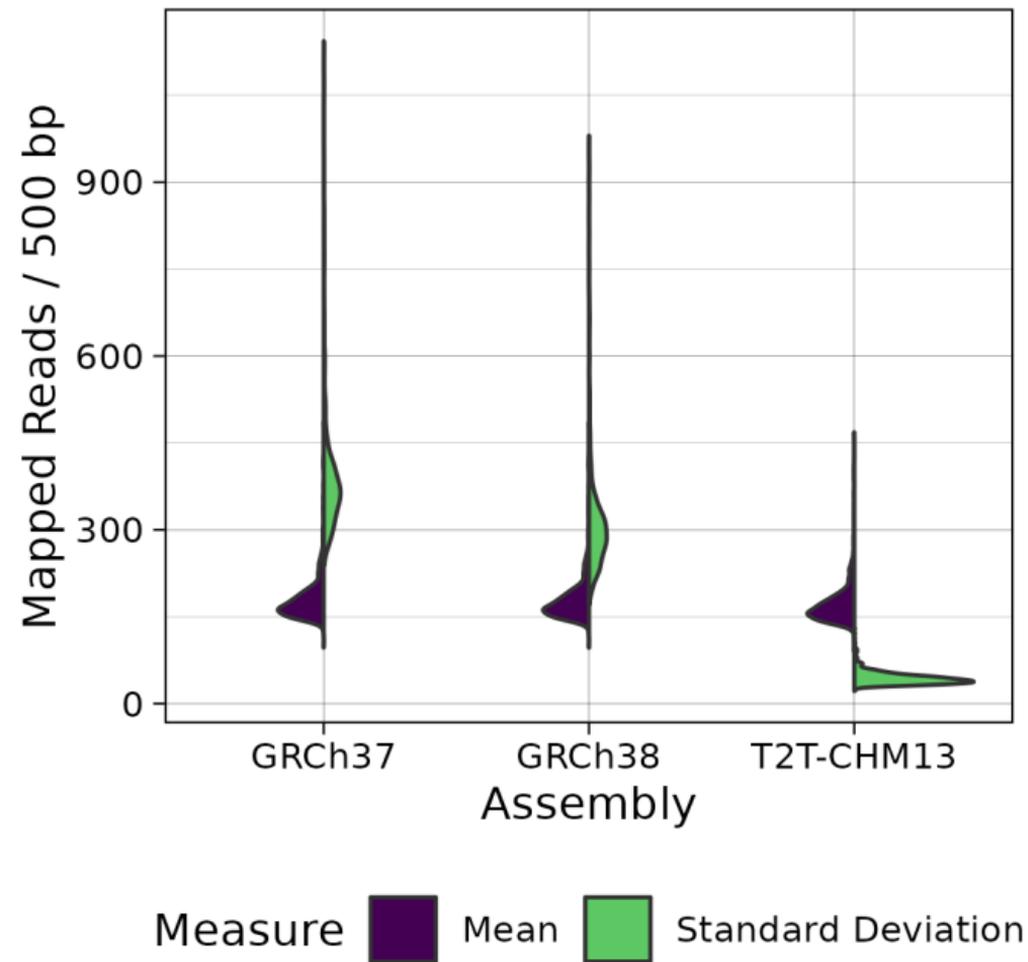


Fig. 3

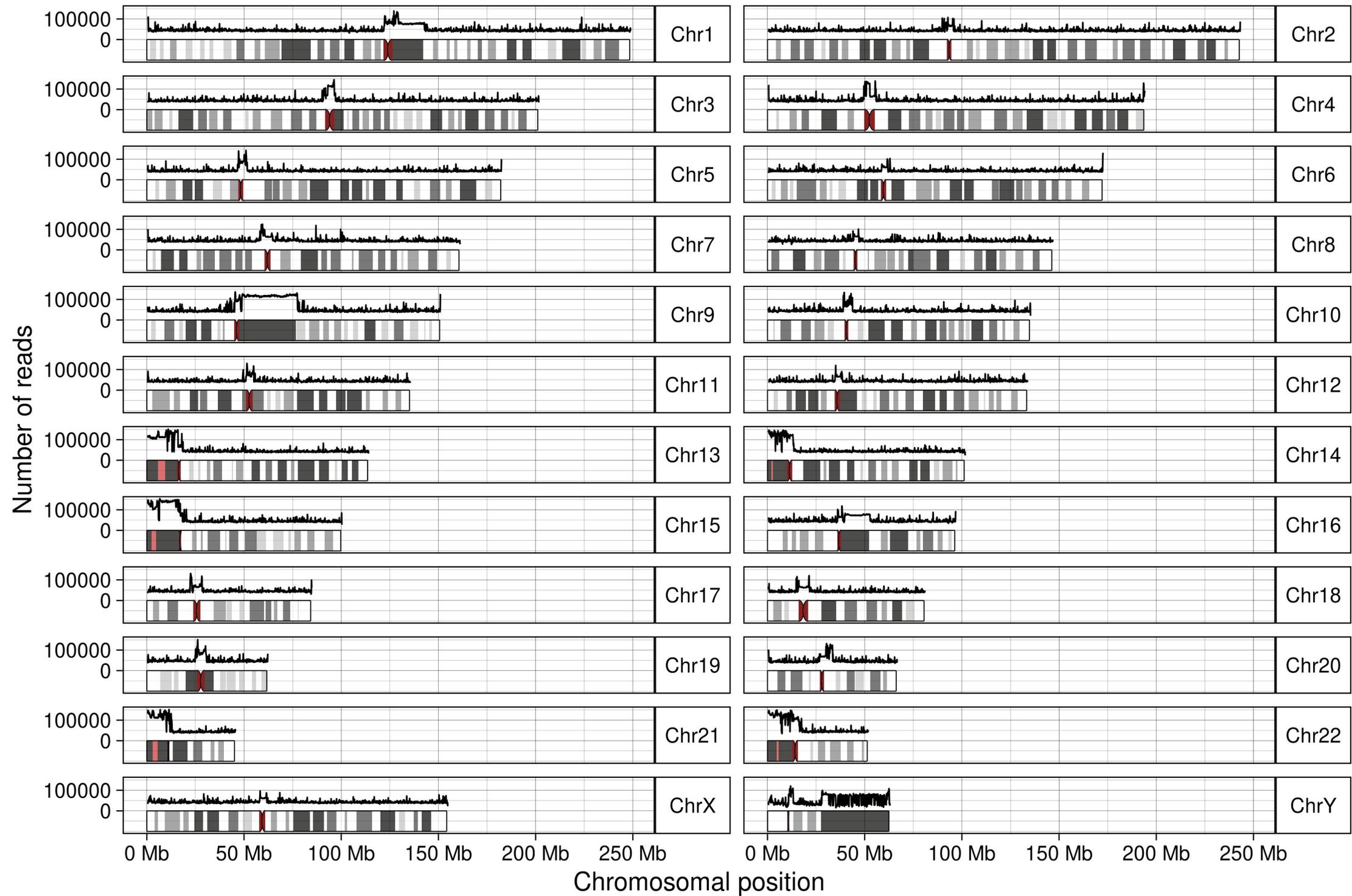
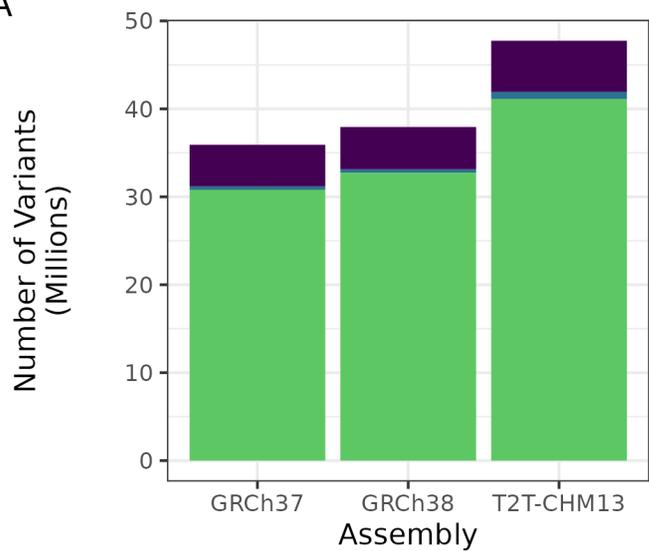
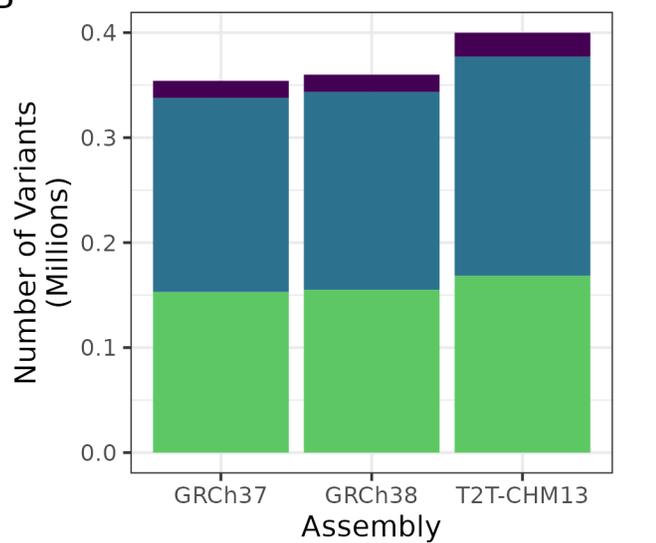


Fig.4**A****B**

Variant Type

- Indel
- Mixed
- SNV

Impact Rating

- HIGH
- MODERATE
- LOW

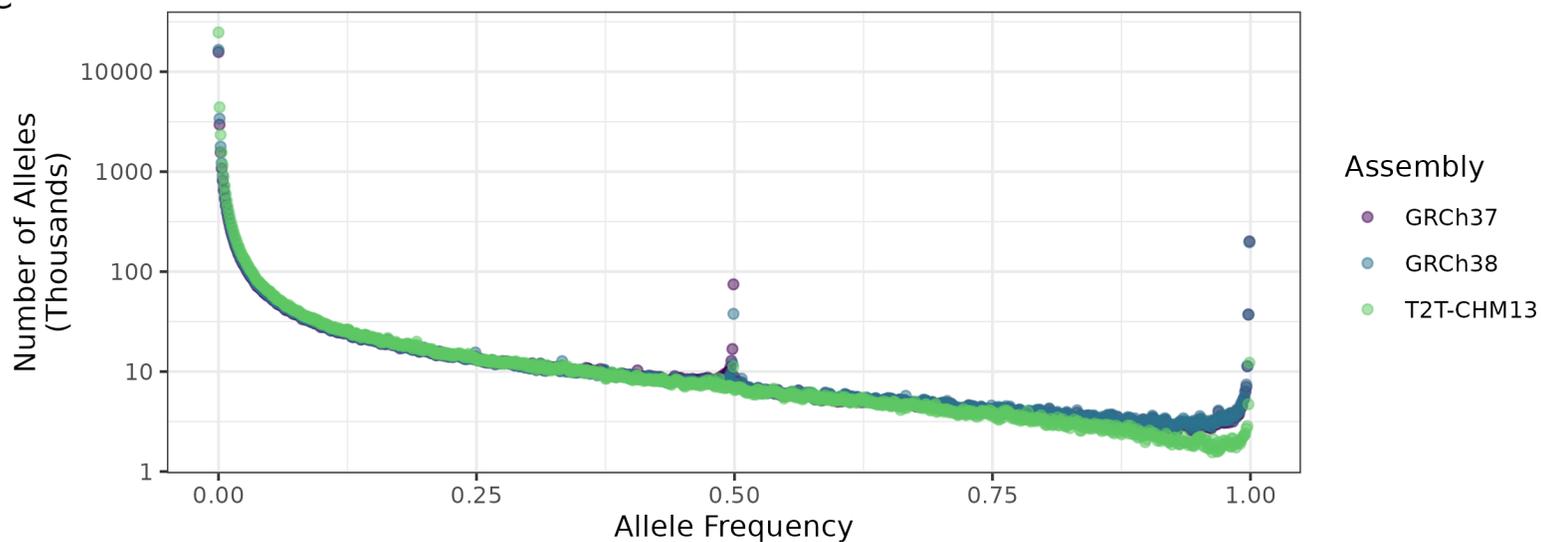
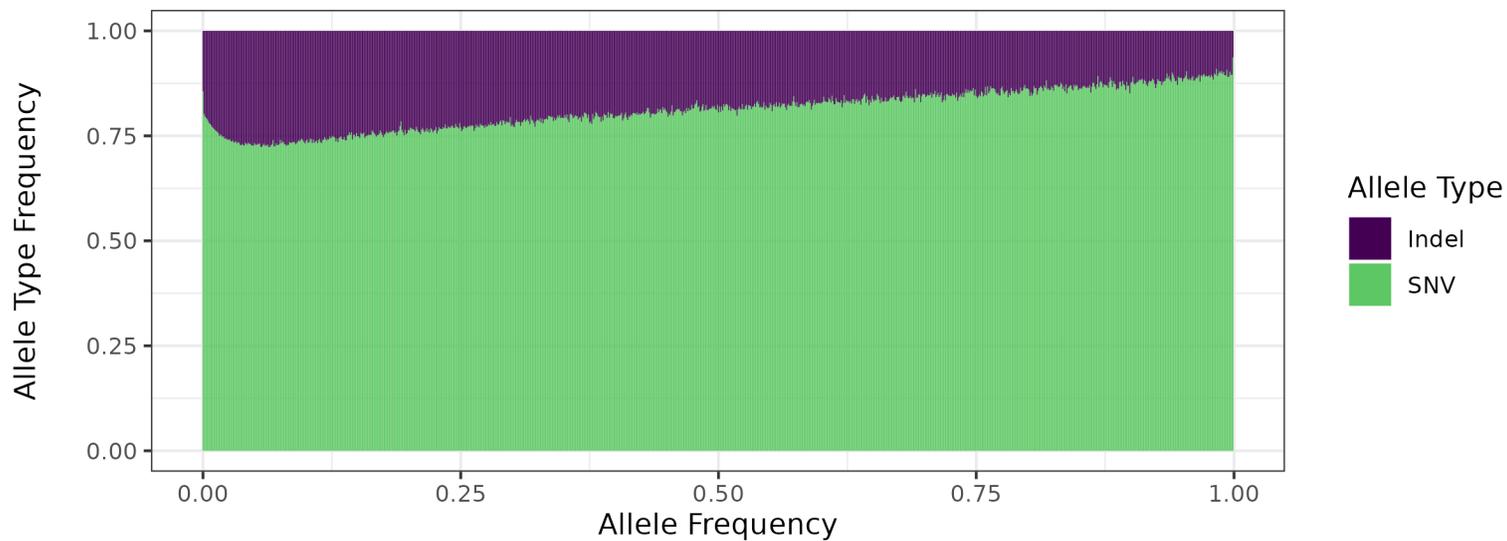
C**D**

Fig. 5

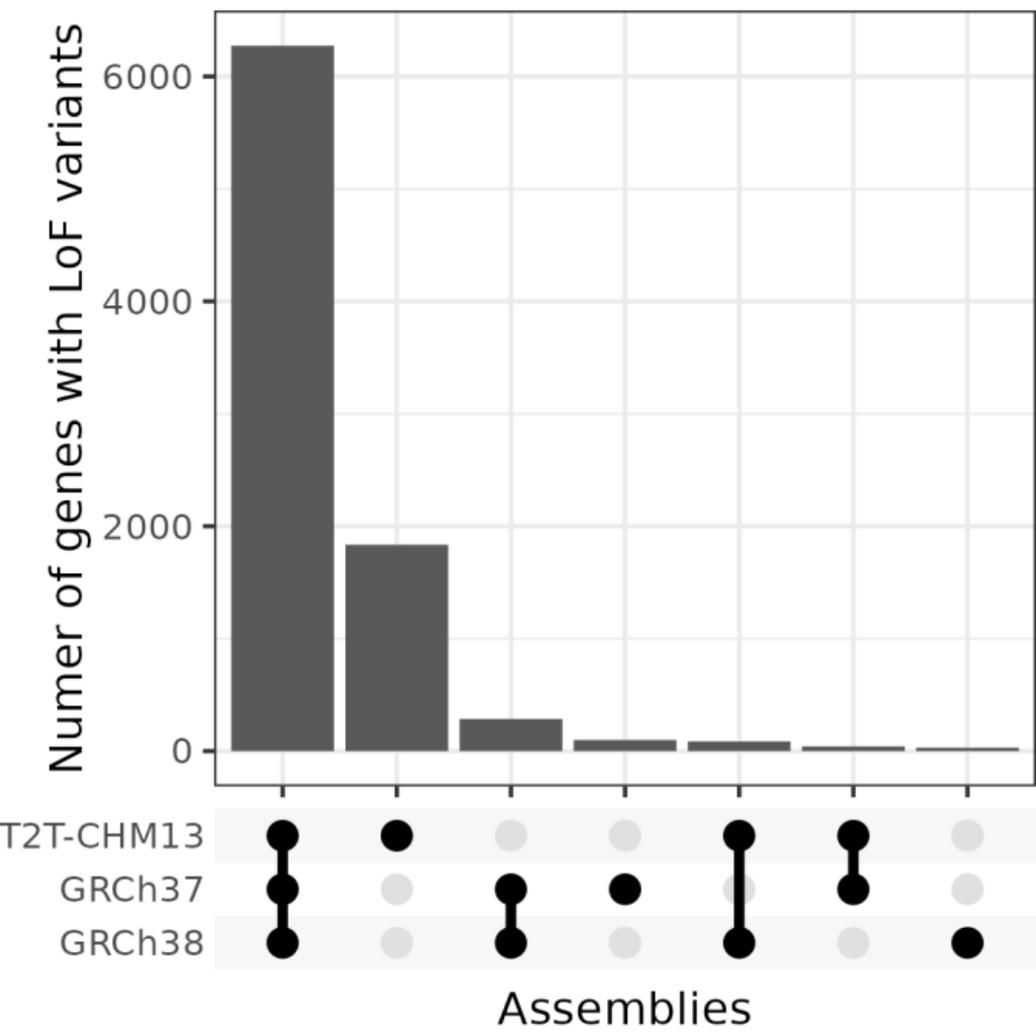


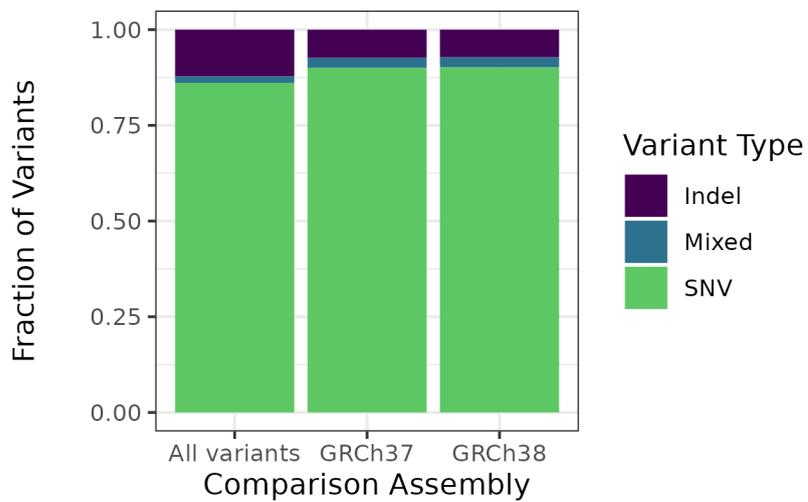
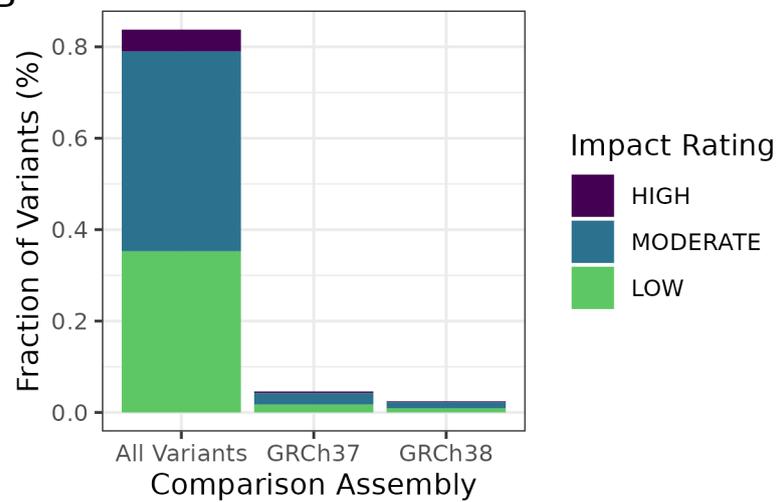
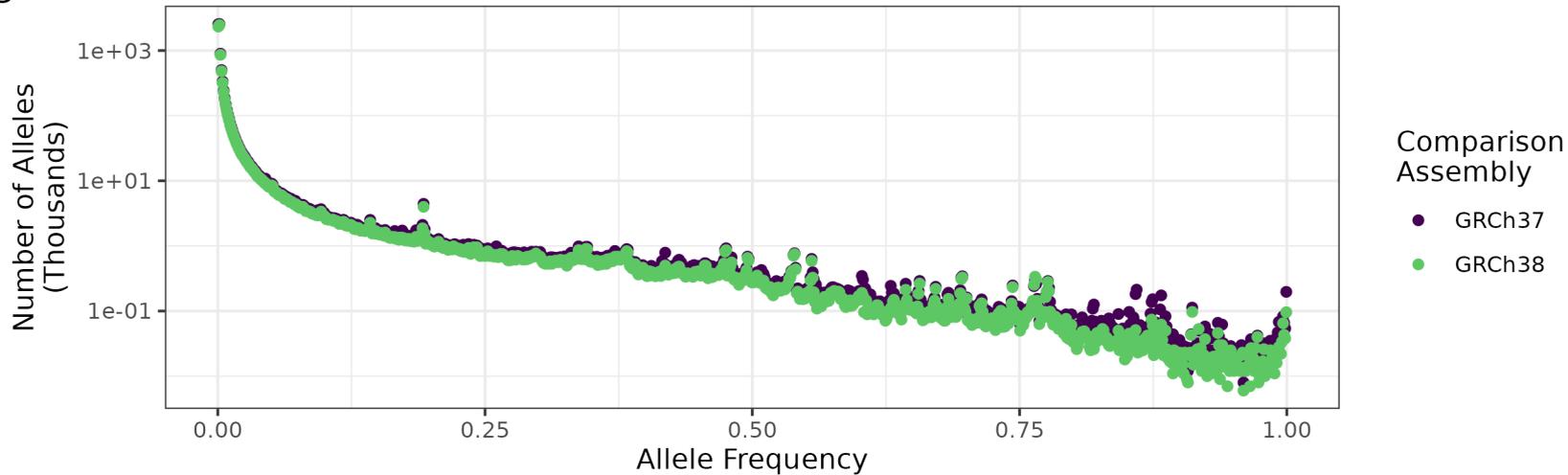
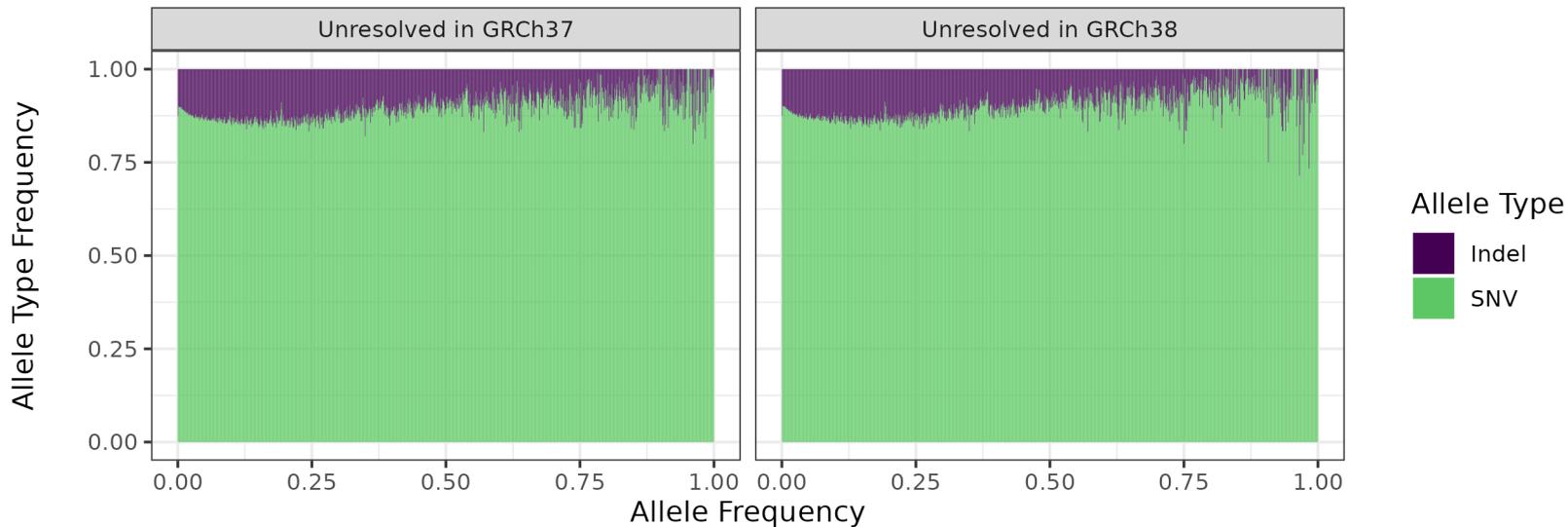
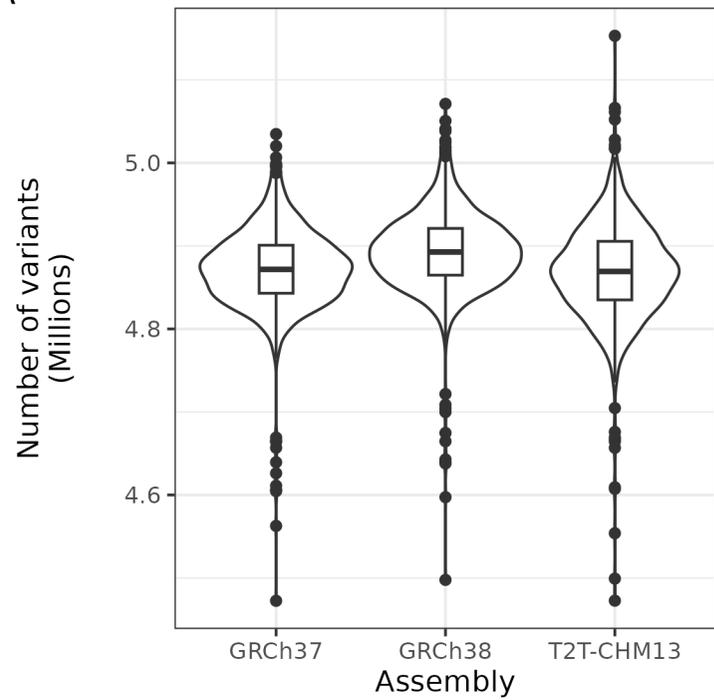
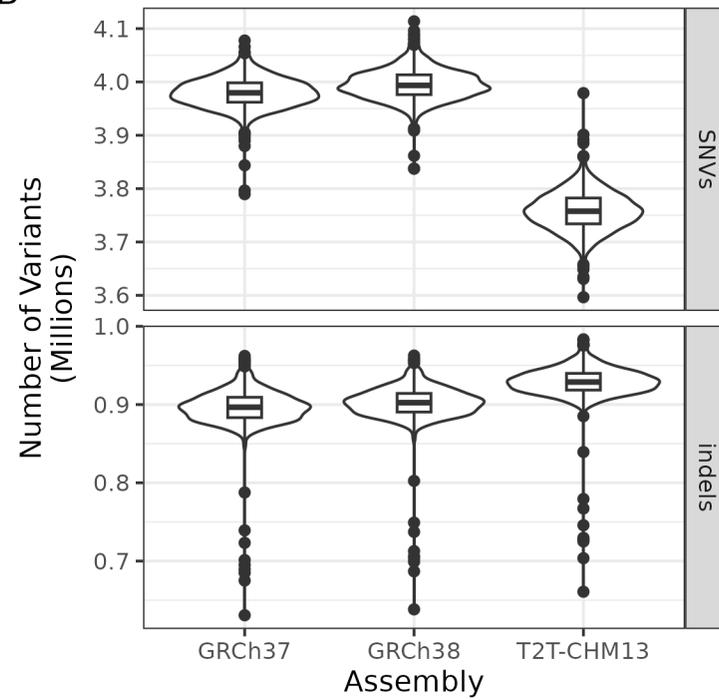
Fig. 6**A****B****C****D**

Fig. 7

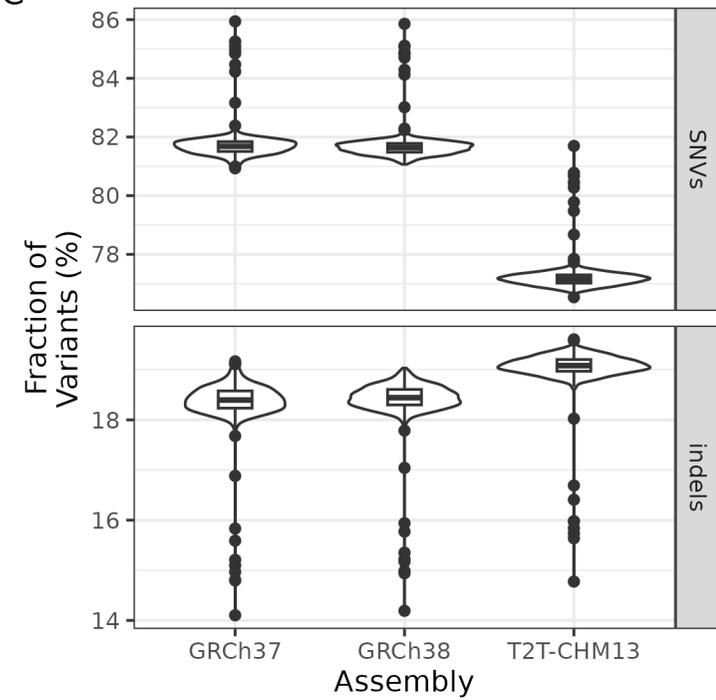
A



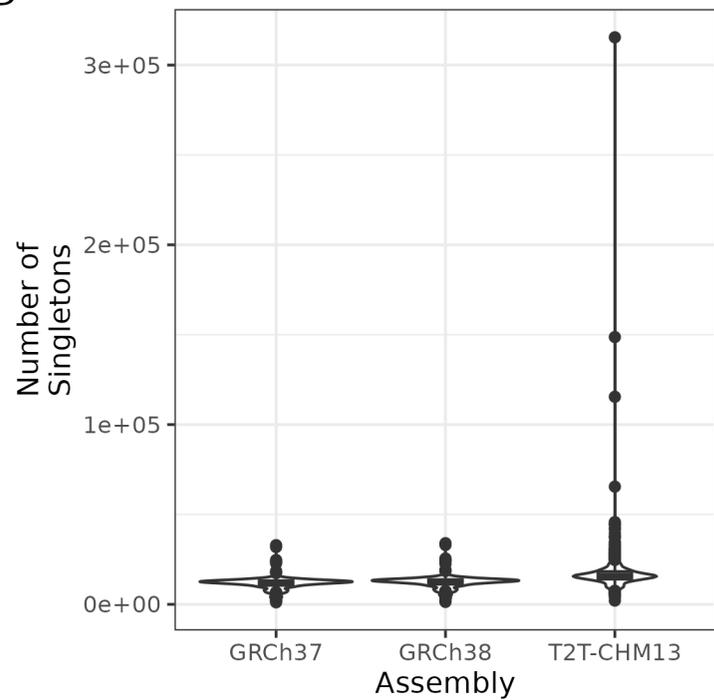
B



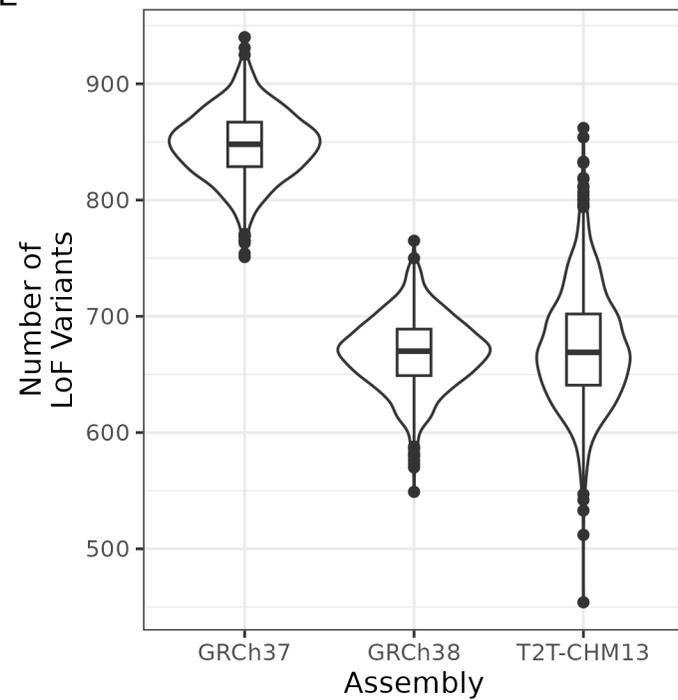
C



D



E



1 Figure Legends

2 Figure 1. General mapping statistics for all 1000 samples. A) Distribution of the total
3 number of unmapped reads. B) Percentage of unmapped reads per individual. C)
4 Percentage of mapped read pairs that were properly paired i.e., both mates were
5 mapped with the expected distance and orientation. D) Mismatch rate, i.e., the
6 fraction of mapped bases that did not match the reference. E) Number of mappings
7 with quality 0 i.e., not mapping uniquely. F) Percentage of mappings that had quality
8 0. MQ0 = mapping quality 0.

9 Figure 2. Summary of read depth in gene regions. Genes were divided into non-
10 overlapping 500-bp bins and the read depth of each bin was calculated. From these,
11 per-individual means (blue, left facing) and standard deviations (green, right facing)
12 were obtained which are shown here.

13 Figure 3. Number of reads from unplaced contigs mapping across T2T-CHM13 in
14 bins of 100 kbp. For each chromosome, an ideogram including the bands is shown;
15 centromeres are dark red and thinner. The x axis corresponds to the position on the
16 chromosome. The y axis is logarithmic and corresponds to the number of reads The
17 line above each chromosome shows the number of reads mapping to a 100 kbp
18 window.

19

1 Figure 4. Summaries of variant characteristics from the whole cohort. A) Comparison
2 of overall variant counts, separated by variant type (SNV, indel or mixed, i.e., both)
3 between variant calls from the different assemblies. B) Overall numbers of variants
4 by predicted impact in the whole call set as determined by SnpEff for each assembly.
5 Variants with impact rating MODIFIER were excluded. C) Distribution of allele
6 frequencies in the three assemblies. D) Fraction of each variant type by allele
7 frequency using T2T-CHM13. Frequencies were rounded to the closest multiple of
8 0.001 (1/1000).

9 Figure 5. Number of genes with predicted LoF variants shared among assemblies.

10 Figure 6. Summaries of variant characteristics from the whole cohort, restricted to
11 newly resolved regions. A) Fraction of variants by type in the novel regions in the
12 T2T-CHM13 assembly that were unresolved in the reference as shown on the x axis
13 and the whole genome. B) Fractions of variants by predicted impact in the whole call
14 set as determined by SnpEff in previously unresolved regions and the whole
15 genome. Variants with impact rating MODIFIER were excluded. C) Distribution of
16 allele frequencies for T2T-CHM13 specific variants mapped to regions not resolved
17 in the two previous assemblies. D) Fraction of each variant type by allele frequency.
18 Frequencies were rounded to the closest multiple of 0.001 (1/1000).

19 Figure 7. Summaries of several per-individual variant call statistics. Statistics
20 included only variants passing filters. A) Number of overall variants. B) Number of
21 SNVs and indels. C) Fraction of SNVs and indels of called variants. D) Number of
22 singletons. E) Number of LoF variants.

Table 1. General cohort-level characteristics of alignments and variant calls with T2T-CHM13 in comparison with previous assemblies

	T2T-CHM13	GRCh37	GRCh38
Ungapped length (bp)	3,117,275,501	2,861,327,195	2,937,639,396
Avg coverage (min /max)	39.0 (20.3/105.9)	39.8 (21.3/107.7)	40.2 (21.5/108.7)
Number of variant sites*	47,744,487	35,571,130	37,938,450
Variant rate (variants / kbp)**	15.3	12.4	12.9
Number of SNV sites (not in dbSNPv155)	41,931,920 (10,806,250)	30,866,176 (14,806,371)	33,133,695 (12,070,889)
Number of indel sites (not in dbSNPv155)	6,650,809 (3,709,821)	5,089,447 (3,170,927)	5,221,969 (2,457,457)
Ts/Tv ratio	1.50	1.95	1.80
Number of common variants (MAF \geq 5%)***	9,468,024	9,800,121	9,832,496
Number of low-frequency variants (MAF < 5%)***	43,151,191	29,714,010	31,895,019
Number of rare variants (MAF < 1%)***	37,402,076	24,807,638	26,741,004
Number of singletons (per kbp) [% of total variants]	19,077,986 (6.12) [40.0]	11,502,698 (4.02) [32.3]	13,547,273 (4.61) [35.7]
Number of non-reference alleles with AF > 50%	1,940,331	2,424,122	2,465,936
Number of LoF variants (% of total variants) [in ClinVar genes]	17,446 (0.037) [17,297]	12,628 (0.036) [12,326]	12,829 (0.034) [12,517]
Number of variants with low/medium/high impact (% of total variants)	168,420/208,625/ 22,848 (0.353/0.437/0.048)	153,001/184,694/ 16,392 (0.426/0.514/0.045)	155,088/188,347/ 16,542 (0.409/0.496/0.044)

Abbreviations: Ts: transition, Tv: transversion, MAF: minor allele frequency, AF: allele frequency, LoF: loss of function

* Multiallelic variants were only counted once.

** based on ungapped length

*** Each alternative allele was counted separately. If an alternative allele had a frequency > 50%, the reference allele was considered the minor allele.

Table 2. Top 10 ClinVar genes by number of loss-of-function variants by assembly

GRCh37		GRCh38		T2T-CHM13	
Gene	Number of LoF variants	Gene	Number of LoF variants	Gene	Number of LoF variants
<i>MUC4</i>	143	<i>MUC4</i>	172	<i>MUC4</i>	245
<i>MUC19</i>	99	<i>MUC3A</i>	127	<i>MUC19</i>	61
<i>MUC3A</i>	71	<i>MUC19</i>	100	<i>MUC2</i>	54
<i>MUC16</i>	48	<i>MUC6</i>	47	<i>MAGEC1</i>	46
<i>MUC6</i>	43	<i>HLA-DRB1</i>	39	<i>MUC5AC</i>	46
<i>HLA-DRB1</i>	37	<i>MUC16</i>	36	<i>HLA-DRB1</i>	35
<i>HLA-DRB5</i>	29	<i>HLA-DRB5</i>	33	<i>AHNAK2</i>	30
<i>ZNF717</i>	23	<i>ZNF717</i>	25	<i>MUC17</i>	29
<i>PDE4DIP</i>	22	<i>ZNF880</i>	19	<i>NACA</i>	28
<i>ANKRD36</i>	21	<i>GOLGA6L6</i>	18	<i>IGFN1</i>	26

Abbreviations: LoF – loss of function

Table 3. Overall characteristics of variant calls in previously unresolved regions

	Regions unresolved in		All of T2T-CHM13
	GRCh37	GRCh38	
Total length (bp)	268,965,814	251,330,203	3,117,275,501
Number of variant sites*	8,321,170	7,826,115	47,744,487
Variant rate (variants / kbp)*	30.9	31.1	15.3
Number of SNV sites	7,710,249	7,265,225	41,931,920
Number of indel sites	829,457	771,429	6,650,809
Number of common variants (MAF \geq 5%)**	695,892	621,512	9,293,456
Number of low-frequency variants (MAF $<$ 5%)**	8,670,575	8,211,562	43,103,217
Number of rare variants (MAF $<$ 1%)**	7,721,885	7,316,843	37,800,790
Number of singletons [% of total variants]	2,540,127 (30.5)	2,335,181 (29.8)	19,077,986 (40.0)
Number of LoF variants	185	99	17,446
Number of variants with low/medium/high impact (% of total variants)	1,435/2,102/286 (0.017/0.025/0.003)	719/1,048/155 (0.009/0.013/0.002)	168,420/208,625/22,848 (0.353/0.437/0.048)

Abbreviations: MAF: minor allele frequency, LoF: loss of function

* Multiallelic variants were only counted once.

** Each alternative allele was counted separately. If an alternative allele had a frequency $>$ 50%, the reference allele was considered the minor allele.

Table 4. Per-individual statistics of variant calls.

	T2T-CHM13	GRCh37	GRCh38
Mean number of variants per individual	4,868,620	4,939,144	4,892,944
Mean number of SNVs (Fraction of total variants)	3,757,819 (80.2%)	4,041,524 (81.8%)	3,995,225 (81.6%)
Mean number of indels (Fraction of total variants)	928,186 (19.8%)	901,550 (18.2%)	901,818 (18.4%)
Mean number of homozygous variants (IQR)	1,218,295 (1,197,992– 1,235,602)	1,491,412 (1,476,977– 1,503,380)	1,491,407 (1,476,992– 1,503,377)
Mean number of heterozygous variants (IQR)	2,479,761 (2,444,590– 2,512,655)	2,490,186 (2,462,633– 2,520,624)	2,490,402 (2,462,660– 2,520,740)
Mean number of singletons (IQR)	17,034 (14,007-18,263)	11,501 (10,550–13,122)	12,145 (11,181– 13,756)
Mean number of LoF variants (IQR)	669 (641-702)	848 (829-867)	668 (649-689)

Abbreviations: IQR: Inter-quartile range, LoF: loss of function



T2T-CHM13 improves read mapping and detection of clinically relevant genetic variation in the Swedish population

Daniel Schmitz, Adam Ameer and Åsa Johansson

Genome Res. published online September 16, 2025

Access the most recent version at doi:[10.1101/gr.279320.124](https://doi.org/10.1101/gr.279320.124)

Supplemental Material <http://genome.cshlp.org/content/suppl/2025/10/10/gr.279320.124.DC1>

P<P Published online September 16, 2025 in advance of the print journal.

Accepted Manuscript Peer-reviewed and accepted for publication but not copyedited or typeset; accepted manuscript is likely to differ from the final, published version.

Open Access Freely available online through the *Genome Research* Open Access option.

Creative Commons License This manuscript is Open Access. This article, published in *Genome Research*, is available under a Creative Commons License (Attribution 4.0 International license), as described at <http://creativecommons.org/licenses/by/4.0/>.

Email Alerting Service Receive free email alerts when new articles cite this article - sign up in the box at the top right corner of the article or [click here](#).



To subscribe to *Genome Research* go to:
<https://genome.cshlp.org/subscriptions>
