

1 **Title:**

2 **Validation of HER2 status in whole genome sequencing data of breast cancers with AI-driven, ploidy-**
3 **corrected approach**

4

5 **Authors: Wojtaszewska Marzena¹, Stępień Rafał², Woźna Alicja³, Piernik Maciej⁴, Dąbrowski Maciej⁵, Gniot**
6 **Michał⁶, Szymański Sławomir⁷, Socha Maciej⁸, Kasprzak Piotr⁹, Matkowski Rafał^{9,10}, Zawadzki Paweł³**

7

8 **Addresses:**

9 1 Department of Hematology, Frederic Chopin Provincial Specialist Hospital in Rzeszów, Poland

10 2 Biostatistics Group, Department of Genetics, Wrocław University of Environmental and Life Sciences, 25

11 Norwida St.50-375 Wrocław, Poland

12 3 Faculty of Physics, Adam Mickiewicz University, Uniwersytetu Poznańskiego 2 Str, 61-614 Poznan, Poland

13 4 Laboratory of Computer Systems, Institute of Computing Science, Poznan University of Technology, ul.

14 Piotrowo 2, 60-965 Poznan, Poland

15 5 Institute of Human Genetics Polish Academy of Sciences. ul. Strzeszyńska 32 60-479 Poznań

16 6 Department of Hematology and Bone Marrow Transplantation, Poznan University of Medical Sciences,

17 Szamarzewskiego 84, 60-569, Poznan, Poland.

18 7 Pomeranian Medical University in Szczecin, Department of surgical gynecology of adults and adolescent girls,

19 Powstańców Wielkopolskich 72, 70-111 Szczecin, Poland

20 8 Department of Perinatology, Gynecology and Gynecologic Oncology, Faculty of Health Sciences, Nicolaus

21 Copernicus University, Collegium Medicum in Bydgoszcz, Łukasiewicza 1, 85-821 Bydgoszcz, Poland

22 9 Department of Oncology, Wrocław Medical University, L. Hirszfelda 12, 50-367, Wrocław, Poland

23 10 Breast Unit, Wrocław Comprehensive Cancer Center, Plac Hirszfelda 12, 53-413 Wrocław, Poland

24

25 **Corresponding author:**

26 Marzena Wojtaszewska, PhD

27 Central Clinical Hospital of the Internal Affairs in Warsaw,

28 Woloska 137 Str, 02-507 Warsaw, Poland

29 **NOTE: This preprint reports new research that has not been certified by peer review and should not be used to guide clinical practice.**
wojtaszewska@gmail.com

30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57

Running title: Utility of WGS in diagnostics of HER2 status.

Keywords: whole genome sequencing, ERBB2 amplification, HER2 overexpression, next generation sequencing

Abbreviations:

- WGS whole genome sequencing,
- HER2 human epidermal growth factor receptor 2
- BC breast cancer
- CN copy number
- FISH fluorescence in situ hybridization
- IHC immunohistochemistry
- ASCO American Society Cancer of Clinical Oncology
- CAP College American Pathologists
- NGS next-generation sequencing
- CEP Centromere enumeration probe
- AI artificial intelligence
- TCGA The Cancer Genome Atlas
- ICGC International Cancer Genome Consortium
- HMF Hartwig Medical Foundation
- FFPE formalin-fixed, paraffin-embedded
- TNBC Triple-negative breast cancer

Abstract

The HER2 protein overexpression is one of the most significant biomarkers for breast cancer diagnostics, prediction, and prognostics. The availability of HER2-inhibitors in routine clinical practice directly translates into the diagnostic need for precise and robust marker identification.

58 At the brink of the genomic era, multigene next-generation sequencing
59 methodologies slowly take over the field of single-biomarker molecular and cytogenetic tests.
60 However, copy number alterations such as amplification of the HER2-coding *ERBB2* gene, are
61 certainly harder to validate as an NGS biomarker than simple SNV mutations. They are
62 characterized by several compound genomic factors i.a. structural heterogeneity,
63 dependence on chromosome count and genomic context of ploidy. In our study, we tested
64 the approach of using whole genome sequencing instead of NGS panels to robustly and
65 accurately determine HER2 status in clinical setup. Based on the large dataset of 877 breast
66 cancer patients' genomes with curated clinical data and a machine learning approach for
67 optimization of an unbiased diagnostic classifier, we provide a reliable algorithm of HER2
68 status assessment.

69

70 **1. Introduction**

71 Human epidermal growth factor receptor 2 (HER2) is an important biomarker for
72 targeted therapy in breast cancer (BC). Patients with overexpression of the receptor were
73 considered the worst prognosis group before HER2 inhibitors were introduced into clinical
74 practice [1]. Nowadays, the first and second generation of these drugs slow down disease
75 progression, improving the outcomes in HER2-positive subgroup of BCs. Therefore, it is crucial
76 to accurately and precisely pinpoint the HER2-overexpression status [2].

77

78 The molecular mechanism of HER2 overexpression is, in most cases, amplification of
79 a 17q12 chromosome region containing the HER2 coding *ERBB2* gene. The reference method
80 for the assessment of *ERBB2* amplification is immunohistochemistry (IHC) coupled with
81 fluorescence in situ hybridization (FISH) [3]. Currently, diagnostic companies and medical

82 services are beginning to offer novel NGS assays, detecting dozens of actionable biomarkers
83 in a single test. They are trying to incorporate the *ERBB2* copy number (*ERBB2* CN) into their
84 portfolio as well. Unfortunately, *ERBB2* amplification status cannot be easily determined by
85 establishing a simple threshold for negative and positive values, as the genomic context of
86 chromosome 17 copy number and tumour ploidy are interrelated with *ERBB2* CN. Firstly,
87 duplication or triplication of the whole chromosome set or just a subset of chromosomes is a
88 common feature of BC. However, changes in ploidy are seldom associated with
89 overexpression of *ERBB2* gene, as average global transcript levels remain unchanged.
90 Secondly, the isolated deletion or duplication events of chromosome 17 may influence the
91 *ERBB2* transcription [4,5]. The gain of an additional copy of chromosome 17, called polysomy,
92 is correlated with tumour ploidy and is considered its surrogate in the FISH test, but
93 discrepancies between these parameters are in part the reason for inaccuracy in *ERBB2*
94 amplification detection [6].

95 As it is not feasible to determine ploidy in conventional FISH, the ratio between *ERBB2*
96 CN and chromosome 17 centromeric probe (CEP17) CN serves as a diagnostic criterion in dual
97 probe assays, recommended by official ASCO/CAP Clinical Practice Guidelines for diagnostics
98 of HER2 in breast cancer patients [3].

99 Whole genome sequencing (WGS) on the other hand is capable of acquiring absolute
100 *ERBB2* copy number, centromere 17 CN and mean ploidy of tumour cells simultaneously.
101 Moreover, WGS can estimate the tumour content of the sample, providing quality control of
102 the material. As WGS is based on PCR-free methodology, it preserves the original proportions
103 of DNA fragments, in contrast to enrichment or PCR-based NGS panels, which may distort the
104 original proportions of DNA fragments and skew the quantification [7].

105 The purpose of this study was to determine the feasibility of accurately distinguishing
106 between HER2-positive and HER2-negative cases of BC based on matched tumour-normal
107 WGS. Up to date, there have been only a few studies evaluating the clinical utility of NGS
108 testing of *ERBB2* gene status, including WGS method [8–12]. Some of them directly address
109 the clinical need to verify the relevance of their findings for patient management, reporting
110 the overall concordance between IHC/FISH and NGS at about 90% level.

111 Our study operates on the large population-based cohort of 877 BCs from publicly
112 available databases, supplied with the final clinical HER2 status based on ASCO/CAP guidelines
113 and targeted treatment information, which serves to validate metastatic samples status. We
114 analyzed the whole cohort of patients, aiming to establish the criteria for WGS *ERBB2* status
115 assessment as close to the golden standard as possible, optimized for both sensitivity and
116 precision with a bias-free machine learning approach. We also provide proof-of-concept that
117 genomic data acquired on different platforms with different chemistry yield sufficiently
118 uniform results for molecular diagnostics of *ERBB2* amplification by WGS.

119

120 **2. Materials and methods**

121

122 **2.1 Sample choice**

123 Matched tumor-normal genomes from 877 breast cancer patients sequenced within
124 three large Genomic Consortia (119- International Cancer Genome Consortium, 70- The
125 Cancer Genome Atlas, 688- Hartwig Medical Foundation; HMF) were downloaded from
126 controlled-access databases after meeting formal criteria [10,13–15]. The samples were
127 sequenced using a low PCR amplification or PCR-free library preparation protocols and paired-
128 end 100-150 base pair Illumina reads with 350- 550 base pair insert size (for details, see the

129 Supplementary Table 1). For analyses of primary tumour samples, we included the datasets
130 with clinical HER2 status described as positive or negative, according to ASCO/CAP guidelines
131 2007-2018 (depending on the year of the original study was conducted, see the
132 Supplementary Table 1). For metastatic/advanced tumour samples from HMF database,
133 metadata on HER2 status were available only for primary tumours, the IHC/FISH status for
134 sequenced sample from second biopsy was not provided. Because of the high rate of
135 conversion from HER2-negative to HER2-positive status (and vice versa) during the cancer
136 evolution [8,9], in metastatic cancers we have taken into consideration also the patients'
137 treatment metadata and discarded all samples, for which treatment history (pre- and post-
138 biopsy) was discordant with initial HER2 status (eg. if trastuzumab was included in any line of
139 treatment even though HER2 status was reported negative). For details on discarded samples
140 see Supplementary Data.

141 As there were no new tissue/DNA/RNA samples processed, the written consent of
142 each subject is in possession of data providers. The primary data were collected in accordance
143 with the standards set by the Declaration of Helsinki and the highest data security standards
144 of ISO 27001. There was no need of acquiring approval from a local ethics committee as no
145 actual tumour samples were used.

146

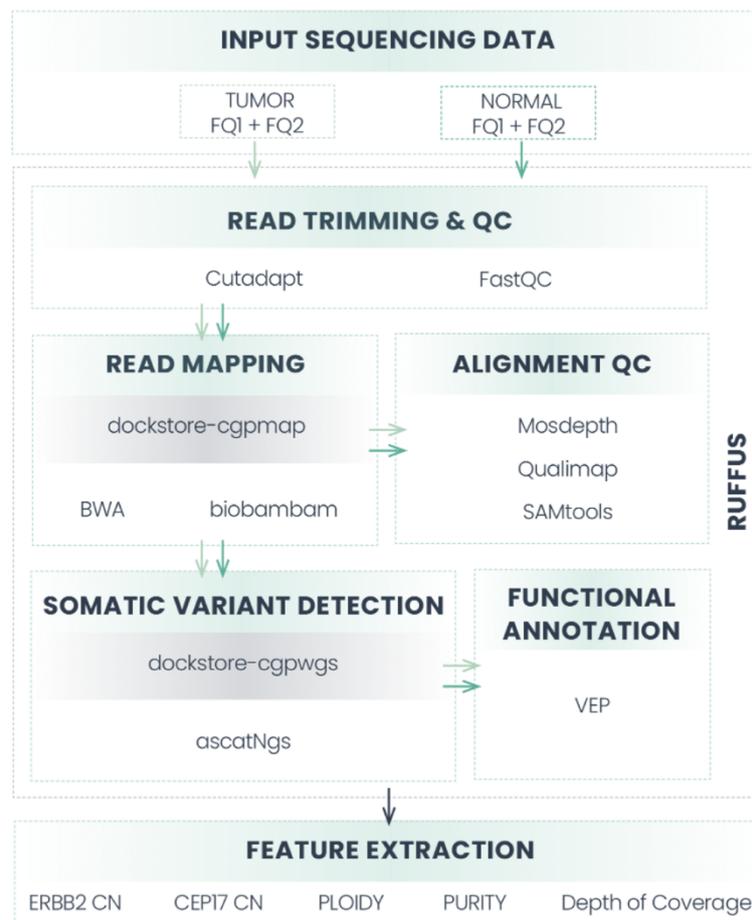
147 **2.2 Whole-genome data processing**

148 The samples were analyzed using publicly available, open-source software embedded
149 within an in-house pipeline (Figure 1) implemented using Ruffus [16]. The analysis started
150 with FASTQ files extraction from the downloaded BAM/CRAM files using Broad Institutes'
151 Picard [17]. Tumour samples with coverage exceeding 75x were downsampled with Seqtk
152 v1.3-r106 [18] to approx. 60x mean coverage. Next, all reads were trimmed using cutadapt

153 v2.10 [19] and mapped to the GRCh37 genome using Sanger's Cancerit CGPMAP pipeline
154 v3.0.0 [20]. Samples with uniquely-mapped read coverage below 20x for either tumor or
155 normal genomes were excluded from the analysis [21,22]. Mean tumour samples' coverage
156 across all datasets after downsampling was 48x, reference blood/EBV-transformed
157 lymphocyte samples' mean coverage was 36x (detailed data are provided in Supplementary
158 Table 2 and Supplementary Figure 1).

159 Variant calling was performed using Sanger's Cancerit CGPWGS pipeline v2.0.1 [20],
160 and specifically copy number variants, purity, and ploidy were identified with ascatNgs [23].
161 Identified variants were annotated using Ensembl VEP v102 [24].

162



163

164

165 Figure 1: The summary of the in-house pipeline used for data extraction and processing.

166

167

168

169 **2.3 Analyzed parameters and method validation**

170 In the study, we used only clinical data on HER2 status according to ASCO/CAP
171 recommendations or, in the case of HMF metastatic/advanced tumours, the presence of
172 targeted treatment with HER2 inhibitors, which was indicative of the confirmed presence of
173 HER2 expression. Based on ASCAT copy number alteration calling, *ERBB2*
174 (NC_000017.10:37844167_37886679) and uniquely mapped 8250 bp sequence adjacent to
175 CEP17 (NC_000017.10:22236000_22244250) copy numbers were extracted along with ploidy
176 and purity estimation for all the tumour samples. The data were used to create 3 features for
177 HER2 status assessment: absolute *ERBB2* CN, *ERBB2_CN*-n (ploidy-adjusted *ERBB2* CN), and
178 *ERBB2_CN/CEP17_CN* ratio. Based on these features, a machine learning-based classifier was
179 constructed, which determined the best approach for HER2 status discrimination. 614
180 samples from the datasets were used as a training set, the remaining 264 samples served as
181 a validation hold out set for the classifier and were not analyzed *a priori*.

182 A decision tree-based classifier was chosen after comparing the effectiveness of
183 logistic regression, random forest and decision tree models. Decision tree outperformed
184 other classifiers in terms of accuracy and interpretability.

185 For the decision-tree-based modelling, the discovery cohort was randomly split into a
186 training (75%) and a test set (25%). The model was constructed on 3 aforementioned features
187 and trained on the training set. Since the number of samples in IHC/FISH HER2-positive and
188 negative groups was unbalanced (there were almost eight times less HER2+ samples than

189 negative), we added class weights (8:1) to minimize the bias. After constructing the model we
190 measured its performance on 264 samples from the validation set. We used accuracy,
191 precision, and recall along with the F1-score. Cohen's Kappa score was estimated to evaluate
192 the non-randomness of classification.

193 To show how each of the 3 features influences the classifiers performance alone, we
194 have established the same parameters independently for each of them as well and compared
195 all the approaches with random data classification methods (Figure 2).

196 To further test the validity of our results, we decided to evaluate whether differences
197 in tumour purity, heterogeneity of ploidy, or differences in mean depth of coverage had any
198 deteriorative effects on the correctness of the results (see Supplementary Tables 2-3 and
199 Supplementary Figures 1-2). For these experiments, we divided the samples into two near-
200 equinumerous groups for each comparison and evaluated the differences in the tests'
201 performance. Finally, we determined the overall predictive value, PPV, and NPV with
202 confidence intervals for the whole dataset of 877 genomes.

203

204 **3. Results**

205 In the analyzed dataset, 159 patients were categorized as triple-negative breast cancer
206 (TNBC) (18%), among HER2-negative patients ER+/HER- accounted for 599 (88%). 110 (13%)
207 of samples were identified by clinical testing as HER2-positive, among them: 74 ER+/HER2+
208 (8%), 36 ER-/HER2+ (4%). For 8 patients' ER status was unavailable.

209 HER2 positivity was slightly underrepresented in favour of TNBC in comparison with
210 statistics for the caucasian population (18%) which may be an accidental or sampling bias
211 related to genomic consortia's sample collection process, or an effect of discarding datasets
212 with incomplete clinical data.

213 The decision tree machine learning approach has demonstrated the best
214 discrimination between HER2-positive and negative cases based on a single parameter, ploidy
215 corrected-ERBB2 CN with a threshold of 2.265 (fig.2). The decision tree algorithm was
216 evaluated in 3-fold cross-validation repeated 10 times to estimate the mean value and
217 standard deviation for each metric. The results were as follows: accuracy = 96,7% (+- 0.87%),
218 precision = 86% (+- 5%) , recall = 89% (+- 6%), Cohen's Kappa = 85% (+- 3.7%) and F1 = 87%
219 (+- 3%). A high value of Cohen's Kappa strongly indicates that our model classifies samples in
220 a non-random fashion.

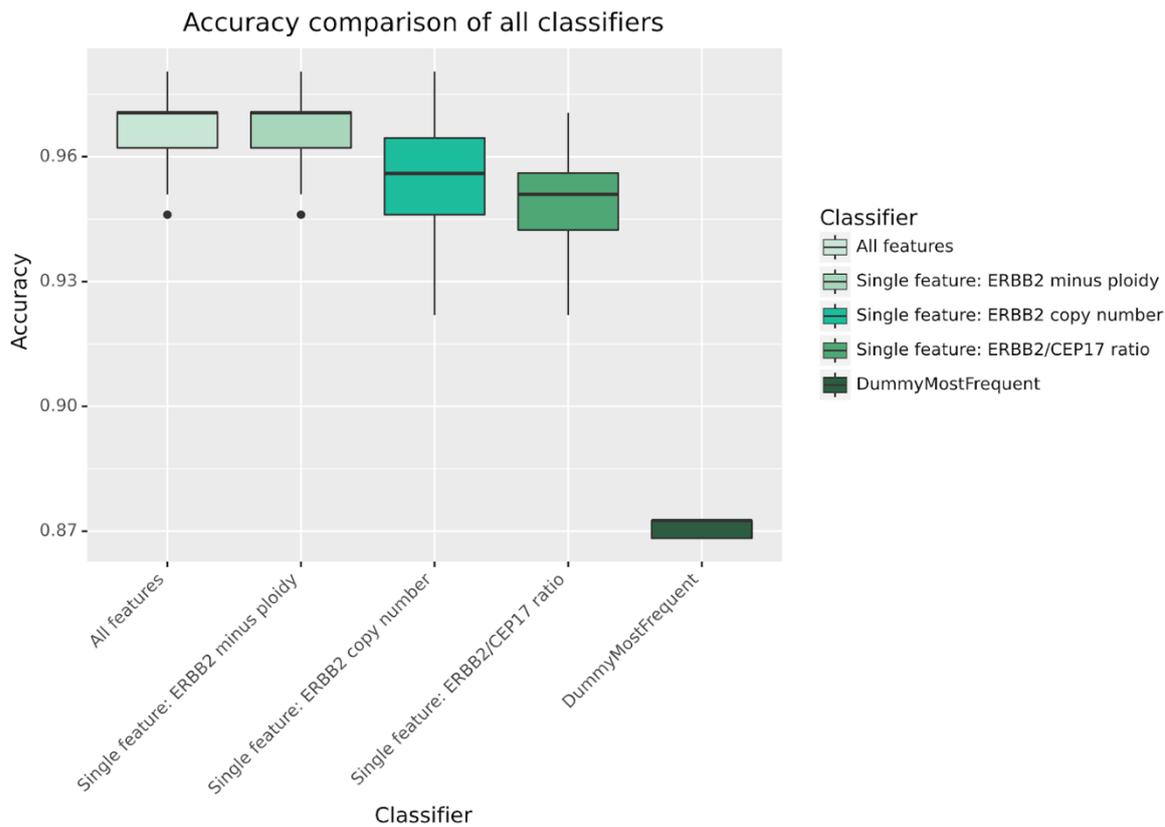
221
222 The learning curve displayed no further improvement with sample numbers
223 exceeding 150 instances, therefore we believe the results display the best reflection of the
224 biological phenomenon of HER2 amplification we could extract from genomic data.
225 Moreover, Principal Component Analysis of the dataset (fig.3) has shown a very good and
226 robust separation of data into two groups, representing differences in HER2 status.

227 As data distribution across depths of coverage, tumour purities, and ploidies were not normal,
228 we decided to compare the accuracy distributions for these parameters with the Wilcoxon
229 signed-rank test. The evaluation of results across data coverages has shown no significant
230 differences ($p>0.05$) between groups.

231 The comparison of low vs high purity also has not yielded significant differences ($p>0.05$).
232 However, there is a significant decrease in mean accuracy of the test from 0.97 to 0.94,
233 dependent on increased tumour ploidy above two ($p=5.1\times 10^{-6}$) (Figure 4).

234

235

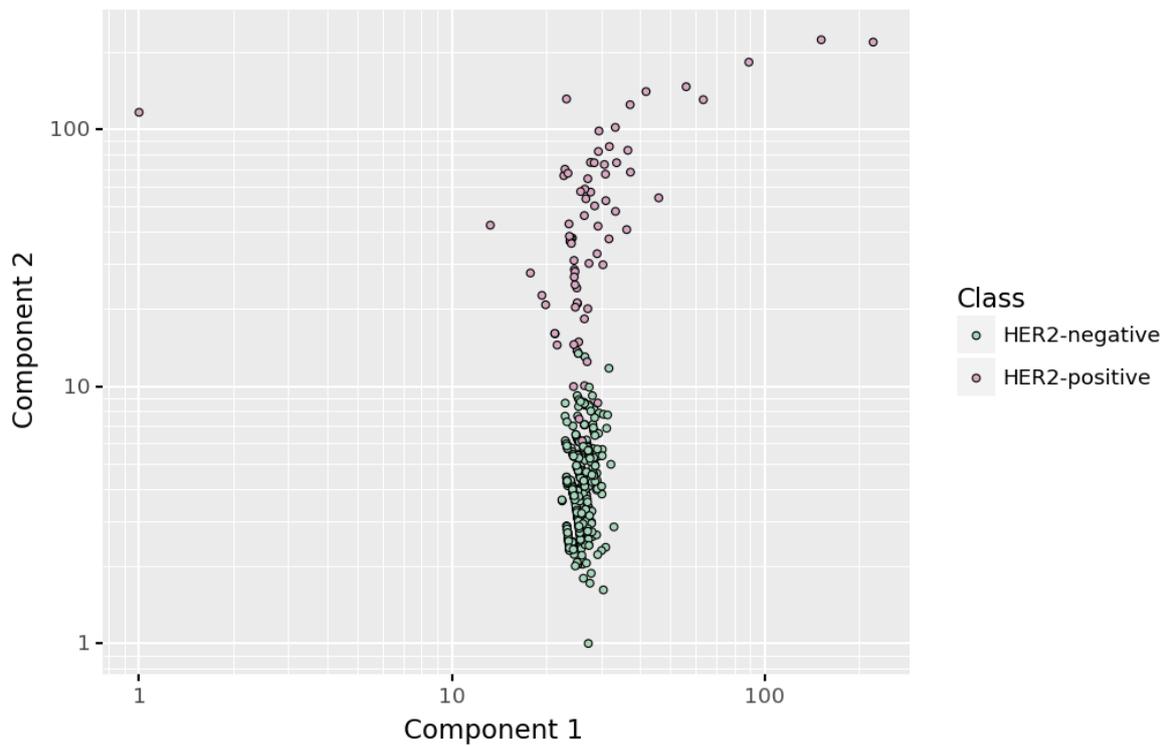


236

237 Figure 2: Accuracy comparison between 3 features used to determine HER2 amplification

238 status in WGS data.

239

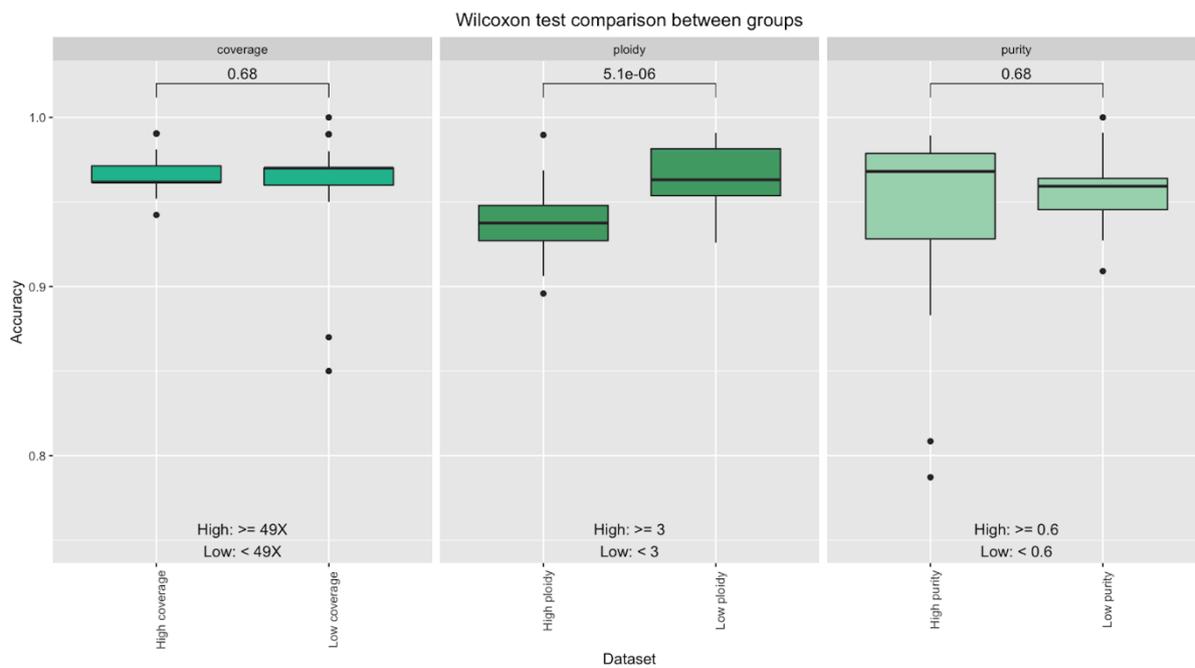


240

241 Figure 3: Principal Component Analysis of the dataset with 6 features: purity, Ploidy, ERBB2

242 CN, CEP17 CN, ERBB2 CN / CEP17 CN ratio, ploidy-corrected ERBB2 CN.

243



244

245 Figure 4: Wilcoxon test comparison of means between distributions of accuracies in: A) High
246 vs low coverage data (threshold 49X), B) High vs low ploidy data (threshold 3), C) High and
247 low purity data (threshold 0.6).

248

249 The analytical validation of the ploidy-corrected ERBB2_CN method gave the overall
250 diagnostic sensitivity of 92.92% (95%CI 86.53-96.89%) and specificity of 97.91% (95%CI 96.62-
251 98.8%) (Table 1).

A

IHC/FISH		
Non-Amplified	Amplified	WGS RESULTS
748	8	Non-Amplified
16	105	Amplified

B

Results	Point estimate	Lower CI	Upper CI
Sensitivity	92.92%	86.53%	96.89%
Specificity	97.91%	96.62%	98.80%
Accuracy	97.16%	95.83%	98.15%
Positive Predictive Value	88.67%	82.78%	92.73%
Negative Predictive Value	98.74%	97.57%	99.35%

252

253 Table 1: Analytical validation of the whole genome sequencing ploidy-corrected ERBB2 copy
254 number. For details on samples used see Supplementary Data.

255

256 4. Discussion

257 Decreasing next-generation sequencing prices and increasing availability of this
258 technology in medical practice have encouraged the transition from conventional cytogenetic
259 and molecular methods to NGS in oncology. However, the evidence on the reliability of NGS
260 for clinical use in copy number detection is still very limited. As the HER2 protein is one of

261 the most significant biomarkers for breast cancer diagnostics, prediction, and prognostics,
262 there were several attempts to show the applicability of NGS techniques in this indication.

263 The largest analytical validation study was conducted by Memorial Sloan Kettering on
264 their proprietary MSK-IMPACT Assay [8]. This hybrid-capture based panel NGS test was
265 analyzed in 213 BC samples and evaluated in a clinical setting on further 599 BCs. The cutoff
266 for positive result was established based solely on *ERBB2* CN, adjusted to background and
267 normal signal of diploid genomes (defined as ‘fold change’, $FC=1.5$). The group reported 95%
268 specificity and 100% sensitivity on >10% of tumour content, with IHC/FISH evaluated by
269 newest, 2018 guidelines and a dual-probe FISH assay [8]. Last year, a continuation of the study
270 exploited the borderline cases with excellent concordance [12]. Several other studies have
271 also proven the clinical value of panel NGS for HER2 testing in breast cancer and other solid
272 tumours, using the same strategy of fold change determination, using either Illumina [1–4] or
273 Ion Torrent methodology [5].

274 On the other hand, data on clinical whole genome sequencing utility is scarce. There
275 have only been two small clinical validation studies with direct comparison to the orthogonal
276 methods. The first, released by Hartwig Medical Foundation, was a part of a WGS pan-cancer
277 validation study. The *ERBB2* status was evaluated on only 16 samples with overall
278 concordance of 93%. HMF group compared ploidy and chromosome 17 CN with absolute CN
279 of *ERBB2* but did not draw any conclusions due to the small sample size [5]. The second,
280 performed by King’s College Hospital in London, was performed on 145 BC samples with only
281 27 positives for HER2. With the 4 samples discrepant, the sensitivity in the UK cohort was 88%
282 and specificity 98% [25].

283

284 We attempted to systematically determine the criteria for whole genome sequencing
285 of ERBB2 CN in matched-normal tumour samples. Our strategy was to gather publicly
286 available breast cancer datasets with reliable clinical metadata and analyze them uniformly
287 with minimal 20x depth of coverage. Our machine learning approach, based on Decision Tree
288 classifier, objectively captured the superiority of ploidy-corrected ERBB2 CN over ERBB2
289 CN/CEP17 CN ratio and absolute ERBB2 CN for HER2 status evaluation in breast cancer by
290 WGS. To measure the test's reliability, we used Cohen's kappa coefficient. The high value of
291 85% rules out the possibility of the data agreement occurring by chance.

292 This is, to the best of our knowledge, the first, the largest, and the most objective
293 study of its kind, utilizing artificial intelligence and machine learning approaches for
294 establishing diagnostic criteria. The optimal CAP/AAP guidelines for IHC/FISH testing took 11
295 years to refine, because it involved a series of consecutive evaluations and quality control
296 rounds of diagnostic parameters which probably could have been done nowadays in an AI-
297 based manner more robustly and quickly [6,7]. First proof-of-concept AI-based solutions for
298 robust FISH and IHC assessment are already tested in clinical setup [26,27].

299 The AI/ML approach is an emerging field of medicine, improving the efficiency of
300 pathomorphological assessment [28] radiology [29] and clinical chemistry [30]. In the field of
301 breast cancer diagnostics, the genomics and transcriptomics is being applied to distinguish
302 between intrinsic BC subtypes with different prognosis [31], identify new potential
303 biomarkers or repurpose the existing. These strategies may only be used in the clinical setting
304 after well-planned validation, showing concordance and stability of the test. Our results prove
305 that WGS is a reliable method for HER2 clinical diagnostics, and it may be implemented as a
306 standalone test or in combination with IHC instead of FISH or other NGS-based methods in
307 routine practice. With diagnostic sensitivity of 92.92% and specificity of 97.91% determined

308 on unselected and heterogeneous groups of patients, we conclude that the technology is
309 mature and ready for prospective, multicenter analytical and clinical validation.

310 Our results do not deviate relevantly from those reported by other groups focused on
311 HER2 NGS testing, however the diagnostic sensitivity is still not optimal. We suspect that
312 heterogenous evaluation of IHC/FISH results, made on the basis of different issues of
313 ASCO/AAP guidelines, may have contributed to the discrepancy. There was also great
314 heterogeneity of the whole genome sequencing raw data, acquired on different equipment
315 by different genomic consortia. There were also serious differences in tumour sample
316 collection, DNA extraction, and library preparation methods. All these preanalytical and
317 analytical factors must have contributed to the greater variation in HER2 results than in the
318 single-facility method with uniform IHC/FISH evaluation methodology and a single laboratory
319 protocol for sample management. Even so, the WGS method exhibits superb robustness and
320 effectiveness, which is a great advantage, allowing for a low-cost external, even world-wide
321 quality control assessment program to be held out in the near future.

322 Other factors contributing to slightly lower analytical sensitivity are changes in HER2
323 status, which could have occurred in metastatic tumours from HMF dataset. In these instances,
324 we couldn't directly evaluate the correctness of IHC/FISH data, because they came from the
325 primary biopsy, not the biopsy corresponding with the sample used for WGS. The shift
326 between IHC positive and negative status is reported in up to 11.5% of HER2-negative cancers
327 (conversion to HER2 positive) and in 37% of those initially positive (conversion to HER2
328 negative in presence of selective pressure of trastuzumab) [8,9].

329 Some of the discrepancies may have come from tumour subclonality, which is a
330 common serious diagnostic issue. The signal from a small proportion of HER2 amplified cells
331 may be below the resolution of whole genome sequencing at 30-60x depth of coverage [10].

332 The spatial intra-tumour heterogeneity may have also contributed to false
333 negative/positive results when there were differences in sampling location between tissue
334 collected for FFPE blocks and WGS (e.g., different distant metastases are sampled).

335 In addition, overexpression of HER2 is not always *ERBB2*-amplification based as about
336 5% of non-amplified tumours exhibit high overexpression. Even though there is currently no
337 genomic background of this phenomenon known, whole genome sequencing could
338 potentially detect alterations in HER2 regulatory pathways leading to overexpression, which
339 could further improve WGS diagnostic power.

340

341

342 **5. Conclusion**

343 We provide evidence that the *ERBB2* status can be reliably determined by WGS
344 methodology which may be included into a comprehensive test for breast cancer diagnostics.
345 The 20% of tumour purity and 30x depth of coverage are sufficient to ensure good quality of
346 genomic data in most instances. Given good concordance of a whole genome sequencing with
347 routinely used methods, we suggest that assessment by a WGS method may be an alternative
348 to other NGS-based methods as well as FISH-based diagnostic tools. Hence, it should be
349 subjected for evaluation by ASCO/CAP in the future updates of the HER2 testing
350 recommendations. In our work, we have also proven that short-reads WGS technology bears
351 great potential for establishing a harmonized global quality assessment program for *ERBB2*
352 detection, as the outputs of heterogeneous data gathered from 4 genomic consortia show a
353 high degree of concordance between methodologies and pipelines.

354

355 **Data accessibility**

356 The data that support the findings of this study are openly available in the following
357 repositories:
358 Hartwig Medical Foundation at <https://www.hartwigmedicalfoundation.nl/> which was
359 acquired under data request number DR-169.
360 International Cancer Genome Consortium at <https://dcc.icgc.org/pcawg/> which was acquired
361 under data request number DACO-6030.
362 The Cancer Genome Atlas data was acquired via dbGaP platform (project
363 phs000178.v11.p8) at [https://www.ncbi.nlm.nih.gov/projects/gap/cgi-](https://www.ncbi.nlm.nih.gov/projects/gap/cgi-bin/study.cgi?study_id=phs000178.v11.p8)
364 [bin/study.cgi?study_id=phs000178.v11.p8](https://www.ncbi.nlm.nih.gov/projects/gap/cgi-bin/study.cgi?study_id=phs000178.v11.p8) , under data request number #86794-3.
365 Secondary data that supports the findings of this study that was generated by the Authors
366 are available in the supplementary material of this article.

367

368 **Conflict of Interest**

369 Alicja Woźna and Paweł Zawadzki are share owners in the company MNM Bioscience Inc.
370 16192 Coastal Highway, Lewes, DE 19958. Other authors declare no conflict of interests.

371

372 **Acknowledgements**

373 We thank the international genomic consortia, their researchers and public funding
374 institutions who have generously donated the genomic data and clinical metadata for this
375 project.

376 This publication and the underlying study have been made possible by the data that the
377 Hartwig Medical Foundation and the Center of Personalized Cancer Treatment (CPCT) have
378 made available. The results published here are also in part based upon data generated by
379 the TCGA Research Network: <https://www.cancer.gov/tcga>.

380 This work was also supported by data obtained from ICGC Breast Cancer Working Group and
381 The Cancer Genome Atlas.

382 Importantly, we thank the patients and their families for their participation in the
383 genomic projects and the opportunity to use their clinical and genomic data for
384 fundamental cancer research. We would also like to acknowledge the work of all clinical
385 staff gathering samples and medical records. Lastly, we thank all members of Adam
386 Mickiewicz University, Poznan Supercomputing and Network Centre and MNM Diagnostics
387 staff for help with genomic data processing and general support.

388

389 **Contribution**

390 Wojtaszewska M- first author, study design, manuscript preparation, data analysis.

391 Wozna A- genomic data administration, study design.

392 Piernik M -pipeline development, bioinformatics support, design of the computational
393 framework for machine learning approach.

394 Dabrowski M- data refinement, supporting in writing and corrections of the manuscript.

395 Gniot M- statistics, genomic data administration, manuscript corrections.

396 Szymański S - curation of medical data.

397 Socha M – statistics, pipeline refinement, data analysis

398 Kasprzak P expertise on molecular HER2 status assessment, clinical data evaluation.

399 Matkowski R- expertise on molecular HER2 status assessment, clinical data evaluation.

400 Zawadzki P^{1,3} - supervision of the project, review of methodology and manuscript.

401

402 **References**

- 403 1. King CR, Kraus MH, Aaronson SA (1985) Amplification of a novel v-erbB-related gene in
404 a human mammary carcinoma. *Science* 229: 974–976.
- 405 2. Kunte S, Abraham J, Montero AJ (2020) Novel HER2–targeted therapies for HER2–
406 positive metastatic breast cancer. *Cancer* 126: 4278–4288.
- 407 3. Wolff AC, Hammond MEH, Allison KH, et al. (2018) Human Epidermal Growth Factor
408 Receptor 2 Testing in Breast Cancer: American Society of Clinical Oncology/College of
409 American Pathologists Clinical Practice Guideline Focused Update. *JCO* 36: 2105–2122.
- 410 4. Hansen T v O, Vikesaa J, Buhl SS, et al. (2015) High-density SNP arrays improve
411 detection of HER2 amplification and polyploidy in breast tumors. *BMC Cancer* 15.
- 412 5. Niu D, Li L, Yu Y, et al. (2020) Evaluation of Next Generation Sequencing for Detecting
413 HER2 Copy Number in Breast and Gastric Cancers. *Pathol Oncol Res* 26: 2577–2585.
- 414 6. Halilovic A, Verweij DI, Simons A, et al. (2019) HER2, chromosome 17 polysomy and
415 DNA ploidy status in breast cancer; a translational study. *Sci Rep* 9: 11679.
- 416 7. Arora K, Shah M, Johnson M, et al. (2019) Deep whole-genome sequencing of 3 cancer
417 cell lines on 2 sequencing platforms. *Sci Rep* 9: 19123.
- 418 8. Ross DS, Zehir A, Cheng DT, et al. (2017) Next-Generation Assessment of Human
419 Epithelial Growth Factor Receptor 2 (ERBB2) Amplification Status. *J Mol Diagn* 19: 244–
420 254.
- 421 9. Roepman P, Bruijn E de, Lieshout S van, et al. (2021) Clinical Validation of Whole
422 Genome Sequencing for Cancer Diagnostics. *The Journal of Molecular Diagnostics* 0.
- 423 10. Nik-Zainal S, Davies H, Staaf J, et al. (2016) Landscape of somatic mutations in 560
424 breast cancer whole-genome sequences. *Nature* 534: 47–54.

- 425 11. Pfarr N, Penzel R, Endris V, et al. (2017) Targeted next-generation sequencing enables
426 reliable detection of HER2 (ERBB2) status in breast cancer and provides ancillary
427 information of clinical relevance. *Genes, Chromosomes and Cancer* 56: 255–265.
- 428 12. Hoda RS, Bowman AS, Zehir A, et al. (2021) Next-generation assessment of human
429 epidermal growth factor receptor 2 gene (ERBB2) amplification status in invasive
430 breast carcinoma: a focus on Group 4 by use of the 2018 American Society of Clinical
431 Oncology/College of American Pathologists HER2 testing guideline. *Histopathology* 78:
432 498–507.
- 433 13. Data Access Compliance Office (DACO) Available from: <https://daco.icgc.org/>.
- 434 14. DNA database van uitgezaaide tumoren, inclusief klinische data Stichting Hartwig
435 Medical Foundation. Available from: [https://www.hartwigmedicalfoundation.nl/data-](https://www.hartwigmedicalfoundation.nl/data-catalogue/)
436 [catalogue/](https://www.hartwigmedicalfoundation.nl/data-catalogue/).
- 437 15. Campbell PJ, Getz G, Korbel JO, et al. (2020) Pan-cancer analysis of whole genomes.
438 *Nature* 578: 82–93.
- 439 16. Goodstadt L (2010) Ruffus: a lightweight Python library for computational pipelines.
440 *Bioinformatics* 26: 2778–2779.
- 441 17. Broad Institute (2019) <http://broadinstitute.github.io/picard/>, Picard toolkit, 2019.
442 Available from: <http://broadinstitute.github.io/picard/>.
- 443 18. SeqTk (2021) <https://github.com/lh3/seqtk>, 2021. Available from:
444 <https://github.com/lh3/seqtk>.
- 445 19. Martin M (2011) Cutadapt removes adapter sequences from high-throughput
446 sequencing reads. *EMBnet j* 17: 10.
- 447 20. Sanger Institute (2021) [cancerit/dockstore-cgppmap](https://github.com/cancerit/dockstore-cgppmap), Cancerit, 2021. Available from:
448 [cancerit/dockstore-cgppmap](https://github.com/cancerit/dockstore-cgppmap).

- 449 21. Cibulskis K, Lawrence MS, Carter SL, et al. (2013) Sensitive detection of somatic point
450 mutations in impure and heterogeneous cancer samples. *Nat Biotechnol* 31: 213–219.
- 451 22. Alioto TS, Buchhalter I, Derdak S, et al. (2015) A comprehensive assessment of somatic
452 mutation detection in cancer using whole-genome sequencing. *Nat Commun* 6: 10001.
- 453 23. Raine KM, Loo P, Wedge DC, et al. (2016) ascatNgs: Identifying Somatically Acquired
454 Copy-Number Alterations from Whole-Genome Sequencing Data. *Current Protocols in*
455 *Bioinformatics* 56.
- 456 24. McLaren W, Gil L, Hunt SE, et al. (2016) The Ensembl Variant Effect Predictor. *Genome*
457 *Biol* 17: 122.
- 458 25. Echejoh G, Liu Y, Chung-Faye G, et al. (2020) Validity of whole genomes sequencing
459 results in neoplasms in precision medicine. *J Clin Pathol*.
- 460 26. Rawat RR, Ortega I, Roy P, et al. (2020) Deep learned tissue “fingerprints” classify
461 breast cancers by ER/PR/Her2 status from H&E images. *Sci Rep* 10.
- 462 27. Zakrzewski F, Back W de, Weigert M, et al. (2019) Automated detection of the HER2
463 gene amplification status in Fluorescence in situ hybridization images for the
464 diagnostics of cancer tissues. *bioRxiv* 490052.
- 465 28. Bera K, Schalper KA, Rimm DL, et al. (2019) Artificial intelligence in digital pathology —
466 new tools for diagnosis and precision oncology. *Nat Rev Clin Oncol* 16: 703–715.
- 467 29. Montagnon E, Cerny M, Cadrin-Chênevert A, et al. (2020) Deep learning workflow in
468 radiology: a primer. *Insights into Imaging* 11: 22.
- 469 30. Machine Learning Takes Laboratory Automation to the Next Level *Journal of Clinical*
470 *Microbiology*. Available from: [https://journals.asm.org/doi/abs/10.1128/JCM.00012-](https://journals.asm.org/doi/abs/10.1128/JCM.00012-20)
471 20.

472 31. Dawson S-J, Rueda OM, Aparicio S, et al. (2013) A new genome-driven integrated
473 classification of breast cancer and its implications. *EMBO J* 32: 617–628.

474

475 **Supporting information**

476 **Supplementary Information:** Detailed methodologies of Whole Genome Sequencing used by main
477 Genomic Consortia and input data heterogeneity.

478

479 Supplementary Table 1: Differences in DNA preparation, library preparation and sequencing
480 technologies across Genomic Consortia.

481 Supplementary Table 2: Coverage heterogeneity across datasets

482 Supplementary Table 3: Purity and ploidy heterogeneity across datasets

483 Supplementary Figure 1: Tumour (A) and Normal (B) coverage heterogeneity across
484 Consortia.

485 Supplementary Figure 2: Purity (A) and ploidy (B) heterogeneity across Consortia.

486

487 **Supplementary Data:** List of samples used in the study classified by IHC/FISH and WGS and
488 additional list of discarded samples.