

1 **Comprehensive, integrated, and phased whole-genome analysis of the primary ENCODE**  
2 **cell line K562**

3 Bo Zhou<sup>1,2</sup>, Steve S. Ho<sup>1,2</sup>, Stephanie U. Greer<sup>3</sup>, Xiaowei Zhu<sup>1,2</sup>, John M. Bell<sup>4</sup>, Joseph G.  
4 Arthur<sup>5</sup>, Noah Spies<sup>2,6,7</sup>, Xianglong Zhang<sup>1,2</sup>, Seunggyu Byeon<sup>8</sup>, Reenal Pattni<sup>1,2</sup>, Noa Ben-  
5 Efraim<sup>1,2</sup>, Michael S. Haney<sup>1,2</sup>, Rajini R. Haraksingh<sup>1,2,9</sup>, Hanlee P. Ji<sup>3,4</sup>, Giltae Song<sup>8</sup>, Dimitri  
6 Perrin<sup>10</sup>, Wing H. Wong<sup>5,11</sup>, Alexej Abyzov<sup>12</sup>, Alexander E. Urban<sup>1,2</sup>

7  
8 <sup>1</sup>Department of Psychiatry and Behavioral Sciences, Stanford University School of Medicine,  
9 Stanford, California 94305, USA

10 <sup>2</sup>Department of Genetics, Stanford University School of Medicine, Stanford, California 94305,  
11 USA

12 <sup>3</sup>Division of Oncology, Department of Medicine, Stanford University School of Medicine,  
13 Stanford, California 94305, USA

14 <sup>4</sup>Stanford Genome Technology Center, Stanford University, Palo Alto, California 94304, USA

15 <sup>5</sup>Department of Statistics, Stanford University, Stanford, California 94305, USA

16 <sup>6</sup>Department of Pathology, Stanford University School of Medicine, Stanford, California 94305,  
17 USA

18 <sup>7</sup>Genome-Scale Measurements Group, National Institute of Standards and Technology,  
19 Gaithersburg, Maryland 20899, USA

20 <sup>8</sup>School of Computer Science and Engineering, College of Engineering, Pusan National  
21 University, Busan 46241, South Korea

22 <sup>9</sup>Current affiliation: Department of Life Sciences, The University of the West Indies, Saint  
23 Augustine, Trinidad and Tobago

24 <sup>10</sup>Science and Engineering Faculty, Queensland University of Technology, Brisbane, QLD 4001,  
25 Australia

26 <sup>11</sup>Department of Biomedical Data Science, Stanford University School of Medicine, Stanford,  
27 California 94305, USA

28 <sup>12</sup>Department of Health Sciences Research, Center for Individualized Medicine, Mayo Clinic,  
29 Rochester, Minnesota 55905, USA

30

31 **Corresponding author:**

32

33 Alexander E. Urban, Ph.D

34 Department of Psychiatry and Behavioral Sciences

35 Department of Genetics

36 Stanford Center for Genomics and Personalized Medicine

37 Tasha and John Morgridge Faculty Scholar, Stanford Child Health Research Institute

38 3165 Porter Drive, Room 2180

39 Palo Alto, CA 94304-1213

40 USA

41 [aeurban@stanford.edu](mailto:aeurban@stanford.edu)

42

43 **Running title:** Comprehensive whole-genome analysis of K562

44

45 **Keywords:** ENCODE, K562, structural variation (SV), haplotype phasing, retrotransposon  
46 insertions, CRISPR, allele-specific expression (ASE), allele-specific methylation (ASM)

1 **ABSTRACT**

2  
3 K562 is widely used in biomedical research. It is one of three tier-one cell lines of ENCODE and  
4 also most commonly used for large-scale CRISPR/Cas9 screens. Although its functional  
5 genomic and epigenomic characteristics have been extensively studied, its genome sequence  
6 and genomic structural features have never been comprehensively analyzed. Such information  
7 is essential for the correct interpretation and understanding of the vast troves of existing  
8 functional genomics and epigenomics data for K562. We performed and integrated deep-  
9 coverage whole-genome (short-insert), mate-pair, and linked-read sequencing as well as  
10 karyotyping and array CGH analysis to identify a wide spectrum of genome characteristics in  
11 K562: copy numbers (CN) of aneuploid chromosome segments at high-resolution, SNVs and  
12 Indels (both corrected for CN in aneuploid regions), loss of heterozygosity, mega-base-scale  
13 phased haplotypes often spanning entire chromosome arms, structural variants (SVs) including  
14 small and large-scale complex SVs and non-reference retrotransposon insertions. Many SVs  
15 were phased, assembled, and experimentally validated. We identified multiple allele-specific  
16 deletions and duplications within the tumor suppressor gene *FHIT*. Taking aneuploidy into  
17 account, we re-analyzed K562 RNA-seq and whole-genome bisulfite sequencing data for allele-  
18 specific expression and allele-specific DNA methylation. We also show examples of how  
19 deeper insights into regulatory complexity are gained by integrating genomic variant information  
20 and structural context with functional genomics and epigenomics data. Furthermore, using  
21 K562 haplotype information, we produced an allele-specific CRISPR targeting map. This  
22 comprehensive whole-genome analysis serves as a resource for future studies that utilize K562  
23 as well as a framework for the analysis of other cancer genomes.

24  
25

## 1 INTRODUCTION

2 K562 is an immortalized chronic myelogenous leukemia (CML) cell line derived from a  
3 53-year-old Caucasian female in 1970 (Lozzio and Lozzio 1975). Since being established, K562  
4 has been widely used in biomedical research as a “work-horse” cell line, resulting in over 17,000  
5 publications to date. In most cases, its use is similar to that of a model organism, contributing to  
6 the understanding of fundamental human biological processes as well as to basic and  
7 translational cancer research (Grzanka et al. 2003; Drexler et al. 2004; Butler and Hirano 2014).  
8 Along with the H1 human embryonic stem cell line and the GM12878 lymphoblastoid cell line,  
9 K562 is one of the three tier-one cell lines of the ENCyclopedia Of DNA Elements Project  
10 (ENCODE) (The ENCODE Project Consortium 2012), forming the basis of over 1,300 ENCODE  
11 datasets to date (Sloan et al. 2016). Furthermore, it is also one of the few cell lines most  
12 commonly used for large-scale CRISPR/Cas9 gene-targeting screens (Wang et al. 2015; Arroyo  
13 et al. 2016; Morgens et al. 2016; Han et al. 2017; Adamson et al. 2016; Liu et al. 2017).

14 Although the functional genomic characteristics of K562 have been extensively studied  
15 and documented, reflected in close to 600 ChIP-seq, 400 RNA-seq, 50 DNase-Seq, and 30  
16 RIP-Seq datasets available through the ENCODE portal (Sloan et al. 2016), the sequence and  
17 structural features of the K562 genome have never been comprehensively characterized, even  
18 though past cytogenetic studies using G-banding, fluorescence *in situ* hybridization (FISH),  
19 multiplex-FISH, and comparative genomic hybridization (CGH) showed that K562 cells contain  
20 pervasive aneuploidy and multiple gross structural abnormalities (Selden et al. 1983; Wu et al.  
21 1995; Naumann et al. 2001; Gribble et al. 2000), not unexpected for a cancer cell line. In other  
22 words, the rich amount of K562 functional genomics and epigenomics work conducted to date,  
23 in particular integrative analyses that have been carried out in various settings using the vast  
24 troves of K562 ENCODE data, were done without taking into account the many differences of  
25 the K562 genome relative to the human reference genome. This leads to skewed interpretations

1 and reduces the amount of knowledge that can be gained from the rich, multi-layered ENCODE  
2 datasets that continue to accumulate.

3 Here, we report for the first time a comprehensive characterization of the K562 genome  
4 that include copy numbers (CN) of chromosome segments at high-resolution, single-nucleotide  
5 variants (SNVs, also including single-nucleotide polymorphisms, i.e. SNPs) and small insertions  
6 and deletions (Indels) with allele-frequencies corrected by CN in aneuploid regions, loss of  
7 heterozygosity, mega-base-scale phased haplotypes often spanning entire chromosome arms,  
8 and structural variants (SVs) including small and large-scale complex SVs with phasing. We  
9 then took first steps into exploring how knowledge about genome sequence and structural  
10 features can influence the interpretation of functional genomics and epigenomics data and show  
11 examples of how deeper insights into genome regulatory complexity can be obtained by  
12 integrating genomic context. These insights also shed light on important questions regarding  
13 cancer evolution.

## 14 **RESULTS**

### 15 **Karyotyping**

16 The K562 cell line exhibits pervasive aneuploidy (Fig. 2A). Analysis of 20 individual K562  
17 cells using GTW banding showed that all cells demonstrated a near-triploid karyotype and are  
18 characterized by multiple structural abnormalities. The karyotype of our line of K562 cells is  
19 overall consistent (although not identical) with previously published karyotypes (Selden et al.  
20 1983; Wu et al. 1995; Naumann et al. 2001; Gribble et al. 2000), suggesting that its near-triploid  
21 state arose during leukemogenesis or early in the establishment of the cell line. It also suggests  
22 that different K562 cell lines kept and passaged in different laboratories may exhibit some  
23 additional karyotypic differences. Although the karyotype for all chromosomes in our K562 cell  
24 line was supported by previous karyotype analyses, slight variations do exist among the various  
25 published analyses (Supplemental Table S1) with chromosomes 10, 12, and 21 showing the  
26 most variability.

## 1 **Identification of Copy Number (CN) by Chromosome Segments**

2 We used read-depth analysis (Abyzov et al. 2011) to assign a CN i.e. ploidy to all  
3 chromosome segments at 10kb-resolution or entire chromosomes in the K562 genome (Fig. 1,  
4 Supplemental Table S2). We first calculated WGS coverage in 10 kb bins across the genome  
5 and plotted it against %GC content where five distinct clusters were clearly observed  
6 (Supplemental Fig. S2). Clusters were designated as corresponding to particular CNs based on  
7 the mean coverage of each cluster (Supplemental Methods). Such designations confirm that the  
8 triploid state is the most common in the K562 genome. The CN assigned to all chromosome  
9 segments using this approach are consistent with array CGH (Supplemental Fig. S3,  
10 Supplemental Data) and also with previous CGH analyses (Gribble et al. 2000; Naumann et al.  
11 2001) with minor differences on chromosomes 7, 10, 11, and 20 (Supplemental Table S3).  
12 While on a general level, the CNs identified based on read-depth analysis tracks the findings  
13 from karyotyping, read-depth analysis reveals the CNs of many chromosome segments that  
14 would not have been apparent from karyotyping alone (Supplemental Fig. S3, Supplemental  
15 Data, Supplemental Table S2). We see that 53.5% of the K562 genome has a baseline CN of  
16 three (consistent with the karyotype, Fig. 2A), 16.9% CN of four, 1.9% CN of five, 2.4% CN of  
17 one, and only 30.0% has remained in a diploid state (Figure 2B). In addition, two large regions  
18 (5.8 Mb and 3.1 Mb in size) on chromosome 9 (20,750,000-26,590,000 and 28,560,000-  
19 31,620,000 respectively) were lost entirely (Supplemental Table S2).

## 20 **SNVs and Indels**

21 We identified SNVs and Indels in the K562 genome. By taking into account the CN of  
22 the chromosomal segments in which they reside, we assigned heterozygous allele frequencies  
23 to these variants, including non-conventional frequencies (e.g. 0.33 and 0.67 in triploid regions;  
24 0.25, 0.50, and 0.75 in tetraploid regions). Using this approach, we detected and genotyped a  
25 total of 3.09 M SNVs (1.45 M heterozygous, 1.64 M homozygous) and 0.70 M Indels (0.39 M  
26 heterozygous, 0.31 M homozygous) (Table 1, Dataset S1). Interestingly, there are 13,471

1 heterozygous SNVs and Indels that have more than two haplotypes in aneuploid regions where  
2 CN is >2 (Dataset S1). Furthermore, chromosomes 3, 9, 13, 14, and X along with large  
3 stretches of chromosomes 2, 10, 12, 17, 20, and 22 show striking loss of heterozygosity (LOH)  
4 (Fig. 1 and Supplemental Table S4). While a normal tissue sample corresponding to K562 is not  
5 available for comparative analysis, we overlapped these SNVs and Indels with those in  
6 dbSNP138 (Sherry et al. 2001) and found the overlap to be 98% and 79% respectively (Fig. 2C,  
7 Dataset S1), suggesting an accumulation of a significant number of K562-specific SNVs and  
8 Indels relative to germline variants present in the population. After filtering for protein-altering  
9 SNVs and Indels in K562 that overlap with those identified from the 1000 Genomes Project or  
10 from the Exome Sequencing Project, we found that 424 SNVs and 148 Indels are private  
11 protein-altering (PPA) (Table 1, Supplemental Table S5). Furthermore, the overlap between the  
12 PPA variants and the Catalogue of Somatic Mutations in Cancer (COSMIC) is 53% and 31% for  
13 SNVs and Indels respectively (Supplemental Table S6). Eighteen genes that acquired PPA  
14 variants overlap with the Sanger Cancer Gene Census; canonical tumor suppressor genes and  
15 oncogenes such as *RAD51B*, *TP53*, *PDGFRA*, *RABEP1*, *EPAS1*, and *WHIS1* are notably  
16 present among them (Supplemental Table S7).

## 17 **Haplotype Phasing**

18 We performed haplotype phasing for the K562 genome by performing 10x Genomics  
19 linked-read library preparation and sequencing (Zheng et al. 2016; Marks et al. 2018). This  
20 library was sequenced ( $2 \times 151$  bp) to 59x genome coverage. Post sequencing quality-control  
21 analysis showed that 1.06 ng, or approximately 320 genome equivalents, of high molecular  
22 weight (HMW) K562 genomic DNA fragments (average fragment size = 59 kb, 95.3% >20kb,  
23 11.9% >100 kb) were partitioned into 1.56 million oil droplets for uniquely barcoding (16 bp)  
24 within each droplet. Half of all reads come from HMW DNA molecules with at least 64 linked-  
25 reads (N50 Linked-Reads per Molecule or LPM) (Table 1). We estimate the actual physical  
26 coverage ( $C_F$ ) to be 191x. The overall sequencing coverage is  $C = C_R \times C_F = 59x$ . The length of

1 sequence coverage per  $2 \times 151$  bp paired-ended read minus 16 bp of “HWM fragment barcode”  
2 is 286 bp, thus coverage ( $C_R$ ) of the average input HMW genomic DNA (59 kb) is 18,304 bp  
3 (286 bp  $\times$  64 linked-reads) or 31.0% of 50 kb. Using Long Ranger (Marks et al. 2018), 1.41 M  
4 (97.2%) of heterozygous SNVs and 0.58 M (83.7%) of Indels (previously identified, Dataset S1)  
5 were successfully phased into 4,987 haplotype blocks (Fig. 1, Table 1, and Dataset S2). The  
6 longest is 11.95 Mb ( $N_{50} = 2.72$  Mb) (Fig. 2D, Table 1, and Dataset S2); however, haplotype  
7 block lengths vary widely across different chromosomes (Supplemental Fig. S4, Fig. 1) with  
8 poorly phased regions corresponding to regions with LOH (Fig. 1, Supplemental Table S4,  
9 Dataset S2).

### 10 **Mega-Haplotypes Encompassing Entire Chromosome Arms**

11 Leveraging the haplotype imbalance in aneuploid regions, we constructed mega-  
12 haplotypes (Table 2, Supplemental Data), often encompassing entire K562 chromosome arms,  
13 by “stitching” the phased haplotype blocks derived from linked-reads using a recently published  
14 method (Bell et al. 2017). Briefly, we counted the number of linked-read barcodes for each  
15 phased heterozygous SNVs assigned to haplotype blocks that contain  $\geq 100$  phased SNVs  
16 (Dataset S2). Since each barcode is specific to a given HMW DNA molecule, the total number  
17 of unique barcodes is directly associated with the number of individual HMW DNA molecules  
18 sequenced. In other words, the counting of unique barcode associated with a particular  
19 sequence gives the fractional representation of that sequence (or genomic locus). Thus, for  
20 each phased haplotype in aneuploid regions with  $CN > 2$ , major and minor haplotypes can be  
21 assigned according to the number of barcodes associated with each haplotype (Fig. 3), where  
22 the major haplotype simply has more associated unique barcodes than the minor. In diploid  
23 regions, the two haplotypes are expected to have similar barcode counts. A matched normal  
24 control genome is required in order to confidently discriminate between the major and minor  
25 haplotypes in a case genome (Bell et al. 2017). Because K562 has no matching normal sample,  
26 we used a female genome (NA12878) for which linked-read data is publicly available and which



1 is of the same ethnicity as K562 (Fig. 3). After verifying aneuploidy (or haplotype imbalance) by  
2 barcode counting and performing the normalization procedures and statistical tests as described  
3 in (Bell et al. 2017), we then “stitched” together contiguous phased haplotype blocks based on  
4 the imbalance between the major and minor haplotypes. Using this approach, a total of 31  
5 autosomal mega-haplotypes were constructed (Table 2, Supplemental Data); 15 of which  
6 encompass entire (or >95%) chromosome arms such as 19p, 19q, 10p, 7p, and 5q (Fig. 3).  
7 The average mega-haplotype is 50.7 Mb or approximately 4 times longer than the longest  
8 phased haplotype block from Long Ranger (Fig. 2D, Table 1, Dataset S2, Table 2). The longest  
9 mega-haplotype is approximately 137 Mb long (4q).

#### 10 **Identification and Reconstruction of Structural Variants (SVs) from Linked-Reads**

11 In addition to phasing, another use for the linked-read sequencing data is to identify  
12 breakpoints of large-scale SVs by searching for the discordant mapping of clusters of linked-  
13 reads carrying the same barcodes. The identified SVs can then also be assigned to specific  
14 haplotypes if the breakpoint-supporting reads contain phased SNVs or Indels (Zheng et al.  
15 2016). Using this approach, which is also implemented by the Long Ranger software from 10x  
16 Genomics, we identified 186 large SVs >30 kb (98% phased) (Dataset S3) and 3,541 deletions  
17 between 50 bp and 30 kb (79% phased) (Dataset S4). The large SVs include deletions,  
18 inversions, duplications, and inter- and intra-chromosomal rearrangements (Dataset S3, Fig.  
19 4A). As expected, we detected the *BCR/ABL1* gene fusion, a hallmark of K562, as one of the  
20 SV calls with highest quality score (Fig. 4A, Dataset S3), along with two other known gene  
21 fusions in K562 (Engreitz et al. 2012): *XKR3/NUP214* between chromosomes 9 and 22 (Fig. 4A)  
22 and *CDC25A/GRID1* between chromosomes 3 and 10 (Dataset S5, Supplemental Data).

23 In addition, we also leveraged the long-range information derived from the linked-reads  
24 to identify, assemble, and reconstruct SV-spanning breakpoints (including those of large-scale  
25 complex rearrangements) in the K562 genome using the recently established method Genome-  
26 wide Reconstruction of Complex Structural Variants (GROC-SVs) (Spies et al. 2017). In this



1 method, long DNA fragments that span breakpoints are statistically inferred and refined by  
2 quantifying barcode similarity between pairs of genomic regions, similar to Long Ranger (Marks  
3 et al. 2018). Sequence reconstruction is then performed by assembling the relevant linked-  
4 reads around the identified breakpoints from which complex SVs are then automatically  
5 reconstructed. The breakpoints that have supporting evidence from the K562 3 kb-mate-pair  
6 dataset (see Supplemental Methods) were determined as high-confidence events (Dataset S5).  
7 GROC-SVs identified a total of 161 high-confidence breakpoints including 12 inter-chromosomal  
8 events (Fig. 1, Dataset S5, Fig. 4B); each event is accompanied with visualization  
9 (Supplemental Data); 138 of the breakpoints were successfully sequence-assembled with  
10 nucleotide-level resolution of breakpoints as well the exact sequence in the cases where  
11 nucleotides have been added or deleted (Dataset S5). A notable example of assembly by  
12 GROC-SVs is a complex intra-chromosomal rearrangement on chromosome 13 (Figure 4B).

13       Using the methods (“gemtools”) as described in (Greer et al. 2017), we identified phased  
14 structural rearrangements (multiple deletions and tandem duplications) within the tumor  
15 suppressor gene *FHIT* on 3p14.2 (Waters et al. 2014) (Fig. 4C). Since K562 exhibits LOH on  
16 chromosome 3, the SVs within *FHIT* were phased using linked-read barcodes instead of  
17 heterozygous SNVs. The hemizygous deletion between (59.74 Mb – 60.08 Mb) of 3p14.2  
18 results in the loss of *FHIT* exons 6, 7, and 8. For the two phased tandem duplications on the  
19 same allele, one is intronic, and the other duplicates exon 5 (Fig. 4C). The two deletions  
20 downstream to the phased duplications are on two different alleles of *FHIT*. Another allele-  
21 specific, complex, intra-chromosomal rearrangement in K562 spans approximately 0.5 Mb on  
22 16q11.2 and 16q12.1 (Fig. 4D), involving two overlapping inversions (62 kb and 125 kb) and a  
23 tandem duplication (163 kb). These events affect *ORC6*, *MYLK3*, *RHBDF1* also known as  
24 *C16orf8*, and *NETO2*, which has recently been identified as a cancer marker gene (Oparina et  
25 al. 2012; Hu et al. 2015). This rearrangement resides on the non-duplicated haplotype of this  
26 triploid region. *ORC6* is located entirely within the more centromeric inversion of this locus on

1 16q11.2 and is “deleted” the by left breakpoint of the more telemetric inversion, which also  
2 “deletes” *C16orf8* and inverts *MYLK3*, possibly disrupting its promoter region or proximal  
3 enhancers or disconnecting *MYLK3* from their regulation (Fig. 4D).

#### 4 **Small-Scale Complex SVs from Deep-Coverage WGS**

5 Small-scale complex SVs (Fig. 5A-E) as well as non-complex SVs were identified using  
6 a novel algorithm called Automated Reconstruction of Complex Structural Variants (ARC-SV)  
7 (Arthur et al. 2017) from deep-coverage WGS data (Dataset S6, Supplemental Data). These  
8 small-scale complex SVs are defined as genomic rearrangements with multiple breakpoints that  
9 cannot be explained by one well-defined (non-complex) SV type such as deletions, insertions,  
10 tandem duplications, or inversions. After filtering out SVs <50 bp or with breakpoints that reside  
11 in simple repeats, low complexity regions, satellite repeats, or segmental duplications, we  
12 identified 122 complex SVs (accompanied with schematic visualizations), 2,235 deletions, 320  
13 tandem duplications, and 6 inversions (Dataset S6). Examples of complex SVs include  
14 dispersed duplications where duplicated sequences are inserted elsewhere in the genome in a  
15 non-tandem fashion (Fig. 5A). These dispersed duplications sometimes involve inversions of  
16 the inserted sequence and deletions at the insertion site (Fig. 5B, C). Other examples include  
17 inversions flanked on one or both sides by deletions (Fig. 5D), duplications that involve multiple  
18 non-exact copies, as well as deletion, inversion, and multiple duplications residing at the same  
19 locus (Fig. 5E). No other published algorithm to date has the capability to automatically identify  
20 and reconstruct these complex SVs. Eight out of ten breakpoints from five complex SVs were  
21 successfully validated by PCR and Sanger sequencing (Supplemental Table S8).

#### 22 **SVs from Mate-Pair Sequencing Analysis**

23 To increase the sensitivity of detecting medium-sized SVs (1 kb-100 kb) in K562, we  
24 constructed a 3 kb-mate-pair library and sequenced ( $2 \times 151$  bp) to 6.9 $\times$  non-duplicate  
25 coverage. The sequence coverage ( $C_R$ ) of each 3 kb insert is 302bp or 10%, which translates to  
26 a physical coverage ( $C_F$ ) of 68.5 $\times$ . From the mate-pair library, SVs (deletions, inversions, and

1 tandem duplications) were identified by clustering discordant read-pairs and split-reads using  
2 LUMPY (Layer et al. 2014). Only SVs that have both discordant read-pair and split-read support  
3 were retained. Overall, we identified 270 deletions, 35 inversions, and 124 tandem duplications  
4 using this approach (Dataset S7). Approximately 83% of these SVs are between 1 kb-10 kb,  
5 and 88% are between 1 kb-100 kb (Dataset S7). Twelve deletions and five tandem duplications  
6 were randomly selected for PCR and Sanger sequencing validation (Supplemental Table S8).  
7 The validation rates were 83% and 80% respectively.

### 8 **Non-Complex SVs from Deep-Coverage WGS**

9 Non-complex SVs (deletions, inversions, insertions, and tandem duplications) in K562  
10 were called from deep-coverage WGS data using a combination of established methods,  
11 namely Pindel (Ye et al. 2009), BreakDancer (Chen et al. 2009), and BreakSeq (Lam et al.  
12 2010). These SVs were combined with those of the same SV type that were identified using  
13 ARC-SV, LUMPY and Long Ranger, where SVs ( $n=2,665$ ) with support from multiple methods  
14 by  $\geq 50\%$  reciprocal overlap were merged. Through this combination of methods, a total of 9,082  
15 non-complex SVs were identified in the K562 genome, including 5,490 deletions, 531  
16 duplications, 436 inversions, and 2,602 insertions (Supplemental Data). (We note that only  
17 BreakDancer (Chen et al. 2009) was designed to call insertions.) Consistent with previous  
18 analyses (e.g. as in (Lam et al. 2012)), deletions show the highest number of concordant calls  
19 across the various methods compared to duplications and inversions ( Supplemental Fig S5,  
20 Supplemental Data). Eighteen deletions ( $>1$  kb) and and 18 tandem duplications, both with split-  
21 read support, were randomly chosen for experimental validation using PCR and Sanger  
22 sequencing. The validation rates were 89% and 72% respectively ( Supplemental Table S8).

### 23 **LINE1 and Alu Insertions**

24 We identified non-reference LINE1 and Alu retrotransposon insertions (REIs) in the  
25 K562 genome from our deep-coverage short-insert WGS data using a modified RetroSeq

1 (Keane et al. 2013) approach (Supplemental Methods). Non-reference REIs were identified from  
2 paired-end reads that have one of the paired reads mapping to the human reference genome  
3 and the other read mapping to either the Alu or LINE1 consensus sequence in a full or split-read  
4 fashion (see Methods). We identified 1,147 non-reference Alu insertions and 85 non-reference  
5 LINE1 insertions in K562 (Supplemental Table S9, Fig. 1). Nine Alu and ten LINE1 insertions  
6 with split-read support were randomly chosen for validation using PCR and Sanger sequencing.  
7 The validation rates were 88% and 100% respectively (Supplemental Table S10). PCR primers  
8 were designed such that one anneals within the retrotransposon sequence and the other  
9 anneals in the unique sequences surrounding the predicted insertion site.

## 10 **Allele-Specific Gene Expression**

11 Integrating CN information (i.e. allele frequencies) of the heterozygous SNVs (Dataset  
12 S1), we re-analyzed two replicates of ENCODE polyA-mRNA RNA-seq data to identify allele-  
13 specific gene expression in K562. We identified 5,053 and 5,149 genes that show allele-specific  
14 expression ( $p < 0.05$ ) in replicates one and two, respectively (Fig. 1, Supplemental Table S11).  
15 We also identified 2,342 and 2,176 genes that would have been falsely identified to have allele-  
16 specific expression and 1,641 and 1,710 genes that would not have been identified to have  
17 allele-specific expression in replicates one and two, respectively, if the allele frequencies of  
18 heterozygous SNVs in aneuploid regions were not taken into consideration (Supplemental Table  
19 S12).

## 20 **Allele-Specific DNA methylation**

21 By integrating CN and phase information of heterozygous SNVs of K562, we identified  
22 110 CpG islands (CGIs) that exhibit allele-specific DNA methylation (Figure 1, Supplementary  
23 Table S13). We obtained K562 whole-genome bisulfite sequencing (WGBS) reads ( $2 \times 100$  bp,  
24 library ENCLB742NWU) from the ENCODE portal (Sloan et al. 2016) and aligned the reads to  
25 hg19 using Bismark (Krueger and Andrews 2011), where 76.9% of reads were uniquely mapped  
26 and 26.2% of cytosines were methylated in a CpG context. We then used reads that overlap

1 both phased heterozygous SNVs (Dataset S2) and CpGs to phase the methylated and  
2 unmethylated CpGs to their respective haplotypes. We then grouped the phased individual  
3 CpGs into CGIs. Fisher's exact test (taking the CN of a given CGI locus into account) was used  
4 to evaluate allele-specific methylation ( $p < 0.05$ ), and significant results were selected using a  
5 target false discovery rate of 10% (Supplemental Methods). Of these 110 CGIs, 35 reside  
6 within promoter regions (here defined as 1 kb upstream of a gene); 83 are intragenic, and 28 lie  
7 within 1 kb downstream of 113 different genes. The following 6 genes are within 1 kb of a  
8 differentially methylated CGI and overlap with the Sanger Cancer Gene Census: *ABL1*, *AXIN2*,  
9 *CCND1*, *HOXD11*, *KDR*, and *PRDM16*.

### 10 **Allele-Specific CRISPR Targets**

11 We identified a total of 28,511 targets in the K562 genome suitable for allele-specific  
12 CRISPR targeting (Fig. 1, Supplemental Table S14). Sequences (including reverse complement)  
13 of phased variants that differ by more than one base pair between the alleles were extracted to  
14 find all possible CRISPR targets by searching for the pattern [G, C, or A]<sub>N</sub>GG (Supplemental  
15 Methods). Using a selection method previously described and validated (Sunagawa et al. 2016),  
16 only conserved high-quality targets were retained. We also took gRNA function and structure  
17 into consideration and performed further filtering of CRISPR targets. Targets with multiple  
18 exact matches, extreme GC content, and those containing TTTT (which might break the  
19 secondary structure of gRNA), were removed. We also used the Vienna RNA-fold package  
20 (Lorenz et al. 2011) to compute gRNA secondary structure and eliminated all targets for which  
21 the stem loop structure (for Cas9 recognition) could not form (Nishimasu et al. 2014). Finally, we  
22 calculated the off-target risk score by using the tool as described in (Ran et al. 2013). To ensure  
23 that all targets are as reliable and as specific as possible, we chose a very strict threshold and  
24 rejected candidates with a score below 75. Of the 28,511 allele-specific CRISPR target sites,  
25 15,488 are within an annotated protein-coding or non-coding RNA transcript, 705 within an exon,

1 and 13 targets are within an experimentally validated enhancer (Visel et al. 2007)  
2 (Supplementary Table S14).

### 3 **Genomic Structural Context Provides Insight into Regulatory Complexity**

4 We show examples of how deeper insights into gene regulation and regulatory  
5 complexity can be obtained by integrating genomic structural contexts with functional genomics  
6 and epigenomics data (Fig. 6A-D). One example is the allele-specific RNA expression and  
7 allele-specific DNA methylation in K562 at the *HOXB7* locus on chromosome 17 (Fig. 6A). By  
8 incorporating the genomic context of *HOXB7* in K562, we see that *HOXB7* exhibit highly  
9 preferential RNA expression from the two copies of Haplotype 1 ( $p = 0.007$ ) in which the CGI  
10 near its promoter is completely unmethylated ( $p = 3.18 \times 10^{-18}$ ) (Fig. 6 A, C). The second  
11 example is allele-specific RNA expression and allele-specific DNA methylation of the *HLX* gene  
12 in K562 (Fig. 6B). The *HLX* locus on chromosome 1 is tetraploid, and we see that *HLX* is only  
13 expressed from Haplotype 1 which has three copies and not expressed in Haplotype 2 ( $p =$   
14  $0.043$ ) (Fig. 6B, D). The CGI of the *HLX* locus is unmethylated in Haplotype 2 but highly  
15 methylated on Haplotype 1 ( $p = 5.14 \times 10^{-15}$ ) (Fig. 6B, C). There is also an allele-specific  
16 CRISPR targeting site for both haplotypes within *HLX* (Fig. 6B). In addition, we performed  
17 Pearson correlation analysis between our deep-coverage K562 WGS data and K562 POLR2A  
18 ChIP-seq data (previously released on the ENCODE data portal) to determine whether changes  
19 in K562 genome CN or ploidy affected binding of the polymerase molecule to genomic DNA in a  
20 large-scale fashion (Supplemental Fig. S6). The two sets of data are very well correlated  
21 ( $r=0.51$ ,  $p<2.2 \times 10^{-16}$ ) suggesting that RNA polymerase activity is generally influenced by ploidy  
22 in the K562 genome. In addition, we also correlated the K562 POLR2A ChIP-seq data with the  
23 FPKM values from four independent K562 polyA RNA-seq experiments (also previously  
24 released on the ENCODE portal) and find that these datasets are also very well correlated  
25 consistently ( $r=0.46$ ,  $p<2.2 \times 10^{-16}$ ;  $r=0.58$ ,  $p<2.2 \times 10^{-16}$ ;  $r=0.47$ ,  $p<2.2 \times 10^{-16}$ ;  $r=0.46$ ,  $p<2.2 \times$   
26  $10^{-16}$ ) (Supplemental Fig. S7A-D).

1           Furthermore, we also find allele-specific RNA expression for the rearranged copy of  
2 *MYLK3* ( $p < 1.93 \times 10^{-17}$ ) and the normal, non-rearranged copies (CN=2) of *ORC6* ( $p < 1.58 \times 10^{-8}$ )  
3 where expression from re-arranged allele (CN=1) of *ORC6* is “depleted” in K562 (Table S11, Fig.  
4 4C, D). These observations made by integrating our K562 linked-read data and ENCODE RNA  
5 expression data provide novel insights into gene regulatory mechanisms in terms of ectopic  
6 expression and dosage compensation, which also raises important questions regarding the  
7 history of the K562 cell line in terms of mutations, selective pressures, and adaption (see  
8 Discussion).

## 9 **DISCUSSION**

10           K562 is one of the most widely used laboratory “work-horse” cell lines in the world.  
11 Among the three tier-one cell lines of ENCODE, K562 has by far the most functional genomics  
12 and epigenomics data generated. Furthermore, K562 is also one of the most commonly used  
13 cell lines for large-scale CRISPR/Cas9 gene-targeting screens (Wang et al. 2015; Arroyo et al.  
14 2016; Morgens et al. 2016; Han et al. 2017; Adamson et al. 2016; Liu et al. 2017). Yet, despite  
15 its wide usage and impact on biomedical research, its genomic sequence and structural  
16 features have never been comprehensively characterized, beyond its karyotype (Selden et al.  
17 1983; Gribble et al. 2000; Wu et al. 1995; Naumann et al. 2001) and SNPs called from 30x-  
18 coverage WGS but without taking aneuploidy or CN into consideration (Cavalli et al. 2016).  
19 Analysis, integration, and interpretation of the extensive collection of functional genomics and  
20 epigenomics datasets for K562 had so far relied solely on the human reference genome. Here,  
21 we present the first detailed and comprehensive characterization of the K562 genome so that  
22 future studies no longer have to rely solely on the human reference genome. By performing  
23 deep-coverage short-insert WGS, 3 kb-insert mate-pair sequencing, deep-coverage linked-  
24 reads sequencing, array CGH, karyotyping, and integrating a compendium of novel and  
25 established analysis methods (Supplemental Fig. S1A), we produced a comprehensive  
26 spectrum of genomic structural features (Fig. 1) for K562 that includes SNVs (Dataset S1),



1 Indels (Dataset S1), ploidy by chromosome segments at 10-kb resolution (Supplemental Table  
2 S2), phased haplotypes (Dataset S2, Supplemental Data)—often of entire chromosome arms  
3 (Table 2, Supplemental Data)—phased CRISPR targets (Supplemental Table S14), non-  
4 reference retrotransposon insertions (Table S9), and SVs (Supplemental Data) including  
5 deletions, duplications, and inversions, and complex SVs (Dataset S6, Dataset S7). Many SVs  
6 were also phased, assembled, and experimentally verified (Dataset S2-S5, Supplemental Table  
7 S8, Supplemental Table S10). Of the 3,784,863 variants that were haplotype-phased in the  
8 K562 genome (Dataset S2-S5), 3,088,185 (81.6%) are SNVs; 692,998 (18.31%) are Indels;  
9 3,451 are deletion SVs (51 bp to 30 kb) (0.1%), and 229 are large SVs.

10 Pervasive aneuploidy is a characteristic of many cancers. Previous studies have  
11 confirmed the near triploid karyotype of K562 (Selden et al. 1983; Gribble et al. 2000; Wu et al.  
12 1995; Naumann et al. 2001). In our analysis, however, we also found considerable portions of  
13 the K562 genome to be much more varied than what had previously been reported. This is  
14 because by leveraging deep-coverage WGS, the CN across different chromosome segments,  
15 as determined by our read-depth analysis, is of much higher resolution than karyotyping.  
16 Furthermore, the identified chromosome segments with aneuploidy (CN>2) and orthogonally  
17 supported by karyotyping and array CGH were further validated, also orthogonally, from a  
18 statistical approach in which significant differences in unique linked-read barcode counts  
19 between the major and minor haplotypes were determined using a one-sided *t*-test ( $p < 0.001$ )  
20 (Bell et al. 2017). In addition, it has to be taken into consideration that for a widely used cell line  
21 with decades of history such as K562, additional genome variation is expected. In light of this, it  
22 is reassuring that the overall karyotype has not changed much over the years. However,  
23 researchers should still keep this aspect in mind when working with a version of K562 that has  
24 been separated from the main ENCODE K562 line used here. We expect that the vast majority  
25 of genomic variants that we describe here to be universal for K562, but for individual variants, it  
26 is possible that different lines of K562 may have slightly diverged from each other

1 (Supplemental Table S1, Supplemental Table S3). When using a different K562 line and  
2 following up on findings for individual loci of interest, a first step should always be to  
3 experimentally validate their presence. When incorporating these genomic variants for global  
4 analyses, such as interrogating network interactions, the vast majority of them will exist, thus  
5 such global analyses are expected to yield substantial insights. Even though the pervasive  
6 aneuploidy in K562 renders the design and interpretation of K562 studies more challenging, the  
7 information we provide here enables researchers to continue the use of this cell line to  
8 investigate the effects of genetic variation on the multiple levels of functional genomics activity  
9 and regulation for which ENCODE data already exists or continues to be produced. Thus,  
10 analysis of K562 data should not only be more complex and challenging, but also potentially  
11 much more insightful and rewarding when taking its complex genome structure into account.

12 Sensitive and accurate identification of SNVs and Indels requires relatively deep WGS  
13 coverage (>33× and >60× respectively) (Bentley et al. 2008; Fang et al. 2014). From our >70×  
14 non-duplicate coverage WGS data, we identified large numbers of SNVs and Indels that we  
15 could subsequently correct for their allele frequencies according to ploidy. In addition to being  
16 essential for correct haplotype identification, these ploidy-corrected variants are also needed for  
17 functional genomics or epigenomics analyses such as the determination of allele-specific gene  
18 expression or of allele-specific transcription factor binding in K562 (Cavalli et al. 2016). From  
19 RNA-seq or ChIP-seq data analysis, a statistically significant increase in transcription or  
20 transcription-factor-binding signal in one allele compared to the other at a heterozygous locus,  
21 may be identified as a case of allele-specific expression or allele-specific transcription-factor  
22 binding which usually suggests allele-specific gene regulation at this locus. However, if  
23 aneuploidy can be taken into consideration and the signals normalized by ploidy, the case  
24 identified might be a result of increased CN rather than the preferential activation of one allele  
25 over the other on the epigenomic level. Indeed, in our re-analysis of two replicates of ENCODE  
26 K562 RNA-seq data, we identified 2,359 and 2,643 genes that would have been falsely

1 identified to have allele-specific expression in addition to 1,808 and 2,063 genes that would not  
2 have been identified to have allele-specific expression in replicates one and two, respectively, if  
3 ploidy was not taken into consideration (Supplemental Table S12).

4         The haplotype phase of genomic sequence variants is an essential aspect of human  
5 genetics, but current standard WGS approaches entirely fail to resolve this aspect. We  
6 performed linked-read sequencing of K562 genomic DNA using the Chromium System from 10x  
7 Genomics (Zheng et al. 2016; Marks et al. 2018). After size-selecting for genomic DNA  
8 fragments >30kb, 300 genomic equivalents of HMW DNA were partitioned into more than one  
9 million oil droplets, uniquely barcoded within each droplet, and subjected to random priming and  
10 amplification. Implemented by Long Ranger (Marks et al. 2018), sequencing reads that originate  
11 from the same HMW DNA molecule can be identified by their respective droplet barcodes and  
12 linked together to produce virtual long reads. Then, by looking for virtual long reads that overlap  
13 a previously called set of heterozygous haplotypes (Dataset S1), the phase information of the  
14 heterozygous haplotypes was determined and the virtual long reads were constructed into  
15 phased haplotype blocks with N50 > 2.72 Mb (Dataset S2, Fig. 2D). Chromosomes 3, 9, 13, 14,  
16 X, and large portions of chromosomes 2, 10, 12, 20, 22 were difficult to phase, resulting in  
17 comparatively shorter phased blocks (Dataset S2, Fig. 1, Supplemental Fig. S4). This is not  
18 surprising since these chromosomes and chromosomal regions exhibit a very high degree of  
19 LOH (Fig. 1 and Supplemental Table S4). Heterozygous loci in aneuploidy regions with more  
20 than two haplotypes were excluded from phasing linked-read analysis due to software and  
21 algorithmic limitations (Zheng et al. 2016). However, the phase information of these loci could  
22 be resolved from our linked-read data in principle, should new algorithms become available.

23         We extended on the already-impressive phasing capabilities of Long Ranger and  
24 constructed mega-haplotypes in K562—often spanning entire chromosome arms (Fig. 3, Table  
25 2, Supplementary Data)—by leveraging the haplotype imbalance in aneuploid chromosomes  
26 using a recently developed method for which its effectiveness in cancer has already been

1 demonstrated (Bell et al. 2017). Since gene dosage is a fundamental component of genome  
2 biology and for which aneuploidy contribute large effects in terms of amplification and reduction,  
3 the ability to haplotype across long stretches of aneuploidy is essential for understanding of the  
4 genetic regulations of cancer and an important component for developing genetically targeted  
5 cancer treatment.

6 It has been shown previously that integrating orthogonal methods and signals improves  
7 SV-calling sensitivity and accuracy (Mohiyuddin et al. 2015; Layer et al. 2014). Here, we  
8 combined deep-coverage short-insert WGS, mate-pair sequencing, linked-read sequencing, and  
9 several SV-calling methods to identify many non-complex SVs. To obtain the union set of non-  
10 complex SV calls from the various methods, the SVs identified by multiple methods were  
11 merged and indicated accordingly (Supplemental Data). For deletions (Supplemental Fig. S5A),  
12 we see strong overlap for the various methods, but this overlap is less pronounced for  
13 duplications (Supplemental Fig. S5B) and inversions (Supplemental Fig. S5C). This is  
14 consistent with previous analysis (Lam et al. 2012) as inversions and duplications are more  
15 difficult in principle to accurately resolve (Lin et al. 2015; Sudmant et al. 2015). We also expect  
16 the detection of many SVs to be method-specific, since each method is designed to utilize  
17 different types of signals and also optimized to identify different classes of SVs (Pabinger et al.  
18 2014; Lin et al. 2015). Again, if particular SVs are of interest for follow-up studies, they should  
19 first be experimentally validated.

20 The complex rearrangements identified by using ARC-SV from short-insert WGS (Fig.  
21 5A-E, Dataset S6, Supplemental Data) and by using GROC-SVs from linked-reads (Fig. 4B,  
22 Dataset S5, Supplemental Data) are classes of SVs that could not be easily identified and  
23 automatically reconstructed using previously existing methods. The small-scale complex SVs  
24 that were identified by using ARC-SV and experimentally validated (Fig. 5, Supplemental Table  
25 S8, Supplemental Data) describe a subtle class of complex rearrangements in cancer genomes  
26 that have never been previously demonstrated by others, but have long speculated to exist

1 (Perry et al. 2008; Quinlan and Hall 2012; Collins et al. 2017). Detecting and automatically  
2 reconstructing these small-scale complex SVs, especially in a “hay” of canonical SVs and in  
3 highly rearranged cancer genomes, has remained an unsolved problem for many years. In  
4 other words, our results reveal a class of previously overlooked complex SVs in cancer that can  
5 now be identified from standard short-insert WGS data and elucidated further. They have clear  
6 implications for the conventional models of cancer evolution which often assume gradual, step-  
7 by-step mutations; however, these complex SVs support a form of punctuated genome  
8 evolution (Davis et al. 2017). A major unsolved question still is how complex SVs arise  
9 mechanistically for which there are general models: template switching during replication (Lee  
10 et al. 2007; Hastings et al. 2009) and chromothripsis (Stephens et al. 2011). These small-scale  
11 complex SVs in K562 implicate that another, yet-undescribed, mechanism might also be  
12 contributing, at least in cancer. Furthermore, the functional consequences of these small-scale  
13 complex SVs are also unknown. These important questions remain unsolved mainly due to the  
14 lack of data and examples. It is possible that this mutational complexity contributes to genome  
15 innovation, at least in cancer, or is just a curious sideshow (Quinlan and Hall 2012). Only the  
16 accumulation of such examples and data will allow researchers in fields such as cancer  
17 evolution to begin to address these important questions.

18         Before the existence of linked-read sequencing, haplotype phasing and resolving large  
19 SVs (>30kb) relied heavily on fosmid libraries (Snyder et al. 2015; Williams et al. 2012; Kitzman  
20 et al. 2011; Cao et al. 2015; Adey et al. 2013) which were laborious, costly, time consuming,  
21 and much less efficient. Using linked-read sequencing and gemtools (Greer et al. 2017), we  
22 phased and resolved complex SVs that are especially compelling on 3p14.2 (within the tumor  
23 suppressor gene *FHIT*) and on 16q11.2 and 16q12.1 of the K562 genome. *FHIT* is frequently  
24 seen to harbor deletions in many types of human cancers, most commonly of epithelial origin,  
25 such as lung, stomach, cervix, head and neck, breast, and kidney (Lubinski et al. 1994; Ohta et  
26 al. 1996; Huebner et al. 1998; Ingvarsson 2001). Reduction or absence in its protein expression

1 occurs in nearly 50% of all cancers (Huebner et al. 1998; Waters et al. 2014). While LOH and  
2 allele-specific deletions within *FHIT* have been previously reported (Wistuba et al. 1997; Li et al.  
3 2016), to our knowledge, this is the first discovery of phased and allele-specific tandem  
4 duplications within *FHIT* and the first report of *FHIT* mutations for CML. Curiously, all previous  
5 reports of deletions within *FHIT* for various cancer types (not including CML) were all centered  
6 on and include exon 5 (Durkin et al. 2008), whereas exon 5 is duplicated in K562. Deletions of  
7 all three *FHIT* exons 6, 7, and 8 (Fig. 4C) are less frequent but have been reported for lung  
8 cancer and esophageal adenocarcinoma (Sozzi et al. 1996; Dagmar et al. 1997).

9 We identified highly allele-specific RNA expression for *ORC6* ( $p < 1.58 \times 10^{-8}$ ) and  
10 *MYLK3* ( $p < 1.93 \times 10^{-17}$ ) in K562 (Supplemental Table S11), which is likely contributed by the  
11 allele-specific complex intra-chromosomal rearrangement residing on the non-duplicated  
12 haplotype of 16q11.2 (triploid) (Fig. 4D). *ORC6* codes for origin recognition protein complex  
13 subunit 6 and is essential for coordinating DNA replication, chromosome segregation, and  
14 cytokinesis (Prasanth et al. 2002). One allele of *ORC6* is “deleted” by one of the inversion  
15 breakpoints of this rearrangement and maintains allele-specific expression on the other allele  
16 (Fig. 4D). The origin recognition complex, in which Orc6 is a subunit, serves as a “landing pad”  
17 for bringing together components of the pre-replicative complex required for DNA replication.  
18 Reduction of Orc6 dosage by small interfering RNA results in decreased DNA replication,  
19 aberrant mitosis, and the formation of multiple nuclei and multipolar spindles in cells; a long  
20 period of this reduction increases cell death (Prasanth et al. 2002). Such effects are not  
21 observed in K562 cells even though RNA expression from one allele of *ORC6* is “depleted” by  
22 this rearrangement ( $p < 1.58 \times 10^{-8}$ , Supplemental Table S11). This is likely because the other  
23 “normal” *ORC6* allele was duplicated, rendering this locus in K562 triploid and thus maintain a  
24 “diploid” gene dosage. Perhaps more importantly, this insight also raises important questions  
25 regarding the history of mutation as well as selective pressures that occurred within K562 cells.  
26 It is possible that one copy of chromosome 16 was first duplicated, freeing this locus from

1 selective pressures, allowing it to acquire new mutations, since “diploid” copies are still  
2 maintained in the genome. It is also conceivable that duplication of this locus is  
3 disadvantageous for cell proliferation, and K562 cells that acquired this rearrangement after the  
4 duplication of chromosome 16 also acquired a selective advantage since they now reverted this  
5 locus back to “diploid” copies. However, it is also possible, though perhaps less likely, that this  
6 rearrangement occurred before chromosome 16 duplication, putting a negative selection  
7 pressure on K562 cells, and that this pressure is released by duplication of the other copy.  
8 Interestingly, this rearrangement also inverts an intact copy of *MYLK3* (Fig. 4D), which was  
9 identified to encode a novel cardiac-specific myosin light chain kinase (Seguchi et al. 2007).  
10 Since its expression is expected to be normally repressed except in heart cells, the allele-  
11 specific RNA expression of *MYLK3* ( $p < 1.93E-17$ , Supplemental Table S11) from this inverted  
12 allele suggests that this inversion activated the ectopic expression of this gene in K562, possibly  
13 by disrupting or disconnecting it from its promoter or proximal enhancer elements that impose  
14 repressive regulatory mechanisms. Finally, this complex rearrangement on chromosome 16 of  
15 K562 also duplicates *NETO2* in a tandem fashion, also on the same allele (Fig. 4D). *NETO2*  
16 codes for a single-pass membrane protein neuropilin and tolloid-like 2 (Stöhr et al. 2002). Its  
17 expression is frequently up-regulated in many types of human cancers including lung, cervical,  
18 colon, and renal carcinomas and has been suggested as a potential genetic marker for cancer  
19 (Oparina et al. 2012). Its up-regulation also correlates with the progression and poor prognosis  
20 of colorectal carcinoma (Hu et al. 2015). It is conceivable that the increase in *NETO2* gene  
21 dosage due to duplication in K562 cells contributes to their efficient proliferation in culture and  
22 that, at least in some cancers, the frequent up-regulations observed for *NETO2* are also  
23 contributed by this similar mechanism of allele-specific tandem duplication.

24 The hallmark of CML is the Philadelphia rearrangement t(9; 22)(q34; q11) which results  
25 in the fusion of *ABL1* and *BCR* (Heisterkamp et al. 1985; de Klein et al. 1982; Groffen et al.  
26 1984). This gene fusion is known to be extensively amplified in the K562 genome by tandem



1 duplication (Wu et al. 1995). FISH analysis showed that fluorescent signals from the *BCR/ABL1*  
2 gene fusion almost always concentrate on a single marker chromosome (Tkachuk et al. 1990;  
3 Wu et al. 1995; Gribble et al. 2000). This is also consistent with our data as the linked-reads that  
4 support the *BCR/ABL1* gene fusion do not share overlapping barcodes with linked-reads that  
5 align elsewhere in the genome, and the *BCR* and *ABL1* gene regions where the fusion occurs  
6 show a >2.8× increase in sequencing coverage relative to average sequencing coverage  
7 across the genome.

8 Data generated from this comprehensive whole-genome analysis of K562 is available  
9 through the ENCODE portal (Sloan et al. 2016) (Supplemental Figure S1B, C). We envision that  
10 this analysis will serve as a valuable resource for further understanding the vast troves of  
11 ENCODE data available for K562, such as determining whether a potential or known regulatory  
12 sequence element has been altered by SNVs or SNPs, Indels, retrotransposon insertions, a  
13 gain or loss of copies of that given element, or allele-specific regulation. As additional examples  
14 of how integrating genomic context can yield further understanding of existing ENCODE data,  
15 we showed, as examples, the complex gene regulatory scenarios uncovered at the *HOXB7* and  
16 *HLX* loci in K562. Hox genes are known to have important roles in hematopoiesis and  
17 oncogenesis (Argiropoulos and Humphries 2007; Shah and Sukumar 2010; Eklund 2011). The  
18 *HOXB7* transcription factor mediates lymphoid development, hematopoietic differentiation and  
19 leukemogenesis (Giampaolo et al. 1995; Carè et al. 1999). *HOXB7* overexpression has been  
20 reported in leukemia (Raval et al. 2007) as well as in many other cancers (Caré et al. 1996; Wu  
21 et al. 2006; Yamashita et al. 2006; Shiraishi et al. 2007; Chen et al. 2008; Storti et al. 2011). It is  
22 directly upstream of *HOXB8*, which is the first Hox gene found to be an oncogene in leukemia  
23 (Blatt et al. 1988). *HLX* has also been suggested to play oncogenic roles in leukemia (Deguchi  
24 et al. 1992; Deguchi and Kehrl 1993; Jawad et al. 2006; Fröhling 2012). By integrating the  
25 genomic context of *HOXB7* and *HLX* in K562 with RNA-seq and WGBS data, we see that the  
26 RNA of both genes are expressed from haplotypes that exhibit aneuploidy and in an allele-

1 specific manner (Fig. 6A, B, D). The allele-specific methylation of the CGIs near these two  
2 genes is associated with active transcription in the case of *HLX* and silencing of transcription in  
3 the case of *HOXB7* (Fig. 6A-C). Such insights into potential oncogene regulation cannot be  
4 obtained by analyzing functional genomics and epigenomics data alone without genome  
5 structural information i.e. correct genomic context. In addition, we also observed that the K562  
6 POLR2A ChIP-seq signal in both replicates is very well correlated with polyA RNA-seq signal  
7 and with WGS coverage, suggesting an association between polymerase binding and active  
8 transcription and between polymerase binding and ploidy (Supplemental Figures S6, S7).

9 Our work here serves to guide future studies that utilize the K562 “workhorse” cell line,  
10 such as CRISPR screens where knowledge of the sequence variants can extend or modify the  
11 number of editing targets (Table S13) while knowledge of aberrant CN will allow for much more  
12 confident data interpretation. To give an example, in a recent study that uses CRISPRi to  
13 screen and elucidate the function of long non-coding RNAs in human cells, out of the seven cell  
14 types studied, the number of gRNA hits varied considerably among the various cell types, with  
15 89.4% of hits unique to only one cell type and none in more than five cell types (Liu et al. 2017).  
16 Although a large portion of the phenomenon are very likely explained by cell-specific effects, it is  
17 still quite possible that many of the gRNA hit differences were the result of differences in  
18 genome sequence or ploidy. Our list of allele-specific CRISPR targets (Table S13) will allow for  
19 the discernment between these two potential reasons for differences in CRISPR effects and  
20 should be particularly valuable for future large-scale CRISPR screens that utilize K562. Lastly,  
21 this study also serves as a technical example for the advanced, integrated, and comprehensive  
22 analyses of other heavily utilized cell lines and genomes in biomedical research such as HepG2.

## 23 **METHODS**

### 24 **Overview**

25 We combined multiple experimental and analysis methods (Supplementary Fig. S1A),  
26 including karyotyping, array CGH, deep (72× non-duplicate coverage ) short-insert whole-

1 genome sequencing (WGS), 3 kb-mate-pair sequencing (Korbel et al. 2007) and 10x Genomics  
2 linked-reads sequencing (Zheng et al. 2016; Marks et al. 2018), to comprehensively  
3 characterize the genome of the primary ENCODE cell line K562 (Fig. 1). The WGS dataset was  
4 used to identify CN i.e. ploidy by chromosome segments, SNVs, Indels, non-reference LINE1  
5 and Alu insertions (Lupski 2010; Sudmant et al. 2015), and SVs such as deletions, duplications,  
6 inversions, insertions, and small-scale complex SVs. These SVs were identified using an  
7 integrated approach that includes BreakDancer (Chen et al. 2009), Pindel (Ye et al. 2009),  
8 BreakSeq (Lam et al. 2010) and ARC-SV (Arthur et al. 2017). The allele frequencies of  
9 heterozygous SNVs and Indels were determined by taking ploidy into account. The linked-reads  
10 were used to phase SNVs and Indels as well as to identify, phase, reconstruct, and assemble  
11 primarily large (>30 kb) and complex SVs (Greer et al. 2017; Spies et al. 2017; Marks et al.  
12 2018), though additional small-scale deletions were also identified and phased (Zheng et al.  
13 2016). SVs and REIs were experimentally validated with PCR and Sanger sequencing. Phased  
14 SNV haplotype blocks in aneuploid regions were “stitched” to mega-haplotypes by leveraging  
15 haplotype imbalance (Bell et al. 2017). The 3 kb-mate-pair data was used to identify additional  
16 SVs and was also used to validate large and complex SVs identified from linked-reads.  
17 Functional genomics and epigenomics datasets from ENCODE were integrated with CN and  
18 phasing information to identify allele-specific RNA expression and allele-specific DNA  
19 methylation. Phased variants were also used to identify allele-specific CRISPR targets in the  
20 K562 genome. For full descriptions of experimental and computational procedures (including  
21 analysis code), see Supplemental Methods.

## 22 **DATA ACCESS**

23 Raw and processed data generated in this study are publicly available on the ENCODE portal  
24 ([encodeproject.org](http://encodeproject.org)) (Sloan et al. 2016) under experiment accessions: ENCSR711UNY,  
25 ENCSR025GPQ, and ENCSR053AXS. ENCODE accessions for individual data files including  
26 Datasets S1-S7 are listed in Supplemental Fig. 1B, C. Analysis code is provided in

1 Supplemental Methods. Overview of all resources generated for K562 in this study is listed with  
2 detailed descriptions in Supplemental Fig. 1B-D.

### 3 **ACKNOWLEDGEMENTS**

4 We thank Aditi Narayanan, Dr. Idan Gabdank, Nathaniel Watson, Dr. Carrie Davis, Kathrina  
5 Onate, and Dr. Cricket Sloan for assistance with data upload to the ENCODE portal. We thank  
6 Dr. Athena Cherry and the Stanford Cytogenetics Laboratory for karyotype analysis and Arineh  
7 Khechaduri for performing genomic DNA preparation. We thank Dr. Minyi Shi for providing K562  
8 cells. A.E.U. was supported by NIH grant P50-HG007735 and the Stanford Medicine Faculty  
9 Innovation Program, and B.Z. was additionally supported by NIH training grant T32-HL110952.  
10 W.H.W. received support from NIH grants HG007834 and HG007735. J.G.A. received funding  
11 from NIH training grant T32-GM096982 and NSF Graduate Fellowship DGE-114747. A.A. was  
12 funded by NIH grant U24CA220242.

### 13 **AUTHOR CONTRIBUTIONS**

14 B.Z. and A.E.U conceived and designed the study. B.Z., R.P., N.B.E, M.S.H, and R.R.H  
15 performed experiments. B.Z., S.S.H., S.U.G., J.M.B., N.S., XW.Z., XL.Z., S.B., J.G.A., G.S.,  
16 D.P., and A.A. performed analysis. H.P.J, W.H.W., and A.E.U contributed resources and  
17 supervised the study. B.Z., S.S.H., and A.E.U. wrote the manuscript.

### 18 **DISCLOSURE DECLARATION**

19 The authors of this manuscript declare no conflicts of interest.

### 20 **REFERENCES**

- 21 Abyzov A, Urban AE, Snyder M, Gerstein M. 2011. CNVnator: An approach to discover,  
22 genotype, and characterize typical and atypical CNVs from family and population genome  
23 sequencing. *Genome Res* **21**: 974–984.  
24 <http://genome.cshlp.org/cgi/doi/10.1101/gr.114876.110>.
- 25 Adamson B, Norman TM, Jost M, Cho MY, Nuñez JK, Chen Y, Villalta JE, Gilbert LA, Horlbeck  
26 MA, Hein MY, et al. 2016. A Multiplexed Single-Cell CRISPR Screening Platform Enables  
27 Systematic Dissection of the Unfolded Protein Response. *Cell* **167**: 1867–1882.e21.  
28 <http://www.ncbi.nlm.nih.gov/pubmed/27984733>.
- 29 Adey A, Burton JN, Kitzman JO, Hiatt JB, Lewis AP, Martin BK, Qiu R, Lee C, Shendure J. 2013.  
30 The haplotype-resolved genome and epigenome of the aneuploid HeLa cancer cell line.  
31 *Nature* **500**: 207–211. <http://www.nature.com/doi/10.1038/nature12064>.

- 1 Argiropoulos B, Humphries RK. 2007. Hox genes in hematopoiesis and leukemogenesis.  
2 *Oncogene* **26**: 6766–6776. <http://www.nature.com/doi/10.1038/sj.onc.1210760>.
- 3 Arroyo JD, Jourdain AA, Calvo SE, Ballarano CA, Doench JG, Root DE, Mootha VK. 2016. A  
4 Genome-wide CRISPR Death Screen Identifies Genes Essential for Oxidative  
5 Phosphorylation. *Cell Metab* **24**: 875–885. <http://www.ncbi.nlm.nih.gov/pubmed/27667664>.
- 6 Arthur JG, Chen X, Zhou B, Urban AE. 2017. Detection of complex structural variation from  
7 paired-end sequencing data. *bioRxiv* 200170.
- 8 Bell JM, Lau BT, Greer SU, Wood-Bouwens C, Xia LC, Connolly ID, Gephart MH, Ji HP. 2017.  
9 Chromosome-scale mega-haplotypes enable digital karyotyping of cancer aneuploidy.  
10 *Nucleic Acids Res* **45**: e162. <http://www.ncbi.nlm.nih.gov/pubmed/28977555>.
- 11 Bentley DR, Balasubramanian S, Swerdlow HP, Smith GP, Milton J, Brown CG, Hall KP, Evers  
12 DJ, Barnes CL, Bignell HR, et al. 2008. Accurate whole human genome sequencing using  
13 reversible terminator chemistry. *Nature* **456**: 53–9.  
14 <http://www.ncbi.nlm.nih.gov/pubmed/18987734>.
- 15 Blatt C, Aberdam D, Schwartz R, Sachs L. 1988. DNA rearrangement of a homeobox gene in  
16 myeloid leukaemic cells. *EMBO J* **7**: 4283–90.  
17 <http://www.ncbi.nlm.nih.gov/pubmed/2907477>.
- 18 Butler MO, Hirano N. 2014. Human cell-based artificial antigen-presenting cells for cancer  
19 immunotherapy. *Immunol Rev* **257**: 191–209.  
20 <http://www.ncbi.nlm.nih.gov/pubmed/24329798>.
- 21 Cao H, Wu H, Luo R, Huang S, Sun Y, Tong X, Xie Y, Liu B, Yang H, Zheng H, et al. 2015. De  
22 novo assembly of a haplotype-resolved human genome. *Nat Biotechnol* **33**: 617–22.  
23 <http://www.ncbi.nlm.nih.gov/pubmed/26006006> (Accessed May 9, 2016).
- 24 Caré A, Silvani A, Meccia E, Mattia G, Stoppacciaro A, Parmiani G, Peschle C, Colombo MP.  
25 1996. HOXB7 constitutively activates basic fibroblast growth factor in melanomas. *Mol Cell*  
26 *Biol* **16**: 4842–51. <http://www.ncbi.nlm.nih.gov/pubmed/8756643>.
- 27 Caré A, Valtieri M, Mattia G, Meccia E, Masella B, Luchetti L, Felicetti F, Colombo MP, Peschle  
28 C. 1999. Enforced expression of HOXB7 promotes hematopoietic stem cell proliferation  
29 and myeloid-restricted progenitor differentiation. *Oncogene* **18**: 1993–2001.  
30 <http://www.ncbi.nlm.nih.gov/pubmed/10208421>.
- 31 Cavalli M, Pan G, Nord H, Wallerman O, Wallén Arzt E, Berggren O, Elvers I, Eloranta M-L,  
32 Rönnblom L, Lindblad Toh K, et al. 2016. Allele-specific transcription factor binding to  
33 common and rare variants associated with disease and gene expression. *Hum Genet* **135**:  
34 485–497. <http://link.springer.com/10.1007/s00439-016-1654-x>.
- 35 Chen H, Lee JS, Liang X, Zhang H, Zhu T, Zhang Z, Taylor ME, Zahnow C, Feigenbaum L,  
36 Rein A, et al. 2008. Hoxb7 inhibits transgenic HER-2/neu-induced mouse mammary tumor  
37 onset but promotes progression and lung metastasis. *Cancer Res* **68**: 3637–44.  
38 <http://www.ncbi.nlm.nih.gov/pubmed/18463397>.
- 39 Chen K, Wallis JW, McLellan MD, Larson DE, Kalicki JM, Pohl CS, McGrath SD, Wendl MC,  
40 Zhang Q, Locke DP, et al. 2009. BreakDancer: an algorithm for high-resolution mapping of  
41 genomic structural variation. *Nat Methods* **6**: 677–681.
- 42 Collins RL, Brand H, Redin CE, Hanscom C, Antolik C, Stone MR, Glessner JT, Mason T,  
43 Pregno G, Dorrani N, et al. 2017. Defining the diverse spectrum of inversions, complex  
44 structural variation, and chromothripsis in the morbid human genome. *Genome Biol* **18**: 36.  
45 <http://genomebiology.biomedcentral.com/articles/10.1186/s13059-017-1158-6>.
- 46 Dagmar M, Beer DG, Wilke CW, Miller DE, Glover TW. 1997. Frequent deletions of FHIT and  
47 FRA3B in Barrett's metaplasia and esophageal adenocarcinomas. *Oncogene* **15**: 1653–  
48 1659. <http://www.nature.com/articles/1201330>.
- 49 Davis A, Gao R, Navin N. 2017. Tumor evolution: Linear, branching, neutral or punctuated?  
50 *Biochim Biophys Acta* **1867**: 151–161. <http://www.ncbi.nlm.nih.gov/pubmed/28110020>.
- 51 de Klein A, van Kessel AG, Grosveld G, Bartram CR, Hagemeijer A, Bootsma D, Spurr NK,



- 1 Heisterkamp N, Groffen J, Stephenson JR. 1982. A cellular oncogene is translocated to the  
2 Philadelphia chromosome in chronic myelocytic leukaemia. *Nature* **300**: 765–7.  
3 <http://www.ncbi.nlm.nih.gov/pubmed/6960256>.
- 4 Deguchi Y, Kehrl JH. 1993. High level expression of the homeobox gene HB24 in a human T-  
5 cell line confers the ability to form tumors in nude mice. *Cancer Res* **53**: 373–7.  
6 <http://www.ncbi.nlm.nih.gov/pubmed/8093351>.
- 7 Deguchi Y, Kirschenbaum A, Kehrl JH. 1992. A diverged homeobox gene is involved in the  
8 proliferation and lineage commitment of human hematopoietic progenitors and highly  
9 expressed in acute myelogenous leukemia. *Blood* **79**: 2841–8.  
10 <http://www.ncbi.nlm.nih.gov/pubmed/1375114>.
- 11 Drexler HG, Matsuo Y, MacLeod RAF. 2004. Malignant hematopoietic cell lines: in vitro models  
12 for the study of erythroleukemia. *Leuk Res* **28**: 1243–51.  
13 <http://www.ncbi.nlm.nih.gov/pubmed/15475063>.
- 14 Durkin SG, Ragland RL, Arlt MF, Mülle JG, Warren ST, Glover TW. 2008. Replication stress  
15 induces tumor-like microdeletions in FHIT/FRA3B. *Proc Natl Acad Sci* **105**: 246–251.  
16 <http://www.pnas.org/cgi/doi/10.1073/pnas.0708097105>.
- 17 Eklund E. 2011. The role of Hox proteins in leukemogenesis: insights into key regulatory events  
18 in hematopoiesis. *Crit Rev Oncog* **16**: 65–76.  
19 <http://www.ncbi.nlm.nih.gov/pubmed/22150308>.
- 20 Engreitz JM, Agarwala V, Mirny LA. 2012. Three-Dimensional Genome Architecture Influences  
21 Partner Selection for Chromosomal Translocations in Human Disease ed. S. Ahmed. *PLoS*  
22 *One* **7**: e44196. <http://dx.plos.org/10.1371/journal.pone.0044196>.
- 23 Fang H, Wu Y, Narzisi G, O’Rawe JA, Barrón LTJ, Rosenbaum J, Ronemus M, lossifov I,  
24 Schatz MC, Lyon GJ. 2014. Reducing INDEL calling errors in whole genome and exome  
25 sequencing data. *Genome Med* **6**: 89. <http://www.ncbi.nlm.nih.gov/pubmed/25426171>.
- 26 Fröhling S. 2012. Widespread over-expression of the non-clustered homeobox gene HLX in  
27 acute myeloid leukemia. *Haematologica* **97**: 1453.  
28 <http://www.ncbi.nlm.nih.gov/pubmed/23053668>.
- 29 Giampaolo A, Pelosi E, Valtieri M, Montesoro E, Sterpetti P, Samoggia P, Camagna A,  
30 Mastroberardino G, Gabbianelli M, Testa U. 1995. HOXB gene expression and function in  
31 differentiating purified hematopoietic progenitors. *Stem Cells* **13 Suppl 1**: 90–105.  
32 <http://www.ncbi.nlm.nih.gov/pubmed/7488973>.
- 33 Greer SU, Nadauld LD, Lau BT, Chen J, Wood-Bouwens C, Ford JM, Kuo CJ, Ji HP. 2017.  
34 Linked read sequencing resolves complex genomic rearrangements in gastric cancer  
35 metastases. *Genome Med* **9**: 57.  
36 <http://genomemedicine.biomedcentral.com/articles/10.1186/s13073-017-0447-8>.
- 37 Gribble SM, Roberts I, Grace C, Andrews KM, Green AR, Nacheva EP. 2000. Cytogenetics of  
38 the Chronic Myeloid Leukemia-Derived Cell Line K562. *Cancer Genet Cytogenet* **118**: 1–8.  
39 <http://linkinghub.elsevier.com/retrieve/pii/S0165460899001697>.
- 40 Groffen J, Stephenson JR, Heisterkamp N, de Klein A, Bartram CR, Grosveld G. 1984.  
41 Philadelphia chromosomal breakpoints are clustered within a limited region, bcr, on  
42 chromosome 22. *Cell* **36**: 93–9. <http://www.ncbi.nlm.nih.gov/pubmed/6319012>.
- 43 Grzanka A, Grzanka D, Orlikowska M. 2003. Cytoskeletal reorganization during process of  
44 apoptosis induced by cytostatic drugs in K-562 and HL-60 leukemia cell lines. *Biochem*  
45 *Pharmacol* **66**: 1611–7. <http://www.ncbi.nlm.nih.gov/pubmed/14555241>.
- 46 Han K, Jeng EE, Hess GT, Morgens DW, Li A, Bassik MC. 2017. Synergistic drug combinations  
47 for cancer identified in a CRISPR screen for pairwise genetic interactions. *Nat Biotechnol*  
48 **35**: 463–474. <http://www.ncbi.nlm.nih.gov/pubmed/28319085>.
- 49 Hastings PJ, Ira G, Lupski JR. 2009. A microhomology-mediated break-induced replication  
50 model for the origin of human copy number variation. *PLoS Genet* **5**: e1000327.  
51 <http://www.ncbi.nlm.nih.gov/pubmed/19180184>.

- 1 Heisterkamp N, Stam K, Groffen J, de Klein A, Grosveld G. 1985. Structural organization of the  
2 bcr gene and its role in the Ph' translocation. *Nature* **315**: 758–61.  
3 <http://www.ncbi.nlm.nih.gov/pubmed/2989703>.
- 4 Hu L, Chen H-Y, Cai J, Yang G-Z, Feng D, Zhai Y-X, Gong H, Qi C-Y, Zhang Y, Fu H, et al.  
5 2015. Upregulation of NETO2 expression correlates with tumor progression and poor  
6 prognosis in colorectal carcinoma. *BMC Cancer* **15**: 1006.  
7 <http://www.ncbi.nlm.nih.gov/pubmed/26699544>.
- 8 Huebner K, Garrison PN, Barnes LD, Croce CM. 1998. The role of the FHIT/FRA3B locus in  
9 cancer. *Annu Rev Genet* **32**: 7–31. <http://www.ncbi.nlm.nih.gov/pubmed/9928473>.
- 10 Ingvarsson S. 2001. FHIT alterations in breast cancer. *Semin Cancer Biol* **11**: 361–366.  
11 <http://linkinghub.elsevier.com/retrieve/pii/S1044579X01903918>.
- 12 Jawad M, Seedhouse CH, Russell N, Plumb M. 2006. Polymorphisms in human homeobox  
13 HLX1 and DNA repair RAD51 genes increase the risk of therapy-related acute myeloid  
14 leukemia. *Blood* **108**: 3916–8. <http://www.ncbi.nlm.nih.gov/pubmed/16902145>.
- 15 Keane TM, Wong K, Adams DJ. 2013. RetroSeq: Transposable element discovery from next-  
16 generation sequencing data. *Bioinformatics* **29**: 389–390.
- 17 Kitzman JO, MacKenzie AP, Adey A, Hiatt JB, Patwardhan RP, Sudmant PH, Ng SB, Alkan C,  
18 Qiu R, Eichler EE, et al. 2011. Haplotype-resolved genome sequencing of a Gujarati Indian  
19 individual. *Nat Biotechnol* **29**: 59–63. <http://www.nature.com/doi/10.1038/nbt.1740>.
- 20 Korbelt JO, Urban AE, Affourtit JP, Godwin B, Grubert F, Simons JF, Kim PM, Palejev D,  
21 Carriero NJ, Du L, et al. 2007. Paired-end mapping reveals extensive structural variation in  
22 the human genome. *Science* **318**: 420–6. <http://www.ncbi.nlm.nih.gov/pubmed/17901297>.
- 23 Krueger F, Andrews SR. 2011. Bismark: A flexible aligner and methylation caller for Bisulfite-  
24 Seq applications. *Bioinformatics* **27**: 1571–1572.
- 25 Krzywinski M, Schein J, Birol I, Connors J, Gascoyne R, Horsman D, Jones SJ, Marra MA. 2009.  
26 Circos: An information aesthetic for comparative genomics. *Genome Res* **19**: 1639–1645.  
27 <http://genome.cshlp.org/cgi/doi/10.1101/gr.092759.109>.
- 28 Lam HYK, Mu XJ, Stütz AM, Tanzer A, Cayting PD, Snyder M, Kim PM, Korbelt JO, Gerstein MB.  
29 2010. Nucleotide-resolution analysis of structural variants using BreakSeq and a  
30 breakpoint library. *Nat Biotechnol* **28**: 47–55.
- 31 Lam HYK, Pan C, Clark MJ, Lacroute P, Chen R, Haraksingh R, O'Huallachain M, Gerstein MB,  
32 Kidd JM, Bustamante CD, et al. 2012. Detecting and annotating genetic variations using  
33 the HugeSeq pipeline. *Nat Biotechnol* **30**: 226–229.  
34 <http://www.nature.com/doi/10.1038/nbt.2134>.
- 35 Layer RM, Chiang C, Quinlan AR, Hall IM. 2014. LUMPY: a probabilistic framework for  
36 structural variant discovery. *Genome Biol* **15**: R84.
- 37 Lee JA, Carvalho CMB, Lupski JR. 2007. A DNA replication mechanism for generating  
38 nonrecurrent rearrangements associated with genomic disorders. *Cell* **131**: 1235–47.  
39 <http://www.ncbi.nlm.nih.gov/pubmed/18160035>.
- 40 Li Y, Zhou S, Schwartz DC, Ma J. 2016. Allele-Specific Quantification of Structural Variations in  
41 Cancer Genomes. *Cell Syst* **3**: 21–34.  
42 <http://linkinghub.elsevier.com/retrieve/pii/S240547121630182X>.
- 43 Lin K, Smit S, Bonnema G, Sanchez-Perez G, de Ridder D. 2015. Making the difference:  
44 integrating structural variation detection tools. *Brief Bioinform* **16**: 852–64.  
45 <http://www.ncbi.nlm.nih.gov/pubmed/25504367>.
- 46 Liu SJ, Horlbeck MA, Cho SW, Birk HS, Malatesta M, He D, Attenello FJ, Villalta JE, Cho MY,  
47 Chen Y, et al. 2017. CRISPRi-based genome-scale identification of functional long  
48 noncoding RNA loci in human cells. *Science (80- )* **355**: eaah7111.  
49 <http://www.sciencemag.org/lookup/doi/10.1126/science.aah7111>.
- 50 Lorenz R, Bernhart SH, Höner Zu Siederdisen C, Tafer H, Flamm C, Stadler PF, Hofacker IL.  
51 2011. ViennaRNA Package 2.0. *Algorithms Mol Biol* **6**: 26.



- 1 <http://www.ncbi.nlm.nih.gov/pubmed/22115189>.
- 2 Lozzio CB, Lozzio BB. 1975. Human chronic myelogenous leukemia cell-line with positive  
3 Philadelphia chromosome. *Blood* **45**: 321–34. <http://www.ncbi.nlm.nih.gov/pubmed/163658>.
- 4 Lubinski J, Hadaczek P, Podolski J, Toloczko A, Sikorski A, McCue P, Druck T, Huebner K.  
5 1994. Common regions of deletion in chromosome regions 3p12 and 3p14.2 in primary  
6 clear cell renal carcinomas. *Cancer Res* **54**: 3710–3.  
7 <http://www.ncbi.nlm.nih.gov/pubmed/8033088>.
- 8 Lupski JR. 2010. Retrotransposition and Structural Variation in the Human Genome. *Cell* **141**:  
9 1110–1112. <http://linkinghub.elsevier.com/retrieve/pii/S0092867410006689>.
- 10 Marks P, Garcia S, Barrio AM, Belhocine K, Bernate J, Bharadwaj R, Bjornson K, Catalanotti C,  
11 Delaney J, Fehr A, et al. 2018. Resolving the Full Spectrum of Human Genome Variation  
12 using Linked-Reads. *bioRxiv* 230946.  
13 <https://www.biorxiv.org/content/early/2018/01/09/230946>.
- 14 Mohiyuddin M, Mu JC, Li J, Bani Asadi N, Gerstein MB, Abyzov A, Wong WH, Lam HYK. 2015.  
15 MetaSV: An accurate and integrative structural-variant caller for next generation  
16 sequencing. *Bioinformatics* **31**: 2741–2744.
- 17 Morgens DW, Deans RM, Li A, Bassik MC. 2016. Systematic comparison of CRISPR/Cas9 and  
18 RNAi screens for essential genes. *Nat Biotechnol* **34**: 634–6.  
19 <http://www.ncbi.nlm.nih.gov/pubmed/27159373>.
- 20 Naumann S, Reutzel D, Speicher M, Decker HJ. 2001. Complete karyotype characterization of  
21 the K562 cell line by combined application of G-banding, multiplex-fluorescence in situ  
22 hybridization, fluorescence in situ hybridization, and comparative genomic hybridization.  
23 *Leuk Res* **25**: 313–22. <http://www.ncbi.nlm.nih.gov/pubmed/11248328>.
- 24 Nishimasu H, Ran FA, Hsu PD, Konermann S, Shehata SI, Dohmae N, Ishitani R, Zhang F,  
25 Nureki O. 2014. Crystal structure of Cas9 in complex with guide RNA and target DNA. *Cell*  
26 **156**: 935–49. <http://www.ncbi.nlm.nih.gov/pubmed/24529477>.
- 27 Ohta M, Inoue H, Cotticelli MG, Kastury K, Baffa R, Palazzo J, Siprashvili Z, Mori M, McCue P,  
28 Druck T, et al. 1996. The FHIT gene, spanning the chromosome 3p14.2 fragile site and  
29 renal carcinoma-associated t(3;8) breakpoint, is abnormal in digestive tract cancers. *Cell*  
30 **84**: 587–97. <http://www.ncbi.nlm.nih.gov/pubmed/8598045>.
- 31 Oparina NY, Sadritdinova AF, Snezhkina A V., Dmitriev AA, Krasnov GS, Senchenko VN,  
32 Melnikova N V., Belenikin MS, Lakunina VA, Veselovsky VA, et al. 2012. Increase in  
33 NETO2 gene expression is a potential molecular genetic marker in renal and lung cancers.  
34 *Russ J Genet* **48**: 506–512. <http://link.springer.com/10.1134/S1022795412050171>.
- 35 Pabinger S, Dander A, Fischer M, Snajder R, Sperk M, Efremova M, Krabichler B, Speicher MR,  
36 Zschocke J, Trajanoski Z. 2014. A survey of tools for variant analysis of next-generation  
37 genome sequencing data. *Brief Bioinform* **15**: 256–278.
- 38 Perry GH, Ben-Dor A, Tsalenko A, Sampas N, Rodriguez-Revenga L, Tran CW, Scheffer A,  
39 Steinfeld I, Tsang P, Yamada NA, et al. 2008. The fine-scale and complex architecture of  
40 human copy-number variation. *Am J Hum Genet* **82**: 685–95.  
41 <http://www.ncbi.nlm.nih.gov/pubmed/18304495>.
- 42 Prasanth SG, Prasanth K V, Stillman B. 2002. Orc6 involved in DNA replication, chromosome  
43 segregation, and cytokinesis. *Science* **297**: 1026–31.  
44 <http://www.ncbi.nlm.nih.gov/pubmed/12169736>.
- 45 Quinlan AR, Hall IM. 2012. Characterizing complex structural variation in germline and somatic  
46 genomes. *Trends Genet* **28**: 43–53.  
47 <http://linkinghub.elsevier.com/retrieve/pii/S0168952511001685>.
- 48 Ran FA, Hsu PD, Wright J, Agarwala V, Scott DA, Zhang F. 2013. Genome engineering using  
49 the CRISPR-Cas9 system. *Nat Protoc* **8**: 2281–2308.  
50 <http://www.nature.com/doifinder/10.1038/nprot.2013.143>.
- 51 Raval A, Tanner SM, Byrd JC, Angerman EB, Perko JD, Chen S-S, Hackanson B, Grever MR,

- 1 Lucas DM, Matkovic JJ, et al. 2007. Downregulation of death-associated protein kinase 1  
2 (DAPK1) in chronic lymphocytic leukemia. *Cell* **129**: 879–90.  
3 <http://www.ncbi.nlm.nih.gov/pubmed/17540169>.
- 4 Seguchi O, Takashima S, Yamazaki S, Asakura M, Asano Y, Shintani Y, Wakeno M, Minamino  
5 T, Kondo H, Furukawa H, et al. 2007. A cardiac myosin light chain kinase regulates  
6 sarcomere assembly in the vertebrate heart. *J Clin Invest* **117**: 2812–24.  
7 <http://www.ncbi.nlm.nih.gov/pubmed/17885681>.
- 8 Selden JR, Emanuel BS, Wang E, Cannizzaro L, Palumbo A, Erikson J, Nowell PC, Rovera G,  
9 Croce CM. 1983. Amplified C lambda and c-abl genes are on the same marker  
10 chromosome in K562 leukemia cells. *Proc Natl Acad Sci U S A* **80**: 7289–92.  
11 <http://www.ncbi.nlm.nih.gov/pubmed/6580644>.
- 12 Shah N, Sukumar S. 2010. The Hox genes and their roles in oncogenesis. *Nat Rev Cancer* **10**:  
13 361–371. <http://www.nature.com/doifinder/10.1038/nrc2826>.
- 14 Sherry ST, Ward MH, Kholodov M, Baker J, Phan L, Smigielski EM, Sirotkin K. 2001. dbSNP:  
15 the NCBI database of genetic variation. *Nucleic Acids Res* **29**: 308–11.
- 16 Shiraishi K, Yamasaki K, Nanba D, Inoue H, Hanakawa Y, Shirakata Y, Hashimoto K,  
17 Higashiyama S. 2007. Pre-B-cell leukemia transcription factor 1 is a major target of  
18 promyelocytic leukemia zinc-finger-mediated melanoma cell growth suppression.  
19 *Oncogene* **26**: 339–48. <http://www.ncbi.nlm.nih.gov/pubmed/16862184>.
- 20 Sloan CA, Chan ET, Davidson JM, Malladi VS, Strattan JS, Hitz BC, Gabdank I, Narayanan AK,  
21 Ho M, Lee BT, et al. 2016. ENCODE data at the ENCODE portal. *Nucleic Acids Res* **44**:  
22 D726–D732. <https://academic.oup.com/nar/article-lookup/doi/10.1093/nar/gkv1160>.
- 23 Snyder MW, Adey A, Kitzman JO, Shendure J. 2015. Haplotype-resolved genome sequencing:  
24 experimental methods and applications. *Nat Rev Genet* **16**: 344–358.  
25 <http://www.nature.com/doifinder/10.1038/nrg3903>.
- 26 Sozzi G, Veronese ML, Negrini M, Baffa R, Cotticelli MG, Inoue H, Torielli S, Pilotti S, De  
27 Gregorio L, Pastorino U, et al. 1996. The FHIT Gene at 3p14.2 Is Abnormal in Lung Cancer.  
28 *Cell* **85**: 17–26. <http://linkinghub.elsevier.com/retrieve/pii/S0092867400810788>.
- 29 Spies N, Weng Z, Bishara A, McDaniel J, Catoe D, Zook JM, Salit M, West RB, Batzoglou S,  
30 Sidow A. 2017. Genome-wide reconstruction of complex structural variants using read  
31 clouds. *Nat Methods*. <http://www.nature.com/doifinder/10.1038/nmeth.4366>.
- 32 Stephens PJ, Greenman CD, Fu B, Yang F, Bignell GR, Mudie LJ, Pleasance ED, Lau KW,  
33 Beare D, Stebbings LA, et al. 2011. Massive genomic rearrangement acquired in a single  
34 catastrophic event during cancer development. *Cell* **144**: 27–40.  
35 <http://www.ncbi.nlm.nih.gov/pubmed/21215367>.
- 36 Stöhr H, Berger C, Fröhlich S, Weber BHF. 2002. A novel gene encoding a putative  
37 transmembrane protein with two extracellular CUB domains and a low-density lipoprotein  
38 class A module: isolation of alternatively spliced isoforms in retina and brain. *Gene* **286**:  
39 223–31. <http://www.ncbi.nlm.nih.gov/pubmed/11943477>.
- 40 Storti P, Donofrio G, Colla S, Airoidi I, Bolzoni M, Agnelli L, Abeltino M, Todoerti K, Lazzaretti M,  
41 Mancini C, et al. 2011. HOXB7 expression by myeloma cells regulates their pro-angiogenic  
42 properties in multiple myeloma patients. *Leukemia* **25**: 527–537.  
43 <http://www.nature.com/doifinder/10.1038/leu.2010.270>.
- 44 Sudmant PH, Rausch T, Gardner EJ, Handsaker RE, Abyzov A, Huddleston J, Zhang Y, Ye K,  
45 Jun G, Hsi-Yang Fritz M, et al. 2015. An integrated map of structural variation in 2,504  
46 human genomes. *Nature* **526**: 75–81.  
47 <http://www.nature.com/doifinder/10.1038/nature15394>.
- 48 Sunagawa GA, Sumiyama K, Ukai-Tadenuma M, Perrin D, Fujishima H, Ukai H, Nishimura O,  
49 Shi S, Ohno R-I, Narumi R, et al. 2016. Mammalian Reverse Genetics without Crossing  
50 Reveals Nr3a as a Short-Sleeper Gene. *Cell Rep* **14**: 662–677.  
51 <http://www.ncbi.nlm.nih.gov/pubmed/26774482>.

- 1 The ENCODE Project Consortium. 2012. An integrated encyclopedia of DNA elements in the  
2 human genome. *Nature* **489**: 57–74. <http://www.ncbi.nlm.nih.gov/pubmed/22955616>.
- 3 Tkachuk DC, Westbrook C a, Andreeff M, Donlon T a, Cleary ML, Suryanarayan K, Homge M,  
4 Redner a, Gray J, Pinkel D. 1990. Detection of bcr-abl fusion in chronic myelogeneous  
5 leukemia by in situ hybridization. *Science* **250**: 559–562.
- 6 Visel A, Minovitsky S, Dubchak I, Pennacchio LA. 2007. VISTA Enhancer Browser--a database  
7 of tissue-specific human enhancers. *Nucleic Acids Res* **35**: D88–92.  
8 <http://www.ncbi.nlm.nih.gov/pubmed/17130149>.
- 9 Wang T, Birsoy K, Hughes NW, Krupczak KM, Post Y, Wei JJ, Lander ES, Sabatini DM. 2015.  
10 Identification and characterization of essential genes in the human genome. *Science* **350**:  
11 1096–101. <http://www.ncbi.nlm.nih.gov/pubmed/26472758>.
- 12 Waters CE, Saldivar JC, Hosseini SA, Huebner K. 2014. The FHIT gene product: tumor  
13 suppressor and genome “caretaker.” *Cell Mol Life Sci* **71**: 4577–4587.  
14 <http://link.springer.com/10.1007/s00018-014-1722-0>.
- 15 Williams LJS, Tabbaa DG, Li N, Berlin AM, Shea TP, MacCallum I, Lawrence MS, Drier Y, Getz  
16 G, Young SK, et al. 2012. Paired-end sequencing of Fosmid libraries by Illumina. *Genome*  
17 *Res* **22**: 2241–2249. <http://genome.cshlp.org/cgi/doi/10.1101/gr.138925.112>.
- 18 Wistuba II, Virmani AK, Gazdar AF, Lam S, LeRiche J, Behrens C, Fong KM, Samet JM,  
19 Srivastava S, Minna JD. 1997. Molecular Damage in the Bronchial Epithelium of Current  
20 and Former Smokers. *JNCI J Natl Cancer Inst* **89**: 1366–1373.  
21 <https://academic.oup.com/jnci/article-lookup/doi/10.1093/jnci/89.18.1366>.
- 22 Wu SQ, Voelkerding K V, Sabatini L, Chen XR, Huang J, Meisner LF. 1995. Extensive  
23 amplification of bcr/abl fusion genes clustered on three marker chromosomes in human  
24 leukemic cell line K-562. *Leukemia* **9**: 858–62.  
25 <http://www.ncbi.nlm.nih.gov/pubmed/7769849>.
- 26 Wu X, Chen H, Parker B, Rubin E, Zhu T, Lee JS, Argani P, Sukumar S. 2006. HOXB7, a  
27 homeodomain protein, is overexpressed in breast cancer and confers epithelial-  
28 mesenchymal transition. *Cancer Res* **66**: 9527–34.  
29 <http://www.ncbi.nlm.nih.gov/pubmed/17018609>.
- 30 Yamashita T, Tazawa S, Yawei Z, Katayama H, Kato Y, Nishiwaki K, Yokohama Y, Ishikawa M.  
31 2006. Suppression of invasive characteristics by antisense introduction of overexpressed  
32 HOX genes in ovarian cancer cells. *Int J Oncol* **28**: 931–8.  
33 <http://www.ncbi.nlm.nih.gov/pubmed/16525643>.
- 34 Ye K, Schulz MH, Long Q, Apweiler R, Ning Z. 2009. Pindel: a pattern growth approach to  
35 detect break points of large deletions and medium sized insertions from paired-end short  
36 reads. *Bioinformatics* **25**: 2865–71. <http://www.ncbi.nlm.nih.gov/pubmed/19561018>  
37 (Accessed June 12, 2017).
- 38 Zheng GXY, Lau BT, Schnall-Levin M, Jarosz M, Bell JM, Hindson CM, Kyriazopoulou-  
39 Panagiotopoulou S, Masquelier DA, Merrill L, Terry JM, et al. 2016. Haplotyping germline  
40 and cancer genomes with high-throughput linked-read sequencing. *Nat Biotechnol* **34**:  
41 303–11. <http://www.ncbi.nlm.nih.gov/pubmed/26829319> (Accessed July 14, 2016).

## 42 43 **FIGURE LEGENDS**

### 44 45 **Figure 1. Comprehensive overview of the K562 genome**

46 Circos (Krzywinski et al. 2009) visualization of K562 genome with the following tracks in inward  
47 concentric order: chromosomes; CN i.e. ploidy by chromosome segment; merged SV density in  
48 1.5 Mb windows of deletions, duplications, and inversions identified using ARC-SV (Arthur et al.  
49 2017), BreakDancer (Chen et al. 2009), BreakSeq (Lam et al. 2010), LUMPY (Layer et al. 2014),  
50 Pindel (Ye et al. 2009), and Long Ranger (Zheng et al. 2016; Marks et al. 2018); phased  
51 haplotype blocks (demarcated with 4 colors for clearer visualization); SNV density in 1 Mb

1 windows; Indel density in 1 Mb windows; dominant zygosity in 1 Mb windows (heterozygous or  
2 homozygous > 50%) with regions exhibiting loss of heterozygosity (LOH) indicated; RNA-seq  
3 reads for loci exhibiting allele-specific expression; CpG islands (CgI) exhibiting allele-specific  
4 methylation; histogram (log-scale) of allele-specifically methylated CpGs in 50 kb windows; non-  
5 reference Alu and LINE-1 insertions; allele-specific CRISPR target sites; large-scale  
6 rearrangements detected by Long Ranger (Zheng et al. 2016; Marks et al. 2018) (light blue:  
7 intrachromosomal; dark blue: interchromosomal); and by GROC-SVs (Spies et al. 2017) (light-  
8 gray: intrachromosomal; dark-gray: interchromosomal).

## 9 10 **Figure 2. K562 ploidy and haplotypes**

11 (A) Representative karyogram of K562 cells produced by GTW banding showing multiple  
12 numerical and structural chromosomal abnormalities and an overall near triploid karyotype.  
13 ISCN 2013 description in relationship to a triploid karyotype [ $<3n>$ ]: 53~70 $<3n>$ ,XX,-X or Y,-  
14 3,?dup(6)(p21p25),+7,?inv(7)(p13p22),add(7)(q32),-9,add(9)(p24),del(9)(p13),add(10)(q22),-  
15 13,add(13)(p11),-14,add(17)(p11.2)x2,add(18)(q23),-20,der(21)t(1;21)(q21;p11),-  
16 22,+4~7mar[cp20]. (B) CN (i.e. ploidy) by percentage across the K562 genome. (C) Percentage  
17 of K562 SNVs and Indels that are novel and known in dbSNP (Sherry et al. 2001). (D) Violin plot,  
18 with overlaid boxplot, of phased haplotype block sizes (Y-axis, log-scaled) where the dashed  
19 line represents the N50 value (2,721,866 bp).

## 20 21 **Figure 3. Mega-Haplotypes of entire K562 chromosome arms**

22  
23 X-axis: chromosome coordinate (Mb). Y-axis: difference in unique linked-read barcode counts  
24 between major and minor haplotypes, normalized for SNV density. Haplotype blocks from of  
25 normal control sample (NA12878) in blue and from K562 in dark gray. Density plots on the right  
26 reflects the distribution of the differences in haplotype-specific barcode counts for control  
27 sample (blue) and K562 (dark gray). These density distributions are used for testing of  
28 significant difference ( $p < 0.001$ ) using one-sided  $t$ -test. Significant difference in haplotype-  
29 specific barcode counts indicate aneuploidy and haplotype imbalance. Haplotype blocks (with  $\geq$   
30 100 phased SNVs) generated from Long Ranger (Dataset S2) for the major and minor  
31 haplotypes were then “stitched” to mega-haplotypes encompassing the entire chromosome arm  
32 of (A) 5q (triploid) and (B) 7p (tetraploid).

## 33 34 **Figure 4. K562 SVs including large complex rearrangements resolved using linked-read 35 sequencing**

36 (A) Heat maps of overlapping barcodes for SVs in K562 resolved from linked-read sequencing  
37 using Long Ranger (Zheng et al. 2016; Marks et al. 2018). *BCR/ABL1* translocation between  
38 chromosomes 9 and 22. *XKR3/NUP214* translocation between chromosomes 9 and 22.  
39 Duplication within *GPHN* on chromosome 14. Deletion that partially overlaps *ZRANB1* and  
40 *CTB2* on chromosome 10. (B) Large complex rearrangement occurring on chromosome 13 with  
41 informative reads from only one haplotype (region with loss-of-heterozygosity). Each line depicts  
42 a fragment inferred from linked-reads based on clustering of identical barcodes (Y-axis) using  
43 GROC-SVs (Spies et al. 2017). Abrupt endings (vertical dashed lines) of fragments indicate  
44 locations of breakpoints of this complex rearrangement. Fragments are phased locally with  
45 respect to surrounding SNVs (colored orange for same haplotype and black when no  
46 informative SNVs are found nearby). Gray lines indicate portions of fragments that do not  
47 support the current breakpoint. Fragments end abruptly at 81.47 Mb, indicating a breakpoint,  
48 picking up again at 81.09 Mb and continuing to 81.11 Mb where they end abruptly, then picking  
49 up again at 90.44 Mb. Coverage from 81.12 Mb to 81.20 Mb are from reads with different sets of  
50 linked-read barcodes and thus not part of this fragment set. (C, D) Complex rearrangements



1 involving multiple haplotype-resolved SVs. Using gemtools (Greer et al. 2017), each SV is  
2 identified from linked-reads grouped by identical barcodes (i.e. SV-specific barcodes, Y-axis)  
3 indicative of single HMW DNA molecules (depicted by each row) that span the breakpoints. SVs  
4 are represented in different colors. X-axis: hg19 genomic coordinate. Dotted lines represent  
5 individual breakpoints with schematic diagram of the rearranged structures drawn below the plot.  
6 (C) Multiple SVs within *FHIT* on 3p14.2. (red) Deletion (DEL) (59.74 Mb – 60.08 Mb) results in  
7 the loss of multiple exons. Two overlapping duplications (DUP) (blue & green) – the presence of  
8 HMW molecules spanning both DUPs indicates a *cis* orientation (same allele of *FHIT*). Two  
9 adjacent DELs (pink & purple) – the spanning HMW molecules for each DEL do not share SV-  
10 specific barcodes, indicating that these DELs are in *trans* (different alleles of *FHIT*). SV  
11 haplotypes analyzed using SV-specific barcodes (not enough informative SNVs due to LOH). (D)  
12 Complex, intra-chromosomal rearrangement spanning approximately 0.5 Mb on 16q11.2 and  
13 16q12.1 that involve two overlapping inversions, 63 kb (red) and 125 kb (blue), and a 163 kb  
14 tandem duplication (green). This rearrangement resides on the non-duplicated haplotype of this  
15 triploid region. *ORC6* is located entirely within the 63 kb inversion on 16q11.2 and is “deleted”  
16 the by the left breakpoint of the 125 kb inversion, which also inverts *MYLK3*. *C16orf8* on the  
17 same haplotype is also partially “deleted” by the 125 kb inversion (blue); *NETO2* is duplicated  
18 by the 163 kb tandem duplication (green). Inset: *MYLK3* and *ORC6* show allele-specific  
19 expression (Supplemental Table S12). *MYLK3* is only expressed from this rearranged allele  
20 (Haplotype 2); *ORC6* is expressed from the non-rearranged “diploid” allele (Haplotype 1).

## 21 22 **Figure 5. Small-scale complex SVs in K562 resolved using ARC-SV**

23 Examples of small-scale complex SVs resolved using ARC-SV (Arthur et al. 2017) from the  
24 K562 WGS dataset. (A) Deletion of Block C and duplication of Block E between Blocks B and D  
25 on chromosome 20 (135,111-136,565). This variant has been validated by PCR. (B) Deletion of  
26 Block B and inverted duplication of Block D between Blocks A and C on chromosome 1  
27 (81,660,347-81,661,554). (C) Duplication and inversion of Blocks B, C, and D between Blocks B  
28 and D on chromosome 3 (158,795,874-158,795,955) overlapping *IQCJ-SCHIP1*. (D) Inversion  
29 of Block C flanked by deletions of Blocks C and D on chromosome 5 (147,553,038-147,554,778)  
30 inside *SPINK14* (coding for a serine peptidase inhibitor). (E) Deletion of Block G, duplications of  
31 blocks I, D, and E, and inverted duplication of Block B between Blocks F and H on chromosome  
32 10 (127,190,417-127,201,193).

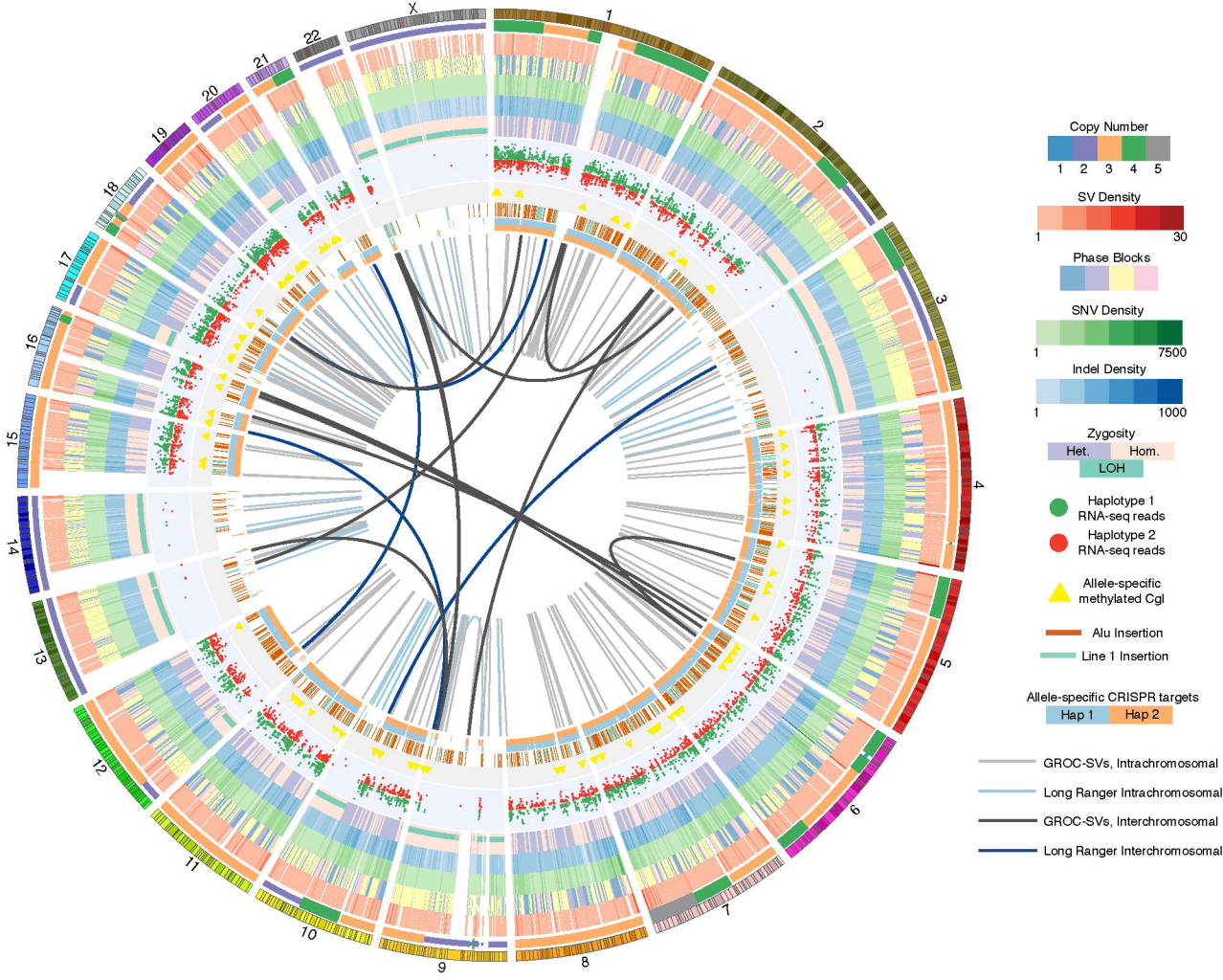
## 33 34 **Figure 6. Genomic structural contexts provide insights into regulatory complexity**

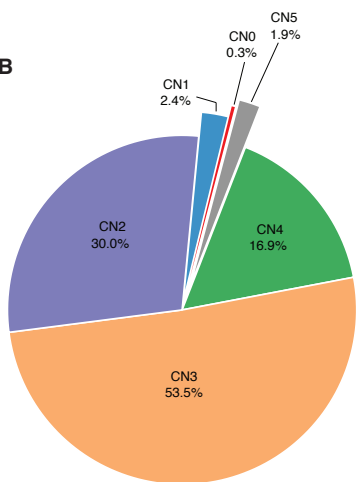
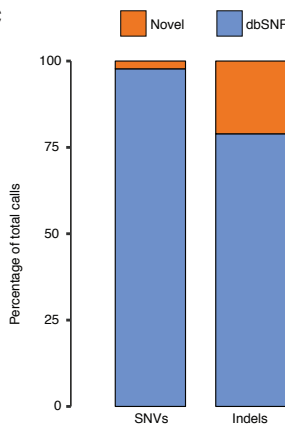
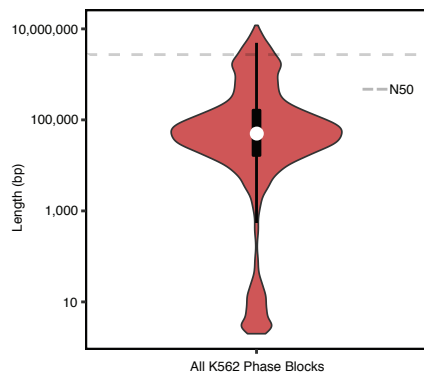
35 (A) Chr17:46,687,000-46,700,000 locus (triploid in K562) containing *HOXB7* and *HOXB8* and  
36 CpG Island (CGI) 22086 (1,203 bp) where phased Haplotype 1 has two copies and Haplotype 2  
37 has one copy. Allele-specific expression of *HOXB7* from Haplotype 1. CpGs in CGI 22086 are  
38 unmethylated in Haplotype 1 and methylated in Haplotype 2. (B) Chr1:221,052,000-221,059,000  
39 locus (tetraploid in K562) containing *HLX* and CGI 2209 (294 bp) where phased Haplotype 1  
40 has three copies and Haplotype 2 has one copy. Allele-specific expression of *HLX* from  
41 Haplotype 1. CpGs in CGI 2209 are unmethylated in Haplotype 2 and highly methylated in  
42 Haplotype 1. Allele-specific CRISPR targeting site 797 bp inside the 5' end of the *HLX* for both  
43 Haplotypes. (C) Number of methylated and unmethylated phased WGBS reads for Haplotypes 1  
44 and 2 in CGI 22086 and CGI 2209 where both CGIs exhibit allele-specific DNA methylation. (D)  
45 Number of RNA-seq reads for Haplotypes 1 and 2 of *HLX* and *HOXB7* where both genes exhibit  
46 allele-specific RNA expression.

## 47 48 **Table 1. Summary of K562 SNVs and Indels**

## 49 50 **Table 2. Haplotypes constructed in aneuploid regions by leveraging haplotype imbalance**

51

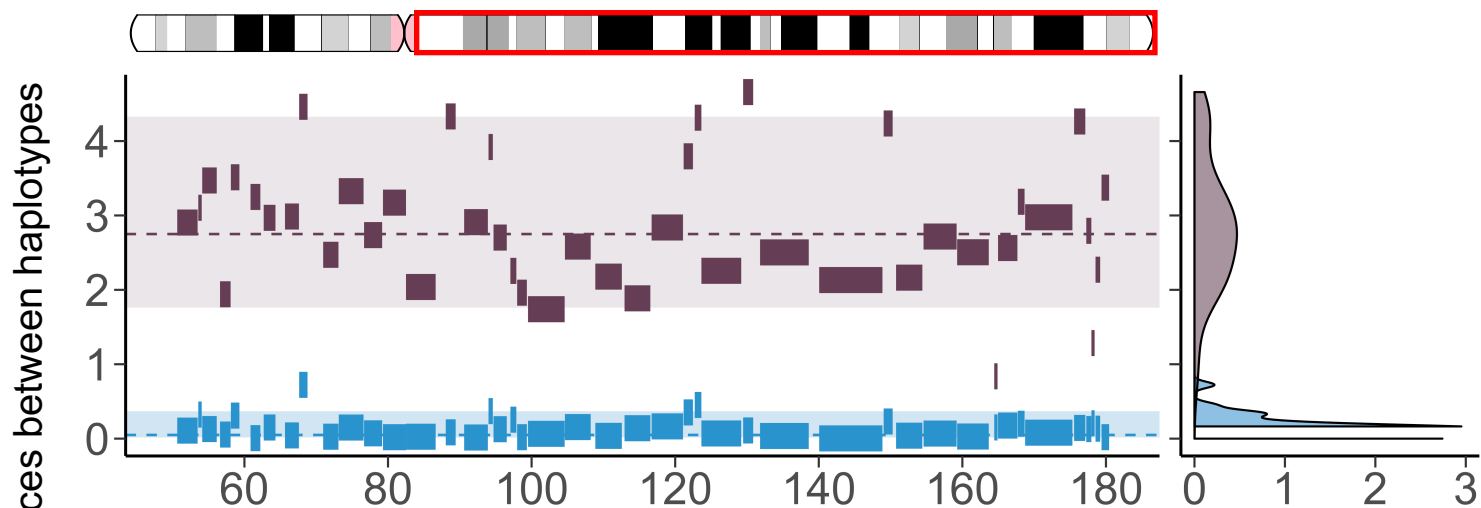


**A****B****C****D**



Control (NA12878) K562

A  
Chr 5q  
Triploid



B  
Chr 7p  
Tetraploid

