# T-CNV: a robust tool for detecting and visualizing copy number variants in targeted sequencing data.

**liu ye**
  Top Gene Tech(Guangzhou) Co., Ltd.

**wu yangming**
  Top Gene Tech(Guangzhou) Co., Ltd.

**zheng zexin**
  Top Gene Tech(Guangzhou) Co., Ltd.

**zhou tianliangwen** ( ✉ tian.zhou@topgene.cn )
  Taylor and Francis Group    https://orcid.org/0000-0003-1269-1489

1    T-CNV: a robust tool for detecting and visualizing copy number variants in targeted

2    sequencing data.

3    Ye Liu[1,2], Yangming Wu[1,2], Zexin Zheng[1,2], Tianliangwen Zhou[1,3]

4    [1] R&D department, Top Gene Tech (Guangzhou) Co., Ltd., Guangzhou, Guangdong, P.R.China

5    [2] These authors contributed equally to this work.

6    [3] Corresponding author. Contact e-mail: tian.zhou@topgene.cn

7    **Abstract**

8    **Background**: Copy number variants (CNVs) are widespread among human genes, causing Mendelian

9    or sporadic traits, or associating with complex diseases. Several tools have been developed for CNV

10    assessment based on next generation sequencing (NGS) data using Read-depth (RD) strategy.

11    However, maintaining high level of sensitivity and specificity is always challenging. Here, we present

12    a novel, powerful, user-friendly and open accessed tool, T-CNV for CNV detection and visualization in

13    targeted NGS panel.

14    **Results**: T-CNV consists of primary CNV detection and CNV candidates confirmation steps. After

15    computing log2 values of normalized read depth ratio of tumor and normal/control sample, T-CNV

16    confirms each possible CNV candidates by bins method, Gaussian Mixture Model (GMM) clustering

17    approach and window-sliding method. We benchmarked its capacity with MLPA-validated dataset.

18    Compared to three other advanced tools, T-CNV presents excellent performance with 95.42%

19    sensitivity, 99.93% specificity and 93.63% positive predict value in MLPA-validated dataset, while

20    achieving satisfactory performance in simulation study (sensitivity 65.95%, positive predict value

21    88.71% at coverage 100X).

22    **Conclusions**: T-CNV is a novel and robust tool for CNV detection and visualization in targeted NGS

23    panel consisting of determination of possible CNV candidates and further confirmation by three

24    different methods. It's publicly available at https://github.com/Top-Gene/T-CNV.

26    **Background**

27    Copy number variants (CNVs) are widespread among human genes, covering 12% of human

28    genome[1]. It can cause Mendelian or sporadic traits or be associated with complex diseases[2, 3],

29    such as neurodevelopmental disorders[4, 5] and cancers[6]. CNV is a large category of structural

30    variants (SVs) that first defined as a segment of DNA that is 1 kb or larger and presents at a variable

31    copy number in comparison with a reference genome[7]. With increasing detections on human

32    genome, CNV has widen to include unbalanced structural variants with >50bp in length[8]. It's

33    considered that the major mechanism for phenotypes derived from CNVs relates to gene dosage

34    effect[9]. CNVs locate either in the dosage-sensitive gene or nearby, which alters or poses an effect on

35    the gene expression level and leads to an abnormal phenotype[10]. Several tools are routinely used

36    for CNVs assessment, including fluorescent *in situ* hybridization (FISH), array comparative genomic

37    hybridization (aCGH), multiplex ligation-dependent probe amplification (MLPA), and recently next

38    generation sequencing (NGS).

39    Although exome- and genome-sequencing techniques are gradually being applied in clinical

40    laboratories, disease-targeted testing still holds a firm place for its high coverage on interested region

41    and cost-effectiveness[11]. To date, MLPA[12] is referred as the gold standard tool for CNV detection,

42    because of its low cost, high sensitivity and specificity, and medium throughput[13]. However, the

43    principle of MLPA mainly encompasses the weakness in lack of sensitivity of regions not directly

44    designed in probe set[14]. In contrast, targeted panel provides a cost-effective and high-throughput

45    way for not only determining SNP and indels, but also identifying common and novel CNVs[14, 15].

46    There are five strategies normally applied in present tools to detect CNVs in NGS data: (1) Paired-end

47    mapping (PEM)[16], (2) read depth (RD)[17], (3) split-read[18], (4) *de novo* assembly[19], and (5)

48    the combination of previously described methods[20, 21]. The concept of RD is similar to that of using

49    density data that a lower region coverage than expected indicates deletion and a higher indicates

50 duplication. Compared with other strategies, RD method has the ability to assess exact copy number

51 of target regions and to detect large insertions and CNVs in complex genomic region classes[17].

52 Although challenges remain, RD has become the routinely used strategy for CNV assessment because

53 of the accumulation of high-coverage NGS data. In order to obtain better prediction, mathematical

54 models were widely used, including hidden Markov chain model (HMM)[22]and Gaussian Mixture

55 Model (GMM)[23, 24].

56 In this study, we develop T-CNV as a novel and powerful tool for CNV detection and visualization in

57 targeted NGS panel. Only autosomal chromosomes were considered in T-CNV to avoid complication

58 of gender. With given interested regions, it utilizes read depth strategy for prior CNV detection and

59 confirmed by three additional steps to achieve high level of sensitivity and specificity. Each possible

60 CNV exons are further confirmed separately by non-overlapped bins method, gaussian mixture model

61 (GMM) clustering approach and overlapped window-sliding method. Later, we conducted

62 performance assessment with three out-of-state targeted sequencing CNV tools, DECoN[25], Atlas-

63 CNV[26] and VisCap[27], on a MLPA-validated dataset ICR96[28] and simulation datasets. DECoN is

64 a modified version of ExomeDepth[22], which applied beta-binomial distribution to capture the

65 variability in read count ratio and combined the likelihood across multiple exons using HMM. The

66 main modifications include compatibility of software and version, and alteration of HMM transition

67 probabilities depends upon distance between exons. Atlas-CNV detects CNVs by measuring log2 ratio

68 and C-score (Z-like-score) to reduce false discovery rate, while VisCap detects CNVs by user-defined

69 log2 ratio threshold and boxplot method. The three tools include visualization of CNV results and are

70 validated on targeted sequencing data. In this study, T-CNV gave satisfactory performance on CNV

71 detection on both validated dataset and simulative datasets. In the meanwhile, T-CNV also provides

72 comprehensive visualization result for better understanding.

3

## Results

### Overview of T-CNV

T-CNV was developed in Python 3.5.6 and R 3.5.1 or higher. It's publicly available at https://github.com/Top-Gene/T-CNV. The overview of T-CNV pipeline is illustrated in Figure 1. The input includes, (1) a bed file containing target regions, (2) DNA sequencing read alignments in BAM format for tumor samples in one pool, (3) if available, DNA sequencing read alignments in BAM format for corresponding normal/control sample. NGS data must be converted to fastq file and aligned to reference genome (hg19) before CNV assessment. Burrows-Wheeler Aligner (BWA)[29] and Genome Analysis Tool Kit (GATK)[30] are recommended for sequence alignment, realignment, recalibration.

### Normalization and GC correction

T-CNV identifies CNV candidates based on read depth strategy. Thus, normalization and bias correction before identification is crucial. GC content was found to have influence on the depth of coverage of NGS data[17, 31]. Since GC content influence may vary between samples even in the same pool, LOESS for GC correction was applied at original read depth of each sample for normalization[32, 33]. We observed the dependency between tumor/normal read depth ratio and GC content (Supplementary Figure S1A). Therefore, we conducted Local Polynomial Regression (LOESS) and smoothing to reduce the influence of GC-bias on log2 values of tumor/normal ratio.

Prior researches corrected GC-content bias by performing LOESS at different level, including fragment read depth[32] and read depth ratios[23, 34]. We compared the outcome for GC correction at different level (Supplementary Figure S1A&B&C&D) and the two strategies showed highly correlated on ICR96 dataset with Pearson correlation r=0.91 in 96 samples (min. 0.76, max. 0.99) (Supplementary Figure S1E&F&G).

## Quality control of ICR96 dataset

In order to reduce false discovery rate, quality control step was introduced. QC poor samples were removed from the pool for better prediction. We conducted a prior noise test on exon log2 value of tumor/normal ratio ($log_2(\frac{T}{N})_{Normalised}$) to derive an index for quality control. We first assumed a sample pool containing 50 pairs of tumor and normal samples with $log_2(\frac{T}{N})_{Normalised} = 1$ targeting 1500 exons. Secondly, random noise was computed in Python 3.5.6 by random function from 5% to 30% in 5% increments and spiked into each exon normalized tumor/normal ratio. The results (Supplementary Figure S2) showed the 50 pairs distribution overview in the left, while in the right showed the fluctuation of $log_2(\frac{T}{N})_{Normalised}$ of each exon in one pair. The worst case that samples were contaminated with 25% noise, may not influence the capacity of CNV calling under threshold of [0.32, -0.42], yielding an overall standard deviation 0.2 of $log_2(\frac{T}{N})_{Normalised}$. Thus, we considered the standard deviation of exon log2 value ($Std_{Exon}$) of tumor and normal ratio, after normalization and GC-bias correction, as a quality control index. A sample with $Std_{Exon} > 0.2$ was regarded as discordant and marked as 'QC poor'. It was removed from the pool in further assessment.

Because the samples in ICR96 dataset were single tumor samples without corresponding normal pair, T-CNV generated the control for the samples in the same pool as illustrated in section "Normalization and quality control" in Method. An overview of the log2 values distribution in pool 1 and pool 2 was illustrated in Figure 2A&B. Nine samples (4 in pool 1, 5 in pool 2) with $Std_{Exon}$ higher than 0.2 were found in ICR96 dataset. Since further CNV assessment was based on log2 value, these nine samples would not be assessed in further performance assessment.

## Optimization of bins method

Normally, RD CNV tools employ segmentation by dividing chromosome into non-overlapped segments to estimate the copy number[17]. T-CNV first identified possible CNV candidates by the exon $log_2(\frac{T}{N})_{Normalised}$. Later, non-overlap bins method was used to verify the candidates. We tuned the bin size in optimization test using ICR96 dataset. The results (Supplementary Figure S3A) indicated

120    that the best solution was setting 30bp bin size, yielding overall sensitivity 95.42%, specificity 99.93%,

121    positive predict value (PPV) 93.63% and negative predict value (NPV) 99.95%. Also, comparing the

122    performance between predicting deletions (Supplementary Figure S3B) and duplications

123    (Supplementary Figure S3C), setting bin size as 30bp presented 95.91%, 94.51% sensitivity and

124    92.66%, 95.56% PPV, respectively.

125    ## Confirmation by GMM clustering

126    Two distinct peaks were observed in the log2 value distribution histogram of genes containing CNV

127    candidates (Supplementary Figure S4A) indicating the distribution of positive CNVs and normal exons.

128    It therefore motivated to detect positive CNVs using GMM clustering method. T-CNV considered the

129    log2 value from three genes (possible CNV candidate gene and other two negative genes) as clustering

130    feature. The specific selection of genes was illustrated in Supplementary Figure S5.

131    Clustering results were given by GMM model and Expectation-Maximization (EM) using standardized

132    log2 value for CNV candidate gene as single feature (Supplementary Figure S4B). Green and red dots

133    in Supplementary Figure S4B stand for two different clusters of log2 value in NF1 in sample 17338,

134    which indicated GMM was able to distinguish positive CNV candidates from normal exons. In

135    accordance with our GMM sampling plan (Supplementary Figure S5), the GMM clustering result for

136    possible CNV candidate exon 46 in gene NF1 in sample 17338 was shown in Supplementary Figure

137    S5C&D&E.

138    ## Optimization of window-sliding method

139    Even though GMM clustering provides a powerful approach to identify true CNVs, the idea behind

140    GMM is soft clustering that each point is assigned to component considering its maximum probability

141    ( $argmax_{component}(P(component|each\ point))$ ). In addition, T-CNV performs GMM clustering

142    approach based on certain sampling strategy (Supplementary Figure S5). False positive would arise

143    when the sampling strategy doesn't fit. During pre-test on ICR96dataset, GMM clustering approach

144    gave 11.35% (32/ (250+32)) false discovery rate (FDR). Thus, other approach was necessary for

145    better performance.

146    Compared to MLPA results, false positive CNVs by bins method and GMM clustering approach

147    presented larger fluctuation (higher standard deviation) at overlapped windows (Supplementary

148    Figure S6). We used overlapped window-sliding method to reduce FDR. Optimization of window-

149    sliding method was conducted by using GMM clustering positive CNV candidates (310 in total), 5bp

150    increments in windows size and 1bp increments in step length. Under 10bp window size and 6bp step

151    length, window-sliding method reached highest area under curve (AUC) 0.82 in receiver operating

152    characteristic curve (ROC) and cut-off value was standard deviation 0.14 (Supplementary Figure

153    S7A&B). The precision-recall graph (Supplementary Figure S7C&D) also demonstrated the same

154    optimal settings as 10bp window size and 6bp step with cut-off value 0.14.

155    Performance comparison on ICR96 dataset

156    To estimate the performance of our tool, we benchmarked T-CNV with a MLPA-validated ICR96

157    dataset, which includes 32 validated genes from prior MLPA test[28]. According to the MLPA validated

158    results, 262 exons (171 deletions, 91 duplications) were determined as true CNV candidates in 59

159    samples, excluding 9 "poor-QC" samples.  In addition, "Normal" exons in MLPA validated results were

160    considered as Non-CNV. T-CNV gave an excellent performance with an overall 95.42% sensitivity

161    (250/262), 95.91% (164/171) sensitivity for deletions and 94.51% (86/91) for duplication,

162    respectively.

163    Furthermore, we compared the performance of T-CNV with other three out-of-state targeted

164    sequencing CNV tools, including DECoN, Atlas-CNV and VisCap. We ran the other three tools under

165    their default setting based on their protocol. Because only autosomal chromosomes are assessed in

166    our study, we fixed the gender as female in tools setting, if necessary.

167    In summary, we compared four parameters, sensitivity, specificity, PPV and NPV among four tools on

168    ICR96 dataset (Figure 3A). In accordance with Roca *et al*[35] findings, DECoN presented best

169    performance among four tools, yielding an overall sensitivity 99.24% (260/262) and PPV 97.01 %

170    (260/268). Compared to Altas and VisCap, T-CNV outperformed with higher sensitivity 95.42%,

171    specificity 99.93%, PPV 93.63% and NPV 99.95% (Figure 3A). Regards to the performance on

172    detecting deletion or duplication, T-CNV and DECoN showed no significance (Figure 3B&C), while,

173    Atlas-CNV and VisCap, gave better prediction on duplications than deletions (Atlas-CNV: false positive

174    rate: deletion 13.87%, duplication: 4.94%; VisCap deletion 24.31%, duplication 6.67%).

175    Performance assessment on simulative dataset

176    We generated simulative dataset containing 100 samples with random size CNVs spiked in each

177    sample using TargetedSim (https://sourceforge.net/projects/targetedsim/files/TargetedSim/,

178    accessed September 10, 2019). We evaluated the performance of T-CNV and DECoN on the simulative

179    dataset at average coverage 50X, 100X and 500X in triplicate, under default settings of both tools. The

180    result for sensitivity and PPV study (Figure 4A) demonstrated steady performance of T-CNV at

181    different depth. While DECoN presented higher sensitivity (Figure 4B), T-CNV achieved higher PPV

182    than DECoN, indicating that less false positive was predicted by T-CNV. T-CNV showed higher sensitive

183    on 100X and 500X, mainly due to the minimum coverage setting 30X in default settings.

184    **Discussion**

185    We reported T-CNV as a powerful tool for CNV detection and visualization in targeted DNA sequencing

186    panels. It allows a set to samples or tumor samples with their normal pair as input. Also, it accepts

187    bam files or locus-depth files generated from SAMtools by command "depth". When locus-depth files

188    are not available, T-CNV converts bam files into locus-depth files automatedly and starts CNV analysis.

189    T-CNV implemented LOESS at log2 value level, since its result was highly correlated to the GC-

190    correction result of LOESS at original depth level. In addition, T-CNV considered a standard deviation

191    of log2 value as an index to separate discordant samples to achieve high sensitivity. Based on the QC

192    criteria, we found 9 samples in ICR96 dataset were marked as poor QC. We also measured the Pearson

193    correlation coefficient of samples and its corresponding pool control (Figure 2C). The QC poor samples

194  also showed less correlation to their controls, indicating the noise and bias in QC poor samples were

195  not effectively removed by previous normalization steps.

196  Since single exon CNV detection is challenging[36], after filtering possible CNV candidates, T-CNV

197  preforms three critical confirmation methods at each candidate.  Thus, it's compromised to give false

198  negative prediction in a CNV covering consecutive exons (Figure 5A&B). T-CNV separates possible

199  CNV candidates in the primary step by measuring whether the log2 value exceeds threshold, where

200  false negative candidates might be yielded, such as NF1 exon 45 in sample 17338 (Figure 5B). Another

201  possibility for false negative prediction results from confirmation methods. A positive candidate is

202  identified only when all confirmation methods determine it as positive.

203  False positive prediction is inevitable due to cut-off setting. T-CNV has optimized thresholds and cut-

204  off value in bins method and window-sliding method using ICR96 dataset to achieve best performance.

205  When comparing false positive prediction with other tools, we observed that Atlas-CNV and T-CNV

206  detected BAP1 exon 1 in sample 17340 as positive duplication (Figure 5C), while all four tools defined

207  a CDH1 exon 1 in sample 17301 as positive deletion (Figure 5D). These false positive CNVs have small

208  size (around 50bp), which might influence the capacity of bins method and window-sliding method

209  due to the default size selection, resulting in false identification.

210  An example of the visualization results (sample 17375 from pool2 in ICR96 dataset) was  shown in

211  Supplementary Figure S8, plotting with the log2 value of all targeted exons. The predicted copy

212  number of positive CNVs was indicated under the ploting, calculated from the log2 value.

213  T-CNV has been benchmarked by ICR96, achieving excellent sensitivity 95.42%, specificity 99.93%,

214  PPV 93.63%, NPV 99.95%. Furthermore, T-CNV achieved satisfactory results in simulation study.

215  Since GMM clustering approach requires specific sampling plan, T-CNV achieve best performance in

216  panel with multiple genes. Thus, it's highly recommended to select appropriate target region when

217  applying T-CNV.

218 **Conclusions**

219 This study reports a robust, novel and open source CNV tool, named T-CNV, consisting of
220 determination of possible CNV candidates and further confirmation by three different methods. It
221 accepted a pair of tumor and normal samples of a set of samples as input. It is able to automatedly
222 convert bam file into SAMtools depth file, while depth files are acceptable. T-CNV is benchmarked by
223 MLPA-validated dataset and also presents satisfactory performance on simulative dataset. As NGS are
224 widely used clinically and academically, we believe T-CNV provides a parallel solution for CNV
225 detection while calling SNPs and indels.

226 **Methods**

227 T-CNV pipeline

228 Normalization and quality control

229 To start CNVs calling, RD for each locus inside target region was determined by SAMtools[37]. Tumor
230 sample and corresponding normal sample were considered as a pair in T-CNV. When normal/control
231 sample was not available, the median RD within the interval $[mean(RD) \pm 2 \cdot$
232 $standard\ deviation(RD)]$ in sample pool served as control sample.

233 Normalization between the RD of a pair sample (tumor and normal/control) is crucial. T-CNV
234 computed reads per thousand bases per million reads sequenced (RPKM) at each exon by

235 $RPKM = RD\ at\ exon \cdot 10^9/(RD\ in\ target\ regions \cdot Exon\ length)$.

236 The normalized RD of tumor sample and normal sample was converted to log2-value
237 $(log_2(RPKM_{tumor}/RPKM_{normal}))$. Similar to Benjamini and Speed's[33] research, log2-value was GC-
238 bias corrected by LOESS in R 3.5.1with 75% smoothing span and degree 2.

10

239      T-CNV considered thresholds of -0.42 for deletion and 0.32 for duplication for diploid. Any exon with

240      log2 value higher than 0.32 or lower than -0.42 was considered as possible CNV candidates. The

241      possible CNV candidates were further confirmed by follow-up steps.

242      Samples with evenly distribution of exon RD were expected in a pool. T-CNV considered the standard

243      deviation of all exon log2-value ($Std_{Exon}$) as an index for quality control. A noise randomly introduced

244      test was conducted in a pool of 50 samples containing 1500 exons. In the test, 5% increment of noise

245      was randomly introduced into all samples on log2 value using uniform function in "random" package

246      in Python 3.5.6. The results (Supplementary Figure S2) illustrated that maximum 25% of noise,

247      corresponding to $Std_{Exon}$ 0.2, under threshold [-0.42, 0.32] may not influence CNV candidate

248      identification. Therefore, to reduce false positive, samples with $Std_{Exon}$ higher than 0.2 were

249      considered as poor quality.

250      Confirmation of CNVs

251      **Bins method**

252      The possible CNV exons are divided into non-overlap bins with 30bp in length as default value.

253      Average RD for each bin in tumor and normal samples were calculated and log2 value was determined.

254      One CNV was confirmed by bins method, if 90% or more bins in the exon presented log2 values higher

255      than 0.32 or lower than -0.42.

256      **GMM clustering method**

257      We first assume that the probability distribution of GC-corrected log2 value can be approximated with

258      a set of gaussian (deletion, duplication and non-CNV). Let $\{\mu_j, \Sigma_j, \phi_j\}$ denote the mean, covariance, the

259      weight of $j^{th}$ component of GMM ($1 \leq j \leq K$). The likelihood of the observed GC-corrected log2 value

260      $x_i$ is

261      $P(x_i) = \sum_{j=1}^{K} P(Z_i = j) \cdot P(x_i | Z_i = j) = \sum_{j=1}^{K} \phi_j \cdot \mathcal{N}(x_i | \mu_j, \Sigma_j),$

262 where $Z_i$ represents the hidden CNV state of $x_i$, $Z_i \in \{1, \dots, K\}$ and $\sum_{j=1}^{K} \phi_j = 1$ ($\phi_j \geq 0, \forall j$). Based

263 on our assumption, observed data (GC-corrected log2 values) are 1-dimension. Therefore, the

264 probability for each set of gaussian is

265 $\mathcal{N}(x_i|\mu_j, \sigma_j^2) = \frac{1}{\sigma_j\sqrt{2\pi}} \exp\left(-\frac{(x_i-\mu_j)^2}{2\sigma_j^2}\right).$

266 To confirm the possible CNV candidate, we select two other genes without CNV candidates, located in

267 the same chromosome as observed data (illustrated in Supplement Figure S5). In this case, two-

268 components (positive and normal) GMM ($K = 2$) is fitted with unknown parameters $\theta =$

269 $\{\mu_1, \dots, \mu_K, \sigma_1, \dots, \sigma_K, \phi_1, \dots, \phi_K\}$. Given observed $X = \{x_1, x_2, \dots, x_i, \dots, x_N\}$, the log-likelihood is

270 $\ell(\theta) = \sum_{i=1}^{N}(\ln \sum_{j=1}^{K} \phi_j \cdot \mathcal{N}(x_i|\mu_j, \Sigma_j)).$

271 Our goal is to maximize the log-likelihood. Thus, we employed EM algorithm to estimate the parameter

272 $\{\mu_j, \sigma_j, \phi_j\}$ of GMM as followed.

273 E-step: compute the conditional probability that $x_i$ belongs to $j^{th}$ component of GMM for observed X

274 $P(Z_i = j|x_i) = \hat{\gamma}_{i,j} = \frac{P(Z_i=j) \cdot P(x_i|Z_i=j)}{P(x_i)} = \frac{\hat{\phi}_j \cdot \mathcal{N}(x_i|\hat{\mu}_j, \hat{\sigma}_j^2)}{\sum_{k=1}^{K} \hat{\phi}_k \cdot \mathcal{N}(x_i|\hat{\mu}_k, \hat{\sigma}_k^2)}.$

275 M-step: use the updated $\hat{\gamma}_{i,j}$ calculate $\{\hat{\mu}_j, \hat{\sigma}_j^2, \hat{\phi}_j\}$ by

276 $\hat{\phi}_j = \sum_{i=1}^{N} \frac{\hat{\gamma}_{i,j}}{N},$

277 $\hat{\mu}_j = \frac{\sum_{i=1}^{N} \hat{\gamma}_{i,j} \cdot x_i}{\sum_{i=1}^{N} \hat{\gamma}_{i,j}},$

278 $\hat{\sigma}_j^2 = \frac{\sum_{i=1}^{N} \hat{\gamma}_{i,j} \cdot (x_i - \hat{\mu}_j)^2}{\sum_{i=1}^{N} \hat{\gamma}_{i,j}}.$

279 Iterate the E-step and M-step until convergence, when $\{\hat{\mu}_j, \hat{\sigma}_j^2, \hat{\phi}_j\}^{t+1} \approx \{\hat{\mu}_j, \hat{\sigma}_j^2, \hat{\phi}_j\}^{t}$. The above steps

280 were completed by using GaussianMixture in scikit-learn[38] package in Python 3.5.6.

**Window-sliding method**

The possible CNV exon was divided into overlapped M bp windows with N bp steps 'slide' away. RD for each window and log2 values in tumor and normal samples were computed. Later, the standard deviation for log2 values of windows ($Std(log_2(windows))$) was calculated. A possible CNV candidate was confirmed as positive when the $Std(log_2(windows))$ was lower than 0.14.

A CNV candidate was confirmed as positive if above all confirmation methods identified as positive. Thus, the output of a sample includes the confirmation results from bins method, GMM clustering method and window-sliding method, quality control results and visualization file. An example of T-CNV result for CNV-positive sample was shown in Supplementary Figure S8.

MLPA validated dataset

The MLPA validated dataset ICR96[28] can be accessed through the European-Genome phenome Archive (EGA) under accession number EGAS00001002428. ICR96 consists of 96 targeted NGS samples (66 MLPA validated CNV-positive samples and 30 CNV-negative samples). Validated positive samples contains CNVs in BRCA1, BRCA2, TP53, MLH1, MSH2, PMS2, EPCAM or PTEN, which were most frequently tested in clinical practice as cancer predisposition genes.

Simulation dataset

We simulated Illumina paired-end reads datasets spiked with deletion and duplication CNVs by TargetedSim (https://sourceforge.net/projects/targetedsim/files/TargetedSim/, accessed September 10, 2019). Homozygous and heterozygous were indicated by 100% and 50% reduce/increase in depth for deletion/duplication, respectively. Simulative dataset consisted of 100 samples, randomly spiked with CNVs (3 heterozygous deletions, 3 heterozygous duplications and 2 homozygous duplications with random length from 1kb to 10kb) in each sample with targeted region of 1634 exons. To compare CNV calling performance with DECoN, we conducted analysis from dataset generation to CNV detection at coverage 50X, 100X and 500X.

305   In this study, true positive (TP) is defined as MLPA-validated or known CNVs, while true negative (TN)

306   is defined as negative exons. Sensitivity (TP/ (TP + false negative (FN))), specificity (TN/ (TN + false

307   positive (FP))), PPV (TP/ (TP+FP)) and negative predict value (NPV) (TN/ (TN+FN)) were calculated

308   in exon basis.

309   **List of abbreviations**

310   **CNVs:** copy number variants

311   **SVs:** structural variants

312   **FISH:** fluorescent *in situ* hybridization

313   **aCGH:** array comparative genomic hybridization

314   **MLPA:** multiplex ligation-dependent probe amplification

315   **NGS:** next generation sequencing

316   **RD:** read depth

317   **PEM:** Paired-end mapping

318   **HMM:** hidden Markov chain model

319   **GMM:** Gaussian Mixture Model

320   **BWA:** Burrows-Wheeler Aligner

321   **GATK:** Genome Analysis Tool Kit

322   **LOESS:** Local Polynomial Regression

323   **PPV:** positive predict value

324   **NPV:** negative predict value

325   **EM:** Expectation-Maximization

326   **FDR:** false discovery rate

327  **AUC:** area under curve

328  **ROC:** receiver operating characteristic curve

329  **RPKM:** reads per thousand bases per million reads sequenced

330  **TP:** true positive

331  **TN:** true negative

332  **FN:** false negative

333  **FP:** false positive

334  **Declarations**

335  Ethics approval and consent to participate

336  Not applicable.

337  Consent for publication

338  Not applicable.

339  Competing interests

340  All authors are employed by Top Gene Tech (Guangzhou) Co., Ltd. by the time of manuscript

341  submission. The authors declare that they have no competing interests.

342  Funding

343  Not applicable.

344  Availability of data and materials

345  The ICR96 dataset is public by Professor Nazneen Rahman's team and can be accessed through the

346  European-Genome phenome Archive (EGA) under the accession number EGAS00001002428[28]. The

15

347 simulative dataset was generated by using TargetedSim

348 ([https://sourceforge.net/projects/targetedsim/files/TargetedSim/](https://sourceforge.net/projects/targetedsim/files/TargetedSim/), accessed September 10, 2019).

## Authors' contributions

350 YL prepared data, generated all figures and wrote the manuscript. YW and ZZ designed the tool and

351 prepared the performance assessment on ICR96 dataset. ZZ and YL prepared the performance

352 assessment on simulative dataset. TZ supervised the project. YL, YW and ZZ contribute equally to this

353 work. All authors have read and approved the final manuscript.

## Author information

355 Affiliations:

356 R&D department, Top Gene Tech (Guangzhou) Co., Ltd., Rm H10F, GT Land Winter Plaza,

357 Zhujiangdong Road, Guangzhou, Guangdong, P.R.China

358 Ye Liu, Yangming Wu, Zexin Zheng, Tianliangwen Zhou

## Corresponding author

360 Tianliangwen Zhou

361 R&D department, Top Gene Tech (Guangzhou) Co., Ltd., Rm H10F, GT Land Winter Plaza,

362 Zhujiangdong Road, Guangzhou, Guangdong, P.R.China

363 tian.zhou@topgene.cn

367     **References**

368     1.     Redon R, Ishikawa S, Fitch KR, Feuk L, Perry GH, Andrews TD, Fiegler H, Shapero MH, Carson
369            AR, Chen W *et al*: **Global variation in copy number in the human genome**. *Nature* 2006,
370            **444**(7118):444-454.

371     2.     Zhang F, Gu W, Hurles ME, Lupski JR: **Copy number variation in human health, disease,**
372            **and evolution**. *Annu Rev Genomics Hum Genet* 2009, **10**:451-481.

373     3.     Mikhail FM: **Copy number variations and human genetic disease**. *Curr Opin Pediatr* 2014,
374            **26**(6):646-652.

375     4.     Coe BP, Girirajan S, Eichler EE: **The genetic variability and commonality of**
376            **neurodevelopmental disease**. *Am J Med Genet C Semin Med Genet* 2012, **160C**(2):118-129.

377     5.     Lee JA, Lupski JR: **Genomic rearrangements and gene copy-number alterations as a cause**
378            **of nervous system disorders**. *Neuron* 2006, **52**(1):103-121.

379     6.     Mamlouk S, Childs LH, Aust D, Heim D, Melching F, Oliveira C, Wolf T, Durek P, Schumacher D,
380            Blaker H *et al*: **DNA copy number changes define spatial patterns of heterogeneity in**
381            **colorectal cancer**. *Nat Commun* 2017, **8**:14093.

382     7.     Feuk L, Carson AR, Scherer SW: **Structural variation in the human genome**. *Nat Rev Genet*
383            2006, **7**(2):85-97.

384     8.     Alkan C, Coe BP, Eichler EE: **Genome structural variation discovery and genotyping**. *Nat*
385            *Rev Genet* 2011, **12**(5):363-376.

386     9.     Lupski JR, Wise CA, Kuwano A, Pentao L, Parke JT, Glaze DG, Ledbetter DH, Greenberg F, Patel
387            PI: **Gene dosage is a mechanism for Charcot-Marie-Tooth disease type 1A**. *Nat Genet*
388            1992, **1**(1):29-33.

389     10.    Stankiewicz P, Lupski JR: **Structural variation in the human genome and its role in**
390            **disease**. *Annu Rev Med* 2010, **61**:437-455.

391     11.    Rehm HL: **Disease-targeted sequencing: a cornerstone in the clinic**. *Nat Rev Genet* 2013,
392            **14**(4):295-300.

393     12.    Schouten JP, McElgunn CJ, Waaijer R, Zwijnenburg D, Diepvens F, Pals G: **Relative**
394            **quantification of 40 nucleic acid sequences by multiplex ligation-dependent probe**
395            **amplification**. *Nucleic Acids Res* 2002, **30**(12):e57.

396     13.    Homig-Holzel C, Savola S: **Multiplex ligation-dependent probe amplification (MLPA) in**
397            **tumor diagnostics and prognostics**. *Diagn Mol Pathol* 2012, **21**(4):189-206.

398     14.    Schenkel LC, Kerkhof J, Stuart A, Reilly J, Eng B, Woodside C, Levstik A, Howlett CJ, Rupar AC,
399            Knoll JHM *et al*: **Clinical Next-Generation Sequencing Pipeline Outperforms a Combined**
400            **Approach Using Sanger Sequencing and Multiplex Ligation-Dependent Probe**
401            **Amplification in Targeted Gene Panel Analysis**. *J Mol Diagn* 2016, **18**(5):657-667.

402     15.    Kerkhof J, Schenkel LC, Reilly J, McRobbie S, Aref-Eshghi E, Stuart A, Rupar CA, Adams P,
403            Hegele RA, Lin H *et al*: **Clinical Validation of Copy Number Variant Detection from**
404            **Targeted Next-Generation Sequencing Panels**. *J Mol Diagn* 2017, **19**(6):905-920.

405   16.   Korbel JO, Urban AE, Affourtit JP, Godwin B, Grubert F, Simons JF, Kim PM, Palejev D, Carriero
406         NJ, Du L *et al*: **Paired-end mapping reveals extensive structural variation in the human**
407         **genome**. *Science* 2007, **318**(5849):420-426.

408   17.   Yoon S, Xuan Z, Makarov V, Ye K, Sebat J: **Sensitive and accurate detection of copy number**
409         **variants using read depth of coverage**. *Genome Res* 2009, **19**(9):1586-1592.

410   18.   Ye K, Schulz MH, Long Q, Apweiler R, Ning Z: **Pindel: a pattern growth approach to detect**
411         **break points of large deletions and medium sized insertions from paired-end short**
412         **reads**. *Bioinformatics* 2009, **25**(21):2865-2871.

413   19.   Li Y, Zheng H, Luo R, Wu H, Zhu H, Li R, Cao H, Wu B, Huang S, Shao H *et al*: **Structural**
414         **variation in two human genomes mapped at single-nucleotide resolution by whole**
415         **genome de novo assembly**. *Nat Biotechnol* 2011, **29**(8):723-730.

416   20.   Zeitouni B, Boeva V, Janoueix-Lerosey I, Loeillet S, Legoix-ne P, Nicolas A, Delattre O, Barillot
417         E: **SVDetect: a tool to identify genomic structural variations from paired-end and mate-**
418         **pair sequencing data**. *Bioinformatics* 2010, **26**(15):1895-1896.

419   21.   Medvedev P, Fiume M, Dzamba M, Smith T, Brudno M: **Detecting copy number variation**
420         **with mated short reads**. *Genome Res* 2010, **20**(11):1613-1622.

421   22.   Plagnol V, Curtis J, Epstein M, Mok KY, Stebbings E, Grigoriadou S, Wood NW, Hambleton S,
422         Burns SO, Thrasher AJ *et al*: **A robust model for read count data in exome sequencing**
423         **experiments and implications for copy number variant calling**. *Bioinformatics* 2012,
424         **28**(21):2747-2754.

425   23.   Gusnanto A, Wood HM, Pawitan Y, Rabbitts P, Berri S: **Correcting for cancer genome size**
426         **and tumour cell content enables better estimation of copy number alterations from**
427         **next-generation sequence data**. *Bioinformatics* 2012, **28**(1):40-47.

428   24.   Li Y, Zhang J, Yuan X: **BagGMM: Calling copy number variation by bagging multiple**
429         **Gaussian mixture models from tumor and matched normal next-generation sequencing**
430         **data**. *Digital Signal Processing* 2019, **88**:90-100.

431   25.   Fowler A, Mahamdallie S, Ruark E, Seal S, Ramsay E, Clarke M, Uddin I, Wylie H, Strydom A,
432         Lunter G *et al*: **Accurate clinical detection of exon copy number variants in a targeted**
433         **NGS panel using DECoN**. *Wellcome Open Res* 2016, **1**:20.

434   26.   Chiang T, Liu X, Wu TJ, Hu J, Sedlazeck FJ, White S, Schaid D, Andrade M, Jarvik GP, Crosslin D
435         *et al*: **Atlas-CNV: a validated approach to call single-exon CNVs in the eMERGESeq gene**
436         **panel**. *Genet Med* 2019, **21**(9):2135-2144.

437   27.   Pugh TJ, Amr SS, Bowser MJ, Gowrisankar S, Hynes E, Mahanta LM, Rehm HL, Funke B, Lebo
438         MS: **VisCap: inference and visualization of germ-line copy-number variants from**
439         **targeted clinical sequencing data**. *Genet Med* 2016, **18**(7):712-719.

440   28.   Mahamdallie S, Ruark E, Yost S, Ramsay E, Uddin I, Wylie H, Elliott A, Strydom A, Renwick A,
441         Seal S *et al*: **The ICR96 exon CNV validation series: a resource for orthogonal assessment**
442         **of exon CNV calling in NGS data**. *Wellcome Open Res* 2017, **2**:35.

443   29.   Li H, Durbin R: **Fast and accurate short read alignment with Burrows-Wheeler**
444         **transform**. *Bioinformatics* 2009, **25**(14):1754-1760.

445   30.   McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernytsky A, Garimella K, Altshuler
446         D, Gabriel S, Daly M *et al*: **The Genome Analysis Toolkit: a MapReduce framework for**
447         **analyzing next-generation DNA sequencing data**. *Genome Res* 2010, **20**(9):1297-1303.

448   31.   Boeva V, Zinovyev A, Bleakley K, Vert JP, Janoueix-Lerosey I, Delattre O, Barillot E: **Control-**
449         **free calling of copy number alterations in deep-sequencing data using GC-content**
450         **normalization**. *Bioinformatics* 2011, **27**(2):268-269.

451   32.   Boeva V, Popova T, Lienard M, Toffoli S, Kamal M, Le Tourneau C, Gentien D, Servant N,
452         Gestraud P, Rio Frio T *et al*: **Multi-factor data normalization enables the detection of copy**
453         **number aberrations in amplicon sequencing data**. *Bioinformatics* 2014, **30**(24):3443-
454         3450.

455   33.   Benjamini Y, Speed TP: **Summarizing and correcting the GC content bias in high-**
456         **throughput sequencing**. *Nucleic Acids Res* 2012, **40**(10):e72.

457   34.   Talevich E, Shain AH, Botton T, Bastian BC: **CNVkit: Genome-Wide Copy Number Detection**
458         **and Visualization from Targeted DNA Sequencing**. *PLoS Comput Biol* 2016,
459         **12**(4):e1004873.

460   35.   Roca I, Gonzalez-Castro L, Fernandez H, Couce ML, Fernandez-Marmiesse A: **Free-access**
461         **copy-number variant detection tools for targeted next-generation sequencing data**.
462         *Mutat Res* 2019, **779**:114-125.

463   36.   de Ligt J, Boone PM, Pfundt R, Vissers LE, Richmond T, Geoghegan J, O'Moore K, de Leeuw N,
464         Shaw C, Brunner HG *et al*: **Detection of clinically relevant copy number variants with**
465         **whole-exome sequencing**. *Hum Mutat* 2013, **34**(10):1439-1448.

466   37.   Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G, Durbin R,
467         Genome Project Data Processing S: **The Sequence Alignment/Map format and SAMtools**.
468         *Bioinformatics* 2009, **25**(16):2078-2079.

469   38.   Pedregosa F, Varoquaux G, Gramfort A, Michel V, Thirion B, Grisel O, Blondel M, Prettenhofer
470         P, Weiss R, Dubourg V *et al*: **Scikit-learn: Machine Learning in Python**. *JMLR* 2011,
471         **12**(85):2825–2830.

472   **Figure Legends**

473   **Figure 1 Pipeline of T-CNV**

474   To start CNVs assessment, read depth at each locus located with the targeted region was determined.

475   Tumor sample and corresponding normal/control sample were considered as a pair in T-CNV, but the

476   latter was not strictly required.Log2 value was calculated and GC- correction was computed. Any exon

477   with log2 value higher than 0.32 or lower than -0.42 was considered as possible CNV candidates. The

478   CNV candidates were further confirmed by the 3 methods including bins-method, GMM-clustering

479   method and window-slicing method.

480    **Figure 2 Box-plot for Log2 value and Pearson correlation coefficient in ICR96 dataset**

481    (A) Four samples with in pool1 were marked as QC poor with $Std_{Exon}$ >0.2. (B) Five samples in pool2

482    were marked as QC poor with $Std_{Exon}$ >0.2. (C) Pearson correlation coefficient of each sample and its

483    corresponding pool control in ICR96 dataset. QC poor samples are indicated with red font.

484    **Figure 3 Performance comparison with DECoN, Atlas-CNV and VisCap on ICR96 dataset**

485     (A) T-CNV, DECoN, Atlas-CNV and VisCap were used to detect CNVs on ICR96 dataset and their

486    performance were presented using four parameters: sensitivity, specificity, positive predict value

487    (PPV) and negative predict value (NPV). (B) Four tools performance on deletion CNVs. (C) Four tools

488    performance on duplication CNVs.

489    **Figure 4 Performance assessment of T-CNV and DECoN on simulative dataset**

490    (A) Performance of T-CNV on simulative dataset under 50X, 100X and 500X was presented as

491    sensitivity and PPV.  "all" represented all simulative CNVs. "del" represented the simulative deletions

492    and "dup" represented the simulative duplications. (B) Performance of DECoN on simulative dataset.

493    **Figure 5 Example of false positive and false negative prediction in ICR96 dataset**

494    (A) A false negative CHEK2 exon 13 in sample 17332 was identified by T-CNV. (B) Four false negative

495    exons located in NF1 gene in 17338 sample were identified by T-CNV. (C) A false positive BAP1 exon

496    1 in sample 17340 was identified by T-CNV. (D) A false positive CDH1 exon 1 in sample 17301 was

497    identified by T-CNV.

498    **Additional file 1**: Supplementary figures. Supplementary Figure S1 LOESS used in GC-content

499    correction. Supplementary Figure S2 Random noise test. Supplementary Figure S3 Optimazation of

500    bins method. Supplementary Figure S4 An example of GMM clustering result on sample 17338 NF1

501    exon 37-57 delietion. Supplementary Figure S5 Sampling plan for GMM clustering in T-CNV.

502    Supplementary Figure S6 True positive and false positive in T-CNV window-sliding method.

503    Supplementary Figure S7 ROC curve and PR curve for optimization of window-sliding method.

504     Supplementary Figure S8 The visulization result of sample 17375 from pool2 in ICR96 dataset. (DOCX
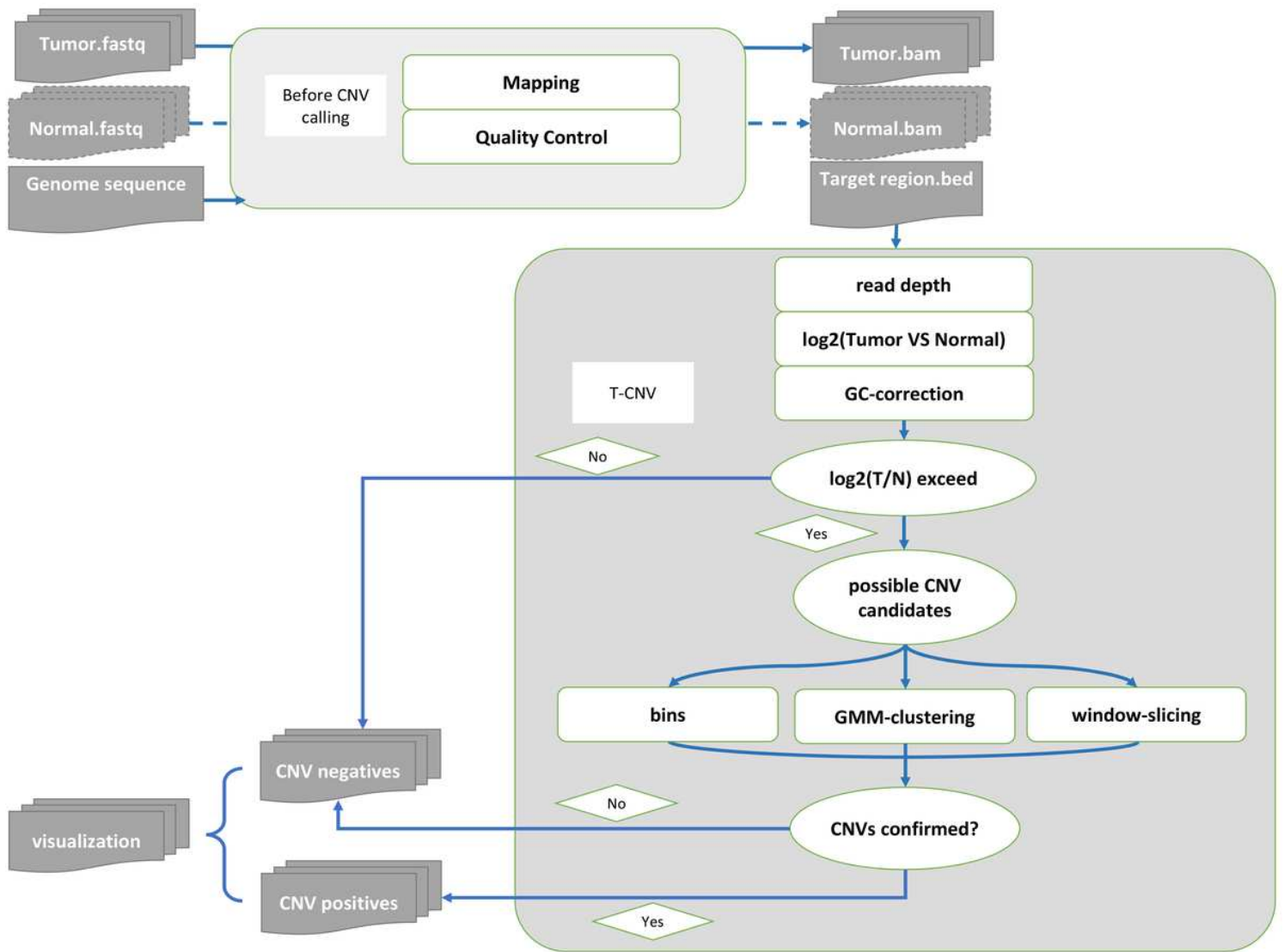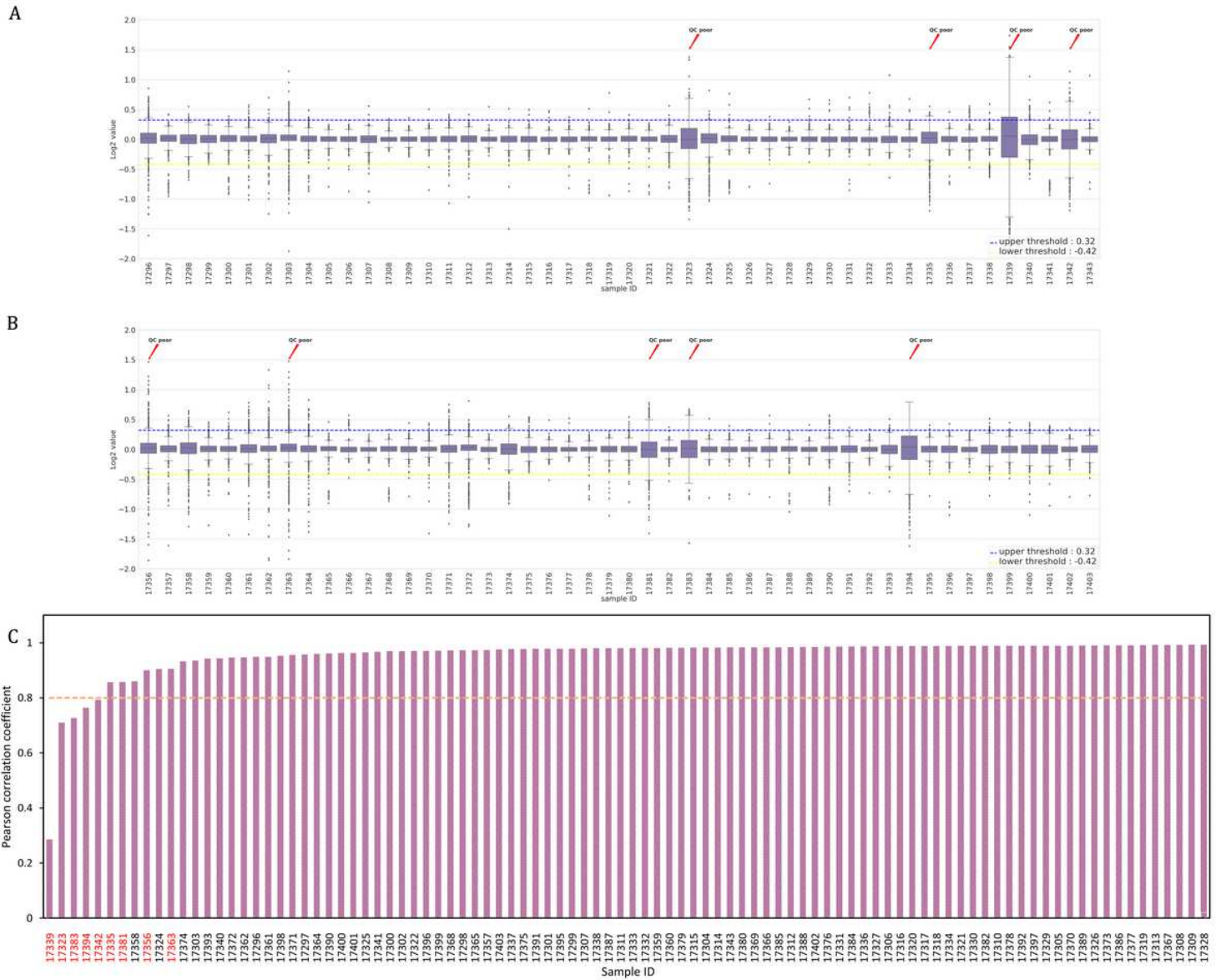
505     2.5MB)

# Figures



## Figure 1

Pipeline of T-CNV To start CNVs assessment, read depth at each locus located with the targeted region was determined. Tumor sample and corresponding normal/control sample were considered as a pair in T-CNV, but the latter was not strictly required.Log2 value was calculated and GC- correction was computed. Any exon with log2 value higher than 0.32 or lower than -0.42 was considered as possible CNV candidates. The CNV candidates were further confirmed by the 3 methods including bins-method, GMM-clustering method and window-slicing method.

## Figure 2

Box-plot for Log2 value and Pearson correlation coefficient in ICR96 dataset (A) Four samples with in pool1 were marked as QC poor with ☐Std☐_Exon >0.2. (B) Five samples in pool2 were marked as QC poor with ☐Std☐_Exon >0.2. (C) Pearson correlation coefficient of each sample and its corresponding pool control in ICR96 dataset. QC poor samples are indicated with red font.
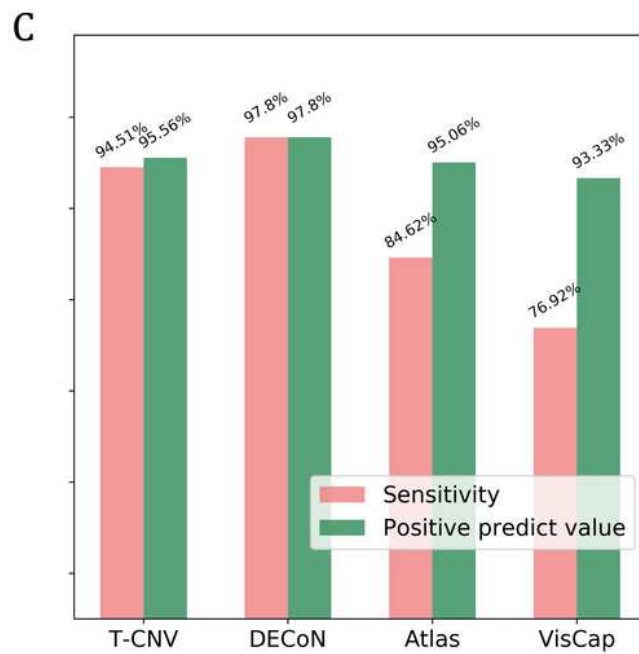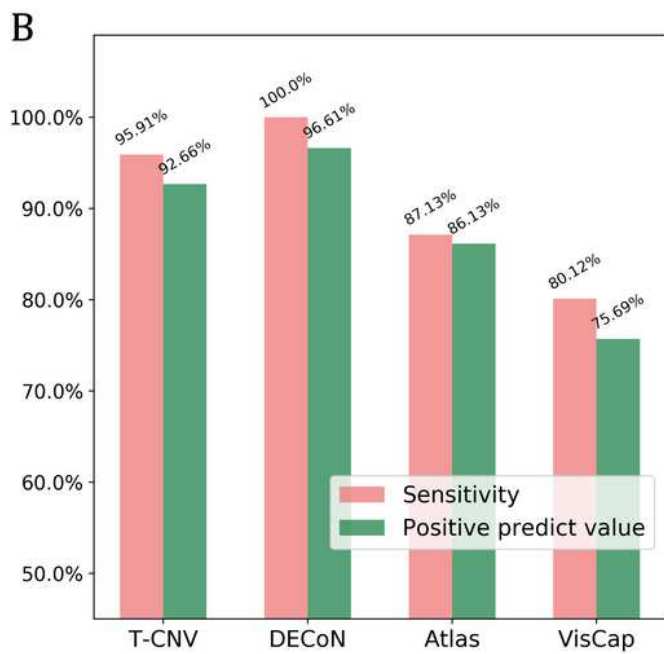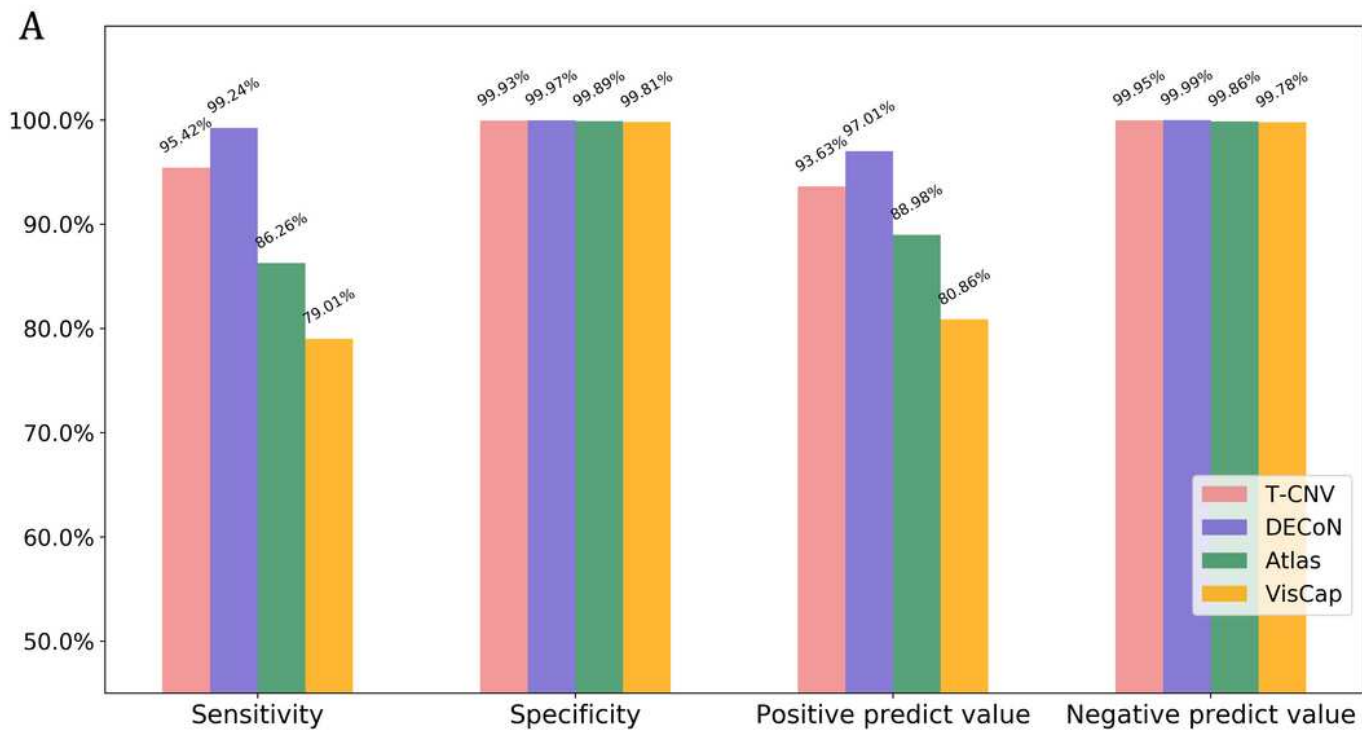
**Figure 3**

Performance comparison with DECoN, Atlas-CNV and VisCap on ICR96 dataset (A) T-CNV, DECoN, Atlas-CNV and VisCap were used to detect CNVs on ICR96 dataset and their performance were presented using four parameters: sensitivity, specificity, positive predict value (PPV) and negative predict value (NPV). (B) Four tools performance on deletion CNVs. (C) Four tools performance on duplication CNVs.
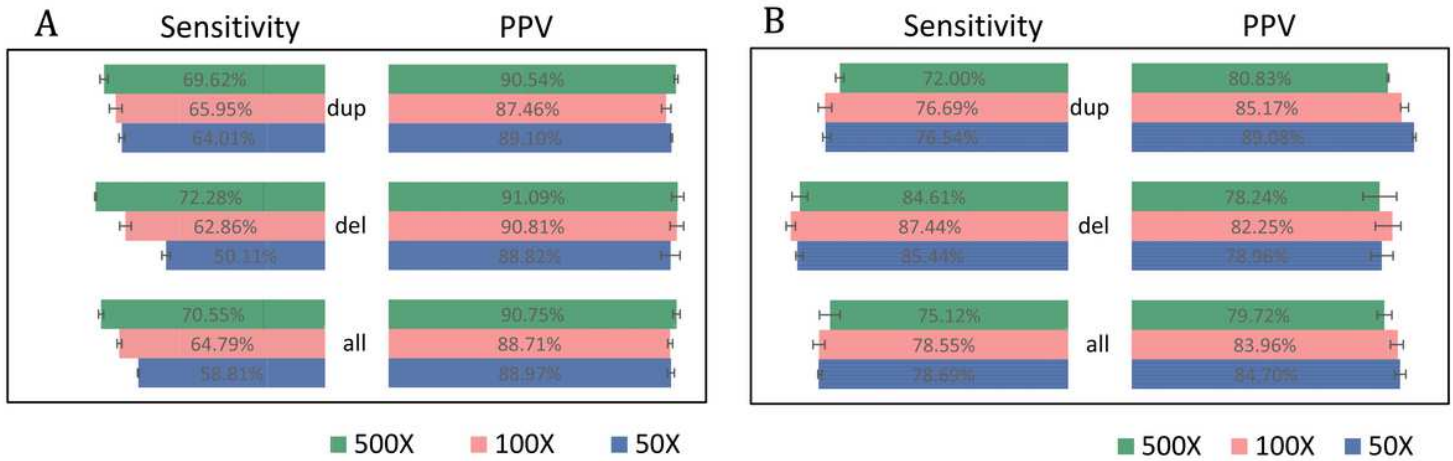
## Figure 4

Performance assessment of T-CNV and DECoN on simulative dataset (A) Performance of T-CNV on simulative dataset under 50X, 100X and 500X was presented as sensitivity and PPV. "all" represented all simulative CNVs. "del" represented the simulative deletions and "dup" represented the simulative duplications. (B) Performance of DECoN on simulative dataset.
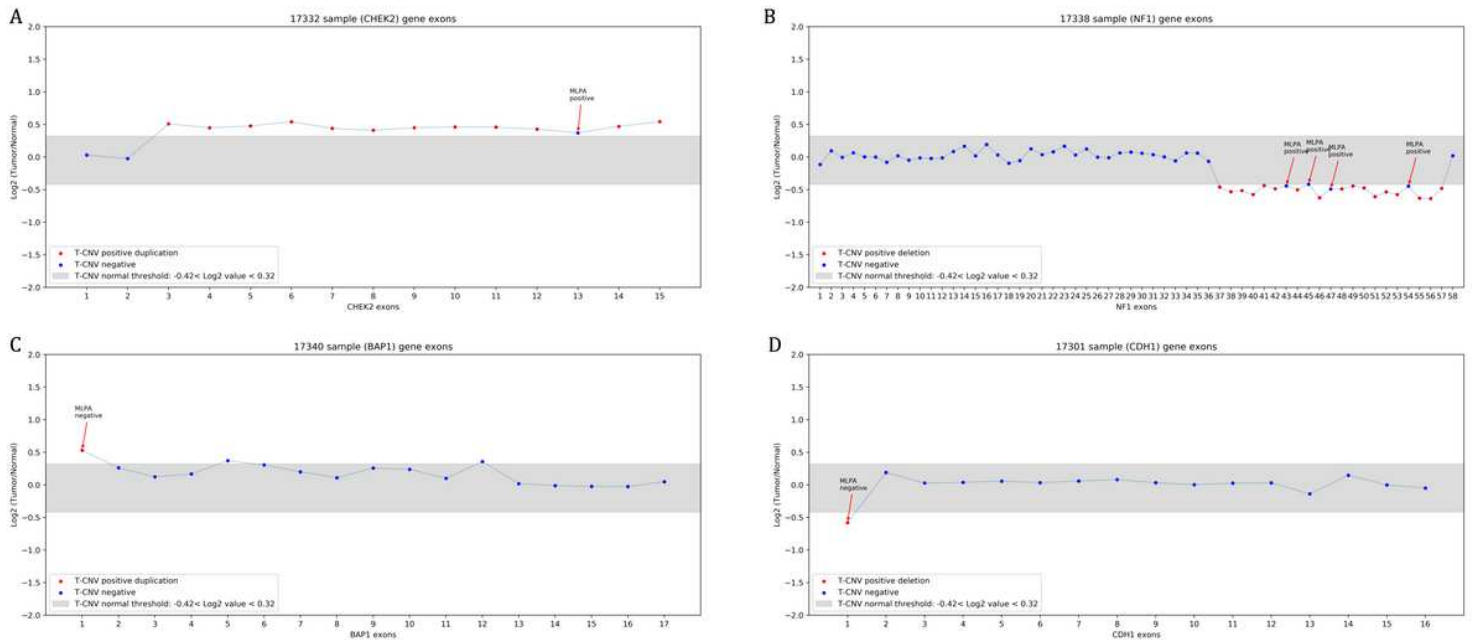


## Figure 5

Example of false positive and false negative prediction in ICR96 dataset (A) A false negative CHEK2 exon 13 in sample 17332 was identified by T-CNV. (B) Four false negative exons located in NF1 gene in 17338 sample were identified by T-CNV. (C) A false positive BAP1 exon 1 in sample 17340 was identified by T-CNV. (D) A false positive CDH1 exon 1 in sample 17301 was identified by T-CNV.

# Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- SupplementaryfigureS3bincomparePSupdate.png
- SupplementaryfigureS1LOESSupdate.png
- SupplementaryfigureS4GMMPSupdate.png
- SupplementaryfigureS6wintpfp.png
- Addtionalfile1.docx
- SupplementaryfigureS2noisetestPS.png
- SupplementaryfigureS817375.png
- SupplementaryfigureS7windowsliding.png
- SupplementaryfigureS5GMMsamplingPS.png