**METHOD**                                                                                    **Open Access**

# GRIDSS2: comprehensive characterisation of somatic structural variation using single breakend variants and structural variant phasing

Daniel L. Cameron[1,2,3*] , Jonathan Baber[3,4], Charles Shale[3,4], Jose Espejo Valle-Inclan[5], Nicolle Besselink[5], Arne van Hoeck[5], Roel Janssen[5], Edwin Cuppen[4,5], Peter Priestley[3,4] and Anthony T. Papenfuss[1,2,6,7*]

* Correspondence: cameron.d@wehi.edu.au; papenfuss@wehi.edu.au
[1]Bioinformatics Division, Walter and Eliza Hall Institute of Medical Research, Parkville, Australia
Full list of author information is available at the end of the article

## Abstract

GRIDSS2 is the first structural variant caller to explicitly report single breakends—breakpoints in which only one side can be unambiguously determined. By treating single breakends as a fundamental genomic rearrangement signal on par with breakpoints, GRIDSS2 can explain 47% of somatic centromere copy number changes using single breakends to non-centromere sequence. On a cohort of 3782 deeply sequenced metastatic cancers, GRIDSS2 achieves an unprecedented 3.1% false negative rate and 3.3% false discovery rate and identifies a novel 32–100 bp duplication signature. GRIDSS2 simplifies complex rearrangement interpretation through phasing of structural variants with 16% of somatic calls phasable using paired-end sequencing.

## Introduction

The reliable detection of structural variants (SVs) is critical to understanding the role genome architecture plays in health and disease. This is especially important in cancer and precision medicine where structural variation can be a key driver mutation [1, 2]. Over the past decade, many tools have been developed for the detection of genomic rearrangements, which have been the subject of recent extensive benchmarks [3, 4]. These tools fall broadly into two camps: those that detect changes in DNA abundance, known as copy number variant or aberration (CNV/CNA) callers, and those that detect non-reference DNA adjacencies, known as structural variant (SV) or breakpoint callers. While CNAs and SVs are merely two different viewpoints of the underlying genomic rearrangements, the methods of detection are fundamentally different. Here, we address the problem of SV detection and show that breakpoint detection alone is insufficient for the comprehensive characterisation of somatic genomic rearrangements that

occur in cancer. A third genomic rearrangement primitive is essential: single breakends.

The variant call format (VCF) [5] defines a single breakend as a breakpoint in which only one side can be unambiguously placed. This can occur due to one of two reasons. Firstly, the sequence on one side of the breakpoint could be absent from the reference. Either non-reference sequence could be present due to the integration of foreign DNA (e.g. provirus) or the reference could lack sequence present in the sample. Secondly, breakpoints into highly repetitive regions cannot be unambiguously placed. Single breakends allow the representation of such breakpoints. Such rearrangements are common in cancer and by reporting single breakends the rearrangement landscape of regions previously considered inaccessible to short read sequence can be explored.

Short read-based SV detection algorithms identify breakpoints by finding clusters of reads that do not support the reference allele. Typically, these use discordant read pairs [6], or split reads [7], with some callers also considering reads with unmapped mates [8] and soft-clipped reads [9]. More sophisticated callers incorporate assembly either through de novo assembly [10], targeted breakpoint assembly [11, or breakend assembly [11]. These callers report breakpoints, that is, novel adjacencies. When reads cannot be unambiguously mapped on either side, a breakpoint call cannot be made and information is lost. Some callers have attempted to address this by considering multiple alignment locations for each read [12], but this only works for regions with a small number of potential alignment locations and has proven impractical for general use. Single breakend calling has the potential to improve short read caller sensitivity above the 50% reported in recent benchmarking [3, 4].

As we move closer to a world in which the CNA and SV primitives can be reliably detected, accurate interpretation of the causative biological events becomes increasingly possible by integrated analysis of this knowledge. While progress has been made on derivative chromosome reconstruction using long reads [13], reconstruction of complex events such as chromothripsis has been problematic for short reads [14, 15]. To date, SV phasing has been used to reduce the complexity of reconstruction for long read based approaches [16] but has not been done by short read callers. The ability of phase somatic structural variants is limited by the read length and, for short read data, by the library fragment size—typically less than 500 bp.

Here, we demonstrate the power of single breakend variant calling using GRIDSS2—a somatic structural variant caller that reports single breakends and phases nearby structural variants. Running GRIDSS2 on 3782 metastatic solid tumours with matched normal samples from the Hartwig cohort, we show that, due to the high prevalence of somatic breakpoints involving low-mappability sequences, GRIDSS2 achieves a false negative rate lower than possible with a traditional breakpoint-only caller. The precision and sensitivity of GRIDSS2 in conjunction with single breakend variant calling and SV phasing lay a strong foundation for downstream tools that enable a deeper understanding of the nature of somatic genomic rearrangements.

## Results

GRIDSS2 utilises the same high-level approach as the first version of GRIDSS, assembling all reads that potentially support a structural variant using a positional de Bruijn graph breakend assembly algorithm [11]. Breakend contigs are then realigned back to

the reference to identify breakpoints and probabilistic structural variant calling is performed based on both the aligned reads and assembled contigs. Single breakend variant calling uses the same probabilistic variant calling approach as breakpoint calling, but instead of split reads, discordant read pairs, and assembly contigs with chimeric alignments support, single breakends are called based on soft-clipped reads, reads with unmapped or ambiguously mapping mates, and assemblies with unmapped or ambiguously mapping breakend sequence (Fig. 1a). SV phasing is performed based on assembly contigs and the presence of transitive calls (Fig. 1). SVs are phased cis if an assembly spans both breaks or a transitive call is found and phased trans if an assembly involves one SV but supports the reference at the other. Since assembly contig length is limited by the library fragment size, only nearby SVs can be phased. GRIDSS2



**Fig. 1** GRIDSS2 overview. **a** contigs are assembled from a single locus of reads mutually supporting the same putative break junction. If the other side cannot be uniquely determined, the contig supports a single breakend call at the break junction position. If different portions of the contig sequence uniquely align to different genomic loci, the assembly supports multiple cis phased breakpoints. **b** Nearby structural variants will have discordant read pairs spanning across multiple breakpoints. These generate spurious transitive calls that are collapsed into the underlying breakpoints, phasing them cis

Cameron *et al. Genome Biology*      (2021) 22:202

Page 4 of 25

includes a 16-step somatic filter specifically tuned for deeply sequenced tumour/normal samples.

## Benchmarking performance

To estimate precision and sensitivity of GRIDSS2, we used a recently generated "gold standard" somatic SV truth set for the COLO829 melanoma cell line and the COLO829BL cell line, which was derived from a normal cell from the same individual, using a combination of Illumina, PacBio, Oxford Nanopore, 10X Genomics linked reads, and optical mapping followed by targeted capture and PCR-based validations and manual curation [17]. To test sensitivity and reproducibility, we ran GRIDSS2, GRIDSS1 [11] Manta [18], svaba [19], and novobreak [20] on 2 independent sequencing replicates of the COLO829T/COLO829BL matched tumour-normal cell lines sequenced to a depth of 100x tumour and 40x normal coverage. With the GRIDSS2 panel of normals (PON) filter applied, GRIDSS2 achieved an average sensitivity/precision of 98%/86% compared to 94%/63% for Manta, 72%/39% for GRIDSS1, 81%/42% for svaba, and 76%/7% for novobreak (Fig. 2a). Without the PON, precision was reduced to 76%, 53%, 13%, 13%, and 7% for GRIDSS2, Manta, GRIDSS1, svaba, and novobreak respectively (Additional file 1: Figure S1).

To evaluate performance at lower sequencing depths and sample purity, we use in-silico downsampling and mixing to simulate a matched normal at 40x and a 60x



**Fig. 2:** Somatic benchmarks. **a** COLO829T/BL tumour and blood cell lines were sequenced in triplicate to 100x/40x. In-silico purity downsampling was performed at 40x normal and 60x tumour coverage. Results are compared against a PCR validated somatic truth set generated from multiple sequencing technologies. **b** Simulation of somatic breakpoints from the CHM13 telomere to telomere assembly against hg38. Single breakend variant calling allows sensitive detection of breakpoints to satellite repeat regions such as centromeres without increasing the false discovery rate. **c** GRIDSS2/Manta validation results on 13 patient samples for 50 bp + events. **d** GRIDSS2/Manta/strelka validation results on 13 patient samples for small (32–100 bp) duplications. **e** GRIDSS2/Strelka validation results for 32–50 bp events. **f** Per sample counts of 32–100 bp somatic tandem duplications in the Hartwig cohort. These mutations are enriched in colorectal cancer and associated with ATM driver mutations. **g** Size distribution of small (32–100 bp) tandem duplications across the Hartwig cohort. This is a distinct signature not associated with microsatellite expansion

tumour sample at 8–100% purity corresponding to 5x, 10x, 15x, 20x, 25x, 30x, 45x, 50x, and 60x effective tumour coverage. Above 10x effective tumour coverage GRIDSS2 achieved higher sensitivity and specificity than the benchmarked callers. At 10x and below, GRIDSS2 retained higher precision, but at lower sensitivity than Manta or svaba (Fig. 2a).

### Single breakend performance

To evaluate the expected performance of single breakend variants in repetitive genomic regions, we ran 2998 simulations each containing a single breakpoint between either the centromeric region, or the region immediately after the centromere of the complete telomere-to-telomere chromosome 8 assembly of CHM13 [21] and a CHM13 region with a unique hg19 chr8 liftover position. For each breakpoint, a simulated normal was generated from the 10 kb sequences flanking the breakpoint and a simulated tumour from the breakpoint and 10 kb flank sequences. Reads were aligned to hg19 chr8 and somatic calling was performed with GRIDSS2, Manta, novobreak, and svaba. Breakpoints were classified based on whether one or both sides of the SV call matched. All callers achieved high sensitivity except in the centromeric satellite repeat regions. In these regions, all callers correctly called both sides of 29% breakpoints with 1%/66%/11%/63% of calls not detected by GRIDSS2/Manta/novobreak/svaba (Fig. 2b). Novobreak's kmer-based approach allows it to detect most of the centromeric breakpoints but at the cost of incorrectly reporting multiple breakpoints 32% of the time as well as potentially incorrect breakpoint location for 66% of breakpoints (c.f. GRIDSS2 9%). This is consistent with our VIRUSBreakend result in which we found that by using GRIDSS2 single breakend variant calls we can reliably detect viral integrations anywhere in the host genome whereas other callers had either low sensitivity or a high false discovery rate in repetitive regions such as centromeres [22].

### Validation on patient samples

To further validate somatic performance, we performed independent validation of GRIDSS2 and manta breakpoint calls from 13 patient tumour samples from the Hartwig cohort [2, 18–20] with a high burden of structural variants. Since the default minimum reported event sizes of GRIDSS2 and Manta are 32 and 50 bp respectively, we compared 32–50 bp events to the short indel caller, Strelka [23]. We used a hybrid capture approach with target probes flanking and overlapping break-junctions to independently validate over 5000 calls identified by any tool. 3403 of 3666 (93%) GRIDSS2 calls were validated compared to 2685 of 4299 (65%) for Manta (Fig. 2c). Of the private Manta calls not found by GRIDSS2, just 230 of 1777 (13%) were validated compared to 836 of 1031 (81%) GRIDSS2 private calls. Imprecise (that is, not base-pair accurate) Manta calls validated at a rate (40/288, 14%) similar to Manta private calls, whereas GRIDSS2 reports only precise somatic SV. No imprecise GRIDSS2 calls passed somatic filtering, whereas all validated imprecise Manta calls were called by GRIDSS2 precisely. In the 32–50 bp range, 329 of 343 (96%) of GRIDSS2 calls validated against 142 of 182 (78%) for Strelka (Fig. 2e). 95% (219 of 232) of 32–50 bp calls private to GRIDSS2 were validated, compared to 47% (35 of 74) for Strelka. Eighty-nine percent (170 of 192) GRIDSS2 single breakend calls were validated, with 17 of the 22 failing validation

occurring in simple repeats. Ninety-eight percent (53 of 54) of centromeric and 88% (61 of 69) of LINE/L1 single breakends passed validation indicating robust single breakend performance outside of simple repeats.

Notably, GRIDSS2 finds many short (32–100 bp) duplications which are largely missed by both Strelka and Manta and does so with 99% precision (Fig. 2d).

### Novel somatic short duplication signature

In addition to reidentifying known kilobase and megabase length duplication signatures, we find a signature consisting of short 32–100 bp non-microsatellite tandem duplications (Fig. 2g). There is a median of 4 of these short (32–100 bp) duplications per sample (Fig. 2f). They are not correlated with larger duplications (R = 0.08), or total breakpoints (R = 0.10). Enrichment of samples with 15 or more short duplications is positively associated with colorectal cancer (q = $1.2 \times 10^{-9}$) and driver mutations in PARK2 (q = 0.0003) and ATM (q = 0.008). Across the Hartwig cohort, 21 patients had driver mutations [2] involving the disruption of a tumour suppressor caused by small duplications (Additional file 2: Table S2).
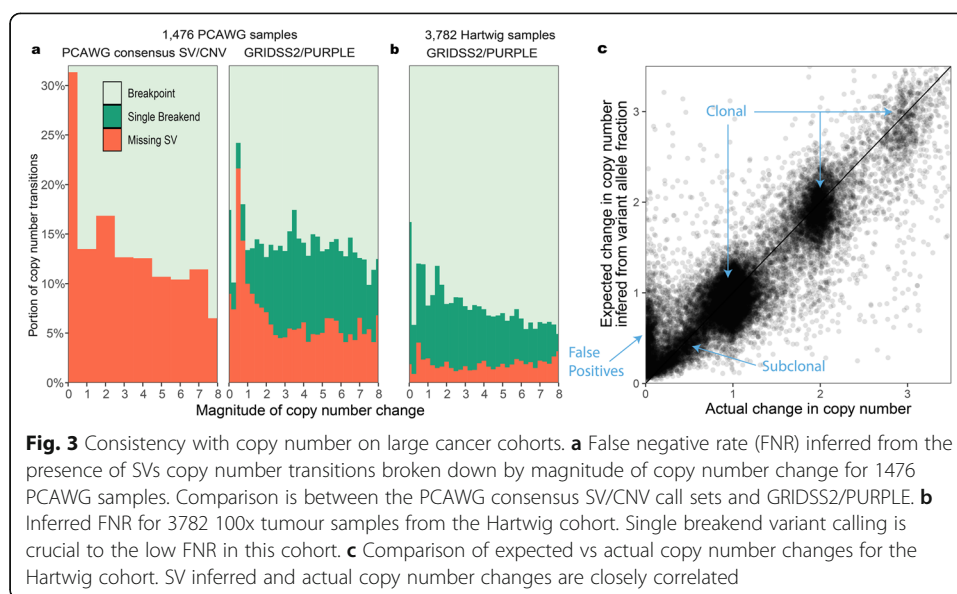
These short tandem duplications are too large to be reliably called by most somatic indel callers, but too short to be reliably called by many SV callers. In part, this is due to the weak read pair signal due to the short variant length, but also since most callers do not report variants shorter than 50 bp threshold used for variant databases such as dbVar. Popular callers such as lumpy [24] and delly [25] do not call duplications shorter than 100 and 300 bp respectively [4], and no duplications shorter than 300 bp were included in the PCAWG consensus call set [1].

### Cohort-level FNR/FDR estimation using copy number consistency

Structural variant and copy number calls are intrinsically related. Any breakpoint must have either a compensating breakpoint (for example, as with inversions) or a copy number change at that SV position. Using this principle, we can estimate a false negative rate (FNR) from the number of unexplained copy number transitions. To generate matching SV and copy number calls, we ran GRIDSS2 and PURPLE [2] on 1476 samples from the ICGC PCAWG WGS cohort and compared results with the state-of-the-art PCAWG consensus call set [26]. Copy number transitions in or within 100 kb of centromeres or a gap in the reference genome were excluded.

Across the 1476 samples, GRIDSS2 identified breakpoints for 83.9% of copy number transitions and single breakends for a further 6.8%, with an estimated 9.3% FNR. The PCAWG consensus call set identified breakpoints for 76.2% of copy number transitions (23.8% FNR). When restricted to clonal copy number transitions, the estimated FNR for the PCAWG consensus dropped to 13.7% and GRIDSS2 to 9.0% (Fig. 3a), indicating robust subclonal GRIDSS2 performance.

As both the PCAWG consensus CNV and PURPLE copy number segmentation are, in part, driven by the PCAWG consensus SV and GRIDSS2 calls, there is the potential for the copy number segmentation to be over-optimised for the corresponding call set. To evaluate the extent of this, we compared the GRIDSS2 SV call sets to the PCAWG consensus CNV, and the PCAWG consensus SV to the PURPLE CNV calls. In both cases, FNR increased with PCAWG/PURPLE at 40.7% and GRIDSS2/PCAWG at 33.8%

**Fig. 3** Consistency with copy number on large cancer cohorts. **a** False negative rate (FNR) inferred from the presence of SVs copy number transitions broken down by magnitude of copy number change for 1476 PCAWG samples. Comparison is between the PCAWG consensus SV/CNV call sets and GRIDSS2/PURPLE. **b** Inferred FNR for 3782 100x tumour samples from the Hartwig cohort. Single breakend variant calling is crucial to the low FNR in this cohort. **c** Comparison of expected vs actual copy number changes for the Hartwig cohort. SV inferred and actual copy number changes are closely correlated

(Additional File 1: Figure S2). Restricting to clonal copy number transitions GRIDSS2/ PCAWG FNR reduced to 19.4% whereas PCAWG/PURPLE remained high at 38.3%. Examining the copy number transitions not explained by their matching SV caller, we found the PCAWG consensus SV call set includes explanatory SVs at 5928 PURPLE copy number transitions that were missed by GRIDSS2 and GRIDSS2 explains 9433 PCAWG copy number transitions (2415 by single breakends) missed by the PCAWG consensus SVs—an indication that the GRIDSS2 call set is likely to be more comprehensive than the PCAWG consensus SV calls.

To evaluate GRIDSS2 on high quality, deeply sequenced samples, GRIDSS2 and PURPLE were run on 3782 40x normal/100x tumour samples from the Hartwig cohort. Excluding those occurring within 1 kb of a gap in the reference genome, 153, 231 of 1,954,548 (7.0%) copy number transitions in the Hartwig cohort were explained only by single breakend variants and 68,171 (3.1%) lacked a corresponding GRIDSS2 SV (Fig. 3b). The higher rate of single breakend calling can be attributed to GRIDSS2 conservatively calling single breakends and the greater sequencing depth in the Hartwig cohort. The 7.0% of copy number transitions in the Hartwig cohort explained by single breakend variant calls represents a lower bound for the FNR of an exclusively breakpoint-based caller. A FNR of 3.1% suggests that, on this cohort, GRIDSS2 achieves a FNR lower than that possible for a breakpoint-based caller.
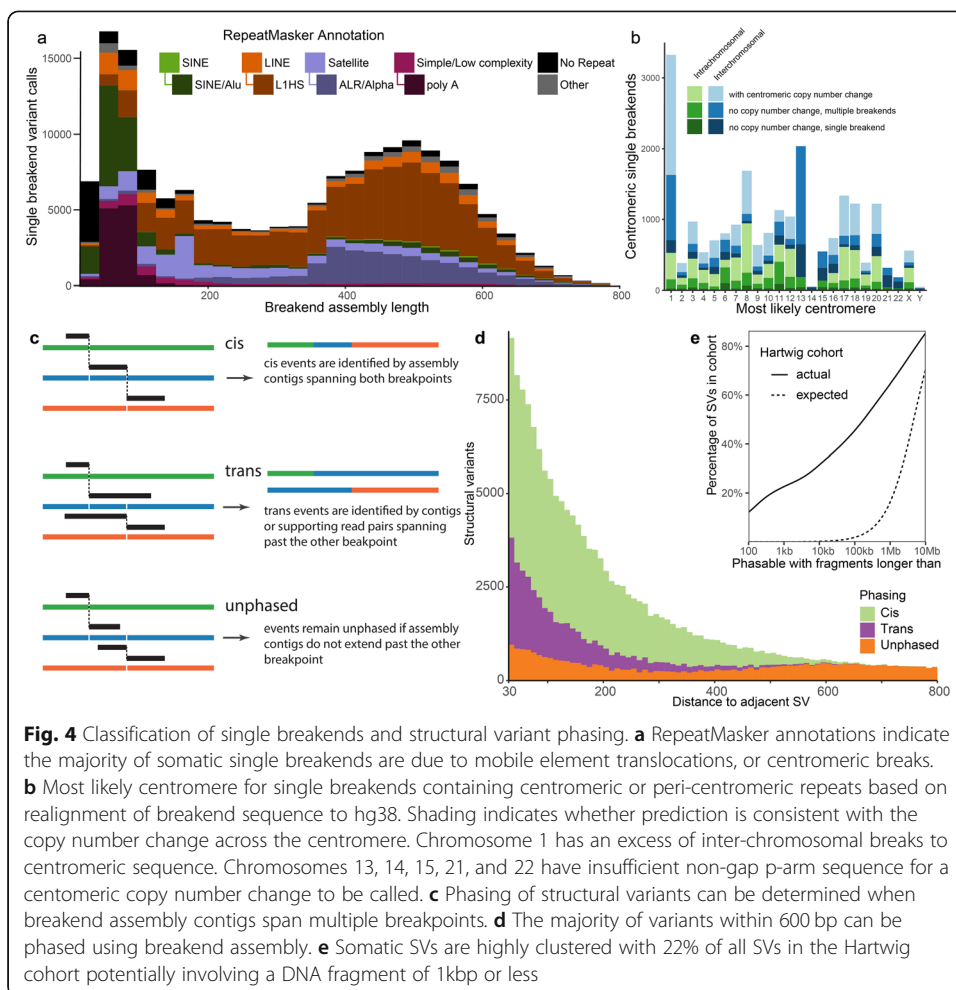
To demonstrate that this reduction in FNR does not come at the cost of a high false discovery rate (FDR), we compared the change in copy number to the change expected based on the variant allele fraction (VAF). For isolated breaks, the change in copy number should match the variant copy number inferred from the variant allele fraction. Using a 3000 bp threshold to ensure at least one full 1kbp copy number bin between SVs, we find that the VAF-inferred SV copy numbers reported by GRIDSS2 are consistent with the copy number changes with no systematic bias in the VAF (Fig. 3c). This

trend remains true for subclonal variants although the false discovery rate does go up. Assuming variants with a copy number change of less than 0.1 and a VAF inferred copy number of at least 0.25 are false positives, GRIDSS2 isolated SV calls have an estimated FDR of 5.4%, with 74% of these subclonal, and single breakends having twice the FDR of breakpoints. Extrapolating these to the rest of the cohort gives an overall estimated FDR of 3.3%.

## Resolving somatic centromeric rearrangements

Although only one side of single breakend variant calls can be uniquely placed, the assembled sequencing flanking the break can be used to classify integrated provirus, mobile element transposition, rearrangements involving centromeric and telomeric sequence, and other events. RepeatMasker annotation reveals that the majority of somatic single breakend calls are caused by SINE Alu, LINE L1HS insertions or rearrangements involving centromeric sequence, a pattern shared between both the Hartwig and PCAWG cohorts (Fig. 4a, Additional File 2: Figure S3). Breakend assembly lengths for SINE single breakends are typically shorter than 150 bp as assemblies longer than this can usually be resolved into breakpoint calls. Similarly, the polyA repeat motif characteristic of LINE translocations [27] is also found in the shorter breakend assemblies. Such assemblies are short as the de Bruijn graph assembler used truncates assemblies at unresolved repeat loops and assemblies able to span the polyA tail are able to be resolved as breakpoints.

Ninety-one percent of the Hartwig cohort samples contain at least one copy number transition occurring in centromeric sequence. Being able to resolve the partners of the centromeric breaks explaining these copy number changes is critical to the accurate reconstruction of the derivative chromosomes. Single breakends into ALR/Alpha and HSATII centromeric repeats are able to give significant insight into the nature of these centromeric breaks. As each human centromere has a slightly different dominant repeat sequence, a mapping between each centromeric single breakend and their most likely centromeric breakpoint partner is possible. To do this, we aligned the single breakend sequences containing a centromeric or peri-centromeric repeat against the hg38 reference genome using BLAT, annotating each with the most likely centromeric partner. Using this approach, we were able to explain 5614 of 11,996 (47%) centromeric copy number changes, implying that approximately half of centromeric rearrangements are centromere to centromere, and the remainder centromere to non-centromeric sequence. Of the 21,587 centromeric single breakends detected, 3148 (15%) had no copy centromeric copy number change, 6850 (32%) had no copy number change but had multiple single breakends linked to the same chromosome, 3358 (16%) had a single breakend associated with a centromere with copy number change, and the remaining 8231 (38%) associated with a centromere with copy number change with multiple breakends mapping to that centromere in that sample. Since chromosomes 13, 14, 15, 21, and 22 have insufficient non-gap p-arm sequence for a centromeric copy number change to be reported, any single breakend predicted to one of these centromeres will have no centromeric copy number change, even if the prediction is correct. These chromosomes make up 31% (3128/9998) of the centromeres with no reported copy number change (Fig. 4b).

**Fig. 4** Classification of single breakends and structural variant phasing. **a** RepeatMasker annotations indicate the majority of somatic single breakends are due to mobile element translocations, or centromeric breaks. **b** Most likely centromere for single breakends containing centromeric or peri-centromeric repeats based on realignment of breakend sequence to hg38. Shading indicates whether prediction is consistent with the copy number change across the centromere. Chromosome 1 has an excess of inter-chromosomal breaks to centromeric sequence. Chromosomes 13, 14, 15, 21, and 22 have insufficient non-gap p-arm sequence for a centromeric copy number change to be called. **c** Phasing of structural variants can be determined when breakend assembly contigs span multiple breakpoints. **d** The majority of variants within 600 bp can be phased using breakend assembly. **e** Somatic SVs are highly clustered with 22% of all SVs in the Hartwig cohort potentially involving a DNA fragment of 1kbp or less

## Somatic phasing

The breakend assembly approach taken by GRIDSS2 also enables the assembly-based phasing of nearby variants. When two structural variants occur in close proximity, they can be phased as cis if the contig aligns across both and trans if the contig aligning across one aligns to the reference sequence at the other (Fig. 4a). Segments shorter than 30 bp are not typically uniquely alignable by BWA and unaligned short DNA segments are treated as insert sequences of an SV connecting the longer flanking segments. Since breakend assembly contig lengths are limited by the fragment size distribution of the DNA library sequenced, only nearby variants can be phased.

For the Hartwig cohort, variants could be phased up to around 500 base pairs. We found that multiple nearby somatic structural variants are frequent, with 22% of all structural variants having an adjacent variant within 1000 bp. This is far in excess of the 0.02% expected if the breakpoints were uniformly randomly distributed (Fig. 4e). Of these, GRIDSS2 was able to phase 70% (Fig. 4d) with 72% cis and 28% trans. This distribution is recapitulated in the 1476 PCAWG samples and LINX classification of these structural variants indicate that that phasable breakpoint clusters occur predominantly in LINE translocations (due to target site duplication, and highly active donor

elements) and highly complex rearrangement events (Additional File 1: Figure S4). This phasing information greatly assists downstream derivative chromosome reconstruction, as it exponentially reduces the number of possible paths through the breakpoint graph.
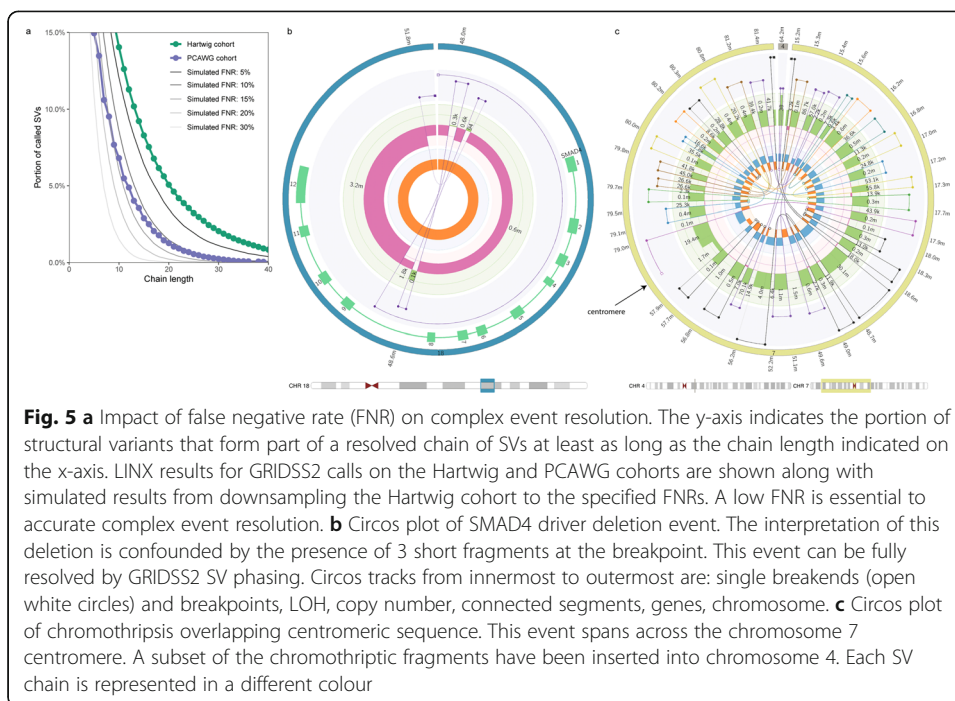
### Impact on complex event resolution

To demonstrate the impact on downstream analysis of complex somatic genomic rearrangements, we ran LINX [28], a rearrangement event interpretation and classification tool, on the Hartwig and PCAWG cohorts. To fully resolve complex rearrangements, structural variants must be chained together to reconstruct the relevant portions of the derivative chromosomes rearranged by the event. If there are errors in the SV call set, it is likely that many complex events will not be able to be fully resolved. To evaluate the impact of FNR on this reconstruction, we evaluated the portion of SVs resolved into long chains for the PCAWG and Hartwig cohorts. In addition, we simulated the effect of increasing FNR by subsampling the Hartwig call set (Fig. 5a). 5.0% of SVs in the Hartwig cohort were reconstructed into chains of 20 SVs or more. Increasing the FNR reduces this to 3.6% of SVs at 5% FNR, 1.5% at 10% FNR, 0.6% at 15% FNR, and 0.25% at 20% FNR. We previously estimated the PCAWG cohort FNR at 11.2%, and we find that the 1.27% of SVs in chains of 20 SVs or more closely match the 1.29% we expect from a simulated downsampling of the Hartwig cohort. This implies that the PCAWG and Hartwig pan-cancer cohorts have a broadly similar composition of complex rearrangements and the differences observed are primarily technical in nature. Small improvements in FNR result in large increases in the ability for downstream tools to resolve complex events. A sub-5% FNR is critical for large event reconstruction.

SV phasing can be critical to the correct interpretation of complex events. For LINX, SV phasing is a critical first step in the chaining of SVs. Of the 486,632 links in the chains resolved by LINX in the Hartwig cohort, 100,007 (21%) were due to GRIDSS2 SV phasing. For the PCAWG cohort, 13,212 of 107,952 links (12%) were resolved by GRIDSS2 SV phasing, with the difference primarily driven by shallower coverage and shorter library fragment sizes resulting in shorter assembly contig lengths (Additional File 1: Figure S2), and the higher FNR. In some cases, apparently complex events can be resolved to simple events containing additional short DNA fragments purely through SV phasing (Fig. 5b).

Finally, we use single breakend repeat annotations to identify instances of chromothripsis overlapping centromeres. In the Hartwig cohort, LINX identifies 270 complex events with at least 10 breakends to centromeric sequence, 17 of which could be fully chained (Fig. 5c). The large number of events with many centromeric single breakends indicates a previously unexplored level of centromeric involvement in complex rearrangements worthy of further investigation.

### Discussion

Through cell line, patient validation, and cohort-level comparisons, we have shown GRIDSS2 has excellent somatic performance above 10x effective tumour coverage. The identification of a small (32–100 bp) duplication signature by GRIDSS2 highlights the importance of robust software tested across a wide range of variant types and sizes. The presence of a signature overlapping the widely accepted but arbitrary 50 bp

**Fig. 5 a** Impact of false negative rate (FNR) on complex event resolution. The y-axis indicates the portion of structural variants that form part of a resolved chain of SVs at least as long as the chain length indicated on the x-axis. LINX results for GRIDSS2 calls on the Hartwig and PCAWG cohorts are shown along with simulated results from downsampling the Hartwig cohort to the specified FNRs. A low FNR is essential to accurate complex event resolution. **b** Circos plot of SMAD4 driver deletion event. The interpretation of this deletion is confounded by the presence of 3 short fragments at the breakpoint. This event can be fully resolved by GRIDSS2 SV phasing. Circos tracks from innermost to outermost are: single breakends (open white circles) and breakpoints, LOH, copy number, connected segments, genes, chromosome. **c** Circos plot of chromothripsis overlapping centromeric sequence. This event spans across the chromosome 7 centromere. A subset of the chromothriptic fragments have been inserted into chromosome 4. Each SV chain is represented in a different colour

threshold separating indels from structural variants suggests it is time to reconsider this threshold as the minimum reported event size for structural variants.

Explicit reporting and handling of single breakend variants represents a significant conceptual advancement in the treatment of genomic rearrangements. Even though only the high-mappability side can be unambiguously placed, sequence classification of the low-mappability side produces useful results and meaningful insights. The identification of frequent somatic centromeric rearrangements demonstrates the utility single breakend variant calling has in regions of the genome traditionally considered inaccessible to short read sequencing. Single breakend variant calling provides a framework for the reliable detection of LINE integration without a specialised caller [29], for the detection of centromeric and telomeric viral integrations [30], and for an entire ecosystem of tools that explicitly model the ambiguity they represent. As single breakends comprise 7.0% of GRIDSS2 calls in the Hartwig cohort, any purely breakpoint-based caller must have a false negative rate of at least 7.0%. GRIDSS2's 3.1% FNR may thus be impossible to achieve for any breakpoint-based caller, at least for this cohort.

One biologically significant finding coming from GRIDSS2's ability to phase structural variants is the degree to which somatic structural variants are clustered. In the Hartwig cohort of metastatic solid tumours, 22% of somatic structural variants potentially involve DNA fragments of less than 1000 bp with GRIDSS2 able to phase 70% of these. Long read sequencing is considered the gold standard for structural variant phasing and phasability is indeed better: 10 kb long reads increase this theoretical phasability to 31%. The high indel error rate of PacBio and ONT sequencing presents a current drawback of long read sequencing: simple long read based SV detection approaches are likely to misidentify complex rearrangements involving nearby cis-phased SVs. Without HiFi sequencing or error correction prior to alignment, the short DNA segments

between the SVs will be unmappable and the long read caller will report a transitive call between the flanking segments as outlined in Fig. 1b. On COLO829, we found 5 instances (of 67 true positives) in which GRIDSS2 based on short-read sequencing was able to correctly place a short DNA segment that the three long-read callers were not. Care must be taken when comparing or combining short and long read variant calls to ensure the different representations of the same event are reconciled and cis phased. GRIDSS2's ability to phase breakpoints involving short DNA fragments is of great utility to downstream rearrangement event classification and karyotype reconstruction as it exponentially reduces the number of possible paths through the breakpoint graph. GRIDSS2's ability to collapse imprecise transitive calls into their corresponding precise breakpoints is similarly essential to complex event reconstruction as these transitive calls result in spurious false positives that are inconsistent with the actual rearrangement structure. The highly clustered nature of somatic SVs means that short read sequencing is surprisingly competitive when it comes to phasing somatic variants. Sophisticated analysis and interpretation of somatic genomic rearrangements does not necessarily require long read sequencing.

Single breakend variant calling enables a sensitivity and specificity unprecedented amongst short read-based somatic structural variant callers, facilitating the resolution of highly complex rearrangements. While breakpoints and copy number segments are widely adopted fundamental genomic rearrangement signals, single breakends have been hitherto unutilised. Their introduction enables the ambiguities present in low mappability regions to be explicitly modelled without compromising FNR or FDR and their potential extends far beyond the examples presented here. GRIDSS2 demonstrates that single breakend variant calling is essential to the comprehensive characterisation of somatic structural variation from short read sequencing data. Combining single breakend variant calling and structural variant phasing with low FNR and FDR, GRIDSS2 represents a foundation upon which sophisticated somatic analysis can be performed.

## Methods

GRIDSS2 extends the GRIDSS [11] software suite with additional features, tools and capabilities. GRIDSS2 is composed of the following 5 phases: (i) preprocessing, (ii) assembly, (iii) variant calling, (iv) annotation, and (v) somatic filtering.

### Preprocessing

GRIDSS2 takes one or more aligned samples in the SAM/BAM [31] file format. These files are pre-processed on a per-file basis and all reads supporting a structural variant are extracted, and all fields or tags referring to another record are corrected. Reads with chimeric alignments (i.e. split reads), reads with a soft or hard clipped alignment CIGAR of at least 5 bp, read pairs in which only one read is mapped, and discordant read pairs are extracted. The library insert size distribution is estimated from the first 10,000,000 reads using picard tools (http://broadinstitute.github.io/picard) and read pairs considered discordant if they fall outside the 99.5% distribution of fragment size lengths. The clipped bases of soft clipped non-chimeric reads are realigned to the reference genome using bwa mem [32] and converted to a chimeric split read alignment if

an alignment is found. Inconsistencies in the mate chromosome and position are corrected (since tools such as GATK indel realignment adjust read alignment positions without updating the mate record), hard clips converted to soft clips, the NM, SA, MC, MQ tags recalculated, and the R2 tag is populated. Improving performance over GRID SS, GRIDSS2 performs this in a two-pass manner with samtools [31] used for name/coordinate sorting the output of the first/second pass respectively.

As with GRIDSS, reads with low alignment sequence entropy and reads with a mapping quality (mapq) less than 20 (c.f. mapq< 10 GRIDSS) are treated as unmapped, soft-clipped reads with clipped sequence having high homology with standard adapter sequences are ignored, reads marked as duplicates, and regions above 50,000x (c.f. 10,000x) coverage are ignored. Read alignments containing an insertion or deletion under 5 bp are considered consistent with the reference.

### Assembly

GRIDSS2 uses the same genome-wide positional de Bruijn graph break-end assembler used by GRIDSS. Reads are split into kmers and associated positions based on the anchoring alignment: kmers from split reads must be assembled only with kmers at the positions inferred by the anchoring alignment, and kmers of unmapped mate reads are assembled at any position consistent with the library fragment size distribution and the anchoring read alignment position. For assembly purpose, split reads and indel alignments are considered multiple soft clipped alignments, and discordant read pairs are treated as multiple read pairs with one read aligned.

The output of the assembly is a set of "soft-clipped" contigs with anchoring bases supporting the reference and non-reference bases supporting a putative breakpoint at a given position. This contig is iteratively realigned to the reference using bwa mem and converted to a split read alignment. Assemblies longer than the 1.5x maximum fragment size distribution, as well as assemblies supporting the reference sequence are ignored. Assembly alignments with a mapq of less than 20 are treated as unmapped. GRIDSS2 introduces a number of refinements to the assembly calling process.

Assembly support is tracked per base pair. Fragments are considered to support a breakpoint only if the fragment support spans at least one base pair beyond any breakpoint homology on both sides. This ensures that when a single contig spans both a germline indel and a somatic SV, the fragments originating from the matched normal sample will not be considered as supporting the somatic breakpoint. This also improves variant allele fraction calculations in regions of complex rearrangement.

GRIDSS2 performs compound realignment of the entire assembly contig. BWA is used to align the entire assembly contig. Assembly contig bases which are soft clipped in the primary alignment reported by BWA are fed back to BWA for realignment. This process is repeated until either all bases are aligned, or no alignment can be found for the remaining bases. Assembly contigs that do not overlap with the locus of origin of the assembly are filtered out. To ensure that valid assemblies are not unnecessarily filtered, GRIDSS2 includes both reads of each fragment in the assembly, and up to 300 bp of anchoring reference-supporting sequence is included in the assembled contig. The remaining contigs are treated as split read alignments. To rectify over-alignment of the primary alignment location in the presence of imperfect breakpoint

Cameron *et al. Genome Biology*    (2021) 22:202

Page 14 of 25

microhomology, the bounds of each split are adjusted to minimise the edit distance to the reference. The originating alignment is tracked using OA SAM tag and contigs that do not partially align to the originating assembly location and strand are filtered.

gridss.SoftClippedToSplitReads invokes bwa, -L 0,0 is added to the command line to remove the soft-clipping alignment penalty. This prevents 1 bp non-template inserted sequences being over-aligned and reported as clean breakpoints with a flanking SNV.

Worse-case assembly performance has been improved by adding an assembly graph path count threshold. Generating 3 assembly contigs with more than 50,000 alternative paths through the assembly graph without advancing the assembly window will flush the assembly window. The maximum assembly window size has been reduced by 2.5x, and more aggressive assembly read downsampling in high coverage regions is performed.

The presence of a contig with at least three non-overlapping alignments results in the breakpoints supported by that assembly being phased cis. If the initially soft-clipped portion of an assembly realigns across one breakpoint but not another, these breakpoints are phased trans.

### Variant calling

Breakpoints are called using a probabilistic model based on the empirical distribution of CIGAR operators, the library fragment size distribution, and mapping rate. Each read/read pair is given a phred-scaled quality score based on the mapping quality and the probability of encountering that read/read pair given the library empirical distribution. Split reads and soft clipped reads use the distribution of soft clipping CIGAR operators. Discordant read pairs use the discordant mapping mate if distal or the library fragment size distribution if falling within the range reported by Picard tools CollectInsertSizeMetrics. Reads with unmapped mates use the unmapped mate fragment mapping rate and indels based on rate of alignments with insertion/deletion CIGAR elements of matching lengths. As with GRIDSS, split reads and breakpoint-supporting assemblies incorporate the mapping quality scores on both sides of the supported break.

The key novel feature of the GRIDSS2 variant calling processing is the reporting of single breakend variants. Single breakends variants are called based on supporting soft clipped reads, assembly contigs, and reads with unmapped mates. Single breakend calling uses the same two-pass approach as breakpoint calling with all maximum cliques first calculated, then evidence uniquely assigned to the highest scoring clique.

In addition to single breakend variant calling, the variant caller has been improved by reducing the default minimum called event to 10 bp, preferentially allocating reads to variants supported by an assembly containing the read; requiring two supporting fragments to call a variant, and excluding inversion-like breakpoints from the minimum variant size filter to prevent filtering of foldback inversions.

### Annotation

GRIDSS2 provides a full per-sample breakdown of all supporting evidence for each variant through the following VCF INFO and FORMAT fields:

- AS, RAS, CAS: counts of assembly contigs supporting a breakpoint originating locally, from the other side of the breakpoint, and from another location respectively. CAS assemblies support multiple variants and provide linking information about those variants.
- ASSR, ASRP: total number of split/soft clipped/indel-containing reads, and discordant read pairs/reads with unmapped mate contributing to any breakpoint-supporting assembly contig at the breakpoint location. Note that read/read pairs that are assembled into a contig but whose interval of support does not span the breakpoint are not counted. The interval of support for a read/read pair is defined as the interval between the first and the last contig base for which that read/read pair contributed to the assembly.
- SR, RP, IC: counts of split reads and discordantly mapped read pairs and indel-containing reads that directly support the breakpoint.
- BA: counts of assembly contigs support a single breakend at this position. Such contigs are aligned only to the local breakend with the breakend sequence either aligning ambiguously or unable to be aligned to the reference genome by bwa.
- BASSR, BASRP: total number of split reads or soft clipped reads and discordant read pairs or reads contributing to any breakend-supporting assembly contig at the variant location.
- BSC, BUM: counts of soft-clipped reads, and reads with unmapped mates at the variant location
- ASQ, RASQ, CASQ, SRQ, RPQ, IQ, BAQ, BSCQ, BUMQ: corresponding quality score contribution for the supporting evidence.
- QUAL, BQ: total contribution to the breakpoint/breakend quality score.
- BANRP, BANSR, BANRPQ, BANSRQ: counts of read pairs/split reads not supporting this breakpoint but assembled into a contig that supports this breakpoint and their corresponding assembly quality score contribution.
- REF/REFPAIR: count of reads/read pairs spanning the local variant position that support the reference allele. Only reads/read pairs that span across the breakpoint microhomology interval (if present) are counted.
- VF/BVF: count of unique fragments supporting the breakpoint/breakend. By tracking unique fragments supporting the variant, a more accurate variant allele fraction can be calculated. This approach prevents double-counting of discordantly mapped fragments for which one of the reads contains a split read alignment. A fragment can support a variant either directly through split read, soft clipped read or discordant alignment of a read pair, indirectly through incorporation of one or both of the constituent reads in an assembly supporting the variant, or both directly and indirectly.
- RF: count of unique fragments supporting the reference allele.
- CQ: variant quality score prior to evidence reallocation.
- BEALN: potential alignment locations of breakend sequence as determined by *gridss.AnnotateInsertedSequence.*
- BEID, BEIDL, BEIDH: identifiers of assembly contigs and the corresponding local and remote alignment base offsets. Single breakend variants do not have a remote breakend, and only breakpoint variants include breakpoint-supporting assemblies. Variants containing the same BEID are phased cis.

- CIPOS: for IMPRECISE variants, CIPOS encodes the interval in which the breakpoint could occur, and for precise variants, CIPOS encodes the homology interval.
- CIRPOS: corresponding CIPOS of the remote breakend.
- IHOMPOS: interval of inexact homology. A Smith-Waterman alignment of the breakpoint sequence against the reference sequence is performed at both breakends. The reference and breakpoint sequence are extended 300 bp from the break on either side with the reference extended an additional 10 bp to account for potential indels in the alignment. The homology length is the length that the sequence alignment could be extended from the common sequence into the breakpoint/reference sequence. Alignments containing a soft clip on the common sequence side are classified as alignment errors and ignored. The SSW library [33] is used for which we implemented a JNI wrapper. Alignment scored 1, -4, 6, 1 for match, mismatch, gap open, and gap extend respectively which correspond to bwa mem alignment scores.
- SC: CIGAR encoding of the anchoring bases that at least one read/read pair/ assembly is aligned to and supports the variant. This is encoded as a CIGAR string with a match for each anchoring base that provides support for the variant call, XNX for the interval over which the breakpoint could occur (due to microhomology or an imprecise call), and a deletion CIGAR element for any intervals over which there is no support (such as a small flanking deletion). Variants with an anchoring SC 10 bp further from the break than a nearby variant are considered to be phased trans.
- SB: Strand bias of the reads supporting the variant. 1 indicates that reads would be aligned to the positive strand if the reference was changed to the variant allele. 0 indicates that reads bases would be aligned to the negative strand if the reference was changed to the variant allele. Strand bias is calculated purely from supporting reads and exclude read pair support since these are intrinsically 100% strand bias. Note that reads both directly supporting the variant and supporting via assembly will be double-counted. Both breakpoint and breakend supporting reads are included.
- IMPRECISE, HOMLEN, HOMSEQ, PARID, EVENT, CIEND, END, and SVTYPE fields carry their usual meaning as per the VCF file format specifications.
- MQ, MQN, MQX, BMQ, BMQN, BMQX mean, min, and max MAPQ score of reads/assembly contigs providing breakpoint/breakend support.

After initial annotation, *gridss.AnnotateInsertedSequence* aligns any single breakend sequences or non-template inserted breakpoint sequences to an arbitrary reference genome and adds an annotation reporting the alignment location. Integrated viral sequence is identified by aligning to a reference of viral sequences. By default, the same reference as the input files were aligned to is used. If a RepeatMasker bed file generated by Bed-Ops [34] rmsk2bed is supplied, inserted sequences will be annotated with the Repeat-Masker class and type corresponding to the BEALN alignments.

### Somatic filtering

By default, GRIDSS2 is a sensitive caller and reports all putative variants supported by at least two well-mapped reads. To generate a set of high and low confidence somatic

call sets, a somatic filter was developed. Variants with 3% of the supporting reads originating from the normal, or deletion or duplication breakpoints under 1000 bp that have any direct split read support in the normal, are hard filtered. Variants are classified as low confidence if any of the following conditions are met: breakend coverage of less than 8 fragments in the normal; allelic fraction of less than 0.5% in the tumour; imprecise variant call; breakend variants without an assembly containing at least one discordant read pair; single breakends with a poly-C or poly-G run of at least 16 bp in the breakend sequence; deletion or duplication breakpoints under 1000 bp with a split read strand bias of 0.95 or greater; breakpoints with a microhomology of over 50 bp; breakpoints with an inexact microhomology of over 50 bp which are not deletion or duplications under 1000 bp; deletion or duplication breakpoints under 1000 bp with no split read support either directly, or through assembly; breakpoints with no discordant read pair support (either directly, or via assembly) which are not deletions or duplications under 1000 bp; deletion or duplication breakpoints under 1000 bp that have any direct split read support in the normal; 100–800 bp deletion breakpoints with an inexact microhomology length of 6 bp or greater; inversion-like breakpoints 40 bp or less that have at least 6 bp of microhomology; deletion-like breakpoints under 1000 bp whose length of sequence inserted at the breakpoint is within 5 bp of the deletion length, except those whose edit distance to the deleted bases is at least 0.5 per base, and less than 0.2 per base to the reverse complement. Breakpoint variants are filtered if either breakend is filtered.

Somatic variants are panel-of-normal (PON) filtered if a match within 2 bp was found in a panel of normals. The default hg19 was constructed from the 40x coverage WGS matched normals for 3972 patients from the Hartwig cohort using the *gridss.GeneratePonBedpe* utility. If multiple samples for a patient existed, only the normal for the first sample was included in the PON. Variants were aggregated across samples using the default setting of ignoring the FILTER field and excluding imprecise calls and breakpoints/single breakends with a QUAL score of less than 75/428.

Viral insertions are annotated using *gridss.AnnotateInsertedSequence*. Single breakend sequences and non-template inserted sequences that do not have an alignment to the reference genome were aligned to a set of human viral reference sequences. Viral reference sequences were obtained from the virus host database [35] and filtered to include only viruses associated with the homo sapiens taxid of 9606. The viral sequences were then masked using RepeatMasker with "-no_is -s -noint -norna -species human" parameters. Generation scripts can be found at https://github.com/hartwigmedical/scripts/tree/master/virus.

Assembly linking: pairs of breakpoints mutually supported by a common assembly contig were annotated as linked by assembly. For assembly contigs spanning more than 2 breakpoints, each adjacent pair was linked with a unique identifier to enable unambiguous traversal of the breakpoint graph.

Transitive linking: chains of precise breakpoint variants were phased trans if an imprecise spanning transitive breakpoint call could be found. To identify transitive calls, a breadth-first search over the breakpoint graph was performed. Variants were considered transitive if the start and end breakends overlapped the start and end breakpoints in a path of precise breakpoint calls. Paths were limited to 1000 bp and 4 segments, with each segment required to be at least 20 bp in length. Paths could not self-

intersect. To prevent exponential runtime in highly rearranged genomes, at most 100, 000 paths and at most 1000 paths per starting breakpoint were considered.

Simple inversion annotation: pairs of breakpoints with orientations consistent with a simple inversion were annotated as simple inversions if the matching breakends were within 35 bp on both sides, no other simple event annotation could be applied, and fragments supporting the constituent variants differed by at most threefold.

Templated insertion annotation: breakend/breakpoint and breakend/breakend pairs were annotated as simple templated insertions if the breakends had opposite orientations, were within 35 bp, no other simple event annotation could be applied, and fragments supporting the constituent variants differed by at most threefold.

Reciprocal translocation: breakpoint/breakpoint pairs were annotated as reciprocal translocations if the breakends on both sides had opposite orientations, were within 35 bp, no other simple event annotation could be applied, and fragments supporting the constituent variants differed by at most threefold.

Equivalent: variants were annotated as equivalent if variants had a breakend within 5 bp of each other and they shared a common breakend sequence. Breakend sequences were truncated to the length of the shorter sequence and were considered matching when the per-base edit distance between breakend sequences was 0.1 or less. For the purposes of this comparison, the nominal breakend sequence was used for single breakends, and the reference sequence of the partner breakend was used for breakpoint variants. For breakpoint variants, the length of the breakend sequence was the maximum of 20 bases, and the width of the interval over which the fragments supporting the partner breakend had anchoring alignments.

Finally, a quality filter was applied to breakpoint variants with a QUAL score of less than 350 and single breakend variants with a QUAL score of less than 1000. Variants linked to a variant passing the qual filter other than through equivalence were rescued from the quality filtering and were considered to have passed regardless of the actual variant quality score. For each input file, two output files were generated: a high confidence call set containing calls passing all filters and a low confidence call set containing all calls except those failing the normal support filter or short events with split read support in the normal.

### Single breakend simulation

To simulate breakpoints to centromeric sequence, we simulated reads from the complete telomere-to-telomere chromosome 8 assembly of CHM13 [21] (chm13.draft_v1.0.chr8.fasta). To ensure an unambiguous breakpoint position on one side, one side of each breakpoint was located in a region with a unique liftover from CHM13 (t2t-chm13-v1.0.hg38.over.chain) to hg38 to hg19 (Bioconductor liftover hg38ToHg19.over.chain) that was at least 20,100 bp in size with the breakpoint positioned 10,000 bp from the CHM13 start of each liftover region. There were 2998 such regions, so 2998 simulations were performed.

The other side of the breakpoints were equally spaced in the CHM13 region chr8:43, 610,000–50,790,000, that is, half in the region defined as centromeric [21] and half in an equally large region immediately after the centromere.

Fasta files were generated with three 20kbp contigs for each breakpoint: one containing the 10kbp preceding the liftover side of the breakpoint concatenated with the 10kbp after

the breakpoint on the other side, one containing the 10kbp either side of the liftover side of the breakpoint, and a third containing the 10kbp either side of the other side of the breakpoint.

Reads were simulated from the 3 contigs using ART version 2.5.8 using parameters "--noALN --paired --seqSys HSXt -l 150 -m 500 -s 100 -rs 1" to a depth of 60x with the tumour and normal reads output to separate fastq files. Novobreak was run directly off the fastq files. The fastq files were aligned with bwa and sorted with samtools. A single pre-generated BAM containing 2,000,000 reads (1,000,000 read pairs) from CHM13 chr8:1000000-1999999 was merged with both the tumour and normal reads to ensure library fragment size distributions could be calculated (svaba requires a minimum of 1,000,000 read pairs to calculate library fragment size distribution). GRIDSS2, manta, and svaba run off the merged BAM files.

Somatic calls were classified as "Matching Breakpoint" if the both sides were within 100 bp of the hg19 liftover position, "Match Breakend" if the call was a single breakend call and the position was within 100 bp of the liftover side of the breakpoint, "Breakpoint (one side correct)" if one side of the breakpoint was within 100 bp of the breakpoint position, and "False positive" if no coordinates were within 100 bp or if another call in that simulation was already classified in one of the 3 true positive categories. Simulations without a match were classified as a "False Negative". Simulations were classified based on the repeat on the non-liftover side of the breakpoint with repeat class annotations taken from chm13.draft_v1.0_chr8_repeatmasker.out.gz.

Simulation code is available from scripts/gridss2_manuscript/*sim* in the GRIDSS github repository. Tools were installed using "conda create -n gridss2_benchmark_sim manta=1.6.0 gridss=2.11.1 svaba=1.1.0 bwa=0.7.17 samtools=1.11 bioconductor-structuralvariantannotation art=2016.06.05", except for novobreak 1.1.3rc which did not use conda due to a dependency on the version of samtool included in the novobreak release and was installed directly from the novobreak github.

### Independent validation of SV calls

Thirteen samples from the Hartwig metastatic cancer cohort were selected for capture panel validation of the structural variant calls. Each variant called in GRIDSS2 was compared with variants called from Manta (for variants longer than 50 bases) and/or Strelka (for variants from 32 to 50 bases in length) to determine if the variant is shared or private. Variants were marked as matching GRIDSS2 if the start and end chromosomes and orientation both matched and start end positions (including confidence intervals) were within 20 bases of each other. Hybrid capture probes were created for each of the shared and private variants. For each breakpoint variant 3 probes of 120 bases each were created: Two reference probes leading up to the breakends from either side, as well as another SV probe going through the structural variant with the break junction close to the middle of the probe. The reference probes were designed to end 20 bases before the start of each structural variant breakend. The SV probe consists of any insert sequence flanked by equal number of bases from each side of the structural variant. For each single breakend variant 2 probes were created: one reference probe as described above and one SV probe which includes no more than 60 bases of the insert sequence with the remainder coming from the reference leading to the break point.

Together, this created a total of 17,125 capture probes of 120 nt in length, targeting 5821 break-junctions (see Additional File 2: Table S1) which were ordered as custom target capture probes from Twist Biosciences (catalog ID 100533). For each of the 13 samples, 50 ng input DNA was used for indexed library construction with enzymatic fragmentation (Twist kit catalog IDs 100253, 100255, and 100401) according to the manufacturer's protocol. A bead-based size selection was performed after PCR to remove the remaining larger fragments (> 700 bp). Multiplexed hybridization was performed using the Twist Hybridization (ID100254), Blockers (PN100856), and Wash Kits (PN100214, 100215, 100216)) using Dynabeads™ MyOne™ Streptavidin T1 (Invitrogen PN65604D) following standard Twist protocol. Enriched library molecules were amplified by PCR for 11 cycles and sequenced on the Illumina NextSeq500 2 × 150 bp High Output run according to manufacturer's standard protocol.

We created a set of 2000 bp predicted alternate contigs from the shared and private structural variant calls using the same technique from above for generating the (non-reference) SV probes and added these to the reference genome. We then mapped each of the reads from the capture panel output with BWA to a hybrid genome including both the GRCH37 reference genome and the novel alternate contigs.

We assessed the viability of the probes by mapping each of the 120 base SV probes to the 2000 base alternate contigs to determine its mapping quality. Of the 5821 SV probes, 80 had 0 mapping quality and 5377 had a perfect mapping quality of 60. Probes with a mapping quality of less than 20 were ignored as well as 77 micro-satellite probes that were unable to be unambiguously validated. Resultant BAM files were examined for evidence of the SV alternate contigs in the SV source sample BAM as well as the BAM files of each of the other samples as controls for systemic effects for each of the predicted variants. Specifically, reads were considering variant-supporting if the read aligned to the alternate contig and the read alignment spanned at least 25 bp over the break position on both sides. The required spanning length was extended by the length of any breakpoint homology and breakpoint inserted sequence. SV calls were marked as validated if all the following criteria were met (and not validated otherwise):

- At least 2 reads were mapped to the alternate contig in the predicted sample
- The support rate for the alternate contig was significantly higher (Poisson model, p = 0.001) in the predicted sample than the maximum of the other 13 samples
- < 40 reads in total across all 13 control samples were mapped to the predicted alternate contig

Single breakend repeat annotations were determined by running RepeatMasker open-4.0.9 against the single breakend sequences reported in the GRIDSS2 VCF using the "-species human" command line with the default DFAM repeat database.

### Comparison to PCAWG

Copy number data was obtained as for the Hartwig cohort running PURPLE [2] with default settings. PCAWG consensus SV and CN calls were obtained from https://dcc.icgc.org/releases/PCAWG/consensus_sv, and https://dcc.icgc.org/releases/PCAWG/consensus_cnv. Copy number transitions were matched with structural variants with a

100 kb margin for PCAWG calls and a 0 bp margin for GRIDSS2/PURPLE. Copy number transitions in or within 100 kb of centromeres or a gap in the reference genome were excluded from analysis. Copy number transitions matched by both a single breakend and a breakpoint, were considered breakpoint matches.

GRIDSS2/PURPLE was run with default settings on 1774 ICGC PCAWG WGS samples. BAMs were converted to fastq and processed using the Hartwig pipeline (https://www.hartwigmedicalfoundation.nl/it-pipeline/). Samples were excluded from analyses if:

– PURPLE quality control failed (QCStatus ! = PASS)
– Multiple specimens for the same donor had WGS data
– Consensus SV or CNV calls were not available

### Hartwig metastatic tumour cohort

GRIDSS2 was run on 3782 paired tumour/normal samples from the Hartwig Medical Foundation cohort of metastatic solid cancers with a 32 bp minimum event size. Samples were aligned with bwa against a GRCH37 reference genome containing only primary contigs. Single breakend RepeatMasker annotations were obtained by running *gridss.AnnotateInsertedSequence* against the UCSC GRCH37 (hg19) RepeatMasker track downloaded from http://hgdownload.cse.ucsc.edu/goldenpath/hg19/bigZips/hg19.fa.out.gz after converting to BED formart using bedops *rmsk2bed*.

Hartwig copy number was determined by PURPLE. Since PURPLE infers the copy number of short segments by the VAF of the flanking SVs, the copy number of these segments does not represent an independent validation of the SV. As such, FDR was segments from only the breakpoints in which the copy number of all 4 flanking segments was determined of depth of coverage and SNP BAF were considered.

Single breakends with a RepeatMasker annotation associated with centromeric or pericentromeric repeats were considered centromeric single breakends. The matching chromosome was considered to be the chromosome for which the BLAT [36] based score score = (1000-((9-floor(Qsize/100))*mismatch+Qcount+Tcount))*min(match/Qsize,1)) is at least 900 and at least 25 higher than the best alignment on a different chromosome when aligning against hg38. A script for annotating likely centromere can be found in example/annotate_most_likely_centromere.R in the GRIDSS repository.

Phasability of the Hartwig cohort was calculated by determining, for each break junction, the length of the DNA segment if it was phased with the first break junction encountered in the appropriate orientation. Known phasing information was ignored for this analysis. Expected phasability was calculated by simulating 3782 randomly fragmented paired genomes with the same number of break junctions as the corresponding Hartwig sample.

For both the PCAWG and Hartwig cohort, rearrangement event classifications were obtained by running LINX [29] 1.12 on the GRIDSS2/PURPLE outputs. Simulated FNR results were obtained by random subsampling of the Hartwig GRIDS2 SV calls and breaking LINX chains whenever a SV was excluded from the subsampling.

### COLO829 somatic benchmark

The COLO829T/COLO829BL cell lines (ATCC® CRL-1974™ and 1980™ respectively) were each sequenced three times to 100x/40x using the HMF workflow [2] and aligned

against GRCH37 without alt contigs using BWA 0.7.15. GRIDSS 2.9.3, GRIDSS 1.3.3, Manta 1.5.0, svaba 1.1.0, and smufin 0.9.3 [37] were run with default parameters. Programs were allocated 8, 16, 16, and 20 cores and 32, 32, 50, and 500Gb of memory respectively. No smufin results were obtained in any replicate as smufin failed to complete in the 100,000 CPU hours/3 months wall time allocated.

Call matching was performed using the StructuralVariantAnnotation BioConductor package (DOI 10.18129/B9.bioc.StructuralVariantAnnotation). A 100 bp matching margin was allowed around the break junction position. Tandem duplication calls matched with insertion calls if the size difference between the duplication and insertion was within 25 bp. False positive calls under 50 bp were filtered after matching so as not to penalise a caller reporting an event slightly larger than 50 bp in the truth set, but slightly smaller than 50 bp in the call set. If multiple calls in a call set matched a single truth set call, all except the highest QUAL call were ignored. GRIDSS1 calls were filtered to call with 0 supporting reads in the normal. The GRIDSS2 PON filter was applied to all callers.

### COLO829 truth set generation

The COLO829 somatic SV truthset was generated using an orthogonal sequencing strategy. We sequenced the COLO829BL and COLO829T cell lines using Illumina HiSeqX (ENA run accessions ERR2752449 and ERR2752450 for COLO829BL and COLO829T, respectively), Oxford Nanopore (ERR2752451 and ERR2752452), Pacific Biosciences (ERR2752447 and ERR2752448), and 10X genomics (ERR2820166 and ERR2820167). All data is grouped under ENA study accession PRJEB27698.

Raw data was analysed for structural variants using the following tools:

– Illumina data was mapped using BWA 0.7.5, SV calling was performed with GRID SS 2.0.1 and somatic SVs were filtered using gridss_somatic_filter.R.
– Nanopore data was mapped using NGMLR 0.2.6 and SV calling was performed with both NanoSV 1.2.0 and Sniffles 1.0.8 separately for COLO829T and COLO829BL. All SVs were merged with an overlap window of 200 base pairs using SURVIVOR 1.0.6 and SVs not present in COLO829BL were kept.
– Pacbio data was mapped using minimap2 2.11-r797 and SVs were called using pbsv 2.1.0 in joint calling mode for COLO829T and COLO829BL. Only SVs with no evidence in COLO829BL were kept.
– 10X genomics data was processed using Longranger 2.2.2 with default settings for COLO829BL and somatic mode for COLO829T. SV calls for both cell lines were merged with an overlap window of 200 base pairs using SURVIVOR 1.0.6 and SVs not present in COLO829BL were kept.

Somatic SV calls for each technology were merged with an overlap window of 200 base pairs using SURVIVOR 1.0.6., and all candidate breakpoints were subjected to independent validation by targeted capture and/or PCR-based approaches. SVs detected with two or more techniques that failed in these validation experiments were curated

by manual inspection of the mapped reads using IGV [38]. A total of 69 SVs were finally considered as true somatic SVs for COLO829T.

## Supplementary Information

The online version contains supplementary material available at https://doi.org/10.1186/s13059-021-02423-x.

**Additional file 1.** Fig.S1-S4. Fig. S1 COLO829 results with/without PON filter. Fig. S2 PCAWG/GRIDSS, and PURPLE/PCAWG SV/CNV consistency. Fig. S3 PCAWG single breakend annotations. Fig. S4 Phasability and LINX classification for Hartwig and PCAWG.

**Additional file 2.** Tables S1-S2. Table S1. Probe validation results. Table S2. Small duplication driver predictions.

**Additional file 3.** Review history.

### Availability of data and materials

GRIDSS2 is freely available as open source software from https://github.com/PapenfussLab/gridss [39] under a GPLv3 license.
Analysis scripts used to generate results are available from https://github.com/PapenfussLab/gridss/tree/master/scripts/gridss2_manuscript [39]. Source code for the GRIDSS2 version used is also available at doi.org/10.5281/zenodo.4739928 [40].
Hartwig cohort data was obtained from the Hartwig Medical Foundation (Data request DR-005 "Development and validation of tumour genome analysis tools"). Standardised procedures and request forms for access to this data can be found at https://www.hartwigmedicalfoundation.nl/en.
Raw and analysed data for the creation of the COLO829T/COLO829BL tumour/normal cell line pair structural variant truth set are available grouped under ENA study accession PRJEB27698 [41]. The COLO829 truth VCF is available from https://github.com/UMCUGenetics/COLO829_somaticSV.
Capture panel validations of 13 patient tumour samples are available under the controlled access dataset accession EGAD00001005525 [42].

## Declarations

### Ethics approval and consent to participate

Patient data was obtained from the Hartwig Medical Foundation through data request DR-005 "Development and validation of tumor genome analysis tools." The data access request was reviewed and approved by the Hartwig Medical Foundation Data Access Board, and all patient data has been used in accordance with informed consent given by the patients.

### Competing interests

The authors declare no competing financial interests.

**Author details**
[1]Bioinformatics Division, Walter and Eliza Hall Institute of Medical Research, Parkville, Australia. [2]Department of Medical Biology, University of Melbourne, Melbourne, Australia. [3]Hartwig Medical Foundation Australia, Sydney, Australia. [4]Hartwig Medical Foundation, Science Park 408, Amsterdam, The Netherlands. [5]Center for Molecular Medicine and Oncode Institute, University Medical Center Utrecht, Heidelberglaan 100, Utrecht, The Netherlands. [6]Peter MacCallum Cancer Centre, Melbourne, Australia. [7]Sir Peter MacCallum Department of Oncology, University of Melbourne, Melbourne, Australia.

**References**
1.  Li Y, et al. Patterns of somatic structural variation in human cancer genomes. Nature. 2020;578(7793):112–21. https://doi.org/10.1038/s41586-019-1913-9.
2.  Priestley P, Baber J. Pan-cancer whole-genome analyses of metastatic solid tumours. Nature. 2019;575(7781):210–6. https://doi.org/10.1038/s41586-019-1689-y.
3.  Kosugi S, Momozawa Y. Comprehensive evaluation of structural variation detection algorithms for whole genome sequencing. Genome Biol. 2019;20(1):117. https://doi.org/10.1186/s13059-019-1720-5.
4.  Cameron, D. L., Di Stefano, L. & Papenfuss, A. T. Comprehensive evaluation and characterisation of short read general-purpose structural variant calling software. Nat Commun 10, 3240 (2019), 1, DOI: 10.1038/s41467-019-11146-4.
5.  Danecek P, Auton A. The variant call format and VCFtools. Bioinformatics. 2011;27(15):2156–8. https://doi.org/10.1093/bioinformatics/btr330.
6.  Fan, X., Abbott, T. E., Larson, D., Chen, K. BreakDancer: identification of genomic structural variation from paired-end read mapping. Current Protocol Bioinformatics 15.6.1–15.6.11 (2014) doi: https://doi.org/10.1002/0471250953.bi1506s45.
7.  Schröder J, Hsu A. Socrates: identification of genomic rearrangements in tumour genomes by re-aligning soft clipped reads. Bioinformatics. 2014;30(8):1064–72. https://doi.org/10.1093/bioinformatics/btt767.
8.  Ye K, Schulz MH, Long Q, Apweiler R, Ning Z. Pindel: a pattern growth approach to detect break points of large deletions and medium sized insertions from paired-end short reads. Bioinformatics. 2009;25(21):2865–71. https://doi.org/10.1093/bioinformatics/btp394.
9.  Wang J, Mullighan CG. CREST maps somatic structural variation in cancer genomes with base-pair resolution. Nat Methods. 2011;8(8):652–4. https://doi.org/10.1038/nmeth.1628.
10. Liu S, et al. Discovery, genotyping and characterization of structural variation and novel sequence at single nucleotide resolution from de novo genome assemblies on a population scale. Gigascience. 2015;4(1):64. https://doi.org/10.1186/s13742-015-0103-4.
11. Cameron DL, Schröder J. GRIDSS: sensitive and specific genomic rearrangement detection using positional de Bruijn graph assembly. Genome Res. 2017;27(12):2050–60. https://doi.org/10.1101/gr.222109.117.
12. Sindi SS, Onal S, Peng LC, Wu H-T, Raphael BJ. An integrative probabilistic model for identification of structural variation in sequencing data. Genome Biol. 2012;13(3):R22. https://doi.org/10.1186/gb-2012-13-3-r22.
13. Aganezov S, Zban I, Aksenov V, Alexeev N, Schatz MC. Recovering rearranged cancer chromosomes from karyotype graphs. BMC Bioinformatics. 2019;20(S20):641. https://doi.org/10.1186/s12859-019-3208-4.
14. Baca SC, Prandi D. Punctuated evolution of prostate cancer genomes. Cell. 2013;153(3):666–77. https://doi.org/10.1016/j.cell.2013.03.021.
15. Cortés-Ciriano I, et al. Comprehensive analysis of chromothripsis in 2,658 human cancers using whole-genome sequencing. Nat Genet. 2020;52(3):331–41. https://doi.org/10.1038/s41588-019-0576-7.
16. Cretu Stancu, M., van Roosmalen M. J. Mapping and phasing of structural variation in patient genomes using nanopore sequencing. Nat Commun 8, 1326 (2017), 1, DOI: https://doi.org/10.1038/s41467-017-01343-4.
17. Valle-Inclan, J. E., Besselink, N. J. M., de Bruijn, E. A multi-platform reference for somatic structural variation detection. bioRxiv (2020).
18. Chen X, Schulz-Trieglaff O. Manta: rapid detection of structural variants and indels for germline and cancer sequencing applications. Bioinformatics. 2016;32(8):1220–2. https://doi.org/10.1093/bioinformatics/btv710.
19. Wala JA, Bandopadhayay P. SvABA: genome-wide detection of structural variants and indels by local assembly. Genome Res. 2018;28(4):581–91. https://doi.org/10.1101/gr.221028.117.
20. Chong Z, Chen K. Structural Variant breakpoint detection with novoBreak. Methods Mol Biol. 2018;1833:129–41. https://doi.org/10.1007/978-1-4939-8666-8_10.
21. Logsdon GA, et al. The structure, function and evolution of a complete human chromosome 8. Nature. 2021. https://doi.org/10.1038/s41586-021-03420-7.
22. Cameron DL, et al. VIRUSBreakend: Viral Integration recognition using single breakends. Bioinformatics. 2021. https://doi.org/10.1093/bioinformatics/btab343.
23. Saunders CT, et al. Strelka: accurate somatic small-variant calling from sequenced tumor-normal sample pairs. Bioinformatics. 2012;28:1811–7.
24. Layer RM, Chiang C, Quinlan AR, Hall IM. LUMPY: a probabilistic framework for structural variant discovery. Genome Biol. 2014;15(6):R84. https://doi.org/10.1186/gb-2014-15-6-r84.
25. Rausch T, Zichner T, Schlattl A, Stutz AM, Benes V, Korbel JO. DELLY: structural variant discovery by integrated paired-end and split-read analysis. Bioinformatics. 2012;28(18):i333–9. https://doi.org/10.1093/bioinformatics/bts378.
26. ICGC/TCGA Pan-Cancer Analysis of Whole Genomes Consortium. Pan-cancer analysis of whole genomes. Nature. 2020;578:82–93.
27. Tubio JMC, et al. Mobile DNA in cancer. Extensive transduction of nonrepetitive DNA mediated by L1 retrotransposition in cancer genomes. Science. 2014;345:1251343.
28. Cameron DL, Baber J, Shale C, Papenfuss AT. GRIDSS, PURPLE, LINX: Unscrambling the tumor genome via integrated analysis of structural variation and copy number. bioRxiv; 2019.

29.   Shale C, Baber J, Cameron DL, Wong M, Cowley MJ, Papenfuss AT, et al. Unscrambling cancer genomes via integrated analysis of structural variation and copy number. bioRxiv; 2020.

30.   Cameron, D. L., Papenfuss, A. T. VIRUSBreakend: viral integration recognition using single breakends. doi: https://doi.org/10.1101/2020.12.09.418731.

31.   Li H, et al. The Sequence alignment/map format and SAMtools. Bioinformatics. 2009;25:2078–9.

32.   Li H, Durbin R. Fast and accurate long-read alignment with Burrows-Wheeler transform. Bioinformatics. 2010;26(5):589–95. https://doi.org/10.1093/bioinformatics/btp698.

33.   Zhao M, Lee W-P, Garrison EP, Marth GT. SSW library: an SIMD Smith-Waterman C/C++ library for use in genomic applications. PLoS One. 2013;8(12):e82138. https://doi.org/10.1371/journal.pone.0082138.

34.   Neph S, et al. BEDOPS: high-performance genomic feature operations. Bioinformatics. 2012;28:1919–20.

35.   Mihara T, Nishimura Y, Shimizu Y, Nishiyama H, Yoshikawa G, Uehara H, et al. Linking virus genomes with host taxonomy. Viruses. 2016;8(3):66. https://doi.org/10.3390/v8030066.

36.   Kent WJ. BLAT--the BLAST-like alignment tool. Genome Res. 2002;12:656–64.

37.   Moncunill V, et al. Comprehensive characterization of complex structural variations in cancer by directly comparing genome sequence reads. Nat Biotechnol. 2014;32:1106–12.

38.   Robinson JT, Thorvaldsdóttir H, Winckler W, Guttman M, Lander ES, Getz G, et al. Integrative genomics viewer. Nature Biotechnol. 2011;29(1):24–6. https://doi.org/10.1038/nbt.1754.

39.   Cameron, D. L., GRIDSS - the Genomic Rearrangement IDentification Software Suite, github, https://github.com/PapenfussLab/gridss, 2021.

40.   Cameron DL. GRIDSS version 2.11.1. zedono. https://zenodo.org/record/4739928. 2021.

41.   Valle-Inclan J, Besselink NJM, de Bruijn E, Cameron DL, Ebler J, Kutzera J, et al. Whole genome sequencing of the COLO829 reference cancer cell line. European Nucleotide Archive. https://www.ebi.ac.uk/ena/browser/view/PRJEB27698. 2021.

42.   Besselink N.J.M. Validation data for the SV analysis package: GRIDSS, PURPLE, LINX, EGAD00001005525, European Genome-Phenome Archive, https://ega-archive.org/datasets/EGAD00001005525, 2021.

## Publisher's Note