1 **Supplementary Notes**

2 **Next-generation sequencing and sequencing data preprocessing**

3 The Agilent (Santa Clara, CA, USA) ClearSeq Inherited Disease panel kit for enrichment was

4 adapted on every enrolled proband. Genomic DNA was fragmented by sonication, ligated to

5 multiplexing paired-end adapters, and amplified by PCR with indexed primers for

6 sequencing, which were then hybridized to biotin-labelled probes, followed by paired-end

7 sequencing (125 bp/150 bp) on the Illumina (San Diego, CA, USA) HiSeq 2000/2500

8 platforms. The average coverage for the yielded data was 124.05X, and the average fraction

9 of target covered with at least 10X and 20X were 99.11% and 97.17%, respectively. Reads

10 were cleaned to pass quality controls and were aligned to the hg19 reference genome by

11 BWA-MEM (V.05.9-r16), sorted by SAMtools (v.1.8), and removed duplicates by Picard

12 (v.2.20.1).

13 **Pipeline for NGS data-based CNV detection**

14 **Data quality assessment**

15 Data quality assessment was performed before further CNV detection. For each BAM file, the

16 number of mapped reads on the annotated exons was calculated by BEDTools with default

17 parameters and standardized into read depth. Mean and coefficient of variance (C. V) values

18 for read depth across all samples were calculated and assessed on each exon. Autosomes and

19 X chromosomes were separately calculated, since the read depth of chromosome X was

1

1  adjusted according to gender information. Gene and exon annotations were based on

2  GENCODE (v19) and downloaded from the UCSC Table Browser

3  (wgEncodeGencodeBasicV19).

4  **CNV calling**

5  Both CANOES[1] and HMZDelFinder[2] were applied for the original detection of CNVs

6  from exome sequencing data. The CANOES approach is optimized to call deletions involving

7  three or more exons and HMZDelFinder is superior in calling single exon-level CNVs.

8  CANOES uses a negative binomial distribution to model read counts and detect CNVs from

9  exome sequencing data. In the published version of CANOES, researchers stated that all CNV

10  callings were restricted to the autosomes due to the complications resulting from the samples

11  being of different sexes[1]. As quite a few DD-associated CNVs are located on chromosome

12  X, we extended the CNV calling to chromosome X by adjusting the target region coverage

13  according to the samples' sex. The detailed extension process included adjusting the X

14  chromosome coverage according to gender and evaluating the baseline for CNV detection on

15  the X chromosome using autosome as well as X chromosome coverage data. Both read

16  coverage and the GC content of the exons were used to call CNVs by CANOES.

17  HMZDelFinder is an algorithm for the detection of rare and intragenic homozygous and

18  hemizygous (HMZ) deletions. BAM files and VCF files of individuals were extracted as input

19  for variant detection. The application was conducted under default-set thresholds.

20

2

1   **CNV annotation and filtering**

2   CNVs called from CANOES and HMZDelFinder were combined and annotated at both the

3   gene and region levels. The annotation and filtration pipeline was called PICNIC (pipeline for

4   clinical NGS-involved CNV detection). The PICNIC is a published automatic process that

5   combines CNV detection results for annotation and filtration[3]. In the published work, 58

6   patients underwent CES-based CNV detection. Compared with the CMA result, PICNIC

7   showed 100% sensitivity (95%CI: 92.13%-100.00%) to pathogenic/likely pathogenic CNVs,

8   demonstrating a reliable analysis result of the CES-based CNV detection[3]. The flowchart of

9   PICNIC is shown in **Supplementary Figure S1.**

10  At the gene level, PICNIC used RefSeq for mRNA annotation and OMIM, HGMD, and

11  UniProt for functional annotation. Then, PICNIC filtered out gene deletions/duplications,

12  which occurred in >10% of the internal samples, as their high frequency made the variation

13  unlikely to cause rare disorders. Next, CNV-influenced genes were compared with clinical

14  phenotypes by using the annotated HPO terms and the standardized clinical records. In this

15  study, to reduce the amount of manual work, we generated a preliminary automatic HPO

16  annotation system to extract HPO terms from the Chinese electronic clinical record. Clinical

17  record sentences in Chinese were translated into English, annotated by MetaMap into UMLS

18  standard phrases and converted to HPO terms. We also established a semantic dictionary

19  supervised by clinicians to correct the annotation process. The final HPO assignment needed

20  to be manually reviewed by clinicians before further analysis. A CNV would be categorized

21  as a candidate variant if any of its affected genes could match the patient's clinical HPO (or

3

1    parental HPO) terms; otherwise, the CNV would be considered unknown or likely benign.

2    Specifically, CNV-affected genes were compared with patient's HPO terms according to the

3    annotations of known disease-causing genes provided by the HPO database

4    (http://www.human-phenotype-ontology.org). Once any CNV-containing gene's HPO

5    annotation matched with the patient's HPO terms, the CNV would be classified as a candidate

6    variant. Considering the hierarchical and specialization differences in HPO terminology, two

7    HPO terms that shared the same parent node in the upper three layers were considered to be

8    consistent.

9    For region-level annotation, we annotated CNVs with DGV, DECIPHER and a prepublished

10   pathogenic database[4]. The detected CNV was classified as a candidate variant if it had any

11   overlap with an established pathogenic region or the known causative genes of an established

12   pathogenic CNV. For other CNVs, variation regions larger than 1 Mb were also be marked as

13   candidate variants, considering that abundant CNV-affected genes would be more likely to be

14   pathogenic. The data of the prepublished pathogenic database are given in the

15   **Supplementary Table S8**.

16

17   **Criteria for variant classification**

18   The variant classification, which is based on the ACMG guidelines but with some adjustments,

19   was first described in our previous work[5] and implemented in this study. Hence, the variant

20   classification would not always be consistent with those in the ACMG guidelines.

21   Specifically, the criteria of pathogenic (P) variants were as follows: 1) the variant would likely

4

1    explain the indication for testing and may be responsible for the patient's clinical presentation;

2    and 2) for SNV and small indel, the variant has the same nucleotide and amino acid change as

3    a previously established pathogenic variant from published studies or the internal database; for

4    CNV, the variant has the same copy-number status as a previously established pathogenic CNV

5    and fully covers the region.

6    Criteria of likely pathogenic (LP) variants: 1) the variant would likely explain the indication

7    for testing and may be responsible for the patient's clinical presentation; and 2) for SNV and

8    small indel, the variant has the same amino acid change as a previously established pathogenic

9    variant regardless of nucleotide change; or Null variant (nonsense, frameshift, canonical +/−1

10   or 2 splice sites, initiation codon) in a gene where loss of function (LOF) is a known mechanism

11   of disease; for CNV, the variant overlaps with a previously established pathogenic CNV

12   (overlapped region is >70% of the reported pathogenic CNV) and the overlapped region

13   included a gene where LOF is a known mechanism of disease or established/predicted to be

14   haploinsufficiency (HI); or both breakpoints of the CNV are within the same established

15   haploinsufficiency gene; and 3) the variant is *de novo* (both maternity and paternity confirmed)

16   in the proband with a negative family history; or is inherited from the affected parents.

17   Specifically, if the parents are not available for the confirmation of *de novo* or compound

18   heterozygous status of pathogenic variants identified in the proband, the variant would be

19   downgraded and classified as 'pathogenic to likely-pathogenic' (P->LP).

20

21   **Process of CNV and SNV integrated (SCI) strategy for diagnosing**

22       The CNV and SNV integrated (SCI) strategy is a diagnosing condition that combines both

5

1    SNV and CNV. Based on the above variant classification criteria, diagnostic results for SNVs

2    and small indels included: 1) one heterozygous P/LP variant in an autosomal dominant, or X-

3    linked dominant gene; 2) one homozygous or two heterozygous P/LP variants (compound

4    heterozygous) on an autosomal recessive gene; and 3) one P/LP variant in an X-linked recessive

5    gene in males. Diagnostic results for CNV included: the P/LP CNV that explains the indication

6    for testing and the inheritance pattern.

7    Moreover, for the SNVs and CNVs that failed to meet the aforementioned conditions but

8    resulted in a rarer event: apparently homozygous (AH) variant caused by overlapping P/LP

9    SNV on one allele and P/LP CNV on the other. These "SNV+CNV AH" cases were also

10    included as one of the diagnosing conditions.

11    In addition, the initial diagnosis was issued without experimental confirmation if the variants

12    were of high confidence[6] (SNV: coverage ≥20×, minor allele fraction ≥35%, and Phred score

13    of variant calling ≥30; CNV: detected variation size >1 Mb).

14

15    **Variant experimental validation**

16    For SNVs, all of the identified P/LP variants were confirmed by Sanger sequencing.

17    For CNVs, referring to a previous study that performed an evaluation of CNV detection from

18    panel-based next-generation sequencing data, researchers found that CNVs covering adequate

19    exons on autosomes can be accurately detected using targeted panel sequencing data same as

20    using CMA, while CNVs detected from sex chromosomes need further evaluation and

21    validation. Additionally, the accuracy of CNV size improves as the size of variants increased[7].

22    Thus, in this study, we considered the CNVs with sizes larger than 1 Mb as variants of high

1    confidence. The quantitative polymerase chain reaction (qPCR), multiplex ligation-dependent

2    probe amplification (MLPA) or CMA validations were performed on CNVs detected from sex

3    chromosomes and CNVs whose detected sizes were less than 1 Mb. Detailed information about

4    the kits used and designed primers for CNV validations are given in **Supplementary Table S1**.

5

6

7

8    **Reference**

9    1.    Backenroth D, Homsy J, Murillo LR, Glessner J, Lin E, Brueckner M, Lifton R, Goldmuntz E,

10   Chung WK, Shen Y. CANOES: detecting rare copy number variants from whole exome sequencing

11   data. *Nucleic Acids Res* 2014;42:e97.

12   2.    Gambin T, Akdemir ZC, Yuan B, Gu S, Chiang T, Carvalho CMB, Shaw C, Jhangiani S, Boone

13   PM, Eldomery MK, Karaca E, Bayram Y, Stray-Pedersen A, Muzny D, Charng WL, Bahrambeigi

14   V, Belmont JW, Boerwinkle E, Beaudet AL, Gibbs RA, Lupski JR. Homozygous and hemizygous

15   CNV detection from exome sequencing data in a Mendelian disease cohort. *Nucleic Acids Research*

16   2017;45:1633-48.

17   3.    Qin Q, Liu B, Yang L, Wu BB, Wang HJ, Dong XR, Lu YL, Zhou WH. Application of copy

18   number variation screening analysis process based on high-throughput sequencing technology.

19   *Chinese Journal of Evidence -Based Pediatric* 2018;13:275-79.

20   4.    Wang HJ, Bi WM, Wu BB, Zhou WH, Cheung SW. Application and explanation of

21   chromosome microarray (CMA) analysis in the diagnosis of clinical genetic disease. *Chinese*

22   *Journal of Evidence -Based Pediatric* 2014;9:227-35.

1    5.    Yang L, Kong YT, Dong XR, Hu LY, Lin YF, Chen X, Ni Q, Lu YL, Wu BB, Wang HJ, Lu

2    QR, Zhou WH. Clinical and genetic spectrum of a large cohort of children with epilepsy in China.

3    *Genet Med* 2019;21:564-71.

4    6.    Gross AM, Ajay SS, Rajan V, Brown C, Bluske K, Burns NJ, Chawla A, Coffey AJ, Malhotra

5    A, Scocchia A, Thorpe E, Dzidic N, Hovanes K, Sahoo T, Dolzhenko E, Lajoie B, Khouzam A,

6    Chowdhury S, Belmont J, Roller E, Ivakhno S, Tanner S, McEachern J, Hambuch T, Eberle M,

7    Hagelstrom RT, Bentley DR, Perry DL, Taft RJ. Copy-number variants in clinical genome

8    sequencing: deployment and interpretation for rare and undiagnosed disease. *Genet Med*

9    2018;21:1121–30.

10   7.    Yao RE, Yu TT, Qing YR, Wang J, Shen YP. Evaluation of copy number variant detection from

11   panel-based next-generation sequencing data. *Mol Genet Genom Med* 2019;7:e00513.

12