

Validation of an Expanded Carrier Screen that Optimizes Sensitivity via Full-Exon Sequencing and Panel-wide Copy Number Variant Identification

Gregory J. Hogan,^{1†} Valentina S. Vysotskaia,^{1†} Kyle A. Beauchamp,¹ Stefanie Seisenberger,¹ Peter V. Grauman,¹ Kevin R. Haas,¹ Sun Hae Hong,¹ Diana Jeon,¹ Shera Kash,¹ Henry H. Lai,¹ Laura M. Melroy,¹ Mark R. Theilmann,¹ Clement S. Chu,¹ Kevin Iori,¹ Jared R. Maguire,¹ Eric A. Evans,¹ Imran S. Haque,^{1‡} Rebecca Mar-Heyming,¹ Hyunseok P. Kang,¹ and Dale Muzzezy^{1*†}

BACKGROUND: By identifying pathogenic variants across hundreds of genes, expanded carrier screening (ECS) enables prospective parents to assess the risk of transmitting an autosomal recessive or X-linked condition. Detection of at-risk couples depends on the number of conditions tested, the prevalence of the respective diseases, and the screen's analytical sensitivity for identifying disease-causing variants. Disease-level analytical sensitivity is often <100% in ECS tests because copy number variants (CNVs) are typically not interrogated because of their technical complexity.

METHODS: We present an analytical validation and preliminary clinical characterization of a 235-gene sequencing-based ECS with full coverage across coding regions, targeted assessment of pathogenic noncoding variants, panel-wide CNV calling, and specialized assays for technically challenging genes. Next-generation sequencing, customized bioinformatics, and expert manual call review were used to identify single-nucleotide variants, short insertions and deletions, and CNVs for all genes except *FMRI* and those whose low disease incidence or high technical complexity precluded novel variant identification or interpretation.

RESULTS: Screening of 36 859 patients' blood or saliva samples revealed the substantial impact on fetal disease-risk detection attributable to novel CNVs (9.19% of risk) and technically challenging conditions (20.2% of risk), such as congenital adrenal hyperplasia. Of the 7498 couples screened, 335 were identified as at risk for an affected pregnancy, underscoring the clinical importance of the test. Validation of our ECS

demonstrated >99% analytical sensitivity and >99% analytical specificity.

CONCLUSIONS: Validated high-fidelity identification of different variant types—especially for diseases with complicated molecular genetics—maximizes at-risk couple detection.

© 2018 American Association for Clinical Chemistry

There are more than 1000 recessive single-gene conditions that vary in severity and age of onset (1). Each is uncommon in the general population, yet collectively these Mendelian diseases account for approximately 20% of infant mortality and 10% of infant hospitalizations (2, 3). Screening for carriers of such conditions in the pre-conception or prenatal period informs couples about the risk of having a child with a serious disease and the available family-planning options. Because of the increasing quality and decreasing cost of genomic technologies, it is now possible to perform pan-ethnic carrier screening for many conditions simultaneously [referred to as expanded carrier screening (ECS)²]. Our previous study of carrier rates in 346 790 patients showed that an ECS panel was expected to identify more pregnancies at risk for serious conditions than ethnic-based panels spanning far fewer genes (4), and the American College of Obstetricians and Gynecologists recently recognized ECS as an acceptable strategy for pre-conception and prenatal carrier screening (5). To the extent that these guidelines increase ECS usage, they will have a large clinical impact because it has recently been shown in clinical-utility studies that approximately 80% of couples found to be at risk for severe conditions pursue alternative reproductive options (6, 7).

¹ Counsyl, Inc., South San Francisco, CA.

* Address correspondence to this author at: Counsyl, Inc., 180 Kimball Way, South San Francisco, CA 94110. Fax 608-541-2450; e-mail dale@counsyl.com.

† These authors contributed equally to this work.

‡ Current affiliation: Freenome, Inc., South San Francisco, CA.

Received January 12, 2018; accepted March 26, 2018.

Previously published online at DOI: 10.1373/clinchem.2018.286823

© 2018 American Association for Clinical Chemistry

² Nonstandard abbreviations: ECS, expanded carrier screening; CNV, copy number variant; NGS, next-generation sequencing; 21-OH CAH, 21-hydroxylase deficiency congenital adrenal hyperplasia; SNV, single nucleotide variant; indels, small insertions and deletions; MFDR, modeled fetal disease risk; ARC, at risk couple; QC, quality control; 1KG, 1000 Genomes; MLPA, multiplex ligation-dependent probe amplification; SNP, single-nucleotide polymorphism.

An ECS must have a high detection rate for each disease on the panel both to identify at-risk couples (i.e., couples wherein both partners are carriers of an autosomal recessive condition or the female partner is a carrier of an X-linked condition) and to minimize the residual risk in couples when only 1 partner has tested positive. The detection rate is particularly important for the recessive diseases that predominate ECS panels because the odds of detecting an at-risk couple scales as the square of the rate for finding an individual carrier (e.g., an 80% detection rate for each parent results in only a 64% detection rate for an at-risk couple). To maximize detection rates, copy number variants (CNVs) must be identified, yet most ECS tests interrogate only a handful of common CNVs with known breakpoints. However, CNVs can vary in size and position, encompassing everything from single exons, which account for 29% of CNVs for Mendelian conditions (8), to the entire gene. A diversity of pathogenic CNVs has been observed in cystic fibrosis carriers, accounting for 1.6% of carriers (9), meaning that the carrier detection rate without CNV detection is <98.4%, which in turn makes the at-risk couple detection rate <96.8%. The inverse is noteworthy: Including novel CNV detection can boost at-risk couple detection for cystic fibrosis to nearly 100%. To our knowledge, the impact of CNVs across all genes on an ECS has not yet been characterized.

Carrier status for a minority of the most prevalent serious conditions is difficult to resolve with standard next-generation sequencing (NGS) and bioinformatics approaches because of the challenging sequence features of the disease genes. Thus, these conditions require special handling. Low complexity sequences (e.g., CGG repeat expansion in *FMR1*³ for fragile X syndrome) and highly homologous regions (e.g., *SMN1* and *SMN2* genes for spinal muscular atrophy) complicate variant identification, yet these hard-to-sequence genes simultaneously contribute substantially to the disease risk. For instance, fragile X syndrome, spinal muscular atrophy, 21-hydroxylase deficiency congenital adrenal hyperplasia (CAH), and α -thalassemia account for 54 affected fetuses per 100 000 pregnancies (10).

Here, we characterized the performance of a 235-gene ECS leveraging NGS to identify single-nucleotide variants (SNVs), small insertions and deletions (indels), CNVs [deletions within nearly all genes and both deletions and duplications for *CFTR* (cystic fibrosis trans-

membrane conductance regulator) and *DMD* (muscular dystrophy, Duchenne and Becker types)], and hard-to-sequence targeted variants. Following medical societies' recommendations (11, 12), we present an analytical validation of the test. Using data from a cohort of 36 859 patients, we modeled the test's clinical impact, focusing on the role of panel-wide CNV calling in detection of at-risk couples. The high analytical sensitivity across genes and variant types that we observed is clinically important because it facilitated identification of 335 couples at high risk out of 7498 couples tested.

Materials and Methods

INSTITUTIONAL REVIEW BOARD APPROVAL

The study protocol was approved by Western Institutional Review Board (number 1145639) and complied with the Health Insurance Portability and Accountability Act (HIPAA). Patient information was deidentified according to the HIPAA Privacy Rule. An informed consent waiver was approved by the Institutional Review Board.

TEST DESCRIPTION

We compiled a panel (Counsyl Foresight Carrier Screen) of 235 genes responsible for 234 clinically important autosomal recessive and X-linked diseases (see Table 1 in the Data Supplement that accompanies the online version of this article at <http://www.clinchem.org/content/vol64/issue7>). This panel consisted of a "Universal" sub-panel (176 diseases) for routine ECS and an opt-in panel (234 diseases) aimed at specific high-risk populations. The design of the Universal panel was as described by Beauchamp et al. (10) and prioritized prevalent diseases with serious and highly penetrant phenotypes that could affect clinical counseling and family planning. Some diseases with moderate but lifelong impact were also included [e.g., familial Mediterranean fever (*MEFV*) and *DFNB1* nonsyndromic hearing loss and deafness (*GJB2*)]. Table 1 summarizes gene-specific methodologies and variant types.

PATIENT COHORT

Modeled fetal disease risk (MFDR) and at-risk couple (ARC) calculations considered patients screened with the 176-disease Universal panel. MFDR was a modeled estimate of the probability that a pregnancy was affected by a condition on the panel (it could be summed across diseases to yield an aggregate MFDR of the panel), whereas an ARC, as used herein, was a couple identified empirically as being at high risk for an affected child. Reported MFDR numbers were calculated as described previously (4, 10) and were US census weighted (calculations outlined in the Methods section of the online Data Supplement). The cohort used in MFDR calculations included

³ Human genes: *FMR1*, fragile X mental retardation 1; *SMN1*, survival of motor neuron 1; *SMN2*, survival of motor neuron 2; *CFTR*, cystic fibrosis transmembrane conductance regulator; *DMD*, dystrophin; *MEFV*, MEFV, pyrin innate immunity regulator; *GJB2*, gap junction protein beta 2; *GBA*, glucosylceramidase beta; *HBA1*, hemoglobin subunit alpha 1; *HBA2*, hemoglobin subunit alpha 2; *CYP21A2*, cytochrome P450 family 21 subfamily A member 2; *GALC*, galactosylceramidase; *SLC12A2*, solute carrier family 12 member 2; *CYP21A1P*, cytochrome P450 family 21 subfamily A member 1, pseudogene; *GBAP1*, glucosylceramidase beta pseudogene 1.

Table 1. Foresight™ ECS panel.

Disease genes	Methodology	Variants reported
<i>General ECS:</i>		
216 genes	NGS	Novel pathogenic SNVs, indels, large deletions
<i>CFTR</i>	NGS	Novel pathogenic SNVs, indels, large deletions, and duplications
<i>DMD</i>	NGS	Novel pathogenic SNVs, indels, large deletions, and duplications
11 genes	NGS	Targeted pathogenic mutations ^a
<i>Technically challenging genes:</i>		
<i>SMN1</i>	NGS	Exon 7 copy number, g.27134T>G SNP
<i>CYP21A2</i>	NGS	<u>Classical:</u> <i>CYP21A2</i> 30-kb deletion, <i>CYP21A2</i> duplication, <i>CYP21A2</i> triplication, c.293-13C>G, p.G111Vfs*21, p.I173N, p.[I237N;V238E;M240K], p.L308Ffs*6, p.Q319*, p.Q319* + <i>CYP21A2</i> dup, p.R357W <u>Nonclassical:</u> p.P31L, p.V281L
<i>HBA1/2</i>	NGS	<u>Single deletions:</u> -alpha3.7, -alpha4.2 <u>Double deletions:</u> -(alpha)20.5, - -BRIT, - -MEDI, - -MEDII, - -SEA, - -THAI, or - -FIL <u>Frequent SNV:</u> Hb constant spring <u>Regulatory deletion:</u> ΔHS-40 (combinations of most variants above with other deleterious variants or duplications can also be detected)
<i>GBA</i>	NGS	p.N370S, p.D409V, p.D448H, IVS2 + 1G>A, p.L444P, p.R463C, p.R463H, p.R496H, p.V394L, p.L29Afs*18
<i>FMR1</i>	PCR/CE ^b	Number of CGG repeats in the 5'-UTR

^a See Table 1 in the online Data Supplement.
^b CE, capillary electrophoresis; UTR, untranslated region.

only patients receiving routine carrier screening and, therefore, excluded those with known family history or infertility. By contrast, the reported couple and ARC counts are raw counts, not necessarily representative of the general US population because of the over-representation of high-risk ethnic groups, infertility patients, and patients with family history of disease. Importantly, the only patients considered to be “at risk” in MFDR and ARC calculations were those with variants that were interpreted as being pathogenic via American College of Medical Genetics and Genomics criteria [described by Beauchamp et al. (10)]. Further, the risk-status calculations accounted for known disease-specific variant combinations that influence pathogenicity (e.g., a couple in which both partners were carriers only for the D444H pathogenic variant were not at risk of a child affected with biotinidase deficiency because homozygosity of D444H—in the absence of other variants—is benign).

NGS WORKFLOW

The molecular workflow of our NGS pipeline was as previously described (13) and is briefly summarized in the Methods section of the online Data Supplement. Sequencing reads were aligned to the human genome ver-

sion 19 using the BWA-MEM algorithm (14). Novel SNVs and indels were identified and genotyped using GATK 1.6 and FreeBayes (15, 16), and 9 known pathogenic sites involving complex indels were detected with simple custom genotyping software that (a) measured the frequency of junction reads in the NGS data that manifest the known mutations and (b) identified samples as carriers if the frequencies exceeded a threshold. CNVs were determined using custom software that leveraged read-depth values, described extensively by Vysotskaia et al. (13). A combination of targeted genotyping and read depth-based copy number analysis identified variants in the technically challenging genes *SMN1*, *GBA*, *HBA1/2*, and *CYP21A2*, as described further in the Methods section of the online Data Supplement. Quality control (QC) metrics and the variant interpretation workflow are also described in the Methods section and Table 2 of the online Data Supplement.

FMR1 CGG REPEAT SIZING

CGG trinucleotide expansions of the *FMR1* promoter were measured by PCR amplification and capillary electrophoresis as previously described (17).

Table 2 Samples and variants included in analytical validation for sensitivity evaluation, assessed for different variant types and technically challenging genes.^a

Variant type	Variant details	Sample source (no. of unique samples)	Number of total variants (no. unique)	Reference data
SNVs and small indels	≤5 bp	Cell lines (92)	41 137 (3364)	1KG and GIAB ^b
Large indels	6-10 bp	Patients (52)	16 (14)	Sanger
	>10 bp		36 (35)	
CNVs	Panel-wide single-exon dels	Patients (26)	5 (4)	MLPA/TaqMan
	Panel-wide multiexon dels		21 (19)	
	<i>CFTR/DMD</i> single-exon dels	Patients (7) and cell lines (11)	3 (2)	MLPA/TaqMan or Coriell
	<i>CFTR/DMD</i> multiexon dels		12 (9)	
<i>CFTR/DMD</i> dups	3 (3)	MLPA/TaqMan or Coriell		
Technically challenging genes:				
<i>SMN1</i>	0-3 copies	Patients (120) and cell lines (8)	234 (2)	TaqMan and/or Coriell
	g.27134T>G SNP		Cell lines (106)	
<i>GBA</i>	p.D448H, IVS2 + 1G>A, p.L444P	Patients (6) and cell lines (1)	7 (3)	TaqMan or Coriell
<i>CYP21A2</i>	Classical, nonclassical, and combination of variants	Patients (15)	15 (14)	Sanger and/or MLPA
<i>FMR1</i>	CGG repeats: normal, intermediate, premutation, and full mutation	Cell lines (40)	40 (37)	Coriell
<i>HBA1/2</i>	Single/double deletions, regulatory deletions, biallelic	Patients (10)	10 (10)	MLPA

^a Sums among number of samples and number of variants might not match because of replicates of individual samples. Samples involved in specificity and reproducibility calculations are in found in Table 3 of the online Data Supplement.

^b GIAB, Genome in a Bottle; dels, deletions; dups, duplications.

ANALYTICAL VALIDATION

Samples and reference data. Samples and reference data were compiled from different sources (Table 2; see also Tables 3 and 4 in the online Data Supplement): purified DNA for 91 cell lines [1000 Genomes (1KG) Project (18)], 70 cell lines with known pathogenic variants in specific genes (Coriell Repository; see Table 5 in the online Data Supplement), and 115 mutation-positive patient blood and saliva samples tested with a previous version of the Counsyl carrier test (a 94-disease panel). Although Coriell and patient samples were selected for inclusion in the validation study to represent a broad range of relevant variants, the 1KG samples were not selected based on any formal criteria. Relevant variants in all mutation-positive patient samples were confirmed orthogonally by PCR/Sanger sequencing, quantitative PCR, or multiplex ligation-dependent probe amplification (MLPA), described in the Methods section and Table 6 of the online Data Supplement. Note that

NA06896 was dropped from *FMR1* accuracy analysis because of inconsistent reference data (19, 20). Validation samples were tested using our standard operating procedure, which included in-process and postprocess QC at the batch, sample, and variant-call level (see Table 2 in the online Data Supplement). Furthermore, consistent with our standard operating procedure for clinical samples, licensed experts, who were blinded to the validation sample set, performed manual review of the sequencing data using our custom review interface. Samples that failed QC and manual review were excluded from further analysis.

Simulation of synthetic CNVs. For every region reportable for CNVs, we simulated a single-copy deletion and tested calling sensitivity; we also simulated single-copy duplications for *DMD* and *CFTR* regions. To create a synthetic CNV in silico, a validation sample that passed QC and manual review was randomly selected, and the depth was adjusted in a specified region such that the normalized

depth ratio was reduced or increased by 50% for a deletion or duplication, respectively. Therefore, the simulations altered copy number but preserved experimental noise. Only 1 deletion or duplication was introduced in a given sample. Variants on chromosome X were only simulated in female samples. Simulated CNVs spanned 1, 2, or 4 exons, and each was simulated independently 5 times. Data from synthetic positive samples and reference samples were analyzed with the CNV calling algorithm described in detail by Vysotskaia et al. (13). A true-positive CNV call matched the intended simulated CNV, ignoring slight differences in precise breakpoints. A false-negative sample had a deletion not identified by the algorithm. Aggregate analytical sensitivity was reported as a weighted sum of the 1-, 2-, and 4-exon sensitivities, for which the weights correspond to empirical frequencies of such CNVs in the literature (8, 21, 22).

Statistical analysis. Validation metrics were defined as follows: Accuracy = $(TP + TN)/(TP + FP + TN + FN)$; Sensitivity = $TP/(TP + FN)$; Specificity = $TN/(TN + FP)$; FDR = $FP/(TP + FP)$, where FDR = false discovery rate. The CIs were calculated by the method of Clopper and Pearson (23). Intraassay and interassay reproducibility was calculated as the ratio of concordant calls to total calls.

Results

Based on published design criteria (10), we developed an NGS-based ECS covering 220 autosomal recessive and 14 X-linked conditions, including technically challenging diseases (Fig. 1 and Table 1; see also Table 1 in the online Data Supplement). For nearly all genes on the 176-disease Universal panel, SNVs and indels were detected via NGS data acquired from regions that could impact gene function (e.g., padded coding exons and known or potentially pathogenic intronic variants). Large CNVs were identified at single-exon resolution panel-wide using relative sample-to-sample changes in sequencing depth. The test was validated as described below and used in a clinical production setting on 36 859 patient samples (tested between November 2016 and February 2018). Of the 7498 couples who underwent testing, 335 were found to be at risk for a condition on the Universal panel. Carrier rates on this cohort enabled calculation of the expected fraction of US pregnancies the Universal panel would identify as affected on a per-disease level (Fig. 1A; see also Table 7 in the online Data Supplement): approximately 1 in 300 pregnancies would be affected by at least 1 serious disease. For a handful of prevalent diseases that comprised 8.1% of the total panel disease risk, high carrier sensitivity required customized CNV analysis because of the genes' complicated techni-

cal features (Fig. 1B; see also the analysis described in the Methods section of the online Data Supplement).

Relative to our previous version of the ECS characterized in a study of 346 790 patients, this updated ECS differed in 2 ways that collectively boosted the assessed MFDR (4, 10). First, the updated Universal panel probed SNVs and indels for 82 more diseases (Fig. 2A); although they have diverse incidence rates, the added diseases collectively accounted for nearly 23% of the updated panel's assessed MFDR (Fig. 2A). Second, the screen additionally detected deletions ranging in size from a single exon to the entire gene (Fig. 2). Panel-wide CNVs contributed approximately 9.2% of the assessed MFDR. CNVs in *DMD* alone represented approximately 6% of the risk, with the remaining genes accounting for approximately 3%, largely consistent with our expectation of CNV-attributable MFDR estimated from the 94-disease panel (10). We have observed hundreds of carriers with exon-level deletions (Fig. 2B); although many spanned multiple exons, 74 deletions encompassed a single exon, demonstrating the need to optimize assay performance to have high analytical sensitivity for short CNVs.

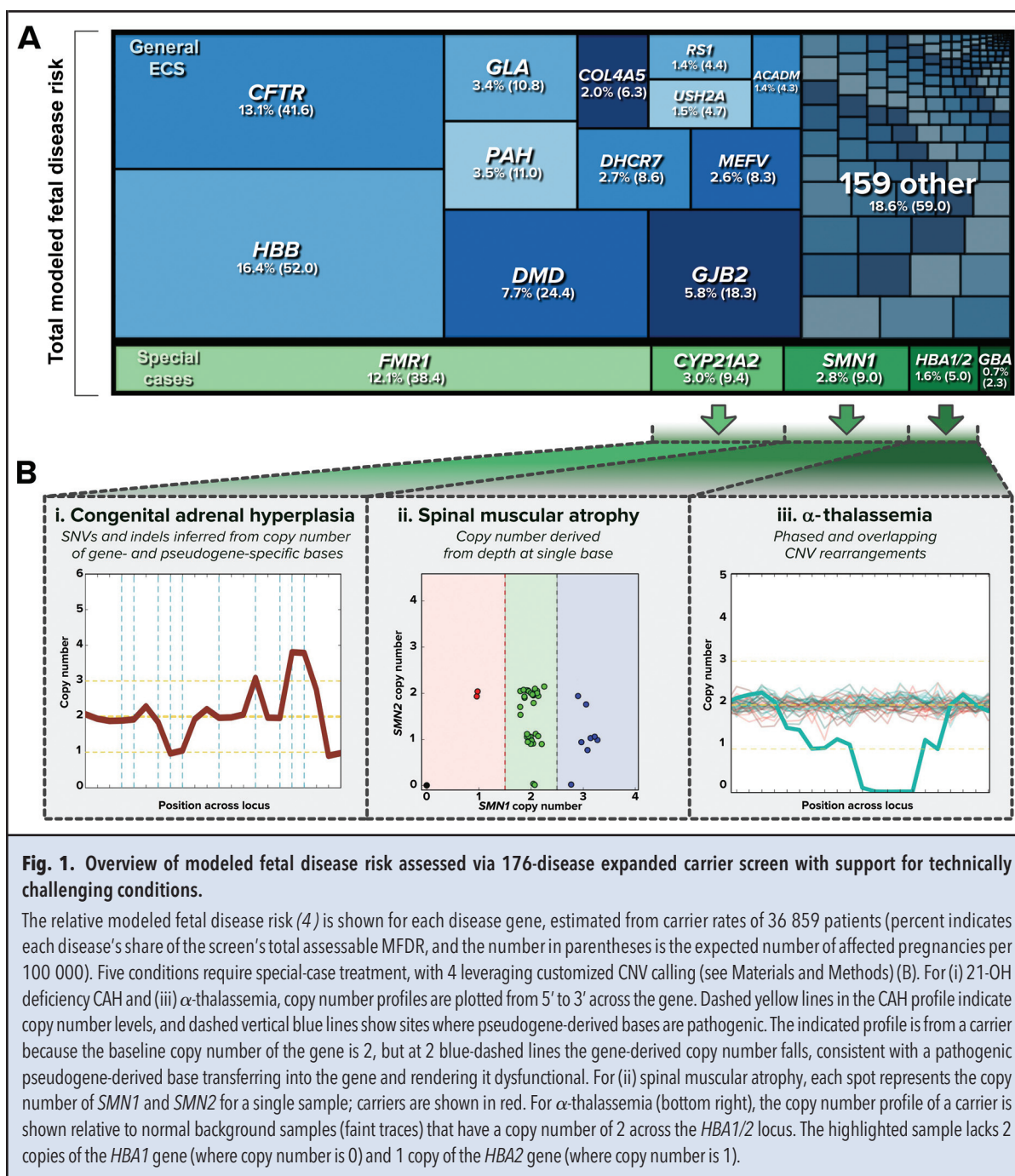
VALIDATION APPROACH

Following guideline recommendations (11, 12) to assess analytical performance of the ECS panel before clinical use, we measured the accuracy of identifying variants that were small (e.g., SNVs and small indels), technically nuanced (e.g., large indels and CNVs), and in hard-to-sequence genes (e.g., *CYP21A2*) (Table 2; see also Table 4 in the online Data Supplement).

ACCURACY AND REPRODUCIBILITY FOR CALLING SINGLE-NUCLEOTIDE POLYMORPHISMS AND SMALL INDELS IN 229 GENES

We compared our ECS data for the reference sample NA12878 with data from the Genome in a Bottle Consortium (24), which includes high-confidence calls for >97.5% of the regions covered by our test. We tested NA12878 across 5 batches and in duplicate within 3 batches for a total of 8 tests, and the results showed an accuracy of >99.99% (see Table 8 in the online Data Supplement). NA12878 is one of our routine controls within every production batch. Since the validation of our test, we have further measured accuracy across 207 batches that spanned reagent lots and instruments, reproducibly observing high calling accuracy across the panel (see Fig. 1 in the online Data Supplement).

To measure SNV and indel calling accuracy across a diverse set of samples, we performed our ECS on reference samples from 1KG. For 90 tested samples that passed QC, we compared genotypes across all exonic regions with sufficient coverage and quality in the 1KG data (248 490 calls in all, 2% of which are indels); 52 discordant calls were adjudicated with Sanger sequencing



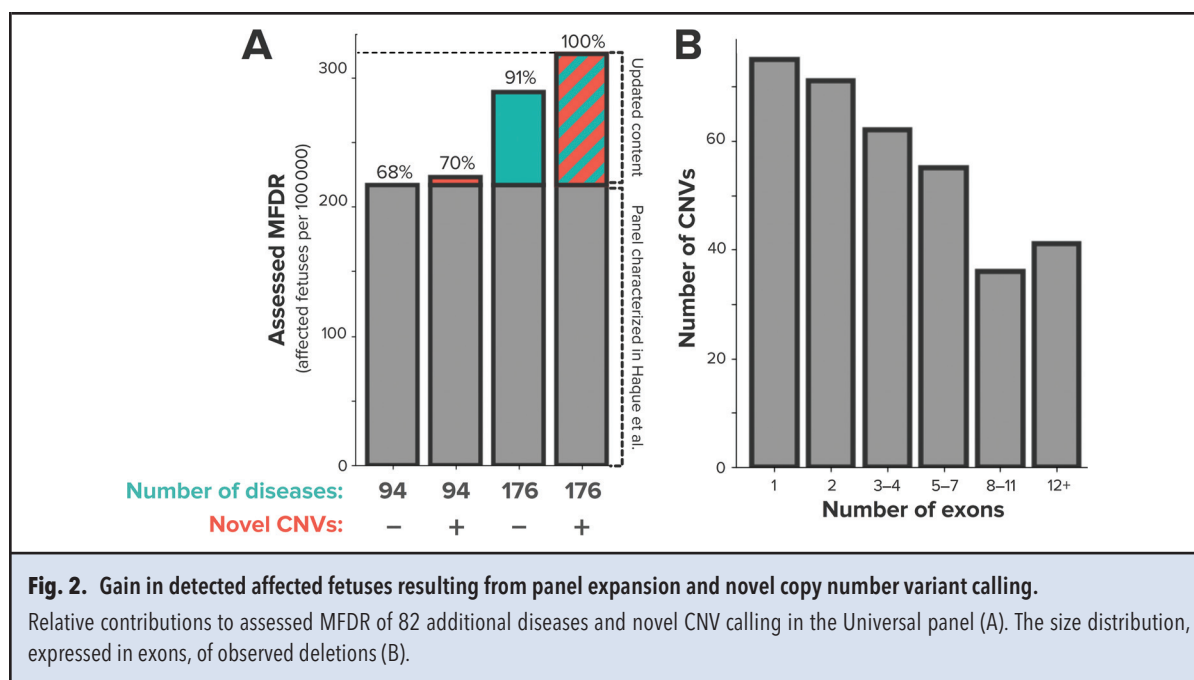
(see Table 9 in the online Data Supplement). Our ECS identified 36 032 true-positive calls and 212 139 true-negative calls, resulting in >99.99% accuracy, analytical sensitivity, and analytical specificity (Fig. 3A).

In addition to establishing the analytical accuracy of the ECS using reference DNA from cell lines, we measured intraassay and interassay reproducibility using different sample types by comparing the equivalence of genotyping

calls starting from separate aliquots of DNA. Overall, the test achieved >99.9% intraassay and interassay reproducibility (see Table 10 in the online Data Supplement).

TECHNICALLY CHALLENGING VARIANTS

Larger indel detection performance. Although only 5% of indels are ≥ 5 bp (25), sensitivity decreases as indel size



increases. Thus, to ensure high analytical sensitivity for detecting indels, we built a cohort of 52 patient samples with 49 unique technically challenging, larger (>5 bp) deletions, insertions, or complex indels in 42 different genes (see Table 11A in the online Data Supplement). All the expected indel calls (52 of 52), including a 33-bp deletion and 21-bp insertion, were observed (Fig. 3B).

CNV detection performance. To overcome the limitation of scarce reference materials for CNV calling, we supplemented available reference material (11 Coriell cell lines; see Table 5B in the online Data Supplement) with orthogonally confirmed positives identified retrospectively (33 clinical samples confirmed by MLPA; see Table 11B in the online Data Supplement). All 44 empirical CNVs—across 13 different genes—were detected (Fig. 3B), demonstrating high analytical sensitivity (100%; 95% CI, 92%–100%; see reproducibility data in Table 12 of the online Data Supplement). Notably, 23 samples had a 1-exon or 2-exon CNV, which can be technically challenging for an NGS-based assay to detect (see Table 4 in the online Data Supplement).

We additionally used >250 000 in silico simulated CNVs to measure analytical sensitivity systematically across the panel. In the in silico CNV simulations, we introduced a synthetic deletion or duplication spanning at least 1 coding exon in the background of an empirical sample's data (Fig. 3, C and D). To assess analytical sensitivity of clinically relevant deletions and duplications in *CFTR* and *DMD*, we scaled the analytical sensitivity for each CNV size by its population frequency cataloged in public databases (21, 22), yielding an aggregate 99.9%

analytical sensitivity for CNVs in each gene. Across the rest of the panel, for which only deletions are reported, our simulations revealed 81.8% analytical sensitivity for single-exon deletions and 98.3% to 100% sensitivity for multiexon deletions.

To assess the analytical specificity of CNV calling, CNV calls in 1KG reference samples were checked against the reference calls. After the standard manual call review of all CNVs, 2 calls were deemed positive (NA12716: *GALC* 7-exon deletion; NA19700: *SLC12A6* 4-exon deletion) and 1 call was flagged as low quality. The NA12716 and NA19700 deletions matched their reference genotypes. Following our production standard operating procedure for low-quality CNV calls, we retested the sample with a flagged call and found it to be a confident negative. Therefore, no false-positive findings were observed in the 1KG reference samples, resulting in an estimated CNV-calling specificity at the sample level of 100% (91 of 91; 95% CI, 96%–100%) and at the gene level of 100% (19 838 of 19 838; 95% CI, 99.98%–100%). Taken together, the empirical and simulation analyses showed the ECS had high analytical sensitivity and specificity for exon-level CNVs.

VARIANT DETECTION PERFORMANCE USING NGS IN THE TECHNICALLY CHALLENGING GENES *CYP21A2*, *HBA1/2*, *GBA*, AND *SMN1*

Several diseases of clinical importance result from mutations in genes that have a paralog or pseudogene that complicates molecular analysis. Such diseases include spinal muscular atrophy (*SMN1* and *SMN2* encode the same protein, but *SMN2* harbors a splicing variant that

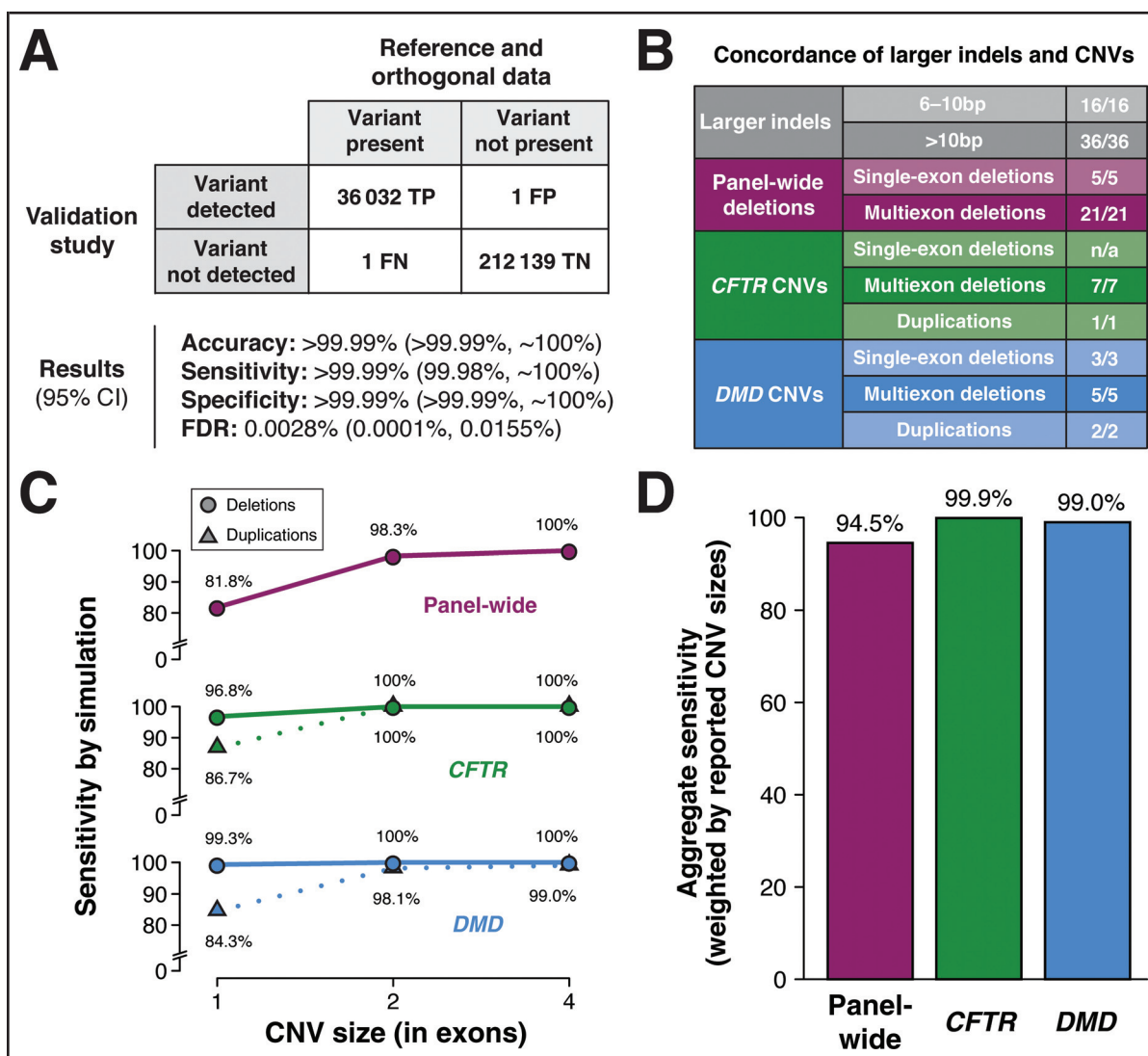


Fig. 3. Analytical performance for calling SNVs, indels, and CNVs.

Contingency table and results for SNV and small indel calling in 229 genes, assessed using 1KG reference material and adjudication by follow-up Sanger (A). For true-negative calculations, all polymorphic positions (positions at which we observed nonreference bases in any sample) across all samples were considered. No-calls were censored from analysis. The no-call rate was 0.13% (317 of 248 490). TP, true positives; TN, true negatives; FP, false positives; FN, false negatives; FDR, false discovery rate. Concordance summary for larger indels and CNVs (B). Sensitivity for CNV calling as measured by simulations, by gene, type, and size (in number of exons) (C). Aggregate sensitivity for CNV calling as measured by simulations (D). Simulation results in (C) were weighted by size and frequency (see Materials and Methods). In (B–D), data reported for “panel-wide” deletions exclude *CFTR* and *DMD*.

results in approximately 10% of functional SMN protein relative to *SMN1*) (26), α -thalassemia (*HBA1* and *HBA2* have identical coding sequences and few distinguishing noncoding bases), 21-OH-deficient CAH (the *CYP21A2* coding sequence is >99% identical to its pseudogene *CYP21A1P*), and Gaucher disease (*GBA* has a nearby pseudogene *GBAP1* with which it shares high sequence identity in certain exons). Recombination and gene con-

version are frequent among these genes and their homologs, which can result in copy number changes (Fig. 1, bottom). We implemented custom variant-calling algorithms combining depth-based copy number and specific mutation analyses for the disease genes and their homologs (see Methods sections here and in the online Data Supplement). For each of these genes, we verified that the ECS identified each of the challenging genotypes correctly

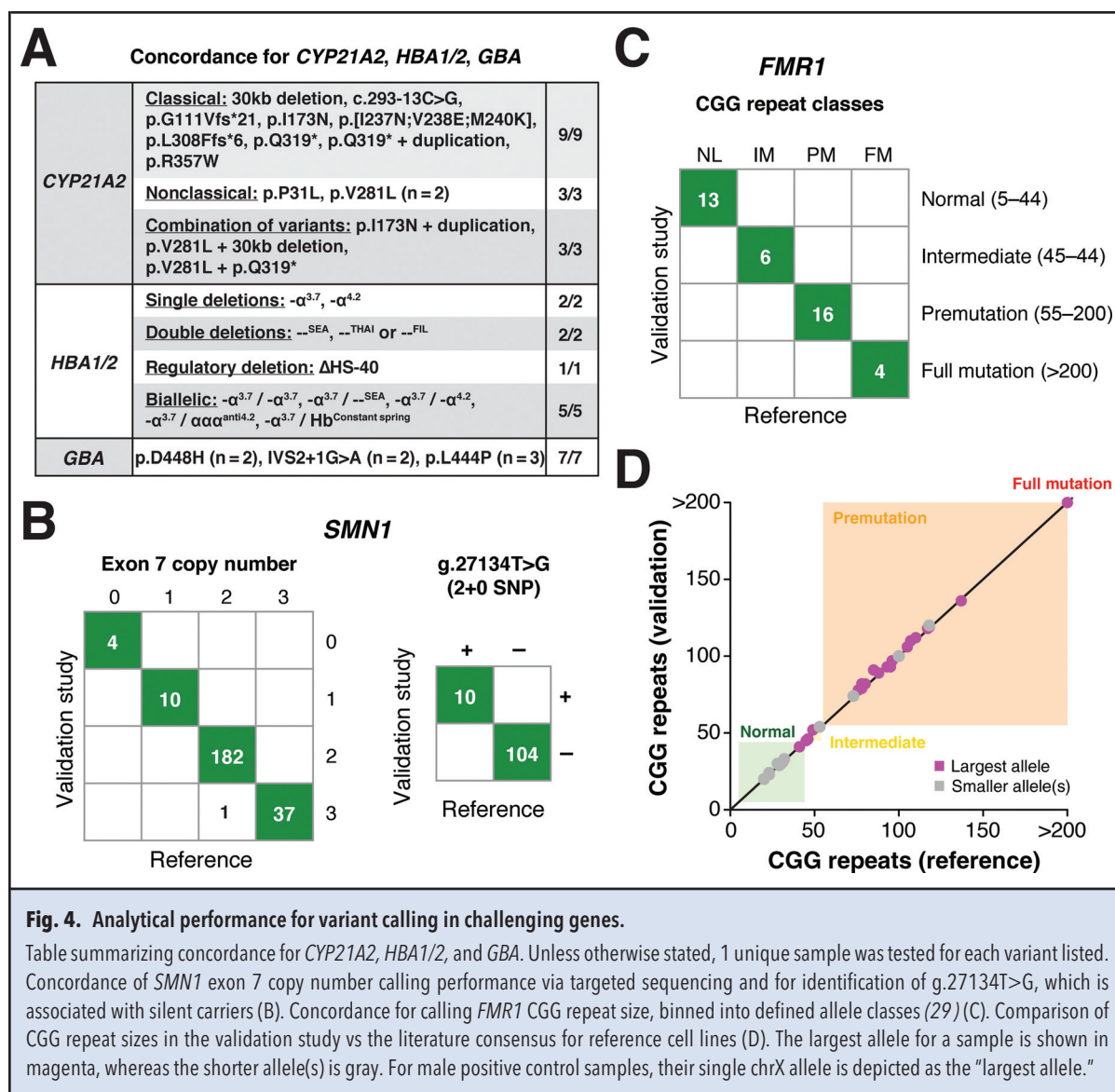


Fig. 4. Analytical performance for variant calling in challenging genes.

Table summarizing concordance for *CYP21A2*, *HBA1/2*, and *GBA*. Unless otherwise stated, 1 unique sample was tested for each variant listed. Concordance of *SMN1* exon 7 copy number calling performance via targeted sequencing and for identification of g.27134T>G, which is associated with silent carriers (B). Concordance for calling *FMR1* CGG repeat size, binned into defined allele classes (29) (C). Comparison of CGG repeat sizes in the validation study vs the literature consensus for reference cell lines (D). The largest allele for a sample is shown in magenta, whereas the shorter allele(s) is gray. For male positive control samples, their single chrX allele is depicted as the “largest allele.”

by testing samples confirmed via orthogonal methods to be carriers of the genotypes of interest. As indicated in Fig. 4A here and Table 11C in the online Data Supplement, 15 variants in *CYP21A2*, 10 variants in the *HBA* locus, and 7 variants in *GBA* were correctly identified.

For spinal muscular atrophy, 128 unique (234 total with replicates) samples with 0, 1, 2, or 3 copies of *SMN1* were analyzed by NGS (see Tables 4 and 5D in the online Data Supplement). Carrier (samples with 0 or 1 copy) vs noncarrier (samples with ≥ 2 copies) identification accuracy by NGS was 100% (95% CI, 98.4%–100%). NGS copy number accuracy was 233 of 234, or 99.6% (95% CI, 97.6%–100%) (Fig. 4B), for which 1 noncarrier patient sample had 3 copies by NGS and 2 by quantitative PCR. We also measured detection of the g.27134T>G

single-nucleotide polymorphism (SNP) associated with 2 + 0 spinal muscular atrophy carrier status (27) in 98 1KG samples and 16 Coriell samples (see Table 5D in the online Data Supplement); all 114 ECS results were concordant with reference data.

FMR1 CGG-REPEAT ANALYSIS

Fragile X syndrome arises from a trinucleotide CGG repeat expansion in the 5' untranslated region of *FMR1* (28). *FMR1* alleles are categorized as normal (5–44 CGG repeats), intermediate (45–54 CGG repeats), premutation (55–200 CGG repeats), and full mutation (>200 CGG repeats) (29). Thirty-nine Coriell samples (see Table 5C in the online Data Supplement) enriched for expansions of various sizes were classified correctly by

our assay (Fig. 4C). Further, the identified CGG repeat allele sizes closely matched the literature consensus sizes (Fig. 4D).

Discussion

ECS provides reliable and affordable risk assessment for many serious recessive and X-linked diseases simultaneously. Genomic technologies like NGS have enabled growth in ECS panel size without incurring a corresponding increase in testing cost, but judicious panel construction and validation are required. For these reasons, we recently published a systematic process for ECS panel design (10) and here present an analytical validation of our updated ECS and a modeling analysis of its clinical impact. Our analytical validation study of variants in hundreds of genes across hundreds of samples demonstrates high analytical sensitivity, analytical specificity, and accuracy of genotype calls. Further, the analysis of clinical impact in a large patient cohort shows that the ECS is expected to identify approximately 1 in 300 pregnancies as being at risk for serious disease in the general US population. Notably, although the incidence of each individual disease is low, the collective frequency of the 176 screened diseases exceeds that of Down syndrome (1 in 800 live births) (30), for which routine screening is offered.

Including panel-wide CNV calling in an NGS-based ECS increases the chance of finding couples at risk for having a child with a serious condition. CNVs can be identified in a clinical workflow via orthogonal technologies (e.g., MLPA) rather than in a single NGS assay, but MLPA testing does not affordably scale to hundreds of genes and incurs additional laboratory handling steps that can introduce operator error. Using known positive samples from biorepositories, retrospectively identified CNV-positive samples, and *in silico* simulations, we demonstrated high sensitivity for novel CNV identification. We expect that simulation analyses, as used here, will become increasingly important during NGS-panel validation, for which performance needs to be evaluated even when clinical samples are rare or nonexistent.

As with other analytical validation studies of NGS-based screens that discover novel variants, a limitation of our study is that analytical sensitivity cannot be established for every variant that the screen could report. The test interrogates many hundreds of kilobases of genomic sequence, and the frequency of variation at many sites is too low to source samples representing all possibilities (i.e., the study involves a finite set of ascertained samples from public repositories and our historical patient cohort). As such, we sought to establish general proficiency for finding different variant types (SNVs, indels, CNVs) and to ensure proper identification of several variants from each of our special-case genes. We have demon-

strated that the screen identifies reference sample genotypes with almost-perfect accuracy. Further, our simulations are a direct attempt to measure analytical sensitivity for the range of possible CNVs the test is designed to identify; aggregate CNV sensitivity was estimated to be >94%. We strongly expect, but cannot formally demonstrate at all sites, that the ECS has high analytical sensitivity for the range of interrogated variant types across the panel. Analytical sensitivity could be reduced at large indels, and clinical sensitivity may be less than analytical sensitivity because of as-yet-undiscovered pathogenic variants in introns or pathogenic CNVs that the test is not designed to identify and report.

The updated ECS contains approximately twice as many genes as the previous version, yet the risk resolved is not twice as great. This phenomenon is driven by the disparate incidence of diseases and highlights the importance of high detection rates for the most common serious conditions, many of which pose challenges because of complicated molecular genetics. For several special cases, we have fine-tuned CNV calling to capture single-base differences (spinal muscular atrophy), phased and overlapping rearrangements (α -thalassemia), and complicated gene conversions (CAH and Gaucher disease). We expect the collective risk of these 4 diseases to rival that of the 100 least-common diseases on the panel.

Ultimately, the clinical value of ECS stems from its ability to identify variants (analytical validity), interpret the pathogenicity of those variants (clinical validity), and impact behavior of couples found to be at risk (clinical utility). The 2-fold aim of this article was to establish the analytical validity of ECS and to quantify how many pregnancies could be impacted by performing ECS in the general US population. Although previous studies have addressed the clinical validity (10) and clinical utility (6, 7) of ECS, we expect future studies to provide further evidence for these separate lines of inquiry.

Author Contributions: *All authors confirmed they have contributed to the intellectual content of this paper and have met the following 3 requirements: (a) significant contributions to the conception and design, acquisition of data, or analysis and interpretation of data; (b) drafting or revising the article for intellectual content; and (c) final approval of the published article.*

Authors' Disclosures or Potential Conflicts of Interest: *Upon manuscript submission, all authors completed the author disclosure form. Disclosures and/or potential conflicts of interest:*

Employment or Leadership: G.J. Hogan, Counsyl; V.S. Vysotskaia, Counsyl; K.A. Beauchamp, Counsyl; S. Seisenberger, Counsyl; P.V. Grauman, Counsyl; K.R. Haas, Counsyl; S.H. Hong, Counsyl; D. Jeon, Counsyl; S. Kash, Counsyl; H.H. Lai, Counsyl; L.M. Melroy, Counsyl; M.R. Theilmann, Counsyl; C.S. Chu, Counsyl; K. Iori, Counsyl; J.R. Maguire, Counsyl; E.A. Evans, Counsyl; I.S. Haque, Counsyl; R. MarHeyming, Counsyl; H.P. Kang, Counsyl; D. Muzzey, Counsyl.

Consultant or Advisory Role: I.S. Haque, Counsyl.

Stock Ownership: G.J. Hogan, Counsyl; V.S. Vysotskaia, Counsyl; K.A. Beauchamp, Counsyl; S. Seisenberger, Counsyl; P.V. Grauman, Counsyl; K.R. Haas, Counsyl; S.H. Hong, Restricted Stock Units; D. Jeon, Counsyl; S. Kash, Counsyl; H.H. Lai, Counsyl; L.M. Melroy, Counsyl; M.R. Theilmann, Counsyl; C.S. Chu, Counsyl; K. Iori, Counsyl; J.R. Maguire, Counsyl; E.A. Evans, Counsyl; I.S. Haque, Counsyl; R. Mar-Heyming, Counsyl; H.P. Kang, Counsyl; D. Muzzey, Counsyl.

Honoraria: None declared.

Research Funding: None declared.

Expert Testimony: None declared.

Patents: D. Muzzey, US20160188793A1.

Role of Sponsor: The funding organizations played no role in the design of study, choice of enrolled patients, review and interpretation of data, or final approval of manuscript.

Acknowledgments: The authors thank Counsyl colleagues for help and advice: members of the CLIA team, including Thi Tran and Jeanette Wong; David Jennions; Kenny Wong; Saurav Guha; Jessica Connor; Dan Davison; Genevieve Haliburton; Andrew Horn; Erik Kase-niit; Brandon Lee; Matt Leggett; Jeff Tratner; and Xin Wang.

References

- Bell CJ, Dinwiddie DL, Miller NA, Hateley SL, Ganusova EE, Mudge J, et al. Carrier testing for severe childhood recessive diseases by next-generation sequencing. *Sci Transl Med* 2011;3:65ra4.
- Costa T, Scriver CR, Childs B. The effect of Mendelian disease on human health: a measurement. *Am J Med Genet* 1985;21:231-42.
- Kumar P, Radhakrishnan J, Chowdhary MA, Giampietro PF. Prevalence and patterns of presentation of genetic disorders in a pediatric emergency department. *Mayo Clin Proc* 2001;76:777-83.
- Haque IS, Lazarin GA, Kang HP, Evans EA, Goldberg JD, Wapner RJ. Modeled fetal risk of genetic diseases identified by expanded carrier screening. *JAMA* 2016;316:734-42.
- Committee opinion no. 690 summary: carrier screening in the age of genomic medicine. *Obstet Gynecol* 2017;129:595-6.
- Ghiossi CE, Goldberg JD, Haque IS, Lazarin GA, Wong KK. Clinical utility of expanded carrier screening: reproductive behaviors of at-risk couples. [Epub ahead of print] *J Genet Couns* September 27, 2017 as doi: 10.1007/s10897-017-0160-1.
- Archibald AD, Smith MJ, Burgess T, Scarff KL, Elliott J, Hunt CE, et al. Reproductive genetic carrier screening for cystic fibrosis, fragile X syndrome, and spinal muscular atrophy in Australia: outcomes of 12,000 tests. [Epub ahead of print] *Genet Med* October 26, 2017 as doi: 10.1038/gim.2017.134.
- Aradhya S, Lewis R, Bonaga T, Nwokekeh N, Stafford A, Boggs B, et al. Exon-level array CGH in a large clinical cohort demonstrates increased sensitivity of diagnostic testing for Mendelian disorders. *Genet Med* 2012;14:594-603.
- Paracchini V, Seia M, Coviello D, Porcaro L, Costantino L, Capasso P, et al. Molecular and clinical features associated with CFTR gene rearrangements in Italian population: identification of a new duplication and recurrent deletions. *Clin Genet* 2008;73:346-52.
- Beauchamp KA, Muzzey D, Wong KK, Hogan GJ, Karimi K, Candille SI, et al. Systematic design and comparison of expanded carrier screening panels. [Epub ahead of print] *Genet Med* June 22, 2017 as doi: 10.1038/gim.2017.134.
- Aziz N, Zhao Q, Bry L, Driscoll DK, Funke B, Gibson JS, et al. College of American Pathologists' laboratory standards for next-generation sequencing clinical tests. *Arch Pathol Lab Med* 2015;139:481-93.
- Rehm HL, Bale SJ, Bayrak-Toydemir P, Berg JS, Brown KK, Deignan JL, et al. ACMG clinical laboratory standards for next-generation sequencing. *Genet Med* 2013;15:733-47.
- Vysotskaia VS, Hogan GJ, Gould GM, Wang X, Robertson AD, Haas KR, et al. Development and validation of a 36-gene sequencing assay for hereditary cancer risk assessment. *PeerJ* 2017;5:e3046.
- Li H. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM. *arXiv* 2013. <http://arxiv.org/abs/1303.3997> (Accessed May 2013).
- Garrison E, Marth G. Haplotype-based variant detection from short-read sequencing. *arXiv* 2012. <http://arxiv.org/abs/1207.3907> (Accessed July 2012).
- McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernysky A, et al. The Genome Analysis Toolkit: a MapReduce framework for analyzing next-generation DNA sequencing data. *Genome Res* 2010;20:1297-303.
- Kaseniit KE, Theilmann MR, Robertson A, Evans EA, Haque IS. Group testing approach for trinucleotide repeat expansion disorder screening. *Clin Chem* 2016;62:1401-8.
- 1000 Genomes Project Consortium, Auton A, Brooks LD, Durbin RM, Garrison EP, Kang HM, et al. A global reference for human genetic variation. *Nature* 2015;526:68-74.
- Chen L, Hadd A, Sah S, Filipovic-Sadic S, Krosting J, Sekinger E, et al. An information-rich CGG repeat primed PCR that detects the full range of fragile X expanded alleles and minimizes the need for southern blot analysis. *J Mol Diagn* 2010;12:589-600.
- Lim GXY, Yeo M, Koh YY, Winarni TI, Rajan-Babu I-S, Chong SS, et al. Validation of a commercially available test that enables the quantification of the numbers of CGG trinucleotide repeat expansion in FMR1 gene. *PLoS One* 2017;12:e0173279.
- The Clinical and Functional Translation of CFTR (CFTR2). <http://cftr2.org> (Accessed June 2017).
- Fokkema IFAC, Taschner PEM, Schaafsma GCP, Celli J, Laros JFJ, den Dunnen JT. LOVD v.2.0: the next generation in gene variant databases. *Hum Mutat* 2011;32:557-63.
- Clopper CJ, Pearson ES. The use of confidence or fiducial limits illustrated in the case of the binomial. *Biometrika* 1934;26:404-13.
- Zook JM, Chapman B, Wang J, Mittelman D, Hofmann O, Hide W, et al. Integrating human sequence data sets provides a resource of benchmark SNP and indel genotype calls. *Nat Biotechnol* 2014;32:246-51.
- Lek M, Karczewski KJ, Minikel EV, Samocha KE, Banks E, Fennell T, et al. Analysis of protein-coding genetic variation in 60,706 humans. *Nature* 2016;536:285-91.
- Monani UR, Lorson CL, Parsons DW, Prior TW, Androphy EJ, Burghes AH, et al. A single nucleotide difference that alters splicing patterns distinguishes the SMA gene SMN1 from the copy gene SMN2. *Hum Mol Genet* 1999;8:1177-83.
- Luo M, Liu L, Peter I, Zhu J, Scott SA, Zhao G, et al. An Ashkenazi Jewish SMN1 haplotype specific to duplication alleles improves pan-ethnic carrier screening for spinal muscular atrophy. *Genet Med* 2014;16:149-56.
- Saul RA, Tarleton JC. FMR1-related disorders. Seattle (WA): University of Washington; 2012.
- Monaghan KG, Lyon E, Spector EB. American College of Medical Genetics and Genomics. ACMG Standards and Guidelines for fragile X testing: a revision to the disease-specific supplements to the Standards and Guidelines for Clinical Genetics Laboratories of the American College of Medical Genetics and Genomics. *Genet Med* 2013;15:575-86.
- de Graaf G, Buckley F, Skotko BG. Estimates of the live births, natural losses, and elective terminations with Down syndrome in the United States. *Am J Med Genet A* 2015;167A:756-67.