

**POLISHING COPY NUMBER VARIANT  
CALLS ON EXOME SEQUENCING DATA  
VIA DEEP LEARNING**

A THESIS SUBMITTED TO  
THE GRADUATE SCHOOL OF ENGINEERING AND SCIENCE  
OF BILKENT UNIVERSITY  
IN PARTIAL FULFILLMENT OF THE REQUIREMENTS FOR  
THE DEGREE OF  
MASTER OF SCIENCE  
IN  
COMPUTER ENGINEERING

By  
Furkan Ozden  
July 2021

Polishing Copy Number Variant Calls on Exome Sequencing Data via  
Deep Learning  
By Furkan Ozden  
July 2021

We certify that we have read this thesis and that in our opinion it is fully adequate,  
in scope and in quality, as a thesis for the degree of Master of Science.

————— A. Ercüment Çiçek (Advisor) —————  
\_\_\_\_\_

\_\_\_\_\_ Can Alkan

\_\_\_\_\_ Tolga Can \_\_\_\_\_

Approved for the Graduate School of Engineering and Science:

\_\_\_\_\_ Ezhan Karaşan  
Director of the Graduate School

# ABSTRACT

## POLISHING COPY NUMBER VARIANT CALLS ON EXOME SEQUENCING DATA VIA DEEP LEARNING

Furkan Ozden

M.S. in Computer Engineering

Advisor: A. Ercument Cicek

July 2021

Accurate and efficient detection of copy number variants (CNVs) is of critical importance due to their significant association with complex genetic diseases. Although algorithms that use whole genome sequencing (WGS) data provide stable results with mostly-valid statistical assumptions, copy number detection on whole exome sequencing (WES) data shows comparatively lower accuracy. This is unfortunate as WES data is cost efficient, compact and is relatively ubiquitous. The bottleneck is primarily due to non-contiguous nature of the targeted capture: biases in targeted genomic hybridization, GC content, targeting probes, and sample batching during sequencing. Here, we present a novel deep learning model, DECoNT, which uses the matched WES and WGS data and learns to correct the copy number variations reported by any off-the-shelf WES-based germline CNV caller. We train DECoNT on the 1000 Genomes Project data, and we show that we can efficiently triple the duplication call precision and double the deletion call precision of the state-of-the-art algorithms. We also show that our model consistently improves the performance independent from (i) sequencing technology, (ii) exome capture kit and (iii) CNV caller. Using DECoNT as a universal exome CNV call polisher has the potential to improve the reliability of germline CNV detection on WES data sets.

*Keywords:* Copy Number Variation, Whole Exome Sequencing, Deep Learning.

## ÖZET

# DERİN ÖĞRENME İLE EKZOM DİZİLEME VERİLERİNDE GEN KOPYA SAYISI ANALİZLERİNİN GELİŞTİRİLMESİ

Furkan Ozden

Bilgisayar Mühendisliği, Yüksek Lisans

Tez Danışmanı: A. Ercument Cicek

July 2021

Gen kopya sayısı varyantlarının (CNV'ler) doğru ve verimli tespiti, karmaşık genetik hastalıklarla önemli ilişkileri nedeniyle kritik öneme sahiptir. Tüm genom dizilimi (WGS) verilerini kullanan algoritmalar, çoğunlukla geçerli istatistiksel varsayımlarla kararlı sonuçlar verse de, tüm ekzom dizileme (WES) verilerinde kopya sayısı tespiti, nispeten daha düşük doğruluk gösterir. WES verileri uygun maliyetli, kompakt ve nispeten her yerde mevcut olduğundan bu talihsiz bir durumdur. Darboğaz öncelikle hedeflenen yakalamamanın bitişik olmayan doğasından kaynaklanmaktadır: hedeflenen genomik hibridizasyondaki önyargılar, GC içeriği, hedefleme problemleri ve sıralama sırasında numune gruplaması. Burada, eşleşen WES ve WGS verilerini kullanan ve kullanıma hazır herhangi bir WES tabanlı germline CNV arayan tarafından bildirilen kopya numarası varyasyonlarını düzeltmeyi öğrenen yeni bir derin öğrenme modeli DECoNT sunuyoruz. DECoNT'u 1000 Genom Projesi verileri üzerinde eğitiyoruz ve son teknoloji algoritmaların çoğaltma çağrısı hassasiyetini verimli bir şekilde üç katına ve silme çağrısı hassasiyetini iki katına çıkarabileceğimizi gösteriyoruz. Ayrıca modelimizin (i) sıralama teknolojilerinden, (ii) ekzom yakalama kitinden ve (iii) CNV arayandan bağımsız olarak performansı sürekli olarak geliştirdiğini gösteriyoruz. DECoNT'u evrensel bir exome CNV çağrı parlatici olarak kullanmak, WES veri setlerinde germline CNV tespitinin güvenilirliğini artırma potansiyeline sahiptir.

*Anahtar sözcükler:* Gen Kopya Sayısı, Tüm Ekzom Sekanslama, Derin Öğrenme.

## Acknowledgement

Firstly, I would like to thank to my supervisor Prof. Ercument Cicek. He is the best supervisor I have ever seen.

Secondly, I would like to express my gratitudes to my parents and my brothers for being supportive throughout my education.

Also, many thanks to jury members, Prof. Can Alkan and Prof Tolga Can, for investing their time to undertake the jury member roles in my defense and evaluate this thesis.

# Contents

- 1 Introduction** 1
  
- 2 Methods** 4
  - 2.1 Data set** . . . . . 4
  
  - 2.2 DECoNT Model** . . . . . 5
    - 2.2.1 Problem Formulation** . . . . . 5
  
    - 2.2.2 DECoNT Architecture** . . . . . 6
  
  - 2.3 Polishing the State-of-the-art WES-based Germline CNV Callers** . 7
    - 2.3.1 Settings for the WES-based CNV Callers** . . . . . 8
  
    - 2.3.2 Training Settings for DECoNT** . . . . . 9
  
    - 2.3.3 Performance Metrics** . . . . . 9
  
    - 2.3.4 Time Performance** . . . . . 10
  
  - 2.4 Polishing Samples from Other Sequencing Platforms** . . . . . 11
  
  - 2.5 Polishing Other WES-based CNV Caller Algorithms** . . . . . 11

**3 Discussion** **13**

**4 Results** **17**

4.1 Bi-LSTM based Neural Network Learns to Correct False Positive  
Germline WES CNV Calls . . . . . 17

4.2 Polishing Performance on a Validated CNV call set . . . . . 22

4.3 Polishing Performance Generalizes to Unseen Sequencing Platforms 23

4.4 Polishing Performance on Calls From Unseen CNV Callers . . . . . 24

**A Supplementary Document** **36**

A.1 Supplementary Figures . . . . . 36

A.2 Supplementary Tables . . . . . 40

A.3 Supplementary Notes . . . . . 41

# List of Figures

4.1	<b>Learning workflow of DECoNT.</b> First, BAM file that corresponds to a WES data set from <b>1000 Genomes data set</b> is used to calculate exome-wide read depth which is input to a third party WES-based CNV caller. The caller generates the calls for various regions which could be (i) a binary prediction like duplication, deletion (e.g., XHMM [1]) as shown in the figure, or (ii) an integer value that indicates the exact copy number (i.e., Control-FREEC [2]). The read depth of the regions for which a call has been made is input to a Bi-LSTM model. Encoded features are passed from a series of fully connected layers along with the original prediction of the caller algorithm. Using the ground truth calls from the WGS data of the same sample the method learns to predict (correct) the calls using cross entropy loss for the binary outputs (as shown in the figure) and using mean squared loss for integral calls. . . . .	18
-----	--	----

<b>4.2 The performance comparison of the WES-based CNV caller’s before and after polishing with DECoNT.</b>	
a) For the tools which predict existence of a CNV event (XHMM, CoNIFER and CODEX2) are evaluated with respect to duplication call precision, deletion call precision and overall precision. DECoNT improves the performance for all tools in all settings and results in drastic improvements. Different shades of gray represent different tools and the attached black bars represent the DECoNT-polished version of those tools.	
b) In this panel, we compare Control-FREEC and the DECoNT-polished with respect to Absolute Error (AE) difference on each sample (i.e., events). Bars to the right indicate the magnitude of the improvement due to polishing of DECoNT. For more than half of the samples, DECoNT results show improvement.	
c) The distribution of the unpolished Control-FREEC predictions in the test samples (pink) is quite different than the ground truth copy number variation distribution. On the other hand, DECoNT polished versions of the same events (dark blue) highly resemble the distribution of the ground truth calls. Black lines across the boxes are median lines for the distributions. Black vertical lines are whiskers and 1.5× Inter Quartile Range is defined with the horizontal lines at the top and bottom of the whiskers. 1000 Genomes data set WGS samples are used as ground truth calls in all analyses.	
Panels d) and e) show the results for CNVkit, similar to b) and c).	27

**4.3 Performance of DECoNT when polishing calls from unseen CNV callers.** DECoNT learns a different set of weights and a different model for each WES based CNV caller. In order to demonstrate the cross-model performance, we used DECoNT to correct CNV calls made by tools other than the ones used for training. We try every pair combination. Tools being pointed by an arrow are call generating tools (i.e., being corrected). Tools at source of the arrow are the tools that are used to train the DECoNT model. Green arrows indicate improvement and red arrows indicate deterioration in the corresponding performance metric. . . . . 29

**A.1** The confusion matrices of the WES-based CNV callers before and after polishing with DECoNT on 1000 Genomes Data test samples. Confusion matrices given with blue borders represent unpolished predictions of corresponding WES-based CNV tools. Since DECoNT only operates on the calls made by a CNV caller, the first column for each unpolished confusion matrix is set as NA (i.e. Not Applicable). The red-bordered confusion matrices are the polished versions of a CNV caller with a DECoNT model trained on the calls made by the same caller (to produce Fig 2). Other confusion matrices are polised version of the CNV caller corrected by a DECoNT model trained on the calls made by a different caller (to produce Fig. 3). Notice the decrease in the number of false positives for both deletion and duplication calls in all platforms. . . . . 37

**A.2** Confusion matrices of the WES-based CNV callers before and after polishing with DECoNT on highly validated CNV callset published in Chaisson et. al. [3]. Similar to Fig. A.1 tool provides great false discovery correction with slight true positive deterioration for both deletion and duplication calls, yielding much better performance metric results. . . . . 38

A.3 Confusion matrices of the WES-based CNV callers before and after polishing with DECoNT on NA12878 data obtained from different sequencing platforms: (i) NovaSeq6000; (ii) HiSeq4000; (iii) BGI500; (iv) MGI2000. Since DECoNT only operates on the calls made by a CNV caller, the first column for each unpolished confusion matrix is set as NA (i.e. Not Applicable). Since CoNIFER does not report any calls on NovaSeq6000 platform, DECoNT has no input to polish and thus the comparison is not applicable. Similar to Figures A.1 and A.2, we observe that DECoNT substantially decreases the number of false discoveries with slight true positive deterioration for both deletion and duplication calls. . . . . 39

A.4 This figure shows the length distribution of true raw XHMM calls and true DECoNT-corrected XHMM calls obtained on the 1000 Genomes WES data set test samples. The ground truth is the CNV calls made by CNVnator on the corresponding WGS samples. We see that for smaller size CNVs XHMM requires more correction by DECoNT. However, again, vast majority of the CNVs cannot be distinguished by the CNV length to decide whether it needs a DECoNT correction. . . . . 39

# List of Tables

4.1	The performances of the WES-based CNV caller algorithms before and after polishing are shown (DEL, DUP and overall precision).	28
4.2	The performance of discrete germline WES-based CNV callers on NA12878 data before and after being polished by DECoNT.	28
4.3	The performance of Control-FREEC on NA12878 data before and after being polished by DECoNT.	29
A.1	This table summarizes the polishing performance of DECoNT on the X chromosome, PAR1 and PAR2 regions of the males in the test split obtained from 1000 Genomes WES samples. The base caller in this analysis isXHMM. The results are obtained on the test samples from the 1000 Genomes dataset and the ground truth is obtained from the CNVnator calls on WGS of the same samples.	40
A.2	In addition to the WES CNV Callers presented in the manuscript, we have also trained a DECoNT model for CNVKit on, again, 1000 Genomes WES samples, using the CNVnator calls obtained on WGS data as ground truth.	40

A.3	The 8 WES CNV calls that DECoNT and CNLearn does not agree	
	are presented. The ground truth CNV calls are obtained through	
	CNVNator WGS CNV calls. Note that, CNLearn samples are pol-	
	ished with a DECoNT model trained with XHMM data. Training	
	a DECoNT model with consensus calls made by CNLearn would	
	increase performance.	41

# Chapter 1

## Introduction

Gene copy number polymorphism in a population due to deletions and duplications of genomic segments substantially derive genetic diversity [4, 5], affecting roughly 7% of the genome [6]. This class of structural variations (SVs), called copy number variations (CNVs), have also been associated with several genetic diseases and disorders such as neurodevelopmental/neurodegenerative disorders [7, 8, 9, 10, 11] and various cancers such as breast, ovary, and pancreas cancers [12, 13, 14, 15]. Karyotyping and microarray analyses have been the standard clinical testing for disease-causing CNVs for many years [16], but High Throughput Sequencing (HTS) has all but replaced these techniques with the ability to theoretically capture all forms of genomic variation. Numerous CNV detection algorithms have enjoyed success by analyzing whole genome sequencing (WGS) data using different sequence signatures such as read depth, discordant paired-end read mappings, and split reads [17]. WGS is a convenient resource for CNV callers as it provides near-Poisson depth of coverage [18]. On the other hand, accurate CNV detection on whole exome sequencing (WES) data has mostly been lacking. The algorithms which call CNVs on the WES data have notoriously high false discovery rates (FDR) reaching up to ~60% which renders them impractical for clinical use [19, 20]. This is mainly due to several problems associated with the WES technology such as non-uniform read-depth distribution among exons caused by biases in (i) sample batches, (ii) GC content, and (iii) targeting

probes [21, 22, 23]. This is unfortunate as WES data size is ten times smaller (i.e.,  $\sim 10\text{GB}$  vs  $\sim 100\text{GB}$ ) and it costs three times less compared to WGS which makes it highly abundant and a common choice for analyzing complex genetic disorders [24, 25, 26, 27]. For instance, the Genome Aggregation Database (gnomAD) contains around 125K WES samples as opposed to 70K WGS samples [28]. Thus, currently, such a rich resource of large scale WES data cannot be fully utilized to investigate the contribution of copy number variation to disease etiology.

Here, we present the first of its kind, exome CNV call *polisher* named *DECoNT* (Deep Exome Copy Number Tuner) to improve the performance of any off-the-shelf WES-based germline CNV detection algorithm. DECoNT is a deep learner that utilizes matched WES and WGS samples present in the 1000 Genomes Project [29] data set to learn the association between (i) calls made by any CNV caller that use WES data and (ii) ground truth calls generated from the WGS data for the same sample. Based on a bidirectional long short-term memory (Bi-LSTM) based architecture, it uses only WES read depth along with the calls from a third party caller and learns to correct noisy predictions (Figure ??). We show that DECoNT can improve the duplication and deletion call precision of the state-of-the-art algorithms by up to 3-fold and 2-fold, respectively. The performance gain is consistent among CNV callers that output integer copy number predictions and categorical predictions (i.e., deletion, duplication, or no call). As the training phase is offline, polishing procedure is memory and time efficient and takes only a few seconds on average per sample. Furthermore, we show that the models learned are universal in the sense that they are (i) sequencing platform, (ii) target capture kit, and (iii) CNV caller independent. For instance, using models learned on 1000 Genomes Project data set that uses Illumina as the sequencing platform and various capture kits such as Agilent and NimbleGen, DECoNT can correct calls made by any of the state-of-the-art CNV caller, for samples obtained from other capture kits like Agilent SureSelect and Illumina Nextera Exome Enrichment Kit or different sequencing platforms like Illumina HiSeq 4000, Illumina NovaSeq 6000, and MGI; that are “unseen” during training. Thus, DECoNT is highly flexible and scalable and makes exome based CNV detection practical by boosting the performance of virtually any WES-based CNV caller algorithm. The

tool and the models are available at <https://github.com/ciceklab/DECoNT>.  
All necessary scripts and data to replicate the results presented in the figures and tables are deposited to <https://zenodo.org/record/3865380>.

# Chapter 2

## Methods

### 2.1 Data set

For training and testing of DECoNT, we used 1,000 samples (i.e., HG00096 to HG02356, when sample IDs are alphabetically ordered) from the 1000 Genomes Project [30]. For these samples we obtain both WES and WGS data. WES samples were captured using the NimbleGen SeqCap EZ Exome v3 as capture kit, and sequenced to an average of 50 $\times$  depth with Illumina Genome Analyzer II and Illumina HiSeq 2000 platforms. The average read length is 76 bps. Reads were aligned to the GRCh38 using the BWA-MEM aligner [31]. WGS samples were also sequenced using the same platforms with an average read length of 100bps. Average depth coverage for this set is 30 $\times$ . For XHMM, CoNIFER and CODEX2, the ground truth CNV calls are obtained using CNVnator [32] tool. For Control-FREEC and CNVkit, the ground truth exact copy number variation events are obtained using mrCaNaVaR [33].

For tools that output a categorical prediction of a CNV, we also use a highly validated CNV call set published in Chaisson *et al.*, [3] as another validation source. The WGS CNV calls in this call set are thoroughly validated. That is,

they were obtained via a consensus of 15 different WGS CNV callers with comparisons against high quality PB-SVs that have single base breakpoint resolution. We obtain WGS CNV calls for these 9 samples from 1000 Genomes data set. (i.e., HG00512, HG00513, HG00514, HG00731, HG00732, HG00733, NA19238, NA19239, NA19240). We also obtained aligned WES reads of these samples, with the exception of HG00514 for which no WES data was available. This data set is only used for testing.

## 2.2 DECoNT Model

### 2.2.1 Problem Formulation

Let  $X$  denote the set of CNV events detected on the WES data set by a WES-based CNV caller, and  $X^{(i)}$  denote the  $i^{th}$  event.  $F_i$  denotes the set of features we use for  $X^{(i)}$  which contains the following information: (i) the chromosome that the CNV event occurred ( $X_{chr}^{(i)}$ ); (ii) the start coordinate of the CNV event ( $X_{start}^{(i)}$ ); (iii) the end coordinate of the CNV event; (iv) the type (e.g., deletion) of the called event ( $X_{call}^{(i)}$ ); and (v) the read depth vector between  $X_{start}^{(i)}$  to  $X_{end}^{(i)}$  ( $X_{RDSeq}^{(i)}$ ). Let  $Y_{gt}^{(i)}$  denote the ground truth label obtained from the WGS CNV call for  $X^{(i)}$ . There are two cases: (i) For the tools that predict the existence of an event  $X_{gt}^{(i)} \in \{0, 1, 2\}$ , denoting no call, deletion or duplication, respectively; and (ii) For the tools that predict the copy number  $Y_{gt}^{(i)} \in \mathbb{Z}^{\geq}$ . Then, the problem at hand is formulated as a classification task for (i), and as a regression task for (ii). That is, our goal is to learn a function  $f$  such that  $f(F_1, \dots, F_n) \rightarrow (Y_{pr}^{(1)}, \dots, Y_{pr}^{(n)})$  such that the difference is between  $(Y_{pr}^{(1)}, \dots, Y_{pr}^{(n)})$  and  $(Y_{gt}^{(1)}, \dots, Y_{gt}^{(n)})$  is minimized with respect to a loss function. Here,  $n = |X|$  and  $Y_{pr}^{(i)}$  is the predicted label for  $X^{(i)}$  and it is in the same domain as  $Y_{gt}^{(i)}$  in respective tasks.

## 2.2.2 DECoNT Architecture

DECoNT is an end-to-end multi-input neural network designed for polishing and improving the performance of the WES-based germline CNV callers. It is capable of improving accuracy of WES CNV calling for both exact CNV prediction (i.e., integer) and categorical CNV prediction cases (i.e., deletion, duplication or no call). For each CNV caller, a distinct network is trained.

DECoNT’s pipeline for the categorical CNV prediction case can be divided into three main building blocks: (i) a data preprocessing step that extracts the read depth for genomic regions of interest (i.e., CNV call regions made by the CNV caller). It also normalizes the read depth sequence and acts as a regularizer for the model. Resulting read depth information is  $-1$  padded to the length of the longest call sequence and masked; (ii) a bidirectional LSTM network (BiLSTM) that inputs the read depth sequence and extracts the required encoded features (i.e., embeddings). This subnetwork has 128 neurons in each direction and is followed by a batch normalization layer; and (iii) a 2-layered fully connected (FC) neural network that inputs the embedding calculated by Bi-LSTM, concatenated with the prior CNV prediction of the CNV caller (a one-hot-encoded vector). The first FC layer has 100 neurons and uses ReLU activation. The output layer has 3 neurons and it calculates the posterior probability of each event via softmax activation: no call, deletion, or duplication. We use weighted cross-entropy as the loss function. This architecture has a total of 160,351 parameters with 159,837 of which are trainable. The rest are the batch normalization parameters.

For a training data set of  $N$  samples, the formulation of DECoNT can be summarized as follows:

$$X_{encoding1}^{(1:N)} = \text{BatchNorm}(\text{BiLSTM}^{(128)}(\text{BatchNorm}(\text{Mask}(X_{RDSeq}^{(1:N)})))) \quad (2.1)$$

$$X_{encoding2}^{(1:N)} = \text{CAT}(X_{encoding1}^{(1:N)}, X_{call}^{(1:N)}) \quad (2.2)$$

$$X_{encoding3}^{(1:N)} = \text{ReLU}(\text{FC}^{(100)}(X_{encoding2}^{(1:N)})) \quad (2.3)$$

$$Y_{pr}^{(1:N)} = \text{Softmax}(\text{FC}^{(3)}(X_{encoding3}^{(1:N)})) \quad (2.4)$$

where BiLSTM<sup>(·)</sup> represents bi-directional LSTM layer with · hidden units in each direction. Similarly, FC<sup>(·)</sup> represents a dense layer with · neurons. ReLU and BatchNorm stand for rectified linear unit activation function and batch normalization respectively.

Using  $Y_{pr}^{(1:N)}$  and  $Y_{gt}^{(1:N)}$  training phase minimizes the categorical cross-entropy loss. We use Adam optimizer [34] with a mini batch size of 128 samples. All weights in the network are initialized using Xavier initialization [35].

DECoNT’s pipeline for the exact (i.e., integer) CNV prediction is almost the same as the one described above. The first difference is instead of taking the one-hot encoded version of the CNV call, it inputs an integer value representing the called copy number. The second difference is at the output layer. Instead of 3 neurons with softmax activation, this version has a single neuron with ReLU activation to perform regression instead of classification. It has a total of 160,149 parameters, 159,635 of which are trainable. Again, the rest are the batch normalization parameters. So, the last layer in the formulation above (Eq. 4) is replaced by the following layer and in this case  $Y_{pr}^{(1:N)} \in \mathbb{Z}^{\geq}$ .

$$X_{pr}^{(1:N)} = \text{ReLU}(\text{FC}^{(1)}(X_{encoding3}^{(1:N)})) \quad (2.5)$$

Using  $Y_{pr}^{(1:N)}$  and  $Y_{gt}^{(1:N)}$  training phase now minimizes the mean absolute error loss. Again, we use Adam optimizer with a mini batch size of 128 samples and we use Xavier initialization for weights.

## 2.3 Polishing the State-of-the-art WES-based Germline CNV Callers

We polish the CNV calls made by four state-of-the-art WES-based germline CNV callers (i) XHMM [1], (ii) CoNIFER [36], (iii) CODEX2 [37], (iv) Control-FREEC

[2], and (v) CNVkit [38]. Tools (i - iii) perform categorical CNV prediction and (iv) and (v) perform exact CNV prediction. We use calls made on the WGS samples by CNVnator [32] as the ground truth call set for discrete predictions and the exact copy number predictions made by mrCaNaVaR as the ground truth call set for integral prediction (i.e., Control-FREEC). First, DECoNT obtains the results of these tools (i - iv). Then, it learns to correct these calls on a portion of the 1000 Genomes data set using ground truth calls. Finally, on the left out test portion of the data, we compare the performance of the CNV callers before and after polishing by DECoNT.

### 2.3.1 Settings for the WES-based CNV Callers

We follow the recommended settings for the WES-based Callers. ForXHMM, the parameters are set as follows: (i)  $Pr(\text{start DEL}) = Pr(\text{start DUP}) = 1e-08$ , (ii) mean number of targets in CNV (geometric distribution) = 6, (iii) mean distance between targets within CNV (exponential decay) =  $70kb$ , and (iv) DEL, DIP, DUP read depth distributions modeled as  $\sim \mathcal{N}(-3, 1)$ ,  $\sim \mathcal{N}(0, 1)$  and  $\sim \mathcal{N}(3, 1)$ , respectively. Also, for XHMM nBins parameter is set as 200 which is the default setting. For CODEX2, minimum read coverage of 20 was enforced at the filtering step. Then, the algorithm automatically chooses its parameter,  $K$ , using BIC (i.e. Bayesian Information Criterion) and AIC (i.e. Akaike Information Criterion). CoNIFER performs SVD on the data matrix and then removes  $n$  singular vectors with  $n$  largest singular values. We set  $n$  to 6. Control-FREEC has 45 parameters which were all set to default values as stated in [2]. CNVkit uses a rolling median technique to recenter each on- or off-target bin with other bins of similar GC content, repetitiveness, target size or distance from other targets, independently of genomic location [38]. We used recommended settings for CNVkit as well where  $\log_2$  read depth below threshold =  $-5$ , above threshold =  $1.0$ .

### 2.3.2 Training Settings for DECoNT

We train a DECoNT model for each of the above-mentioned tools. The set  $X$  of CNV calls per tool is shuffled and divided into training, validation and testing sets which contain 70%, 20%, and 10% of the data, respectively. The number of events in the test sets are 6, 832 (3101 no-calls, 2098 duplications, 1633 deletions); 81, 761 (67885 no-calls, 3042 duplications, 10834 deletions); 180 (85 no-calls, 43 duplications, 52 deletions); 20, 482 (minimum copy number is 0, maximum copy number is 585); 39720 (minimum copy number is 34, maximum copy number is 0) forXHMM, CODEX2, CoNIFER, Control-FREEC, and CNVkit, respectively. The second input of the algorithm is the read depth for the CNV-associated regions on the WES data. We calculate it using the Sambamba tool [39]. For all tools other than CODEX2, DECoNT is trained up to 30 epochs with early stopping by checking the loss on the validation fold. Training for CODEX2 has a maximum epoch number 60. For training, DECoNT uses final CNV calls (i.e., concatenated bin-level calls) made by the CNV callers.

### 2.3.3 Performance Metrics

Tools (i - iii) predict CNVs either as deletion or duplication. The main problem of these callers are false discovery rates [19, 20]. Given a deletion or duplication call by tools (i - iii), DECoNT outputs a probability for the call to be deletion, duplication or no call (i.e., false discovery). The option with the highest probability is returned as the prediction.

In order to assess the performance of tools (i - iii) before and after being polished, we calculate the following performance metrics using  $Y_{pr}^{(1:N)}$  and  $Y_{gt}^{(1:N)}$ : (i) duplication call precision; (ii) deletion call precision, and (iii) overall precision. We first define the following variables:  $TP_1 :=$  number of duplications correctly identified;  $TP_2 :=$  number of deletions correctly identified;  $FP_1 :=$  number of duplications incorrectly identified;  $FP_2 :=$  number of deletions incorrectly identified.

Then, the performance metrics are defined as follows:

$$\text{Micro precision} = \text{Micro recall} = \text{Accuracy} = \frac{\text{TP}_1 + \text{TP}_2 + \text{TP}_3}{\text{TP}_1 + \text{FP}_1 + \text{TP}_2 + \text{FP}_2 + \text{TP}_3 + \text{FP}_3} \quad (2.6)$$

$$\text{Macro precision} = \frac{(\text{TP}_1/(\text{TP}_1 + \text{FP}_1)) + (\text{TP}_2/(\text{TP}_2 + \text{FP}_2)) + (\text{TP}_3/(\text{TP}_3 + \text{FP}_3))}{3} \quad (2.7)$$

$$\text{Macro recall} = \frac{(\text{TP}_1/(\text{TP}_1 + \text{FN}_1)) + (\text{TP}_2/(\text{TP}_2 + \text{FN}_2)) + (\text{TP}_3/(\text{TP}_3 + \text{FN}_3))}{3} \quad (2.8)$$

$$\text{Duplication call precision} = \frac{\text{TP}_1}{\text{TP}_1 + \text{FP}_1} \quad (2.9)$$

$$\text{Deletion call precision} = \frac{\text{TP}_2}{\text{TP}_2 + \text{FP}_2} \quad (2.10)$$

$$\text{Overall precision} = \frac{\text{TP}_1 + \text{TP}_2}{\text{TP}_1 + \text{TP}_2 + \text{FP}_1 + \text{FP}_2} \quad (2.11)$$

In order to test DECoNT’s performance on exact CNV prediction problem, which is a regression task, we use Absolute Error ( $AE$ ) between the predicted and ground truth copy number values. For an event  $X_i$ ,  $AE^{(i)}$  is defined as follows:

$$AE^{(i)} = Y_{pr}^{(i)} - Y_{gt}^{(i)} \quad (2.12)$$

### 2.3.4 Time Performance

All models are trained on a SuperMicro SuperServer 4029GP-TRT with 2 Intel Xeon Gold 6140 Processors (2.3GHz, 24.75M cache), 251GB RAM, 3 NVIDIA GeForce RTX 2080 Ti (11GB, 352Bit) and 1 NVIDIA TITAN RTX GPUs (24GB, 384Bit). We used 4 GPUs in parallel to train all 5 models and total training times were approximately as follows:  $\sim 70, 12, 95, 50$ , and 20 hours for XHMM, CoNIFER, CODEX2, Control-FREEC, and CNVkit, respectively. Note that

training is performed offline. The average polishing time per sample is in the order of seconds, for all models.

## 2.4 Polishing Samples from Other Sequencing Platforms

The training data we use is obtained using Illumina Genome Analyzer II and Illumina HiSeq 2000 machines. We check if models trained on these 1000 Genomes data can be used to polish CNV calls made on WES samples obtained using other sequencing platforms or capture kits that have not been seen by DECoNT.

We obtain the WES data for the sample NA12878, sequenced using four different platforms: (i) Illumina NovaSeq 6000; (ii) Illumina HiSeq 4000; (iii) BGISEQ-500; and (iv) MGISEQ-2000. Reads are aligned to the reference genome (GRCh38) using BWA [\[31\]](#) with *-mem* option and default parameters. Average depth coverage for these samples are  $241\times$ ,  $395\times$ ,  $328\times$ , and  $129\times$ , respectively. We use these four samples only for testing. All considered WES-based CNV callers are used to call CNV events on these four WES samples with default parameters. Using the CNVnator calls obtained on the WGS sample for NA12878 as the ground truth, we measure the performance the CNV callers before and after polishing with DECoNT. Note that NA12878 data is not included in the training data set in any form.

## 2.5 Polishing Other WES-based CNV Caller Algorithms

In our framework, a separate DECoNT model is trained for every WES-based germline CNV caller. We check if a DECoNT model trained using calls made by one algorithm can be used to polish the calls made by others in the absence of a

trained model.

We use the same three models trained with the settings described in Section 2.3 for XHMM, CoNIFER, and CODEX2. For each tool-specific DECoNT model, we polish the calls made by others. Here the training and testing folds are again exclusive. For testing, we use the same test folds for each tool as described in Section 2.3. This experiment results in 6 tests (i.e., for two-way comparison among every tool pair). We measure the performance of the polishing procedure using duplication precision, deletion precision and overall precision to obtain 18 performance results in total.

# Chapter 3

## Discussion

High throughput sequencing platforms, since their inception in 2007, have now become the dominant source of data generation for biological and medical research and on their way to be routinely used for diagnosis and treatment guidance. Although whole human genome sequencing cost is now reduced below the \$1,000 mark, whole exome sequencing will likely remain the main workhorse in clinical settings due to i) lower cost, ii) capturing almost all actionable genetic defects within exons, and iii) smaller data size that reduces computational burden for analysis. However, the main drawback of WES has been discovery and genotyping of CNVs. First, depth of coverage among exons are not uniform, making it very difficult to apply read depth based methods. Second, the reads often do not span CNV breakpoints which is a must for read pair and split read based approaches. Therefore it is often necessary to complement WES studies with alternative approaches such as array comparative genomic hybridization or quantitative RT-PCR.

We specifically designed our new algorithm, DECoNT, to address this limitation as a CNV call *polisher*. Using a deep learning approach, we were able to boost both the precision of several widely-used state-of-the-art algorithms that use WES data for CNV discovery. Although we trained DECoNT using matched WGS and WES samples from the 1000 Genomes Project, we also demonstrated

that the performance gain is independent from the training data, the capture kit, and the sequencing platform. Therefore, the trained model is portable, and it can be applied to data sets regardless of the data generation protocol without requiring new samples to train DECoNT.

Copy number variation is an important cause of genetic diseases that may be difficult to characterize in clinical settings without specific assays. WES is a powerful method to genotype small mutations but so far it has been unsuccessful to discover large CNVs that have a more direct effect in gene losses. DECoNT aims to help ameliorate high false discovery rate problems related to CNV characterization using WES, also including integer copy number prediction. Therefore DECoNT adds an important type of genomic variation discovery to the capabilities of WES and enhances the genome analysis arsenal in the clinic.

DECoNT uses WGS-derived CNV calls as labels for training. Note that these labels cannot serve as the ground truth but rather as the semi-ground truth. Unfortunately, there does not exist sufficiently large hand-curated labeled data for training a model. Chaison et al. provides hand-curated CNV calls on 9 samples which let us only perform validation. While a larger sample set was used, the latest release by HGSC [30] contains only 674 CNVs, which is a very small number for training DECoNT. This call set can also be regarded as a consensus call set (i.e., result of a consensus caller) because it is generated by using three different calling pipelines, and supported by long-read sequencing, StrandSeq, and optical mapping analysis of a subset of these genomes. We attempted training a model with this data set but unsurprisingly the model did not converge. For comparison, for training the DECoNT model for XHMM, we used  $\sim 68k$  calls. While these hand curated high quality data sets and consensus callers are certainly going to help DECoNT to achieve higher precision, they are quite limited in size which currently prohibits training. We foresee that with increasing size of call sets, DECoNT's performance will also increase. Yet, it is evident that CNV calling on WGS data is more accurate compared to CNV calling on WES data even when using a single over-the-shelf WGS-based CNV caller. Thus, DECoNT transfers these higher confidence labels into the WES domain and is limited by the precision of the underlying WGS-based CNV caller (CNVnator in this case)

[32]. Abyzov et al. report that CNVnator has high sensitivity (86%-96%), and low false-discovery rate (3%-20%). This corresponds to a precision range of 80% to 97% (mean 11.8%). Note that [32] obtains these performance results on two high coverage (20 $\times$  - 32 $\times$ ) trios in the 1000 Genomes Pilot Project data set [40]. DECoNT uses the newer 1000 Genomes data set at 30 $\times$  coverage generated using NovaSeq [30]. Thus, we expect the error tolerance of CNVnator to be better or similar in our analyses.

One issue with the polisher is to come up with a recipe to set the parameters of the base caller to achieve the best performance. In this study, we mostly used the suggested parameters and had to relax CoNIFER’s parameters as it hardly returned any calls. One other option is to run the base caller in the most liberal setting to improve sensitivity and let DECoNT correct the likely higher number of false positives. We tested if this is feasible using XHMM which is the best performing algorithm in our benchmarks. As detailed in Supplementary Note 3, we ran it also in a liberal and a conservative setting and then polished the calls using DECoNT. The precision improvement of tool was stable at  $\sim 22\%$  in all three settings. The polished precision of the liberal setting was worse by 10% compared to the suggested setting. Conservative and suggested setting precision values differed by only  $\sim 1\%$ . Thus, we suggest using the *suggested* parameter settings of the base caller unless it hardly makes calls and provides an insufficient number of calls for training.

The next challenge will be relieving DECoNT from dependence on existing variation callers, and make it a standalone, highly accurate CNV discovery tool using whole exome sequencing. One possible such direction for DECoNT could be redesigning the architecture to work with bin-level data. That is, most baseline callers first analyze read depth in small bins and then, adjacent bins are smoothed/combined via a segmentation algorithm which are returned as final calls if enough evidence exists along multiple neighboring bins. Currently, DECoNT-the-polisher, is working with this final decision which is limiting. Working with bin level data will require an architecture that can handle one basepair resolution whereas now DECoNT works with kilobase-sized windows which are averaged. As is, the Bi-LSTM component is prone to vanishing/exploding gradients due to

possibly very long read depth sequences.

# Chapter 4

## Results

### 4.1 Bi-LSTM based Neural Network Learns to Correct False Positive Germline WES CNV Calls

A Bidirectional Long Short-Term Memory (Bi-LSTM) network [41] is a type of recurrent neural network which learns a representation (i.e., embedding) of a sequence by processing it for each time-step (i.e., each read depth value in the CNV region in our case) in the forward and the backward directions. While doing so, it remembers a summary of the sequence observed so far to capture the context for each time-step. RNNs and LSTM-based architectures have been widely and successfully used in natural language processing domain to process sequence data [42].

DECoNT uses a single hidden layered Bi-LSTM architecture with 128 hidden neurons in each direction to process the read depth signal (Methods). First, WES-based germline CNV caller result is obtained along with the read depth signal in those event regions. Bi-LSTM subnetwork learns a transformed representation for the read depth sequence (Figure ??). This embedding and the corresponding

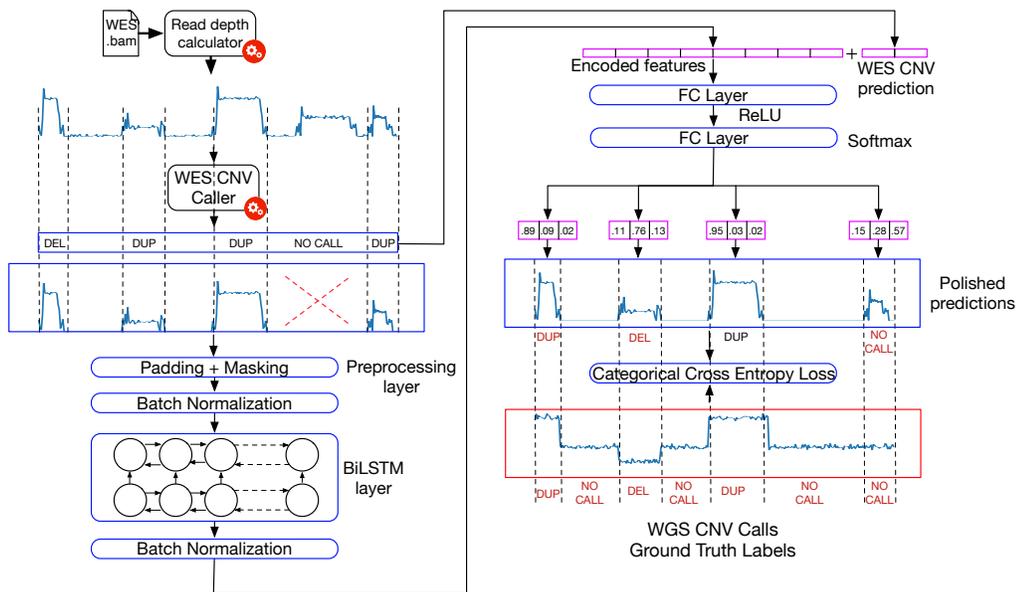


Figure 4.1: **Learning workflow of DECoNT.** First, BAM file that corresponds to a WES data set from **1000 Genomes data set** is used to calculate exome-wide read depth which is input to a third party WES-based CNV caller. The caller generates the calls for various regions which could be (i) a binary prediction like duplication, deletion (e.g., XHMM [1]) as shown in the figure, or (ii) an integer value that indicates the exact copy number (i.e., Control-FREEC [2]). The read depth of the regions for which a call has been made is input to a Bi-LSTM model. Encoded features are passed from a series of fully connected layers along with the original prediction of the caller algorithm. Using the ground truth calls from the WGS data of the same sample the method learns to predict (correct) the calls using cross entropy loss for the binary outputs (as shown in the figure) and using mean squared loss for integral calls.

CNV call are input to a fully connected (FC) layer feed forward neural network. The FC layers predict the polished result for call. DECoNT makes use of the calls made on the WGS data of the same sample as the ground truth for the learning procedure. We use matched WGS data to obtain the ground truth calls for the CNV events called on the WES samples of the same individuals in the 1000 Genomes data set [29].

We polish state-of-the-art WES-based germline CNV callers. There are two types of such algorithms. The first type makes discrete predictions for CNVs (i.e., deletion and duplication). We consider three methods in this category: (i) XHMM [1], (ii) CoNIFER [36], and (iii) CODEX2 [37]. The second type predicts the exact copy number as an integer value. The examples we consider of this type is Control-FREEC [2] and CNVkit [38]. DECoNT architecture is flexible and can be easily modified to polish both types of algorithms (Methods). We train a DECoNT model for every above-mentioned tool using 3 NVIDIA GeForce RTX 2080 Ti and 1 NVIDIA TITAN RTX GPUs in parallel with training times ranging from  $\sim 1$  to  $\sim 4$  days (Methods).

We find that DECoNT is able to substantially improve the performance of all algorithms in almost all comparisons. For algorithms that make discrete predictions, we observe improvements in both duplication and deletion call precision (Figure 4.2a). The largest gain is in duplication call precision for CoNIFER which is improved by 3-fold (i.e., 24.68% to 75%). The largest gain in deletion call precision is again obtained for CoNIFER, which is improved by 1.5-fold (i.e., 45.45% to 68.51%). Also, the overall precision is improved by 2.6-fold (i.e., 27.22% to 71.11%) for CoNIFER. This improvement is especially striking as CoNIFER is relatively conservative compared to other algorithms and seldom make calls despite relaxation of its parameters. For XHMM, we observe 1.4, 1.7, and 1.5 fold increases which correspond to 20%, 29%, and 24% improvements in duplication, deletion and overall precision, respectively. We see a similar trend for CODEX2. Before polishing with DECoNT, CODEX2 achieves 12% duplication precision, 45% deletion precision, and 27% overall precision. DECoNT provides 1.9 fold increase in duplication call precision, 1.5 fold increase in deletion call precision, and 1.75 fold increase in overall precision, respectively. These correspond to 11%,

23%, and 20% improvements in each respective metric. Confusion matrices before and after polishing are shown in Supplementary Figure 1 for all tools. We would like to note that these improvements are obtained in seconds per sample at the test time. Increased precision is an important result for the life scientists who work with these calls as the reliability of the calls are substantially increased as the number of false positives are substantially decreased.

As for the Control-FREEC and CNVkit, which output exact copy number values, we evaluate their performance (i.e., absolute error; *AE*). For Control-FREEC we consider 20,482 CNV events called (Figure 4.2b, Methods). DECoNT improves the absolute error in 74.58% of the test samples for an average *AE* improvement of 47.39 and deteriorates the performance in 25.35% of the test samples for an average *AE* deterioration of only 1.2. While unpolished Control-FREEC predictions have a Spearman correlation coefficient of 0.227 with matched ground truth copy numbers, DECoNT-polished predictions have a Spearman correlation coefficient of 0.568 (Figure 4.2c). DECoNT-polished predictions highly resemble the distribution of the ground truth calls. In order to mimic the discrete prediction case, we also discretized the CNV calls of Control-FREEC to Deletion ( $CN < 2$ ), Duplication ( $CN > 2$ ) and No-Call ( $CN = 2$ ) categories to measure precision as defined in Section 2.3.1. Again, DECoNT was able to improve the DEL and DUP precision up to 3-fold (see Supplementary Note 1 for details). The other tool that outputs exact copy number values is CNVkit. We evaluate its performance on 3,972 CNV events called (Figure 4.2d, Methods). DECoNT improves the absolute error in 86.78% of the test samples for an average *AE* improvement of 1.82 and deteriorates the performance in 13.21% of the test samples for an average *AE* deterioration of only 0.66. Raw CNVkit predictions have a Spearman correlation coefficient of 0.0156 and DECoNT-polished predictions have a Spearman correlation coefficient of 0.122 (Figure 4.2e). Similar to Control-FREEC, DECoNT-polished CNVkit predictions highly resemble the distribution of the ground truth calls.

To show the need for a complex model like DECoNT in this application, we used standard machine learning algorithms for polishing and compared the performance. We used SVM and logistic regression for the discrete prediction case and

polynomial regression for the exact prediction (rounded). We show that these models actually deteriorate the baseline caller performance and we need more complex models like DECoNT for this task. Details are given in Supplementary Note 2.

We also investigated DECoNT’s polishing performance on a consensus WES-based germline CNV caller, CNLearn [43]. CNLearn first runs 4 WES-based callers (CANOES, CODEX, CLAMMS, XHMM) and then using a random forest classifier, learns to aggregate the results of these programs. We obtained the 39 CNV predictions of CNLearn on 4 samples from the 1000 Genomes Project (personal communication with S. Girirajan; default settings are used). The list of these samples are given in Supplementary Note 4. Using the CNVnator calls obtained on the WGS data of the samples as the ground truth, we observed that CNLearn achieved a precision of 0.79. Using the DECoNT model trained using XHMM calls, we polished the results of CNLearn and improved the precision to 0.889. Note that CNLearn requires more computation as it employs many models, yet, DECoNT was able to improve the performance, even when using a cross-model polisher. More specifically, CNLearn and polished-CNLearn do not agree on 8 calls, out of which, polished version is correct in 4, unpolished version is correct in 2 and both are incorrect in 2. The list of these calls is also given in Supplementary Note 5.

We further analyzed the results of DECoNT in various other settings. First, we evaluated DECoNT for calls made on the X chromosome and pseudo-autosomal regions. We observe that the performance of the most reliable base caller XHMM is the lowest in these regions. Yet, DECoNT is still able to improve the performance in 8 out of 9 categories as shown in Supplementary Table 1. We also analyzed the length distribution of the corrected calls. In Supplementary Figure 4, we show that true XHMM calls have a large variance in size (i.e., up to 1 Kbp). Yet, the DECoNT-corrected true calls range up to 500 bps, showing that XHMM shows lower accuracy in shorter calls and needs polishing.

## 4.2 Polishing Performance on a Validated CNV call set

In order to further test the polishing performance of DECoNT, we also use a highly validated CNV call set published by Chaisson et. al., [3]. This data set contains the WGS-based CNV calls of 9 individuals selected from the 1000 Genomes Project samples, for which a consensus call set is obtained using 15 different WGS-based CNV callers with comparisons against high quality SVs generated using long read Pacific Biosciences data with a single basepair breakpoint resolution (Methods). We use data from 8 samples with matched WES data set.

Using the same models explained in Section ??, we correct the CNV calls made on WES data of the 8 samples made byXHMM, CoNIFER and CODEX2. Note that none of the DECoNT models have seen the data of these individuals during training. We validate the performance using this call set. Table 4.1 summarizes the performances before and after polishing with DECoNT, with respect to WGS validated calls.

Similar to analysis above, DECoNT improves the performance of all three algorithms in all comparisons. The most substantial improvements are observed for CoNIFER. 7%, 31.4% and 16% improvements are observed for duplication, deletion, and overall precision, respectively. It is noteworthy that while CoNIFER does not report any deletion events, DECoNT was able to correct incorrect duplication calls into correct deletion calls and increase the precision to 31.4% in this category. See Supplementary Figure 2 for the confusion matrices obtained before and after polishing by DECoNT. For XHMM and CODEX2, we see consistent improvements reaching up to nearly 2-fold for CODEX2.

### 4.3 Polishing Performance Generalizes to Unseen Sequencing Platforms

We obtained the training data from the 1000 Genomes data set, in which the WES component is produced using Illumina Genome Analyzer II and Illumina HiSeq 2000. While data from these platforms is abundant and sufficient training data set size can be met, for users using other sequencing platforms, it might not be possible to train DECoNT due to lack of matched WES and WGS samples. We therefore evaluated whether models trained on the available 1000 Genomes data can be used to polish CNV calls made on WES samples obtained using other sequencing platforms or capture kits that have not been seen by DECoNT (Methods).

We obtain the WES data for the sample NA12878, sequenced using four different platforms: (i) Illumina NovaSeq 6000; (ii) Illumina HiSeq 4000; (iii) BGISEQ-500; and (iv) MGISEQ-2000. We use these four samples only for testing. All considered WES-based CNV callers are used to call CNV events on these four WES samples.

Even though DECoNT has not seen the read depth information or the CNV events on these sequencing platforms, it still can generalize from the training on the 1000 Genomes data and still can substantially improve the performances of XHMM, CoNIFER, and CODEX2 (Table [4.2](#) and Supplementary Figure 3). We observe improvements in 34 out of 36 tests.

The most substantial improvement is observed for CODEX2 that corresponds to an average 2.6-fold increase in performance. This even exceeds testing performance on the same platform as training (i.e.,  $\sim 2$ -fold improvement). For XHMM, the performance is improved 10 out of 12 tests, reaching up to doubling the performance in overall precision performance for BGISEQ and MGISEQ platforms. For NovaSeq 6000 and HiSeq 4000, the performance deteriorates in duplication precision. However, XHMM makes a few duplication calls: 3 and 2, respectively. While DECoNT keeps the true positives, it adds a few false positives and this

results in the performance decrease in these settings. CoNIFER does not report any events on the NovaSeq 6000 platform despite tuning its parameters to more relaxed settings. On BGISEQ-500 and MGISEQ-2000 platforms, even though CoNIFER does not report any duplication calls, DECoNT finds some deletion calls and increases duplication precision from 0 to 1%. While it does not report any duplication calls for the HiSeq 4000 platform, DECoNT is able to increase the precision to 50% but by reporting a true positive and a false positive. The trend in deletion precision performance is similar. Finally, overall precision performance is consistently increased in all tests and the improvement ranges from 0.3% to 2.3%.

For Control-FREEC,  $\sim 65\%$  to  $\sim 74\%$  of the CNV calls have been improved as opposed to only  $\sim 7\%$  to  $\sim 8\%$  of the calls have been deteriorated by DECoNT. We observe a decrease in average absolute error after polishing in all four platforms which ranges from 0.94 to 1.0 (Table [4.3](#)).

We note that the improvements provided by DECoNT on the BGI and MGI platforms are important as these systems belong to a completely different manufacturer. Since these platforms are expected to have different systematic biases in read depth distributions compared to the training data of DECoNT, we would also expect a lower testing performance. Yet, DECoNT is able to generalize well and consistently proves to be useful across a diverse set of technologies. Overall, the performance is on par with the tests obtained on Illumina Genome Analyzer II and Illumina HiSeq 2000. Polishing procedure consistently improves the performance in a platform-independent manner.

## 4.4 Polishing Performance on Calls From Unseen CNV Callers

A distinct DECoNT model is trained for every WES-based germline CNV caller. This makes sense as the call regions and numbers substantially differ among

algorithms in their recommended settings (e.g., CODEX2 calls 10× more events than XHMM). We check if a DECoNT model trained using calls made by one algorithm can be used to polish the calls made by others in the absence of a trained model (e.g., due to time constraints in training).

We use the same DECoNT models trained for XHMM, CoNIFER, and CODEX2 on 1000 Genomes Project data. For each tool-specific DECoNT model, we polish the calls made by others on samples not seen during training. For instance, we polish the calls made by CODEX2, using the DECoNT model trained on XHMM calls. This experiment results in 6 tests (i.e., for two-way comparison among every tool pair). We measure the performance of the polishing procedure using duplication call precision, deletion call precision and overall precision to obtain 18 performance results in total (Methods).

We observe that DECoNT improves the performance metric in 10 out of the 18 comparisons, XHMM-trained DECoNT consistently improves the other tools' performance in all metrics, except DEL precision when polishing calls reported by CoNIFER - ranging from 2% to 13% (Figure 4.3 and Supplementary Figure 1). Duplication precision is improved in most of the cases with the exception of CoNIFER-trained and CODEX2-trained DECoNT models deteriorating the performance of XHMM by 11% and 8% respectively. For deletion precision, this is not the case. Deletion precision is improved for CODEX2 for both DECoNT models. However, for CoNIFER, deletion precision is deteriorated by 13% and 45% when polished with XHMM-trained and CODEX2-trained DECoNT models respectively. This is due to very limited number of deletion CNV predictions of CoNIFER as even a small perturbation to the true positives of deletion calls yield large differences in precision. Also, CoNIFER-trained DECoNT model very slightly deteriorates deletion precision of XHMM calls by 5%. While XHMM improves overall precision for both other methods, in half of the overall precision comparisons, the performance is decreased.

Overall, DECoNT is still somewhat effective despite being trained using a different call set. The training process uses the read depth information for the event regions which enables DECoNT to generalize to polish other tools. While,

arguably, it can be used to polish calls generated by other tools, a DECoNT model trained on the calls of the to-be-polished WES-based caller is suggested, as the improvements in the discussed performance metrics are larger in this case, which is expected.

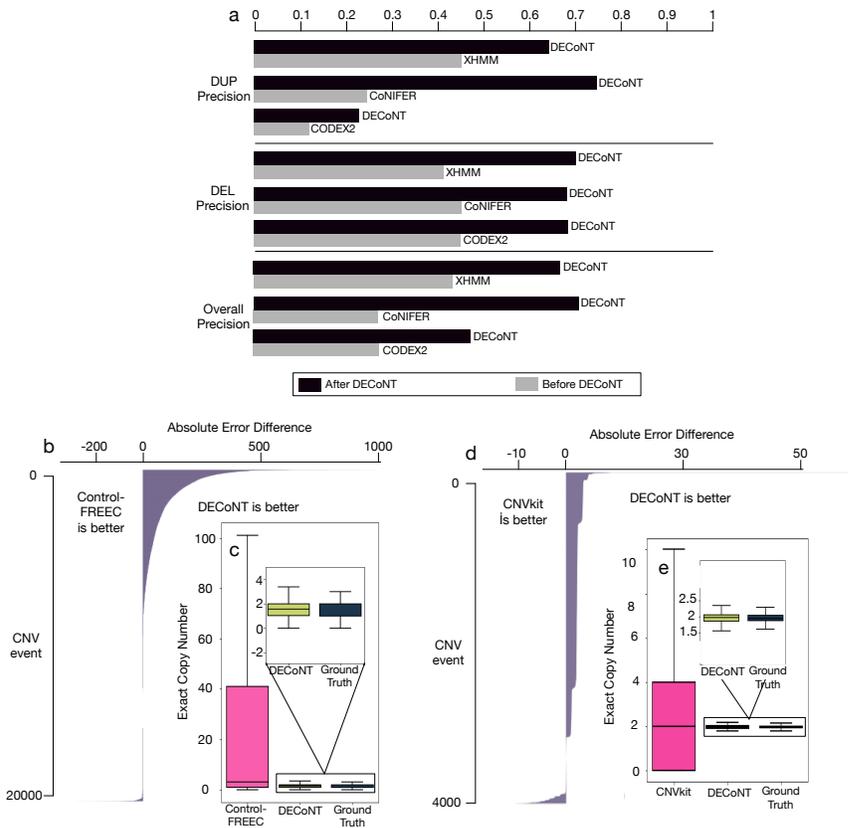


Figure 4.2: **The performance comparison of the WES-based CNV caller’s before and after polishing with DECoNT.** a) For the tools which predict existence of a CNV event (XHMM, CoNIFER and CODEX2) are evaluated with respect to duplication call precision, deletion call precision and overall precision. DECoNT improves the performance for all tools in all settings and results in drastic improvements. Different shades of gray represent different tools and the attached black bars represent the DECoNT-polished version of those tools. b) In this panel, we compare Control-FREEC and the DECoNT-polished with respect to Absolute Error (AE) difference on each sample (i.e., events). Bars to the right indicate the magnitude of the improvement due to polishing of DECoNT. For more than half of the samples, DECoNT results show improvement. c) The distribution of the unpolished Control-FREEC predictions in the test samples (pink) is quite different than the ground truth copy number variation distribution. On the other hand, DECoNT polished versions of the same events (dark blue) highly resemble the distribution of the ground truth calls. Black lines across the boxes are median lines for the distributions. Black vertical lines are whiskers and  $1.5 \times$  Inter Quartile Range is defined with the horizontal lines at the top and bottom of the whiskers. 1000 Genomes data set WGS samples are used as ground truth calls in all analyses. Panels d) and e) show the results for CNVkit, similar to b) and c).

Table 4.1: The performances of the WES-based CNV caller algorithms before and after polishing are shown (DEL, DUP and overall precision).

2*Tool	DUP Precision		DEL Precision		Overall Precision	
	default	polished	default	polished	default	polished
XHMM	0.064	<b>0.071</b>	0.257	<b>0.387</b>	0.135	<b>0.170</b>
CoNIFER	0.090	<b>0.160</b>	0.0*	<b>0.314</b>	0.090	<b>0.250</b>
CODEX2	0.027	<b>0.046</b>	0.387	<b>0.685</b>	0.185	<b>0.350</b>

ValidatedWGS CNV call set of Chaisson *et al.* [3] is used as the ground truth CNV call set. We first use matched WES reads to call WES CNVs using CoNIFER, CODEX2, and XHMM. Then, we use DECoNT to polish obtained CNV calls. Table shows the DEL, DUP and overall precision of the methods. \*CoNIFER does not report any deletion events on this set of WES samples.

Table 4.2: The performance of discrete germline WES-based CNV callers on NA12878 data before and after being polished by DECoNT.

2*Platform	2*Tool	DUP Precision		DEL Precision		Overall Precision	
		default	polished	default	polished	default	polished
NovaSeq 6000	XHMM	<b>0.660</b>	0.330	0.078	<b>0.111</b>	0.097	<b>0.133</b>
	CoNIFER	NA*	NA*	NA*	NA*	NA*	NA*
	CODEX2	0.043	<b>0.139</b>	0.198	<b>0.398</b>	0.112	<b>0.266</b>
HiSeq 4000	XHMM	<b>0.500</b>	0.125	0.093	<b>0.156</b>	0.100	<b>0.152</b>
	CoNIFER	0.0**	<b>0.500</b>	0.191	<b>0.192</b>	0.191	<b>0.214</b>
	CODEX2	0.032	<b>0.075</b>	0.188	<b>0.389</b>	0.099	<b>0.212</b>
BGISEQ-500	XHMM	0.045	<b>0.076</b>	0.157	<b>0.176</b>	0.088	<b>0.200</b>
	CoNIFER	0.0**	<b>0.010</b>	0.052	<b>0.082</b>	0.052	<b>0.055</b>
	CODEX2	0.051	<b>0.156</b>	0.214	<b>0.492</b>	0.125	<b>0.364</b>
MGISEQ-2000	XHMM	0.045	<b>0.076</b>	0.157	<b>0.176</b>	0.088	<b>0.200</b>
	CoNIFER	0.0**	<b>0.010</b>	0.052	<b>0.082</b>	0.052	<b>0.055</b>
	CODEX2	0.051	<b>0.156</b>	0.214	<b>0.492</b>	0.125	<b>0.364</b>

We evaluated caller performance on NA12878 data obtained using four different sequencing platforms: (i) NovaSeq 6000; (ii) HiSeq 4000; (iii) BGISEQ-500; (iv) MGISEQ-2000. Note that DECoNT models did not train on the data sequenced with any of these sequencing platforms or on NA12878 sequencing data of any form. DUP precision, DEL precision and Overall precision results are shown. In all comparisons, DECoNT provides substantial improvements showing the generalizability of our models trained on 1000 Genomes data set. \*CoNIFER does not report any CNV calls on NA12878 WES data sequenced with NovaSeq 6000. For that reason DECoNT has no input to correct and thus that comparison is not applicable. \*\*CoNIFER does not report any duplication events in the unpolished case. 1000 Genomes data set WGS samples are used as ground truth calls.

Table 4.3: The performance of Control-FREEC on NA12878 data before and after being polished by DECoNT.

2*Platform	2*# of Events	% of Improved Events	% of Deteriorated Events	Mean Absolute Error (MAE) Difference Decreased by
NovaSeq 6000	329	<b>73.85%</b>	7.59%	<b>0.9392</b>
HiSeq 4000	437	<b>70.94%</b>	6.86%	<b>1.0022</b>
BGISEQ-500	367	<b>64.57%</b>	8.17%	<b>0.9809</b>
MGISEQ-2000	367	<b>64.57%</b>	8.17%	<b>0.9809</b>

We evaluate caller performance on NA12878 data obtained using four different sequencing platforms: (i) NovaSeq 6000; (ii) HiSeq 4000; (iii) BGISEQ-500; (iv) MGISEQ-2000. Note that DECoNT models did not train on the data sequenced with any of these sequencing platforms or on NA12878 sequencing data of any form. Table shows the number of CNVs reported on each sample, the percentage of improved and deteriorated events, the average decrease in absolute error after being polished by DECoNT. In all comparisons, DECoNT provides substantial improvements showing the generalizability of our models trained on 1000 Genomes data set. 1000 Genomes data set WGS samples are used as ground truth calls.

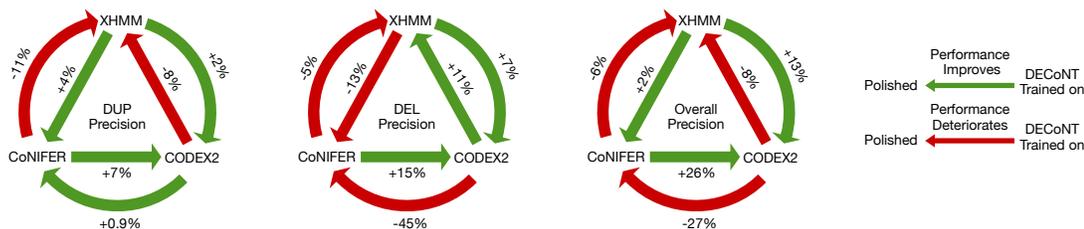


Figure 4.3: **Performance of DECoNT when polishing calls from unseen CNV callers.** DECoNT learns a different set of weights and a different model for each WES based CNV caller. In order to demonstrate the cross-model performance, we used DECoNT to correct CNV calls made by tools other than the ones used for training. We try every pair combination. Tools being pointed by an arrow are call generating tools (i.e., being corrected). Tools at source of the arrow are the tools that are used to train the DECoNT model. Green arrows indicate improvement and red arrows indicate deterioration in the corresponding performance metric.

# Bibliography

- [1] M. Fromer, J. L. Moran, K. Chambert, E. Banks, S. E. Bergen, D. M. Ruderfer, R. E. Handsaker, S. A. Mccarroll, M. C. O'Donovan, M. J. Owen, and et al., “Discovery and statistical genotyping of copy-number variation from whole-exome sequencing depth,” *The American Journal of Human Genetics*, vol. 91, no. 4, p. 597–607, 2012.
- [2] V. Boeva, T. Popova, K. Bleakley, P. Chiche, J. Cappo, G. Schleiermacher, I. Janoueix-Lerosey, O. Delattre, and E. Barillot, “Control-FREEC: a tool for assessing copy number and allelic content using next-generation sequencing data,” *Bioinformatics*, vol. 28, p. 423–425, Jun 2011.
- [3] M. J. Chaisson, A. D. Sanders, X. Zhao, A. Malhotra, D. Porubsky, T. Rausch, E. J. Gardner, O. L. Rodriguez, L. Guo, R. L. Collins, *et al.*, “Multi-platform discovery of haplotype-resolved structural variation in human genomes,” *Nature communications*, vol. 10, no. 1, pp. 1–16, 2019.
- [4] J. Sebat, B. Lakshmi, J. Troge, J. Alexander, J. Young, P. Lundin, S. Månér, H. Massa, M. Walker, M. Chi, *et al.*, “Large-scale copy number polymorphism in the human genome,” *Science*, vol. 305, no. 5683, pp. 525–528, 2004.
- [5] A. J. Iafrate, L. Feuk, M. N. Rivera, M. L. Listewnik, P. K. Donahoe, Y. Qi, S. W. Scherer, and C. Lee, “Detection of large-scale variation in the human genome,” *Nature genetics*, vol. 36, no. 9, pp. 949–951, 2004.
- [6] P. H. Sudmant, S. Mallick, B. J. Nelson, F. Hormozdiari, N. Krumm, J. Huddleston, B. P. Coe, C. Baker, S. Nordenfelt, M. Bamshad, L. B.

- Jorde, O. L. Posukh, H. Sahakyan, W. S. Watkins, L. Yepiskoposyan, M. S. Abdullah, C. M. Bravi, C. Capelli, T. Hervig, J. T. S. Wee, C. Tyler-Smith, G. van Driem, I. G. Romero, A. R. Jha, S. Karachanak-Yankova, D. Toncheva, D. Comas, B. Henn, T. Kivisild, A. Ruiz-Linares, A. Sajantila, E. Metspalu, J. Parik, R. Villems, E. B. Starikovskaya, G. Ayodo, C. M. Beall, A. Di Rienzo, M. F. Hammer, R. Khusainova, E. Khusnutdinova, W. Klitz, C. Winkler, D. Labuda, M. Metspalu, S. A. Tishkoff, S. Dryomov, R. Sukernik, N. Patterson, D. Reich, and E. E. Eichler, “Global diversity, population stratification, and selection of human copy-number variation.,” *Science*, vol. 349, p. aab3761, Sept. 2015.
- [7] N. Pankratz, A. Dumitriu, K. N. Hetrick, M. Sun, J. C. Latourelle, J. B. Wilk, C. Halter, K. F. Doheny, J. F. Gusella, W. C. Nichols, *et al.*, “Copy number variation in familial parkinson disease,” *PloS one*, vol. 6, no. 8, 2011.
- [8] E. L. Heinzen, A. C. Need, K. M. Hayden, O. Chiba-Falek, A. D. Roses, W. J. Strittmatter, J. R. Burke, C. M. Hulette, K. A. Welsh-Bohmer, and D. B. Goldstein, “Genome-wide scan of copy number variation in late-onset alzheimer’s disease,” *Journal of Alzheimer’s Disease*, vol. 19, no. 1, pp. 69–77, 2010.
- [9] D. Levy, M. Ronemus, B. Yamrom, Y.-h. Lee, A. Leotta, J. Kendall, S. Marks, B. Lakshmi, D. Pai, K. Ye, *et al.*, “Rare de novo and transmitted copy-number variation in autistic spectrum disorders,” *Neuron*, vol. 70, no. 5, pp. 886–897, 2011.
- [10] G. M. Cooper, B. P. Coe, S. Girirajan, J. A. Rosenfeld, T. H. Vu, C. Baker, C. Williams, H. Stalker, R. Hamid, V. Hannig, *et al.*, “A copy number variation morbidity map of developmental delay,” *Nature genetics*, vol. 43, no. 9, p. 838, 2011.
- [11] M. Zarrei, C. L. Burton, W. Engchuan, E. J. Young, E. J. Higginbotham, J. R. MacDonald, B. Trost, A. J. Chan, S. Walker, S. Lamoureux, *et al.*, “A large data resource of genomic copy number variation across neurodevelopmental disorders,” *NPJ genomic medicine*, vol. 4, no. 1, pp. 1–13, 2019.

- [12] H. Hieronymus, R. Murali, A. Tin, K. Yadav, W. Abida, H. Moller, D. Berney, H. Scher, B. Carver, P. Scardino, *et al.*, “Tumor copy number alteration burden is a pan-cancer prognostic factor associated with recurrence and death,” *Elife*, vol. 7, p. e37294, 2018.
- [13] M. Kumaran, C. E. Cass, K. Graham, J. R. Mackey, R. Hubaux, W. Lam, Y. Yasui, and S. Damaraju, “Germline copy number variations are associated with breast cancer risk and prognosis,” *Scientific reports*, vol. 7, no. 1, pp. 1–15, 2017.
- [14] B. M. Reid, J. B. Permeth, Y. A. Chen, B. L. Fridley, E. S. Iversen, Z. Chen, H. Jim, R. A. Vierkant, J. M. Cunningham, J. S. Barnholtz-Sloan, *et al.*, “Genome-wide analysis of common copy number variation and epithelial ovarian cancer risk,” *Cancer Epidemiology and Prevention Biomarkers*, vol. 28, no. 7, pp. 1117–1126, 2019.
- [15] G. Macintyre, T. E. Goranova, D. De Silva, D. Ennis, A. M. Piskorz, M. Eldridge, D. Sie, L.-A. Lewsley, A. Hanif, C. Wilson, *et al.*, “Copy number signatures and mutational processes in ovarian carcinoma,” *Nature genetics*, vol. 50, no. 9, pp. 1262–1270, 2018.
- [16] B. Trost, S. Walker, Z. Wang, B. Thiruvahindrapuram, J. R. MacDonald, W. W. Sung, S. L. Pereira, J. Whitney, A. J. Chan, G. Pellecchia, *et al.*, “A comprehensive workflow for read depth-based identification of copy-number variation from whole-genome sequence data,” *The American Journal of Human Genetics*, vol. 102, no. 1, pp. 142–155, 2018.
- [17] S. S. Ho, A. E. Urban, and R. E. Mills, “Structural variation in the sequencing era,” *Nature Reviews Genetics*, pp. 1–19, 2019.
- [18] A. Belkadi, A. Bolze, Y. Itan, A. Cobat, Q. B. Vincent, A. Antipenko, L. Shang, B. Boisson, J.-L. Casanova, and L. Abel, “Whole-genome sequencing is more powerful than whole-exome sequencing for detecting exome variants,” *Proceedings of the National Academy of Sciences*, vol. 112, no. 17, pp. 5473–5478, 2015.

- [19] F. Zare, M. Dow, N. Monteleone, A. Hosny, and S. Nabavi, “An evaluation of copy number variation detection tools for cancer using whole exome sequencing data,” *BMC bioinformatics*, vol. 18, no. 1, p. 286, 2017.
- [20] R. Tan, Y. Wang, S. E. Kleinstein, Y. Liu, X. Zhu, H. Guo, Q. Jiang, A. S. Allen, and M. Zhu, “An evaluation of copy number variation detection tools from whole-exome sequencing data,” *Human mutation*, vol. 35, no. 7, pp. 899–907, 2014.
- [21] L. Kadalayil, S. Rafiq, M. J. Rose-Zerilli, R. J. Pengelly, H. Parker, D. Oscier, J. C. Strefford, W. J. Tapper, J. Gibson, S. Ennis, *et al.*, “Exome sequence read depth methods for identifying copy number changes,” *Briefings in bioinformatics*, vol. 16, no. 3, pp. 380–392, 2015.
- [22] N. Krumm, P. H. Sudmant, A. Ko, B. J. O’Roak, M. Malig, B. P. Coe, A. R. Quinlan, D. A. Nickerson, E. E. Eichler, N. E. S. Project, *et al.*, “Copy number variation detection and genotyping from exome sequence data,” *Genome research*, vol. 22, no. 8, pp. 1525–1532, 2012.
- [23] J. M. Kebschull and A. M. Zador, “Sources of pcr-induced distortions in high-throughput sequencing data sets,” *Nucleic acids research*, vol. 43, no. 21, pp. e143–e143, 2015.
- [24] T. D. D. D. Study, T. Fitzgerald, S. Gerety, W. Jones, M. van Kogelenberg, D. King, J. McRae, K. Morley, V. Parthiban, S. Al-Turki, *et al.*, “Large-scale discovery of novel genetic causes of developmental disorders,” *Nature*, vol. 519, no. 7542, pp. 223–228, 2015.
- [25] S. De Rubeis, X. He, A. P. Goldberg, C. S. Poultney, K. Samocha, A. E. Cicek, Y. Kou, L. Liu, M. Fromer, S. Walker, *et al.*, “Synaptic, transcriptional and chromatin genes disrupted in autism,” *Nature*, vol. 515, no. 7526, pp. 209–215, 2014.
- [26] F. K. Satterstrom, J. A. Kosmicki, J. Wang, M. S. Breen, S. De Rubeis, J.-Y. An, M. Peng, R. Collins, J. Grove, L. Klei, *et al.*, “Large-scale exome sequencing study implicates both developmental and functional changes in the neurobiology of autism,” *Cell*, vol. 180, no. 3, pp. 568–584, 2020.

- [27] T. Singh, B. Neale, M. Daly, S. E. M.-A. Consortium, *et al.*, “Initial results from the meta-analysis of the whole-exomes of over 20,000 schizophrenia cases and 45,000 controls,” *European Neuropsychopharmacology*, vol. 29, pp. S813–S814, 2019.
- [28] K. J. Karczewski, L. C. Francioli, G. Tiao, B. B. Cummings, J. Alföldi, Q. Wang, R. L. Collins, K. M. Laricchia, A. Ganna, D. P. Birnbaum, *et al.*, “Variation across 141,456 human exomes and genomes reveals the spectrum of loss-of-function intolerance across human protein-coding genes,” *BioRxiv*, p. 531210, 2019.
- [29] The 1000 Genomes Project Consortium, “A global reference for human genetic variation,” *Nature*, vol. 526, pp. 68–74, Sep 2015.
- [30] M. Byrska-Bishop, U. S. Evani, X. Zhao, A. O. Basile, H. J. Abel, A. A. Regier, A. Corvelo, W. E. Clarke, R. Musunuri, K. Nagulapalli, *et al.*, “High coverage whole genome sequencing of the expanded 1000 genomes project cohort including 602 trios,” *bioRxiv*, 2021.
- [31] H. Li, “Aligning sequence reads, clone sequences and assembly contigs with bwa-mem,” 2013.
- [32] A. Abyzov, A. E. Urban, M. Snyder, and M. Gerstein, “CNVnator: An approach to discover, genotype, and characterize typical and atypical cnvs from family and population genome sequencing,” *Genome Research*, vol. 21, p. 974–984, Jul 2011.
- [33] C. Alkan, J. M. Kidd, T. Marques-Bonet, G. Aksay, F. Antonacci, F. Hormozdiari, J. O. Kitzman, C. Baker, M. Malig, O. Mutlu, *et al.*, “Personalized copy number and segmental duplication maps using next-generation sequencing,” *Nature genetics*, vol. 41, no. 10, p. 1061, 2009.
- [34] D. P. Kingma and J. Ba, “Adam: A method for stochastic optimization,” *arXiv preprint arXiv:1412.6980*, 2014.
- [35] X. Glorot and Y. Bengio, “Understanding the difficulty of training deep feedforward neural networks,” in *Proceedings of the thirteenth international conference on artificial intelligence and statistics*, pp. 249–256, 2010.

- [36] N. Krumm, P. H. Sudmant, A. Ko, B. J. Oroak, M. Malig, B. P. Coe, A. R. Quinlan, D. A. Nickerson, and E. E. Eichler, “Copy number variation detection and genotyping from exome sequence data,” *Genome Research*, vol. 22, no. 8, p. 1525–1532, 2012.
- [37] Y. Jiang, R. Wang, E. Urrutia, I. N. Anastopoulos, K. L. Nathanson, and N. R. Zhang, “Codex2: full-spectrum copy number variation detection by high-throughput dna sequencing,” *Genome Biology*, vol. 19, no. 1, 2018.
- [38] E. Talevich, A. H. Shain, T. Botton, and B. C. Bastian, “Cnvkit: genome-wide copy number detection and visualization from targeted dna sequencing,” *PLoS computational biology*, vol. 12, no. 4, p. e1004873, 2016.
- [39] A. Tarasov, A. J. Vilella, E. Cuppen, I. J. Nijman, and P. Prins, “Sambamba: fast processing of ngs alignment formats,” *Bioinformatics*, vol. 31, no. 12, p. 2032–2034, 2015.
- [40] The 1000 Genomes Project Consortium, “A map of human genome variation from population-scale sequencing,” *Nature*, vol. 467, no. 7319, p. 1061, 2010.
- [41] S. Hochreiter and J. Schmidhuber, “Long short-term memory,” *Neural computation*, vol. 9, no. 8, pp. 1735–1780, 1997.
- [42] Y. Yu, X. Si, C. Hu, and J. Zhang, “A review of recurrent neural networks: Lstm cells and network architectures,” *Neural computation*, vol. 31, no. 7, pp. 1235–1270, 2019.
- [43] V. K. Pounraja, G. Jayakar, M. Jensen, N. Kelkar, and S. Girirajan, “A machine-learning approach for accurate detection of copy number variants from exome sequencing,” *Genome research*, vol. 29, no. 7, pp. 1134–1143, 2019.

# Appendix A

## Supplementary Document

### A.1 Supplementary Figures

		Unpolished			DECoNT-XHMM Polished			DECoNT-CoNIFER Polished			DECoNT-CODEX2 Polished			
		NO CALL	DUP	DEL	NO CALL	DUP	DEL	NO CALL	DUP	DEL	NO CALL	DEL	DEL	
XHMM	Ground Truth	NO CALL	NA	1708	1447	1974	790	391	2104	310	741	2647	505	3
		DUP	NA	1587	508	360	1589	146	1384	263	448	1666	422	7
		DEL	NA	198	1384	217	84	1281	715	208	659	1359	212	11
CoNIFER	Ground Truth	NO CALL	NA	77	8	15	53	17	76	3	6	69	15	1
		DUP	NA	39	4	2	31	4	5	27	11	30	12	1
		DEL	NA	42	10	4	24	10	9	6	37	32	20	0
CODEX2	Ground Truth	NO CALL	NA	7955	6413	5337	3028	6003	1248	166	958	10974	1773	1621
		DUP	NA	2081	1786	368	1145	2354	339	130	561	896	1794	1177
		DEL	NA	7240	6774	1015	3723	9276	945	371	2282	3662	4216	6136

Figure A.1: The confusion matrices of the WES-based CNV callers before and after polishing with DECoNT on 1000 Genomes Data test samples. Confusion matrices given with blue borders represent unpolished predictions of corresponding WES-based CNV tools. Since DECoNT only operates on the calls made by a CNV caller, the first column for each unpolished confusion matrix is set as NA (i.e. Not Applicable). The red-bordered confusion matrices are the polished versions of a CNV caller with a DECoNT model trained on the calls made by the same caller (to produce Fig 2). Other confusion matrices are polished version of the CNV caller corrected by a DECoNT model trained on the calls made by a different caller (to produce Fig. 3). Notice the decrease in the number of false positives for both deletion and duplication calls in all platforms.

		Unpolished			Polished		
		NO CALL	DUP	DEL	NO CALL	DUP	DEL
		NO CALL	DUP	DEL	NO CALL	DUP	DEL
XHMM Ground Truth	NO CALL	NA	352	210	195	264	103
	DUP	NA	34	18	18	28	6
	DEL	NA	144	79	55	99	69
CoNIFER Ground Truth	NO CALL	NA	79	0	46	13	20
	DUP	NA	11	0	3	4	4
	DEL	NA	32	0	13	8	11
CODEX2 Ground Truth	NO CALL	NA	3251	2472	4643	588	492
	DUP	NA	145	118	128	87	48
	DEL	NA	1973	1636	1225	1209	1175

Figure A.2: Confusion matrices of the WES-based CNV callers before and after polishing with DECoNT on highly validated CNV callset published in Chaisson et. al. [3]. Similar to Fig. A.1 tool provides great false discovery correction with slight true positive deterioration for both deletion and duplication calls, yielding much better performance metric results.

		NovaSeq6000						HiSeq4000						BGI500						MGI2000					
		Unpolished			Polished			Unpolished			Polished			Unpolished			Polished			Unpolished			Polished		
		NO CALL	DUP	DEL	NO CALL	DUP	DEL	NO CALL	DUP	DEL	NO CALL	DUP	DEL	NO CALL	DUP	DEL	NO CALL	DUP	DEL	NO CALL	DUP	DEL	NO CALL	DUP	DEL
XHMM	NO CALL	NA	0	65	24	4	37	NA	0	79	27	6	46	NA	21	10	10	12	9	NA	21	10	10	12	9
	DUP	NA	2	17	6	2	11	NA	1	18	10	1	8	NA	1	6	1	1	5	NA	1	6	1	1	5
	DEL	NA	1	7	2	0	6	NA	1	10	0	1	10	NA	0	3	0	0	3	NA	0	3	0	0	3
CoNIFER	NO CALL	NA	0	0	0	0	0	NA	0	24	9	1	14	NA	0	5299	3088	846	1365	NA	0	5299	3088	846	1365
	DUP	NA	0	0	0	0	0	NA	0	14	6	1	7	NA	0	67	23	10	34	NA	0	67	23	10	34
	DEL	NA	0	0	0	0	0	NA	0	9	4	0	5	NA	0	299	115	58	126	NA	0	299	115	58	126
CODEX2	NO CALL	NA	613	442	840	105	110	NA	503	377	607	164	109	NA	320	250	478	33	59	NA	320	250	478	33	59
	DUP	NA	32	33	14	31	20	NA	19	19	9	19	10	NA	21	13	12	13	9	NA	21	13	12	13	9
	DEL	NA	98	118	43	87	86	NA	107	92	54	69	76	NA	63	72	32	37	66	NA	63	72	32	37	66

Figure A.3: Confusion matrices of the WES-based CNV callers before and after polishing with DECoNT on NA12878 data obtained from different sequencing platforms: (i) NovaSeq6000; (ii) HiSeq4000; (iii) BGI500; (iv) MGI2000. Since DECoNT only operates on the calls made by a CNV caller, the first column for each unpolished confusion matrix is set as NA (i.e. Not Applicable). Since CoNIFER does not report any calls on NovaSeq6000 platform, DECoNT has no input to polish and thus the comparison is not applicable. Similar to Figures [A.1](#) and [A.2](#), we observe that DECoNT substantially decreases the number of false discoveries with slight true positive deterioration for both deletion and duplication calls.

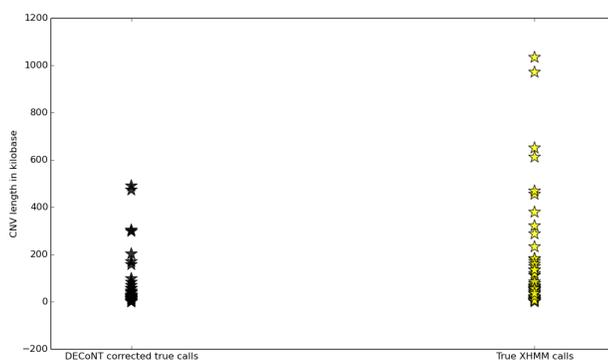


Figure A.4: This figure shows the length distribution of true raw XHMM calls and true DECoNT-corrected XHMM calls obtained on the 1000 Genomes WES data set test samples. The ground truth is the CNV calls made by CNVnator on the corresponding WGS samples. We see that for smaller size CNVs XHMM requires more correction by DECoNT. However, again, vast majority of the CNVs cannot be distinguished by the CNV length to decide whether it needs a DECoNT correction.

## A.2 Supplementary Tables

Table A.1: This table summarizes the polishing performance of DECoNT on the X chromosome, PAR1 and PAR2 regions of the males in the test split obtained from 1000 Genomes WES samples. The base caller in this analysis isXHMM. The results are obtained on the test samples from the 1000 Genomes dataset and the ground truth is obtained from the CNVnator calls on WGS of the same samples.

	Deletion Precision (Unpolished-Polished)	Duplication Precision (Unpolished-Polished)	Overall Precision (Unpolished-Polished)
X Chromosome	0.0350 - 0.1153	0.3018 - 0.4753	0.1702 - 0.4072
PAR1	0.1667 - 0.2500	0.7083 - 0.6112	0.5278 - 0.5455
PAR2	0.3334 - 0.6667	0.25 - 0.50	0.2871 - 0.4

Table A.2: In addition to the WES CNV Callers presented in the manuscript, we have also trained a DECoNT model for CNVKit on, again, 1000 Genomes WES samples, using the CNVnator calls obtained on WGS data as ground truth.

CNVKit 1000 Genomes Ground truth call-set is 1000 Genomes WGS samples.	Deletion Precision	Duplication Precision	Overall Precision
Before DECoNT	0.0940	0.1525	0.1234
After DECoNT	0.1234	0.5497	0.2527

Table A.3: The 8 WES CNV calls that DECoNT and CNLearn does not agree are presented. The ground truth CNV calls are obtained through CNVNator WGS CNV calls. Note that, CNLearn samples are polished with a DECoNT model trained with XHMM data. Training a DECoNT model with consensus calls made by CNLearn would increase performance.

Sample	Chromosome	CNV Start	CNV End	CNLearn Prediction	DECoNT Prediction	Ground Truth (CNVNator WGS Calls)
NA19144	11	6128771	6170380	DEL	NO-CALL	NO-CALL
NA19144	chr14	73541473	73573608	DUP	DEL	NO-CALL
NA11832	chr6	32519300	32666612	DUP	DEL	DEL
NA11832	chr15	34386562	34528116	DUP	DEL	DUP
NA18968	chr6	29889285	29945317	DUP	DEL	DEL
NA18968	chr6	32519300	32579157	DUP	DEL	DEL
NA18968	chr6	32584060	32665112	DUP	NO-CALL	DEL
NA12249	chr16	55810440	55826342	DUP	NO-CALL	DUP

## A.3 Supplementary Notes

**Supplementary Note 1** For the integer CNV calls of Control-FREEC, we have categorized the calls such that Copy Number  $> 2$  is Duplication, Copy Number  $< 2$  is Deletion and Copy Number  $= 2$  is No-Call. Then, we evaluated the polishing performance with the performance metrics defined in Section 4.3 Performance Metrics with respect to the 1000 Genomes WGS CNV calls of CNVNator:

- Duplication Precision was increased from 0.1063 to 0.3932
- Deletion Precision was increased from 0.2578 to 0.5936
- Overall Precision was increased from 0.1277 to 0.4432

**Supplementary Note 2** In order to show the need for a complex machine learning model like DECoNT for this polishing task, we also experimented with traditional machine learning methods such as Support Vector Machines (SVM), Logistic Regression and Polynomial Regression (degree = 2) as polishers. We used the scikit-learn implementations and the default parameters. These algorithms are run with the same settings we used for DECoNT. We worked on the 1000 Genomes dataset samples and same the train-test split. We input the same features into these models as we input to DECoNT: read depth and the call of the baseline caller. We used the corresponding WGS calls by CNVNator as the

ground truth as we do for DECoNT. We used XHMM and FREEC as the baseline callers for this experiment.

Below, we show that these models actually cannot polish the calls and deteriorate the results. See the notes below:

- XHMM predictions result in 0.4541 and 0.4144 precision for duplication and deletion calls, respectively.
- When correcting XHMM calls, SVM based model predictions result in 0.3562 and 0.3321 in precision for duplication and deletion calls, respectively.
- When correcting XHMM calls, Logistic Regression based model predictions result in 0.3334 and 0.2174 in precision for duplication and deletion calls, respectively.
- Control-FREEC predictions result in a MSE of 37.17 with standard deviation of 75.89
- When correcting Control-FREEC calls, Polynomial Regression polished model predictions result in a MSE of 58.10 with standard deviation of 18.91

**Supplementary Note 3** In order to test our assumption that running the base callers in their suggested parameter settings is sound we performed an experiment with XHMM which is the best performing method in our benchmarks. We ran it in also conservative and liberal settings in addition to the suggested setting. The parameter values that correspond to these settings are given in the table below.

XHMM	minTarget Size	maxTarget Size	minMean TargetRD	maxMean TargetRD	minMean SampleRD	maxMean SampleRD	maxSd SampleRD
Conservative	10	1000	10	5000	25	2000	1500
Suggested	5	10000	5	5000	5	2000	1500
Liberal	0	100000	0	50000	0	20000	15000

The precision values before and after polishing with DECoNT are given in the table below.

XHMM	Dup Precision (Unpolished-Polished)	Del Precision (Unpolished-Polished)	Overall Precision (Unpolished-Polished)
Conservative	0.4758 - 0.6548	0.4572 - 0.7120	0.4665 - 0.6834
Suggested	0.4541 - 0.6451	0.4144 - 0.7046	0.4348 - 0.6704
Liberal	0.3785 - 0.5543	0.3028 - 0.5921	0.3406 - 0.5732

We observe that the liberal setting results in a worse polished precision  $\sim 10\%$ . Conservative and suggested setting results are similar. The improvement in precision values are stable across all runs. Thus, we suggest using the default parameter settings for the base callers unless they return insufficient number of calls which prohibit DECoNT training. Then, the parameter choices can be relaxed.

**Supplementary Note 4** The 28 sample taken from 1000 Genomes dataset that were used in CNLearn analysis is given below:

- NA11832, NA12249, NA18968, NA19144

**Supplementary Note 5, Competing Interests:** Authors declare no competing interests.

**Supplementary Note 6, Author Contributions:** AEC and CA designed and supervised the study. AEC and FO designed the model. FO implemented the software and performed the experiments. AEC, CA and FO wrote the manuscript.

**Supplementary Note 7, Data Availability:** The training data used for DECoNT are obtained from 1000 Genomes Project. The WES bam files are available at [ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/data\\_collections/1000\\_genomes\\_project/data/](ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/data_collections/1000_genomes_project/data/). WGS CNV calls are obtained from CNVnator tool, they are already processed and made available at: [ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/data\\_collections/1000G\\_2504\\_high\\_coverage\\_SV/working/20190825\\_Yale\\_CNVnator/](ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/data_collections/1000G_2504_high_coverage_SV/working/20190825_Yale_CNVnator/). WES reads for NA12878 sample were obtained from: <https://www.ncbi.nlm.nih.gov/sra> with accession codes SRX5191370, SRX5191369, SRX5180030, and SRX5180221 for NovaSeq 6000, HiSeq 4000, BGISEQ-500, and MGISEQ-2000, respectively. Highly validated WGS CNV calls presented in Chaisson et. al. obtained from dbVar <https://www.ncbi.nlm.nih.gov/dbvar/> with accession nstd152. The inputs we use (i) CNV calls made by third party CNV callers and (ii) the calculated read depth data are also available at <https://zenodo.org/record/3865380> inside respective folders of the analysis. All other data that support the key findings of this paper can be found in the article and corresponding supplementary tables as referenced in the text.

**Supplementary Note 8, Code Availability:** DECoNT is implemented and released at <https://github.com/ciceklab/DECoNT> under [CC-BY-NC-ND 4.0 International license](https://creativecommons.org/licenses/by-nc-nd/4.0/). All custom python scripts that were used to generate matched WGS and WES CNV data are also available on the GitHub page. The scripts used to generate the data for all figures and tables in the manuscript are provided at <https://zenodo.org/record/3865380>.