

Development of an evidence-based algorithm that optimizes sensitivity and specificity in WES-based diagnostics of a clinically heterogeneous patient population¹

Peter Bauer¹, Maximilian E.R. Weiss¹, Omid Paknia¹, Martin Werber¹, Aida M. Bertoli-Avella¹, Zafer Yüksel¹, Krishna Kumar Kandaswamy¹, Malgorzata Bochinska¹, Gabriela E. Opera¹, Shivendra Kishore¹, Volkmar Weckesser¹, Ellen Karges¹, Arndt Rolfs^{1,2}

¹ CENTOGENE AG, Rostock, Germany
² Albrecht-Kossel-Institute for Neuroregeneration, Rostock, Germany

Sensitivity of whole exome sequencing (WES) is not well-defined. We applied very low thresholds in WES-associated variant calling to also enable investigation of candidate variants that are commonly neglected. As Sanger sequencing revealed ~5% of these to be true positives (Figure 1), we considered numerous variant-specific features (Tables 1 and 2) for the development of a robust predictor for true and false positives. Iterative rounds of receiver operating characteristic (ROC) curve generation identified features and corresponding thresholds with high predictive value (Figure 2). In a corresponding workflow for our data, 91.3% of variants can be pre-classified with 100% specificity and 99.8% sensitivity, while the remaining 8.7% of variants require confirmatory Sanger sequencing (Figure 3).

1. Stringent thresholds during variant calling sacrifice analytical sensitivity

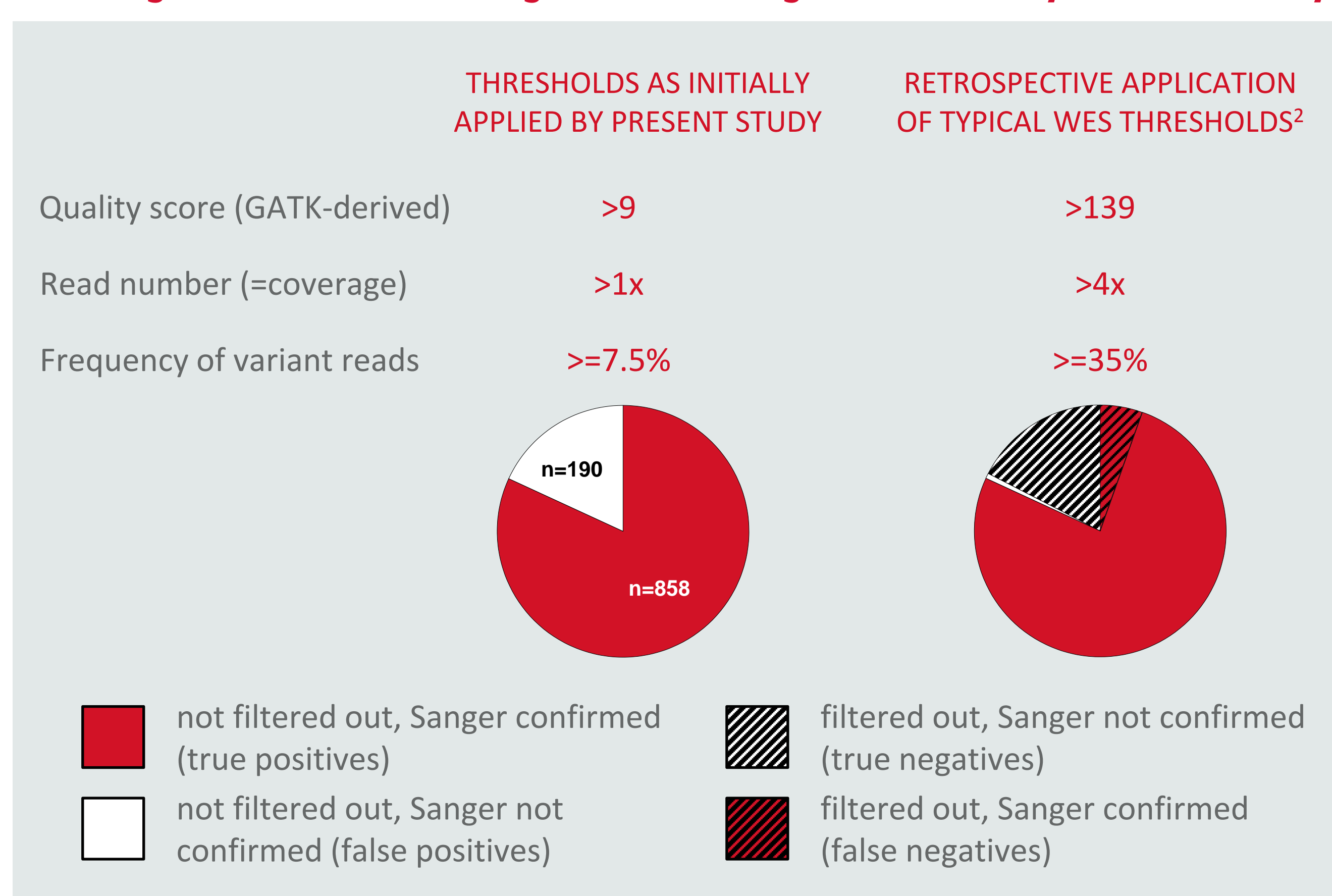


Figure 1: Impact of filtering parameters used for variant calling.

2. Numerous variant-specific features have potential value for predicting true and false positives

FEATURE	BACKGROUND	MEAN (+/- SEM) FOR TRUE POSITIVES	MEAN (+/- SEM) FOR FALSE POSITIVES	P-VALUE (TWO-SIDED STUDENT'S T-TEST)
quality	unit-less score (GATK-derived)	2415 (+/-78.8)	43 (+/-3.1)	5.2 x 10 ⁻¹⁸
read number	count (=coverage)	128 (+/-3.1)	26 (+/-3.1)	3.9 x 10 ⁻²¹
variant reads	count	86 (+/-2.4)	7 (+/-0.8)	8.6 x 10 ⁻²¹
reference reads	count	42 (+/-1.7)	19 (+/-2.4)	2.3 x 10 ⁻⁴
frequency	variant reads divided by read number	0.68 (+/-0.009)	0.32 (+/-0.012)	2.6 x 10 ⁻²⁵
GC-content	GC-content in +/- 100 bp neighbourhood	0.52 (+/-0.004)	0.059 (+/-0.010)	1.2 x 10 ⁻⁵

Table 1: Candidate analogous features.

FEATURE	BACKGROUND	COUNT (FRACTION) AMONGST 858 TRUE POSITIVES	COUNT (FRACTION) AMONGST 190 FALSE POSITIVES	P-VALUE (TWO-SIDED FISHER'S EXACT TEST)
zygosity	homo-/hemizygous (vs. heterozygous)	347 (40%)	5 (3%)	2.2 x 10 ⁻¹³
location	exonic (vs. intronic)	791 (92%)	153 (81%)	6.4 x 10 ⁻⁶
type	indels (vs. SNVs)	314 (37%)	33 (17%)	1.5 x 10 ⁻³
prediction	pathogenic (vs. likely pathogenic or VUS)	202 (24%)	10 (5%)	2.1 x 10 ⁻⁴
homopolymer	from n>3 homopolymer (vs. not)	78 (9%)	11 (6%)	1.5 x 10 ⁻¹

Table 2: Candidate digital features.

3. Iterative analyses reveal robust predictors for Sanger confirmation status

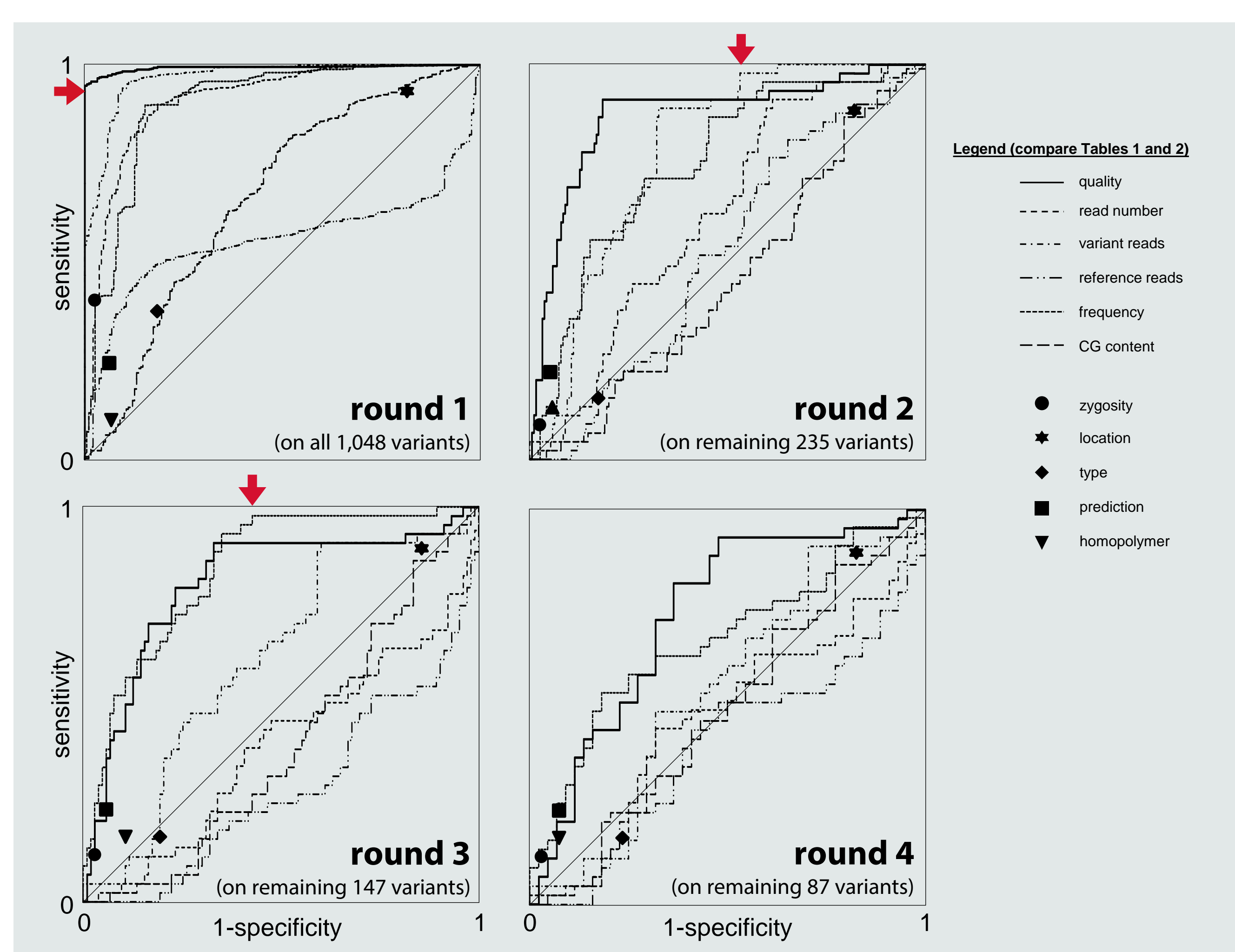


Figure 2: Optimized workflow based on our set of data. Arrows in round 1-3 ROC analyses indicate parameters with highest predictive value for true positives (round 1) and true negatives (rounds 2 and 3).

4. Sanger sequencing load can strongly be reduced without compromising specificity or sensitivity, but remains necessary for a fraction of variants

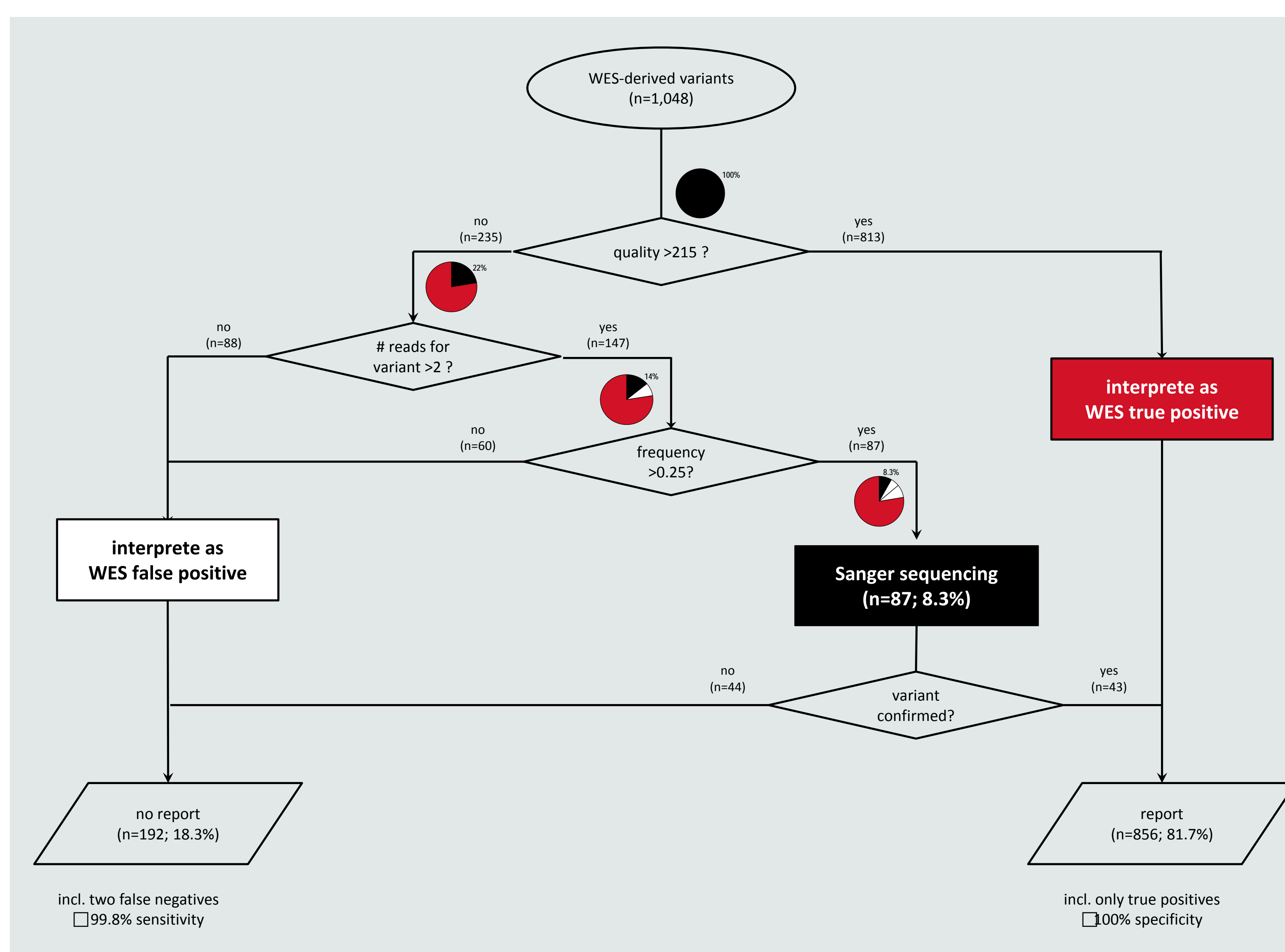


Figure 3: Optimized workflow (compare Figure 2). Our WES-based findings corroborate recent statements on panel-based NGS.³

References

- in press as Bauer et al. in Genet Med
- Strom et al., 2014, Genet Med 16:510-5
- Mu et al., 2016, J Mol Diagn 18:923-32

Disclosure of conflict of interest

This study was sustained in part by Centogene AG, Rostock. All authors of the presentation are employees of CENTOGENE AG, Rostock, Germany.