

1 **Comparison of three bioinformatics tools in the detection of ASD candidate variants from whole**  
2 **exome sequencing data**

3

4 **Apurba Shil**<sup>1,2,3</sup>, **Liron Levin**<sup>4</sup>, **Hava Golan**<sup>2,3,5</sup>, **Gal Meiri**<sup>2,6</sup>, **Analya Michaelovski**<sup>2,7</sup>, **Yair**  
5 **Tsadaka**<sup>2,8</sup>, **Adi Aran**<sup>9</sup>, **Ilan Dinstein**<sup>2,3,10</sup> and **Idan Menashe**<sup>1,2,3\*</sup>

6 <sup>1</sup>Department of Epidemiology, Biostatistics, and Health Community Sciences, Faculty of Health  
7 Sciences, Ben-Gurion University of the Negev, Beer-Sheva, Israel

8 <sup>2</sup>Azrieli National Centre for Autism and Neurodevelopment Research, Ben-Gurion University of the  
9 Negev, Beer-Sheva, Israel;

10 <sup>3</sup>The School of Brain Sciences and Cognition, Ben-Gurion University of the Negev, Beer-Sheva,  
11 Israel;

12 <sup>4</sup>Bioinformatics Core Facility, Ben-Gurion University of the Negev, Beer-Sheva, Israel

13 <sup>5</sup>Department of Physiology and Cell Biology, Faculty of Health Sciences, Ben-Gurion University of  
14 the Negev, Beer-Sheva, Israel

15 <sup>6</sup>Preschool Psychiatric Unit, Soroka University Medical Center, Beer-Sheva, Israel

16 <sup>7</sup>Child Development Center, Soroka University Medical Center, Beer-Sheva, Israel

17 <sup>8</sup>Child Development Center, Ministry of Health, Beer-Sheva, Israel.

18 <sup>9</sup>Psychology Neuropediatric Unit, Shaare Zedek Medical Center, Jerusalem, Israel.

19 <sup>10</sup>Psychology Department, Ben-Gurion University of the Negev, Beer-Sheva, Israel

20

21 \* Correspondence: [idanmen@bgu.ac.il](mailto:idanmen@bgu.ac.il); tel.: +972-8-6477456

22

23

24 **Keywords:** Whole-exome sequencing, Autism Spectrum disorder, Bioinformatics, Genetics,

25 **Genomics**

26

## 27 **Abstract**

## 28 **Background**

29 Autism spectrum disorder (ASD) is a heterogenous multifactorial neurodevelopmental  
30 condition with a significant genetic susceptibility component. Thus, identifying genetic  
31 variations associated with ASD is a complex task. Whole-exome sequencing (WES) is an  
32 effective approach for detecting extremely rare protein-coding single-nucleotide variants  
33 (SNVs) and short insertions/deletions (INDELs). However, interpreting these variants'  
34 functional and clinical consequences requires integrating multifaceted genomic information.

## 35 **Methods**

36 We compared the concordance and effectiveness of three bioinformatics tools in detecting  
37 ASD candidate variants (SNVs and short INDELs) from WES data of 220 ASD family trios  
38 registered in the National Autism Database of Israel. We studied only rare (<1% population  
39 frequency) proband-specific variants. According to the American College of Medical  
40 Genetics (ACMG) guidelines, the pathogenicity of variants was evaluated by the *InterVar*  
41 and *TAPES* tools. In addition, likely gene-disrupting (LGD) variants were detected based on  
42 an in-house bioinformatics tool, *Psi-Variant*, that integrates results from seven in-silico  
43 prediction tools.

## 44 **Results**

45 Overall, 605 variants in 499 genes distributed in 193 probands were detected by these tools.  
46 The overlap between the tools was 64.1%, 17.0%, and 21.6% for *InterVar–TAPES*, *InterVar–*  
47 *Psi-Variant*, and *TAPES–Psi-Variant*, respectively. The intersection between *InterVar* and  
48 *Psi-Variant* ( $I \cap P$ ) was the most effective approach in detecting variants in known ASD genes  
49 (OR = 5.38, 95% C.I. = 3.25–8.53), while the union of *InterVar* and *Psi-Variant* ( $I \cup P$ )  
50 achieved the highest diagnostic yield (30.9%).

## 51 **Conclusions**

52 Our results suggest that integrating different variant interpretation approaches in detecting  
53 ASD candidate variants from WES data is superior to each approach alone. The inclusion of  
54 additional criteria could further improve the detection of ASD candidate variants.

55

## 56 **Background**

57 Autism spectrum disorder (ASD) comprises a collection of heterogeneous  
58 neurodevelopmental disorders that share two behavioral characteristics—difficulties in social  
59 communication and restricted, repetitive behaviors and interests<sup>1,2</sup>. The etiology of ASD has  
60 a significant genetic component, as is evident from multiple twin and family studies<sup>3-6</sup>. Yet,  
61 over the years, very few genetic causes of ASD have been discovered; thus, today, despite  
62 extensive research, an understanding of the overall genetic architecture of ASD remains  
63 obscure<sup>6,7</sup>.

64 The emergence of next-generation sequencing (NGS) approaches in the past decade has  
65 transformed the genetic research of complex traits<sup>8</sup>. These NGS technologies have facilitated  
66 high-throughput DNA sequencing for large cohorts of patients, allowing the comparison of  
67 multiple variants that includes single-nucleotide variants (SNVs) and short  
68 insertions/deletions (INDELs) between large groups of patients<sup>9-12</sup>. In this realm, whole-  
69 exome sequencing (WES) is particularly suitable for studying the genetics of heterogenous  
70 traits such as ASD, as it focuses on a relatively limited number of protein-coding  
71 variants<sup>9,10,13-18</sup>.

72 However, understanding the functional consequences of coding Variants is not a trivial  
73 task. In 2015, the American College of Medical Genetics and Genomics (ACMG) and the  
74 Association for Molecular Pathology (AMP) published standards and guidelines to generalize  
75 sequence variant interpretation and to address the issue of inconsistent interpretation across  
76 laboratories<sup>8</sup>. The resulting system for classifying variants recommends 28 criteria (16 for

77 pathogenic and 12 for benign variants) and provides a set of scoring rules based on variant  
78 population allele frequency, variant functional annotation, variant familial segregation,  
79 etc.<sup>8,19</sup>; Variants are classified as pathogenic (P), likely pathogenic (LP), variants of uncertain  
80 significance (VUS), likely benign (LB) or benign (B). Subsequently, multiple in-silico tools  
81 were developed to implement these ACMG/AMP criteria for annotating the prospective  
82 pathogenicity of variants detected in WES studies.

83 While the ACMG/AMP scoring approach is highly effective for detecting de-novo highly  
84 penetrant mutations for rare Mendelian disorders, it is less suitable for detecting inherited  
85 partially penetrant variants<sup>20</sup>. Such variants, usually annotated as VUS in terms of the  
86 ACMG/AMP criteria, are expected to contribute significantly to the risk of developing  
87 neurodevelopmental conditions, including ASD<sup>9,17,18,21,22</sup>. Thus, relying solely on the  
88 ACMG/AMP criteria for variant annotation in WES studies of ASD may result in an under-  
89 representation of susceptibility variants, which will lead to a lower diagnostic yield for ASD.  
90 To overcome this potential limitation, we have developed “*Psi-Variant*,” a pipeline to detect  
91 different types of likely gene-disrupting (LGD) variants, including protein truncating and  
92 deleterious missense variants. We applied *Psi-Variant* – in comparison with *InterVar* and  
93 *TAPES*, two variant interpretation tools that use the ACMG/AMP criteria – to a large WES  
94 dataset of ASD to evaluate the concordance between these tools to detect variants and to  
95 assess their effectiveness in detecting ASD susceptibility variants.

96

## 97 **Methods**

### 98 **Study sample**

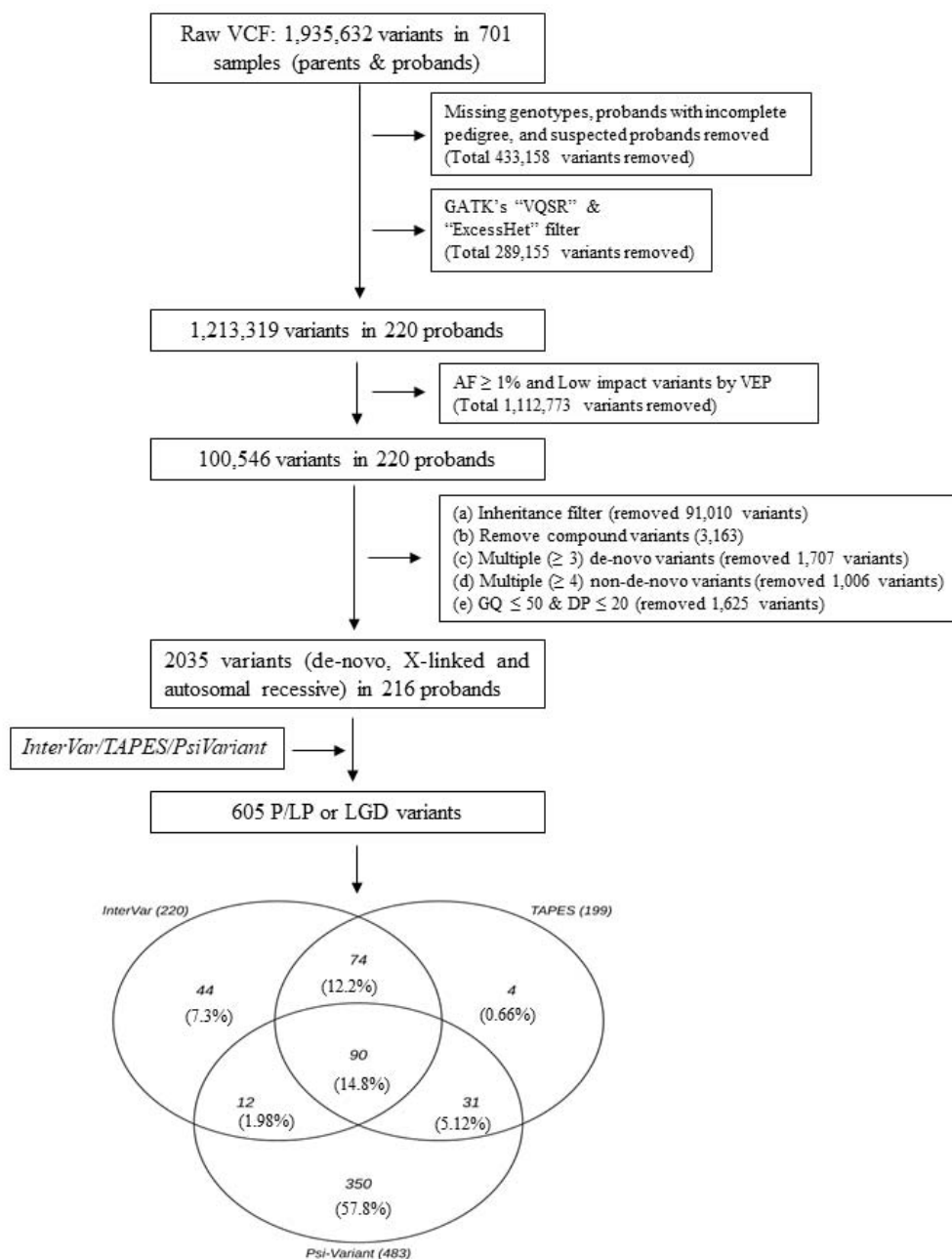
99 Initially, the study sample comprised 250 children diagnosed with ASD who are registered in  
100 the National Autism Database of Israel (NADI)<sup>23,24</sup> and whose parents gave consent for  
101 participation in this study. Based on our clinical records, none of the parents in the study has

102 been diagnosed with ASD, intellectual disability, or any other type of neurodevelopmental  
103 disorder. Genomic DNA was extracted from saliva samples from participating children and  
104 their parents using Oragene®•DNA (OG-500/575) collection kits (DNA Genotek, Canada).

105

#### 106 **Whole exome sequencing**

107 WES analysis was performed on the above-mentioned samples with Illumina HiSeq  
108 sequencers, followed by the Illumina Nextera exome capture kit at the Broad Institute as part  
109 of the Autism Sequencing Consortium, described previously<sup>11</sup>. Sequencing reads aligned to  
110 Genome Reference Consortium Human Build 38 and aggregated into BAM/CRAM files were  
111 analyzed using the Genome Analysis Toolkit (GATK)<sup>25</sup> to generate a joint variant calling  
112 format (vcf) file for all subjects in the study. We excluded data for 30 probands from the raw  
113 vcf file due to incomplete pedigree information or low-quality WES data. Thus, WES data for  
114 220 ASD trios were analyzed in this study (Fig. 1).



115

116 **Fig. 1** Analysis workflow for detecting LP/P/LGD Variants from the WES data. InterVar and  
 117 TAPES detected LP/P Variants by implementing ACMG/AMP criteria. *Psi-Variant* detected  
 118 LGD Variants by utilizing in-house criteria.

119

120

## 121 **Data analysis**

122 The SNV detection process in this study is outlined in Fig. 1. and explained below.

123

### 124 *Data cleaning*

125 The raw vcf file contained 1,935,632 variants. From this file, we removed variants with  
126 missing genotypes and/or variants in regions with low read coverage ( $\leq 20$  reads) and/or with  
127 low genotype quality ( $GQ \leq 50$ ). In addition, we removed all common variants (i.e., those  
128 with a population minor allele frequency  $>1\%$ )<sup>26</sup> as well as those that did not pass the  
129 GATK's "VQSR" and "ExcessHet" filters. Thereafter, we used an in-house machine learning  
130 (ML) algorithm to remove other potentially false-positive variants. The details of this ML  
131 algorithm and its efficiency in classifying true positive and false positive variants are  
132 summarized in the supplementary file S1. Finally, we used the pedigree structure of the  
133 families to identify proband-specific genotypes, including de-novo variants, recessively  
134 inherited variants, and X-linked variants (in males). Recessively inherited variants occur in  
135 the same loci of both copies of a gene in autosomes (where both the parents are carriers).  
136 Whereas one altered copy of the gene in chromosome X among males is defined as X-linked  
137 (males). We removed multiallelic variants from these genotypes and those classified as "de-  
138 novo" that appeared in more than two individuals in the sample. In this study, we haven't  
139 considered compound heterozygote variants (in cis/trans).

140

### 141 *Identifying ASD candidate variants*

142 We searched for candidate ASD Variants using three complementary approaches. First, we  
143 applied *InterVar*<sup>19</sup> and *TAPES*<sup>27</sup>, two commonly used publicly available command-line tools  
144 that use ACMG/AMP criteria<sup>8</sup>, to detect LP/P Variants. In addition, we assigned the  
145 ACMG/AMP PS2 criterion to all the de-novo Variants to detect additional LP/P Variants

146 from the list of VUS. Since *InterVar* and *TAPES* are less sensitive tools for detecting  
147 recessive possible gene disrupting (LGD) variants<sup>20</sup>, we developed an integrated in-house  
148 tool, *Psi-Variant*, to detect LGD variants. The *Psi-Variant* workflow starts using Ensembl's  
149 Variant Effect Predictor (VEP)<sup>26</sup> to annotate the functional consequences for each variant in a  
150 multi-sample vcf file. Then, all frameshift indels, nonsense, and splice acceptor/donor variants  
151 are further analyzed by the LoFtool<sup>28</sup> with scores of  $< 0.25$  are annotated as intolerant  
152 variants. In addition, it applies six different in-silico tools to all missense substitutions and  
153 annotates them as “deleterious/damaging” if at least three ( $\geq 50\%$ ) of them exceed the  
154 following cutoffs: SIFT<sup>29</sup> ( $< 0.05$ ), PolyPhen-2<sup>30</sup> ( $\geq 0.15$ ), CADD<sup>31</sup> ( $> 20$ ), REVEL<sup>32</sup> ( $> 0.50$ ),  
155 M\_CAP<sup>33</sup> ( $> 0.025$ ) and MPC<sup>34</sup> ( $\geq 2$ ). These scores were extracted by utilizing the dbNSFP  
156 database<sup>35</sup>.

157

#### 158 *Comparison between InterVar, TAPES, and Psi-Variant*

159 We compared the number of variants detected by the three tools and the percentages of  
160 variants detected by different combinations. Thereafter, we used the list of ASD genes ( $n =$   
161 1031) from the SFARI Gene database<sup>36</sup> (accessed on 11 January 2022) as the gold standard to  
162 compute the odds ratio (OR) and positive predictive value (PPV) for detecting candidate ASD  
163 variants in SFARI genes. In addition, we assessed the detection yield for each tool  
164 combination by computing the proportion of children with detected candidate ASD variants in  
165 SFARI genes.

166

#### 167 *Software*

168 Data storage, management, and analysis were conducted on a high-performing computer  
169 cluster in a Linux environment using Python version 3.5 and R Studio version 1.1.456. All the  
170 statistical analyses and data visualizations were incorporated into R Studio.



171

## 172 **Results**

### 173 **Detection of candidate variants by the different tools**

174 A total of 605 variants in 193 probands (highlighted in the supplementary Table S2) were  
175 detected by at least one of *InterVar* (n = 220), *TAPES* (n = 199), or *Psi-Variant* (n = 483) from  
176 a dataset of 2,035 high-quality, ultra-rare variants with proband-specific genotypes (Fig. 1).  
177 Of these, 90 variants (14.9%) were detected by all three tools. The highest concordance in  
178 detected variants was observed between *InterVar* and *TAPES* (64.3%), followed by *TAPES*  
179 and *Psi-Variant* (21.6%) and *InterVar* and *Psi-Variant* (17.0%).

180 The characteristics of the detected variants are shown in Table 1. Significantly higher  
181 rates of LoF and missense variants were detected by all three tools compared to the rates of  
182 these variants in the clean vcf file ( $P < 0.001$ ). As expected, missense variants comprised the  
183 majority of detected variants, with 81.6%, 58.8%, and 51.4% of the variants detected by *Psi-*  
184 *Variant*, *TAPES*, and *InterVar*, respectively. Notably, a higher number of frameshift variants  
185 were detected by *Psi-Variant* than by *InterVar* and *TAPES* (58 vs. 39 and 22, respectively),  
186 but the percentages of these variants out of the total number of detected variants were lower  
187 due to the markedly higher number of variants detected by *Psi-Variant*.

188 Almost all ( $\geq 95\%$ ) variants detected by either *InterVar* or *TAPES* were de-novo variants,  
189 while de-novo variants comprised only 36.2% of the variants detected by *Psi-Variant*, which  
190 also detected a high portion of X-linked and autosomal recessive variants (37.1% and 26.7%,  
191 respectively). Examination of the distribution of the detected variants in genes associated  
192 with ASD according to the SFARI Gene database<sup>36</sup> revealed a two-fold enrichment of  
193 variants distributed in ASD genes (for all detection tools) compared to their portion in the  
194 clean vcf file and even a higher enrichment of Variants distributed in high-confidence ASD  
195 genes ( $P < 0.001$ ).

196

197

**Table 1** Characteristics of the detected Variants by *InterVar*, *TAPES*, and *Psi-Variant* from the WES data

Characteristics	Preliminary output (n = 1213319)	<i>InterVar</i> (n = 220)	<i>TAPES</i> (n = 199)	<i>Psi-Variant</i> (n = 483)
<b>Functional consequence</b>				
Frameshift (insertions/deletions)	4232 (0.349%)	39 (17.7%) *	22 (11.1%) *	58 (12.01%) *
Missense	95919 (7.91%)	113 (51.4%) *	117 (58.8%) *	394 (81.6%) *
Stop Gain/Loss/retain, Start Gain/Loss	2105 (0.17%)	16 (7.27%) *	13 (6.53%) *	16 (3.31%) *
Non-frameshift/in-frame	4062 (0.33%)	42 (19.1%) *	43 (21.61%) *	--
Splice acceptor/donor/region	18817 (1.55%)	4 (1.82%)	4 (2.01%)	12 (2.48%)
Synonymous, downstream/upstream gene, intron variant	871205 (71.8%)	6 (2.73%)	0 (0%)	3 (0.62%)
Other	216979 (17.9%)	--	--	--
<b>Inheritance pattern wise</b>				
De-novo	43052 (3.55%)	209 (95%) *	193 (97%) *	175 (36.2%) *
Autosomal recessive	70948 (5.85%)	9 (4.09%)	5 (2.51%)	179 (37.1%) *
X-linked	9103 (0.75%)	2 (0.91%) *	1 (0.5%) *	129 (26.7%) *
Other	1090216 (89.8%)	--	--	--
<b>Gene type wise</b>				
SFARI genes with a score 1	19236 (1.58%)	15 (6.82%) *	12 (6.03%) *	21 (4.35%) *
All SFARI genes (with scores 1-3)	93681 (7.72%)	32 (14.5%) *	24 (12.1%) *	75 (15.5%) *
Other genes	1119638 (92.28%)	188 (85.4%) *	175 (87.9%) *	408 (84.5%) *

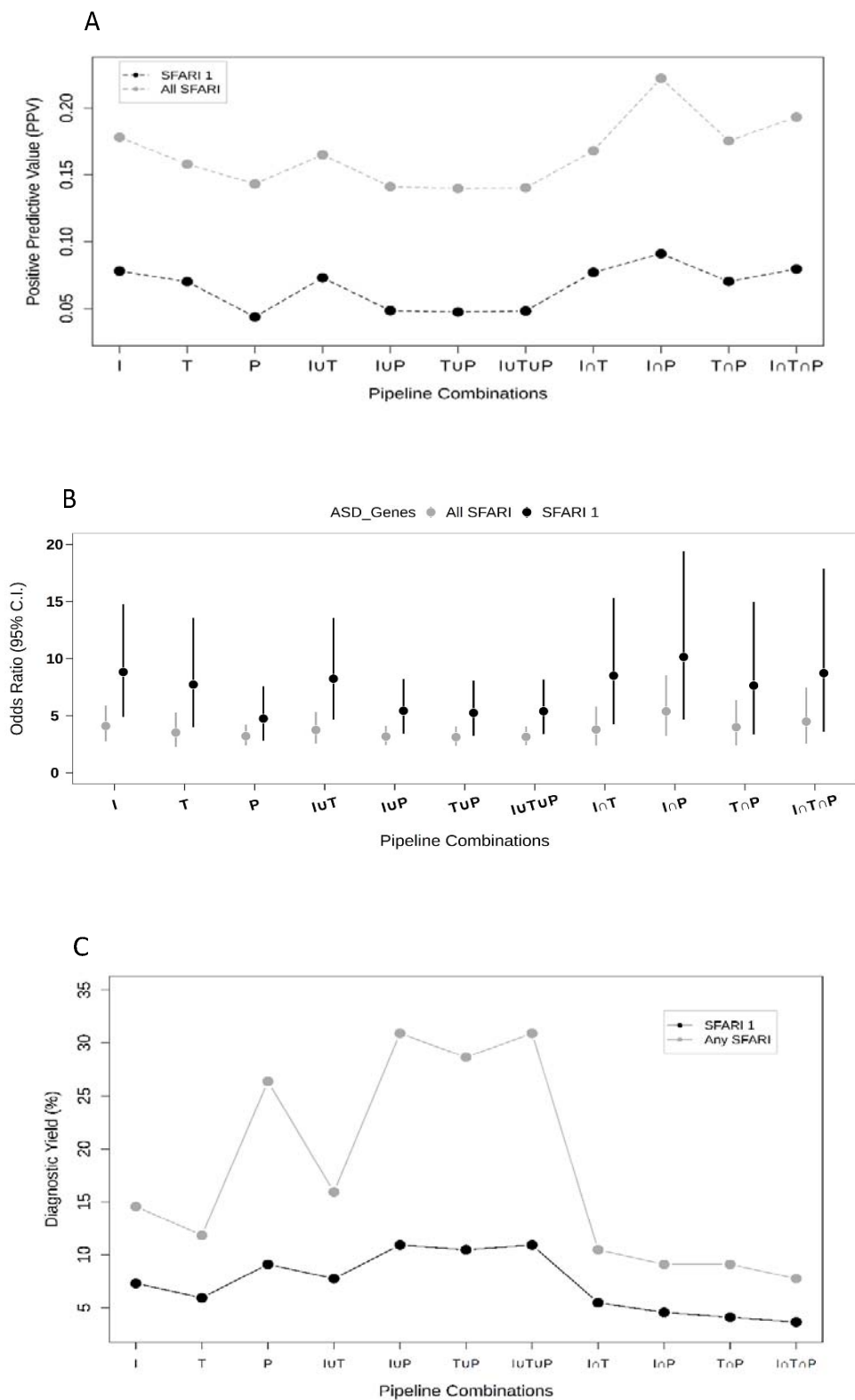
\* <0.05 level of significance; two-sided two proportions Z test

198

### 199 **Effectiveness of ASD candidate Variants detection**

200 To assess the effectiveness of the different tools in detecting ASD candidate SNVs, we  
 201 calculated the PPV and the OR for detecting ASD genes (i.e., those listed in the SFARI Gene  
 202 database<sup>36</sup>) for different combinations of utilization of the three tools. The results of these  
 203 analyses are shown in Fig. 2. Utilization of any of the three tools resulted in a significant  
 204 enrichment of ASD genes, with the highest enrichment being observed in SNVs detected by  
 205 *InterVar* (PPV = 0.178; OR = 4.10, 95% confidence interval (C.I.) = 2.77–5.90) followed by  
 206 *TAPES* (PPV = 0.158; OR = 3.53, 95% C.I. = 2.28–5.27) and *Psi-Variant* (PPV = 0.143; OR  
 207 = 3.21, 95% C.I. = 2.39–4.22). Notably, better performance in detecting ASD candidate  
 208 SNVs was obtained at the intersection of the detected SNVs between *InterVar* and *Psi-*

209 *Variant* ( $I \cap P$ ) (PPV = 0.222; OR = 5.38, 95% CI = 3.25–8.53). The  $I \cap P$  combination was  
210 also the most effective in detecting SNVs in high-confidence ASD genes (i.e., those with a  
211 score of 1 in the SFARI Gene database<sup>36</sup> (Fig. 2A -2B). However, the  $I \cap P$  combination had  
212 a relatively low diagnostic yield of 9.1% for SFARI genes. On the other hand, the union of  
213 *InterVar* and *Psi-Variant* ( $I \cup P$ ) achieved a diagnostic yield of 30.9% (Fig. 2C) (three times  
214 more than  $I \cap P$ ) but had a reduced effectiveness in detecting variants in SFARI genes (PPV  
215 = 0.141; OR = 3.18, 95% C.I. = 2.43–4.10) (Fig. 2A - 2B).



216 **Fig. 2** Effectiveness of *InterVar* (*I*), *TAPES* (*T*), *Psi-Variant* (*P*), and their combinations in  
 217 detecting candidate variants in ASD genes. **A** Positive predictive value (PPV) of detecting

218 candidate variants in SFARI 1 and all SFARI genes. **B** Odds Ratios (ORs) of detecting  
219 candidate variants in SFARI 1 and all SFARI genes. **C** Diagnostic yield (%) achieved by the  
220 different tool combinations for detecting candidate variants in SFARI 1 and all SFARI genes.

221

## 222 **Discussion**

223 In this study, we assessed the concordance and effectiveness of three bioinformatics tools in  
224 the interpretation of variants detected in the WES of children with ASD. There was better  
225 agreement in variant detection between *InterVar* and *TAPES* than between *Psi-Variant* and  
226 each of these two tools, probably because both *InterVar* and *TAPES* are based on the  
227 ACMG/AMP guidelines<sup>8</sup>, while *Psi-Variant* uses the interpretation of seven in-silico tools in  
228 assessing the functional consequences of LGD variants. In addition, most (94%) of the  
229 variants detected by either *InterVar* or *TAPES* were de-novo variants, compared to only 36%  
230 of the variants detected by *Psi-Variant*. This difference may be attributed to the fact that  
231 ACMG/AMP guidelines are particularly designed to detect de-novo highly penetrant variants,  
232 while inherited variants (autosomal recessive and X-linked) are usually classified as VUS<sup>20</sup>.  
233 Importantly, such rare inherited variants have been found to be associated with a variety of  
234 neurodevelopmental conditions, including ASD<sup>9,17,18,21,22</sup>. Another major difference between  
235 these tools lies in the detection of in-frame insertions/deletions that comprised ~20% of the  
236 variants detected by either *InterVar* or *TAPES*, while such SNVs were discarded by *Psi-*  
237 *Variant*. We decided to exclude these variants from *Psi-Variant* because their clinical  
238 relevance has been demonstrated in several genetic disorders<sup>37,38</sup> but not in ASD<sup>39-41</sup>.

239 Another important factor that could affect the concordance between the three tools is the  
240 annotation tools they use. Specifically, both *InterVar* and *TAPES* use AnnoVar<sup>42</sup> for their  
241 variant annotation, while *Psi-Variant* uses Ensembl's VEP<sup>26</sup>. It has already been shown that  
242 AnnoVar and VEP have a low concordance in the classification of LoF variants<sup>43</sup>. In

243 addition, each tool, *InterVar*, *TAPES*, and *Psi-Variant*, utilizes a different set of in-silico tools  
244 for the classification of missense variants, with SIFT<sup>29</sup> alone being shared by all three tools.  
245 These differences are probably the reason for the large differences in the detection of  
246 missense variants between the three tools (Table 1).

247 Today, there are no accepted guidelines for the detection of ASD susceptibility variants  
248 from WES data. Many genetic labs use the ACMG/AMP guidelines<sup>8</sup>, leading to a relatively  
249 low diagnostic yield<sup>44,45</sup>. Our findings suggest that different combinations of bioinformatics  
250 tools for variant interpretation may improve the detection of ASD susceptibility variants.  
251 Furthermore, combining these tools provides more flexibility in selecting the desired  
252 proportion between the detection yield and false positives. Thus, future guidelines for the  
253 detection of ASD susceptibility variants should consider the integration of different variant  
254 interpretation criteria.

255 Of note, many of the variants detected by the integrative pipeline affect genes with no  
256 known association with ASD, according to the SFARI Gene database<sup>36</sup>. This finding  
257 highlights the capability of the integrative pipeline to detect novel ASD genes. Obviously, the  
258 association of these genes and variants with ASD susceptibility needs to be validated in  
259 additional studies.

260 The results of this study should be considered under the following limitations. First, the  
261 effectiveness assessments of the different tools and their combinations were based on ASD  
262 genes from the SFARI Gene database<sup>36</sup>. While this is the most commonly used database for  
263 ASD genes and is continuously updated, it is based on data curated from the literature and  
264 may thus include genes falsely associated with ASD. Second, the variant detection analyses  
265 were performed on WES data of a cohort from the Israeli population, which may not  
266 necessarily be representative of the genetic architecture of ASD. Third, the tools used in this  
267 study were designed to detect only extremely rare variants with relatively large functional

268 effects. Thus, a more effective approach for the detection of ASD susceptibility variants  
269 should also include the interpretation of other types of genomic variations, such as copy-  
270 number and compound heterozygote variants<sup>46-51</sup>, as well as other variants with milder  
271 functional effects<sup>17,52,53</sup>. Finally, it should be noted that there are many other approaches for  
272 variant interpretation from WES data. Thus, it is possible that combinations of other  
273 approaches would be more effective in the detection of ASD susceptibility variants from  
274 WES data than the approaches investigated in this study.

275

## 276 **Conclusions**

277 Our findings suggest that combination of different bioinformatics tools is more effective in  
278 the detection of ASD candidate variants from WES data than each of the examined tools  
279 alone. Future guidelines for the detection of ASD susceptibility variants should consider  
280 integrating different variant interpretation approaches to improve the effectiveness of ASD  
281 candidate variants detection from whole exome sequencing data.

## **List of abbreviations**

ACMG/AMP: American College of Medical Genetics and Genomics/Association of Molecular Pathology; ASD: autism spectrum disorder; C.I.: confidence interval; GATK: Genome Analysis Toolkit; LGD: likely gene disrupting; LoF: loss of function; LP: likely pathogenic; ML: machine learning; NADI: National Autism Database in Israel; NGS: next-generation sequencing; OR: odds ratio; P: pathogenic; PPV: positive predictive value; SNV: single nucleotide variants; VEP: Variant Effect Predictor; vcf: variant calling format; VUS: variants of uncertain significance; WES: whole exome sequencing.

## **Declarations**

### **Institutional Review Board Statement**

The study was conducted according to the guidelines of the Declaration of Helsinki and approved by the Ethics Committee of Soroka University Medical Center (SOR-076-15; 17 April 2016).

### **Ethics approval and consent to participate**

Informed consent was obtained from all the families involved in the study.

### **Consent to publication**

Not applicable.

### **Availability of data and materials**

WES data were generated as part of the ASC and are available in dbGaP with study accession: phs000298.v4.p3. The generated results and codes are available in a GitHub public repository: <https://github.com/AppWick-hub/Psi-Variant> or available upon reasonable request to the corresponding author Prof. Idan Menashe ([idanmen@bgu.ac.il](mailto:idanmen@bgu.ac.il)).

### **Competing interests**

The authors declare no competing interests.

### **Funding**

This study was supported by a grant from the Israel Science Foundation (1092/21).

### **Authors' contributions**

*Conceptualization:* A.S. and I.M.; *methodology:* A.S. and I.M.; *software:* A.S. and L.L.; *validation:* A.S. and I.M.; *formal analysis:* A.S.; *resources:* H.G., G.M., A.M., Y.T., A.A.



and I.D.; *data curation*: A.S.; *writing—original draft preparation*: A.S. and I.M.; *writing—review and editing*: I.M., and A.S.; *supervision*: I.M.; *project administration*: I.M.; *funding acquisition*: I.M. All the authors have read and agreed to the published version of the manuscript.

## Acknowledgments

We thank the families who participated in this research, without their contributions, genetic studies would be impossible.

## References

1. American Psychiatric Association. Diagnostic and Statistical Manual of Mental Disorders, 5th edn. *American Psychiatric Publishing* (2013).
2. Meng-Chuan Lai, Michael V Lombardo, S. B.-C. Autism. *Lancet* (2014).
3. Yoo, H. Genetics of Autism Spectrum Disorder: Current Status and Possible Clinical Applications. *Exp Neurobiol* **24**, 257 (2015).
4. Lord, C. *et al.* Autism spectrum disorder. *The Lancet* **392**, 508–520 (2018).
5. Ronald, A. & Hoekstra, R. A. Autism spectrum disorders and autistic traits: A decade of new twin studies. *American Journal of Medical Genetics, Part B: Neuropsychiatric Genetics* **156**, 255–274 (2011).
6. Hallmayer, J. *et al.* Genetic Heritability and Shared Environmental Factors Among Twin Pairs With Autism. *Arch Gen Psychiatry* **68**, 1095–1102 (2011).
7. Devlin, B. *et al.* Patterns and rates of exonic de novo mutations in autism spectrum disorders. *Nature* **485**, 242–246 (2012).
8. Richards, S. *et al.* Standards and guidelines for the interpretation of sequence variants: A joint consensus recommendation of the American College of Medical Genetics and Genomics and the Association for Molecular Pathology. *Genetics in Medicine* **17**, 405–424 (2015).
9. Satterstrom, F. K. *et al.* Autism spectrum disorder and attention deficit hyperactivity disorder have a similar burden of rare protein-truncating variants. *Nat Neurosci* **22**, 1961–1965 (2019).
10. Fu, J. M. *et al.* Rare coding variation provides insight into the genetic architecture and phenotypic context of autism. *Nat Genet* **54**, (2022).
11. Satterstrom, F. K. *et al.* Large-Scale Exome Sequencing Study Implicates Both Developmental and Functional Changes in the Neurobiology of Autism Article Large-Scale Exome Sequencing Study Implicates Both Developmental and Functional Changes in the Neurobiology of Autism. *Cell* 1–17 (2020) doi:10.1016/j.cell.2019.12.036.
12. Wu, D. *et al.* Large-Scale Whole-Genome Sequencing of Three Diverse Asian Populations in Singapore. *Cell* **179**, 736-749.e15 (2019).
13. Satterstrom, F. K. *et al.* Large-Scale Exome Sequencing Study Implicates Both Developmental and Functional Changes in the Neurobiology of Autism Article Large-Scale Exome Sequencing Study Implicates Both Developmental and Functional Changes in the Neurobiology of Autism. *Cell* 1–17 (2020) doi:10.1016/j.cell.2019.12.036.
14. Feliciano, P. *et al.* Exome sequencing of 457 autism families recruited online provides evidence for autism risk genes. *NPJ Genom Med* **4**, (2019).
15. De Rubeis, S. *et al.* Synaptic, transcriptional and chromatin genes disrupted in autism. *Nature* **515**, 209–215 (2014).

16. Ishay, R. T. *et al.* Diagnostic Yield and Economic Implications of Whole-Exome Sequencing for ASD Diagnosis in Israel. (2022).
17. Wang, T., Zhao, P. A. & Eichler, E. E. Rare variants and the oligogenic architecture of autism. *Trends in Genetics* 1–9 (2022) doi:10.1016/j.tig.2022.03.009.
18. Doan, R. N. *et al.* Recessive gene disruptions in autism spectrum disorder. *Nat Genet* **51**, (2019).
19. Li, Q. & Wang, K. InterVar: Clinical Interpretation of Genetic Variants by the 2015 ACMG-AMP Guidelines. *Am J Hum Genet* **100**, 267–280 (2017).
20. Houge, G. *et al.* Stepwise ABC system for classification of any type of genetic variant. *European Journal of Human Genetics* **30**, 150–159 (2022).
21. Wilfert, A. B. *et al.* Recent ultra-rare inherited variants implicate new autism candidate risk genes. *Nat Genet* **53**, 1125–1134 (2021).
22. Halvorsen, M. *et al.* Exome sequencing in obsessive–compulsive disorder reveals a burden of rare damaging coding variants. *Nat Neurosci* **24**, 1071–1076 (2021).
23. Dinstein, I. *et al.* The National Autism Database of Israel: a Resource for Studying Autism Risk Factors, Biomarkers, Outcome Measures, and Treatment Efficacy. *Journal of Molecular Neuroscience* **70**, 1303–1312 (2020).
24. Meiri, G. *et al.* Brief Report: The Negev Hospital-University-Based (HUB) Autism Database. *J Autism Dev Disord* **47**, 2918–2926 (2017).
25. McKenna, A. *et al.* The Genome Analysis Toolkit: A MapReduce framework for analyzing next-generation DNA sequencing data. *Genome Res* **20**, 1297 (2010).
26. McLaren, W. *et al.* The Ensembl Variant Effect Predictor. *Genome Biol* **17**, 1–14 (2016).
27. Xavier, A., Scott, R. J. & Talseth-Palmer, B. A. TAPES: A tool for assessment and prioritisation in exome studies. *PLoS Comput Biol* **15**, 1–9 (2019).
28. Fadista, J., Oskolkov, N., Hansson, O. & Groop, L. LoFtool: A gene intolerance score based on loss-of-function variants in 60 706 individuals. *Bioinformatics* **33**, 471–474 (2017).
29. Ng, P. C. & Henikoff, S. SIFT: Predicting amino acid changes that affect protein function. *Nucleic Acids Res* **31**, 3812–3814 (2003).
30. Adzhubei, I., Jordan, D. M. & Sunyaev, S. R. *Predicting functional effect of human missense mutations using PolyPhen-2. Current Protocols in Human Genetics* vol. 2 (2013).
31. Rentzsch, P., Witten, D., Cooper, G. M., Shendure, J. & Kircher, M. CADD: Predicting the deleteriousness of variants throughout the human genome. *Nucleic Acids Res* **47**, D886–D894 (2019).
32. Ioannidis, N. M. *et al.* REVEL: An Ensemble Method for Predicting the Pathogenicity of Rare Missense Variants. *Am J Hum Genet* **99**, 877–885 (2016).
33. Jagadeesh, K. A. *et al.* M-CAP eliminates a majority of variants of uncertain significance in clinical exomes at high sensitivity. *Nat Genet* **48**, 1581–1586 (2016).
34. Samocha, K. E. *et al.* Regional missense constraint improves variant deleteriousness prediction. *bioRxiv* (2017) doi:10.1101/148353.
35. Liu, X., Li, C., Mou, C., Dong, Y. & Tu, Y. dbNSFP v4: a comprehensive database of transcript-specific functional predictions and annotations for human nonsynonymous and splice-site SNVs. *Genome Med* **12**, 1–8 (2020).
36. Abrahams, B. S. *et al.* SFARI Gene 2.0: A community-driven knowledgebase for the autism spectrum disorders (ASDs). *Mol Autism* **4**, 2–4 (2013).
37. Sergouniotis, P. I. *et al.* The role of small in-frame insertions/deletions in inherited eye disorders and how structural modelling can help estimate their pathogenicity. *Orphanet J Rare Dis* **11**, 1–8 (2016).
38. Sallah, S. R. *et al.* Assessing the Pathogenicity of In-Frame CACNA1F Indel Variants Using Structural Modeling. *Journal of Molecular Diagnostics* **24**, 1232–1239 (2022).
39. Iossifov, I. *et al.* De Novo Gene Disruptions in Children on the Autistic Spectrum. *Neuron* **74**, 285–299 (2012).
40. Dong, S. *et al.* De novo insertions and deletions of predominantly paternal origin are associated with autism spectrum disorder. *Cell Rep* **9**, 16–23 (2014).
41. Kopp, N., Amarillo, I., Martinez-Agosto, J. & Quintero-Rivera, F. Pathogenic paternally inherited NLGN4X deletion in a female with autism spectrum disorder: Clinical, cytogenetic,

- and molecular characterization. *Am J Med Genet A* **185**, 894–900 (2021).
42. Wang, K., Li, M. & Hakonarson, H. ANNOVAR: Functional annotation of genetic variants from high-throughput sequencing data. *Nucleic Acids Res* **38**, 1–7 (2010).
  43. McCarthy, D. J. *et al.* Choice of transcripts and software has a large effect on variant annotation. *Genome Med* **6**, (2014).
  44. Trost, B. *et al.* Genomic architecture of autism from comprehensive whole-genome sequence annotation. *Cell* **185**, 4409–4427.e18 (2022).
  45. Tammimies, K. *et al.* Molecular diagnostic yield of chromosomal microarray analysis and whole-exome sequencing in children with autism spectrum disorder. *JAMA - Journal of the American Medical Association* **314**, 595–903 (2015).
  46. Husson, T. *et al.* Rare genetic susceptibility variants assessment in autism spectrum disorder: detection rate and practical use. *Transl Psychiatry* **10**, (2020).
  47. Turner, T. N. *et al.* Genomic Patterns of De Novo Mutation in Simplex Autism. *Cell* **171**, 710–722.e12 (2017).
  48. Leppa, V. M. M. *et al.* Rare Inherited and De Novo CNVs Reveal Complex Contributions to ASD Risk in Multiplex Families. *Am J Hum Genet* **99**, 540–554 (2016).
  49. Krumm, N. *et al.* Excess of rare, inherited truncating mutations in autism. *Nat Genet* **47**, 582–588 (2015).
  50. Lin, B. D. *et al.* The role of rare compound heterozygous events in autism spectrum disorder. doi:10.1038/s41398-020-00866-7.
  51. Tuncay, I. O. *et al.* The genetics of autism spectrum disorder in an East African familial cohort. *Cell Genomics* **3**, 100322 (2023).
  52. Du, Y. *et al.* Nonrandom occurrence of multiple de novo coding variants in a proband indicates the existence of an oligogenic model in autism. *Genetics in Medicine* **22**, 170–180 (2020).
  53. Guo, H. *et al.* Genome sequencing identifies multiple deleterious variants in autism patients with more severe phenotypes. *Genetics in Medicine* **21**, 1611–1620 (2019).