

# Genetic Sex Validation for Sample Tracking in Clinical Testing

**Jianhong Hu**

Baylor College of Medicine, Human Genome Sequencing Center (HGSC)

**Viktoriya Korchina**

Baylor College of Medicine, Human Genome Sequencing Center (HGSC)

**Hana Zouk**

Laboratory for Molecular Medicine (LMM), Mass General Brigham

**Maegan V. Harden**

Broad Institute of MIT and Harvard

**David Murdock**

Baylor College of Medicine, Human Genome Sequencing Center (HGSC)

**Alyssa Macbeth**

Broad Institute of MIT and Harvard

**Steven M. Harrison**

Laboratory for Molecular Medicine (LMM), Mass General Brigham

**Niall Lennon**

Broad Institute of MIT and Harvard

**Christie Kovar**

Baylor College of Medicine, Human Genome Sequencing Center (HGSC)

**Adithya Balasubramanian**

Baylor College of Medicine, Human Genome Sequencing Center (HGSC)

**Lan Zhang**

Baylor College of Medicine, Human Genome Sequencing Center (HGSC)

**Gauthami Chandanavelli**

Baylor College of Medicine, Human Genome Sequencing Center (HGSC)

**Divya Pasham**

Baylor College of Medicine, Human Genome Sequencing Center (HGSC)

**Robb Rowley**

National Human Genome Research Institute

**Ken Wiley**

National Human Genome Research Institute

**Maureen E. Smith**

Northwestern University Feinberg School of Medicine

**Adam Gordon**

Northwestern University Feinberg School of Medicine

**Gail P. Jarvik**

University of Washington Medical Center

**Patrick Sleiman**

Children's Hospital of Philadelphia

**Melissa A Kelly**

Genomic Medicine Institute, Geisinger

**Harris T. Bland**

Vanderbilt University Medical Center

**Mullai Murugan**

Baylor College of Medicine, Human Genome Sequencing Center (HGSC)

**Eric Venner**

Baylor College of Medicine, Human Genome Sequencing Center (HGSC)

**Eric Boerwinkle**

Baylor College of Medicine, Human Genome Sequencing Center (HGSC)

**Cynthia Prows**

Cincinnati Children's Hospital Medical Center

**Lisa Mahanta**

Laboratory for Molecular Medicine (LMM), Mass General Brigham

**Heidi L. Rehm**

Broad Institute of MIT and Harvard

**Richard A. Gibbs**

Baylor College of Medicine, Human Genome Sequencing Center (HGSC)

**Donna M. Muzny** (✉ [donnam@bcm.edu](mailto:donnam@bcm.edu))

Baylor College of Medicine, Human Genome Sequencing Center (HGSC)


---

**Short Report**

**Keywords:** Next Generation Sequencing (NGS), Clinical testing, Sex concordance, SNP genotyping

**Posted Date:** September 11th, 2023

**DOI:** <https://doi.org/10.21203/rs.3.rs-3304685/v1>

**License:**  This work is licensed under a Creative Commons Attribution 4.0 International License. [Read Full License](#)

---

# Abstract

## Objective

Data from DNA genotyping via a 96-SNP panel in a study of 25,015 clinical samples were utilized for quality control and tracking of sample identity in a clinical sequencing network. The study aimed to demonstrate the value of both the precise SNP tracking and the utility of the panel for predicting the sex-by-genotype of the participants, to identify possible sample mix-ups.

## Results

Precise SNP tracking showed no sample swap errors within the clinical testing laboratories. In contrast, when comparing predicted sex-by-genotype to the provided sex on the test requisition, we identified 110 inconsistencies from 25,015 clinical samples (0.44%), that had occurred during sample collection or accessioning. The genetic sex predictions were confirmed using additional SNP sites in the sequencing data or high-density genotyping arrays. It was determined that discrepancies resulted from clerical errors, samples from transgender participants and stem cell or bone marrow transplant patients along with undetermined sample mix-ups.

## Introduction

The implementation of next generation sequencing (NGS) technologies in clinical laboratories [1–3] typically involves three phases: (i) the pre-analytic phase including sample collection, DNA extraction and shipment; (ii) the analytic phase of NGS library preparation, DNA sequencing, bioinformatics analysis; and (iii) a post-analytic phase including clinical report generation and delivery. Each phase is inherently subject to sample tracking and identification errors, with prior reports of more than 46% of errors occurring during the pre-analytical phase, caused by inappropriate test requests, order entry errors, patient misidentification, and labelling errors [4]. Validation and tracking of sample identity therefore is a basic and important aspect of effective clinical NGS testing.

DNA-based methods for sample tracking include genotyping of short tandem repeats (STRs) or single nucleotide polymorphisms (SNPs) [5]. STRs are generally located in non-coding regions, prone to high sequencing error rates, and often require longer than typical sequencing read lengths to precisely define the number of repeats, limiting their application. In contrast, SNPs are ubiquitous in the genome and simple to assay [6, 7]. In this study, a 96-SNP panel was used to track samples through the clinical NGS workflow in the National Institute of Health's Electronic Medical Records and Genomics Phase III (eMERGE) program [8]. The network linked together 11 sample collection sites and 2 clinical genetic testing laboratories, the Human Genome Sequencing Center Clinical Laboratory at Baylor College of Medicine (BCM-HGSC-CL) and the Mass General Brigham Laboratory for Molecular Medicine (LMM) in partnership with the Clinical Research Sequencing Platform (CRSP) at the Broad Institute of MIT and Harvard. A total of 25,015 clinical DNA samples were processed. The 96-SNP panel-based procedure provided a robust method for sample tracking in the clinical NGS workflow and showed that the testing of sex can provide a valuable quality control tool.

## Methods

### Fluidigm SNP genotyping assay

Two clinical laboratories harmonized methods for the program[8] and utilized a 96-SNP panel but incorporated different selected SNPs to track samples and determine ancestry. Each 96-SNP panel contained one subset of SNPs on the sex-chromosomes. The rest autosome SNPs are within the target region of the capture design used in the eMERGE program (Supplementary material) [8]. Assays were performed according to the manufacturer's recommendations.

The BCM-HGSC-CL's 96-SNP panel replaced 19 of the original Fluidigm SNPtrace 96 sites to match genomic regions specifically targeted in eMERGE III. The remaining sites included 3 SNPs on Chromosome X and 3 on Chromosome Y[9, 10]. At the Broad Institute, the chosen SNPs included 95 autosomal SNPs and 1 sex determining assay SNP, covering the AMELX and AMELY gene (AMG\_3B) with a sex-specific 6 base-pair insertion/deletion.

## **Illumina Infinium SNP array assays and NGS**

The HumanCoreExome v1-3 BeadChips containing 500K variant sites, including more than 12,900 located on the X chromosome, that are informative for genetic sex prediction, were utilized according to manufacturer's specifications. DNA sequencing for the eMERGE phase III program has been described previously[8].

## **Results**

The BCM-HGSC-CL and LMM/Broad laboratories utilized the same analytical platform foundation, employing slightly different SNP sites for the assays, but generally similar workflows (Fig. 1), to test for concordance between data generated from the 96-SNP panel genotyping and the DNA sequence data. The average SNP call rates were 97.3% and 97.5% for the 25,015 samples processed at the BCM-HGSC-CL and the LMM/Broad, respectively. When comparing the 96-SNP panel genotype-based sex to reported sex at the time of sample accessioning, a total of 110 (0.44%) non-concordant cases from two testing laboratories were identified. The two testing laboratories utilized slightly different workflows to technically validate the sex discrepancies.

At the BCM-HGSC-CL, of the 14,515 samples processed, 73 samples with sex discrepancies were re-tested with the same 96-SNP panel. Identical results were obtained for 70 of the re-tested samples (Table 1). For the remaining 3 cases, where the sex provided on test requisition was male, non-concordant or ambiguous data were observed between the initial and the repeated assays. For two of these samples, the automated software calls from one of each duplicate assays indicated that the DNA source was from individuals with Klinefelter Syndrome (47, XXY). However, further review of the SNP scatter plots for autosome and sex SNPs indicated that the inconsistent sex calls most likely resulted from sample contamination involving a mixture of male and female DNAs (Fig. 2). The third sample was called as female with lower confidence initially. In the repeated assay, one of the X SNPs failed to call due to localization in between clusters in plot analysis. This is most likely due to the female sample mixed up with some DNA sample from another female.

Table 1  
Comparison of genetic sex determined in various assays and reported sex on test requisition

Sequencing site	Total	Sample providing site	Sex on test requisition	Sex from 1st Fluidigm array	Sex from 2nd Fluidigm array	Sex from Illumina array	Sex from sequencing data	Sample number		
BCM-HGSC-CL	73	Site 1	Male	Female	Female	Female	-	5		
			Female	Male	Male	Male	-	5		
		Site 2	Male	Female	Female	Female	-	13		
			Male	Female	Klinefelter	Female	-	1		
			Female	Male	Male	Male	-	7		
		Site 3	Male	Female	Female	Female	-	3		
			Female	Male	Male	Male	-	3		
			Female	Male	Male	NA*	-	1		
		Site 4	Male	Female	Female	Female	-	7		
			Female	Male	Male	Male	-	9		
			Male	Klinefelter	Female	NA*	-	1		
		Site 5	Male	Female	Female	Female	-	6		
			Male	Female	No Call	Female	-	1		
			Female	Male	Male	Male	-	4		
		Site 6	Male	Female	Female	Female	-	4		
			Female	Male	Male	Male	-	3		
		LMM/Broad	37	Site 7	NA **	Female	-	-	Female	1
					Female	NA***	-	-	Male	1
					Male	Female	-	-	Female	16
					Female	Male	-	-	Male	13
Site 8	Male			Female	-	-	Female	1		
Site 9	Male			Female	-	-	Female	1		
	Female			Male	-	-	Male	2		
Site 10	Male			Female	-	-	Female	1		
	Female			Male	-	-	Male	1		

\*: Insufficient gDNA for Illumina array; \*\*: Sex not reported on requisition form; \*\*\*: sex not called in assay; NA: not available

Next, Illumina HumanCore Exome Arrays were utilized as an orthogonal high-density hybridization genotyping assay to further test 71 of the 73 samples with sex inconsistencies except two samples which had insufficient

genomic DNA (Table 1). HumanCore Exome Array results confirmed 96-SNP panel genotyping sex data, including the suspected two contaminated female samples with additional male or other female DNA.

At the Broad/LMM, the reported sex from the test requisition was compared with the genetic sex determined by both the Fluidigm genotyping assay and the data from the eMERGE III sequencing panel. Of the 10,500 samples processed, 151 were initially either identified as discordant or had no sex determination. For 95 samples, the Fluidigm assay data could not return a sex determination, however the sequencing sex matched the reported sex for each and no further action was taken. For 19 of the remaining 56 samples, the sequencing and reported sex were concordant, but did not match the genotyping determined sex. Further review of these 19 samples showed that the genotyping assay calls were generally borderline or low confidence calls, suggesting sub-optimal performance of the single sex determining SNP as the reason for the data discrepancy, rather than either a sex reporting error at accession or sample mix-up in the testing laboratory. The remaining 37 samples had highly confident sex determination calls from both the SNP assay and the subsequent DNA sequencing that were concordant, but did not match the site reported sex (Table 1).

Internal tracking showed that none of the 110 confidently identified sex discrepant samples occur within the clinical DNA sequencing laboratories and that most errors were likely introduced prior to shipment of samples. Sampling sites identified handling errors from test requisitions, sample extraction, and sample handling procedures for 54 cases. Forty-six of these had information that was incorrectly or incompletely entered on the test requisitions and were resolved by examination of other records. In 6 other cases, it was determined that incorrect samples had been shipped from the sampling sites to the genome centers. Biological explanations for the discrepant tracking data were identified for an additional 12 cases. In 4 of these 12 cases, further examination of records revealed that the samples were provided by transgender participants. In addition, 8 sex discrepant samples were determined to be from individuals who had received stem cell or bone marrow transplants. Causes of the sample genetic vs. reported sex discrepancy are listed in Table 2.

Table 2  
Causes of sample sex discrepancy

Sex Discrepant Categories		BCM-HGSC-CL Samples	LMM/Broad samples	Total
Sampling site errors	Incorrect/incomplete information on Test Requisition	45	1	46
	Error during DNA extraction	0	2	2
	Incorrect sample shipped	6	0	6
Transgender		2	2	4
Stem cell/bone marrow transplant recipient		8	0	8
Not solved/under investigation		12	32	44
<b>Total Sex Discrepancies</b>		<b>73</b>	<b>37</b>	<b>110</b>

Where possible, the information on test requisition forms was amended and correct clinical reports were issued for 45 cases processed at the BCM-HGSC-CL, or the incorrect samples were replaced and re-processed. Twelve cases sequenced at the BCM-HGSC-CL with sample-mix ups due to unknown causes were withdrawn from the study. Similarly, 32 unsolved cases sequenced at LMM/Broad were either withdrawn or remain under investigation.

## Discussion

To identify sample swaps during the processing of 25,015 clinical samples in the NIH eMERGE III program, two clinical DNA sequencing laboratories first utilized a Fluidigm-based 96-SNP panel assay to track internal processes. These analyses indicated no sample swaps had occurred in the time interval between sample arrival at the testing laboratories and the delivery of the final DNA sequencing data. In contrast, when the test was expanded to predict the concordance between the self-reported sex of participants at the time of their initial enrollment, with a predicted sex-by-genotype, there were 110 discordant samples. A battery of follow-up tests indicated that these likely arose before the materials were received at the clinical DNA sequencing laboratories. The bases of the sample tracking errors at sample collection sites were determined in 66 of the 110 cases (60%), while leaving the remaining 44 cases unsolved and under investigation. Of these 66 resolved cases, the largest source for the initial discordance occurring in 54 cases (81%) arose from clerical or shipping errors (81%). The remaining 12 cases (18% of the 66 solved) had biological underpinnings that explained the discordant results, as 8 were due to stem cell/bone marrow transplants while 4 were from transgender individuals. Future sample collecting procedures should be modified to ensure that participants are invited to note these types of events at the time of collection, so that this information is available for quality control.

The 96-SNP panel has proven value for precise sample tracking [11]. In general, 20 informative SNP loci are sufficient for unique individual sample identification[12, 13]. Other SNP panels have been used for identification of human samples[6, 14, 15]. A low-density QC genotyping array launched by Illumina which includes 15,949 markers has been utilized in genomic-based clinical diagnostics[16]. Our studies showed that these two different SNP platforms exhibited consistent results when applied for sex identification. In comparison to the use of the Illumina Infinium array platform, the workflow for the 96-SNP panel assay is faster (1-day workflow vs 3-day workflow) and more cost-effective. However, the Illumina Infinium array platform provides more information on linkage analysis, HLA haplotyping, ethnicity determination and other genetic information in addition to fingerprinting and thus may be preferred in some scenarios. Other commercial systems are also available to substitute for the platforms described here if they provide cost-effective and precise data with similar qualities.

This level of tracking error is unacceptable for ongoing clinical practice, but the study does not represent the levels that will be expected in further clinical programs. At least one laboratory declared their initial sample enrollments as 'research samples' and thus committed to later repeat assays under a fully compliant protocol, to verify any findings that may impact care. Others were able to quickly identify points of error and rectify their protocols to ensure faithful future sample handling. All sites committed to rechecking of records and reconciling actionable findings with orthogonal data, including family histories and biochemical tests, before returning results. The 'lessons learned' from these analyses ensure that a repeat of the same program would likely minimize any similar errors.

## Limitations

While false positive rates are low for this application of SNP trace, false negative rates will be high. Here, the overall level of genetic and reported sex discordance of 0.44% is likely an underestimate of the true error rate in this study, as the misclassification of genetic sex from a random sample swap would be expected to result in incorrect, erroneous assignment, only 50% of the time. The true ratio may be skewed by factors introducing a sex-bias in the direction of misclassification. This could be caused by skewed phenotypes of individuals with sex chromosome anomalies or that gender obfuscation may be socially driven in an unequal manner, depending on the gender

identity of the individual. Overall, the rate is likely higher than the 0.44% identified here, but not anticipated to be higher than twice that level.

## Abbreviations

HGSC

Human Genome Sequencing Center

LMM

Laboratory for Molecular Medicine

NGS

Next generation DNA sequencing

STR

short tandem repeat

SNP

single nucleotide polymorphism

PCR

polymerase chain reaction

eMERGE

Electronic Medical Records and Genomics

EMR

electronic medical record

HGSC-CL

Human Genome Sequencing Center Clinical Laboratory

BCM

Baylor College of Medicine

CRSP

Clinical Research Sequencing Platform

STA

specific target amplification

LSP

locus specific primer

IFC

Integrated Fluidic Circuit

ASP

Allele Specific Primers

ACMG

American College of Medical Genetics

NHGRI

National Human Genome Research Institute

IRB

institutional review board

## Declarations



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

The Electronic Medical Records and Genomics (eMERGE) Network is a National Human Genome Research Institute (NHGRI)-funded consortium tasked with developing methods and best practices for utilization of electronic medical record (EMR) as a tool for genomic research. All 11 sample collection sites consented participants under institutional review board (IRB)-approved protocols and the two sequencing centers had IRB-approved protocols that deferred consent to the participating sites. The protocol number for Baylor College of Medicine was (#H-40455).

## **Consent for publication**

Not applicable.

## **Availability of data and materials**

Data are available in dbGaP for controlled public access (phs001616.v1.p1).

## **Competing interests**

JH, DM, MM, RAG, DMM disclose that the Baylor Genetics Laboratory is co-owned by Baylor College of Medicine. EV is cofounder of Codified Genomics, which provides variant interpretation services. DM has received consulting fees from Illumina. The remaining authors disclose they have no competing interests.

## **Funding**

The eMERGE Phase III Network was initiated and funded by the National Human Genome Research Institute (NHGRI) through the following grants: U01HG8657 (Kaiser Permanente Washington Health Research Institute/University of Washington), U01HG8685 (Brigham and Women's Hospital), U01HG8672 (Vanderbilt University Medical Center), U01HG8666 (Cincinnati Children's Hospital Medical Center), U01HG6379 (Mayo Clinic), U01HG8679 (Geisinger Clinic), U01HG8680 (Columbia University Health Sciences), U01HG8684 (Children's Hospital of Philadelphia), U01HG8673 (Northwestern University), MD007593 (Meharry Medical College), U01HG8701 (Vanderbilt University Medical Center serving as the Coordinating Center), U01HG8676 (Partners HealthCare/Broad Institute), and U01HG8664 (Baylor College of Medicine).

## **Authors' contributions**

JH, HLR, RAG, DMM contributed to the study concept and design; JH, VK, HZ, MVH, CK, MES annotated and compiled information regarding sample accessioning; HZ, MVH, DM, EV performed NGS data analysis; NL, MES, GJ, HLR, RAG, DMM provided funding support for the project; Investigation: JH, VK, HZ, MVH, AM, SMH, CK, MES, AG, PS, MK, HB, LM, HLR, RAG, DMM conducted the research and investigation process of sample verification; AB, LZ, GC, DP performed the 96-SNP panel and Illumina array genotyping assay; VK, CK, RR, KW, MM participated in the project administration; MES, AG, GJ, PS, MK, HB, CP provided eMERGE sample collections; JH, MM, EV, HLR, RAG, DMM supervised the studies; JH, HZ, MVH, HLR, RAG, DMM were the major contributors in original draft writing; JH, HZ, MVH, DM, AM, SMH, NL, RR, KW, AG, GJ, PS, MK, HB, MM, EV, EB, CP, LM, HLR, RAG, DMM participated in manuscript revision. All authors read and approved the final manuscript.

## **Acknowledgments**

We thank all eMERGE Phase III Network participants for their engagement in this research effort.

## References

1. Norton N, Li D, Hershberger RE. Next-generation sequencing to identify genetic causes of cardiomyopathies. *Curr Opin Cardiol*. 2012;27(3):214–20.
2. Ku CS, et al. Exome sequencing: dual role as a discovery and diagnostic tool. *Ann Neurol*. 2012;71(1):5–14.
3. Yang Y, et al. Clinical whole-exome sequencing for the diagnosis of mendelian disorders. *N Engl J Med*. 2013;369(16):1502–11.
4. Hammerling JA. A Review of Medical Errors in Laboratory Diagnostics and Where We Are Today. *Lab Med*. 2012;43(2):41–4.
5. Butler JM. *Short tandem repeat analysis for human identity testing*. *Curr Protoc Hum Genet*, 2004. Chapter 14: p. Unit 14 8.
6. Pengelly RJ, et al. A SNP profiling panel for sample tracking in whole-exome sequencing studies. *Genome Med*. 2013;5(9):89.
7. Yousefi S, et al. A SNP panel for identification of DNA and RNA specimens. *BMC Genomics*. 2018;19(1):90.
8. eMerge C. Harmonizing Clinical Sequencing and Interpretation for the eMERGE III Network. *Am J Hum Genet*. 2019;105(3):588–605.
9. Pakstis AJ, et al. SNPs for a universal individual identification panel. *Hum Genet*. 2010;127(3):315–24.
10. Nassir R, et al. An ancestry informative marker set for determining continental origin: validation and extension using human genome diversity panels. *BMC Genet*. 2009;10:39.
11. Liang-Chu MM, et al. Human biosample authentication using the high-throughput, cost-effective SNPtrace(TM) system. *PLoS ONE*. 2015;10(2):e0116218.
12. McGuire AL, Gibbs RA. Genetics. No longer de-identified. *Science*. 2006;312(5772):370–1.
13. Lin Z, Altman RB, Owen AB. Confidentiality in genome research. *Science*. 2006;313(5786):441–2.
14. Miller JK, et al. Use of Sequenom sample ID Plus(R) SNP genotyping in identification of FFPE tumor samples. *PLoS ONE*. 2014;9(2):e88163.
15. Castro F, et al. High-throughput SNP-based authentication of human cell lines. *Int J Cancer*. 2013;132(2):308–14.
16. Ponomarenko P, et al. Clinical utility of the low-density Infinium QC genotyping Array in a genomics-based diagnostics laboratory. *BMC Med Genomics*. 2017;10(1):57.

## Figures

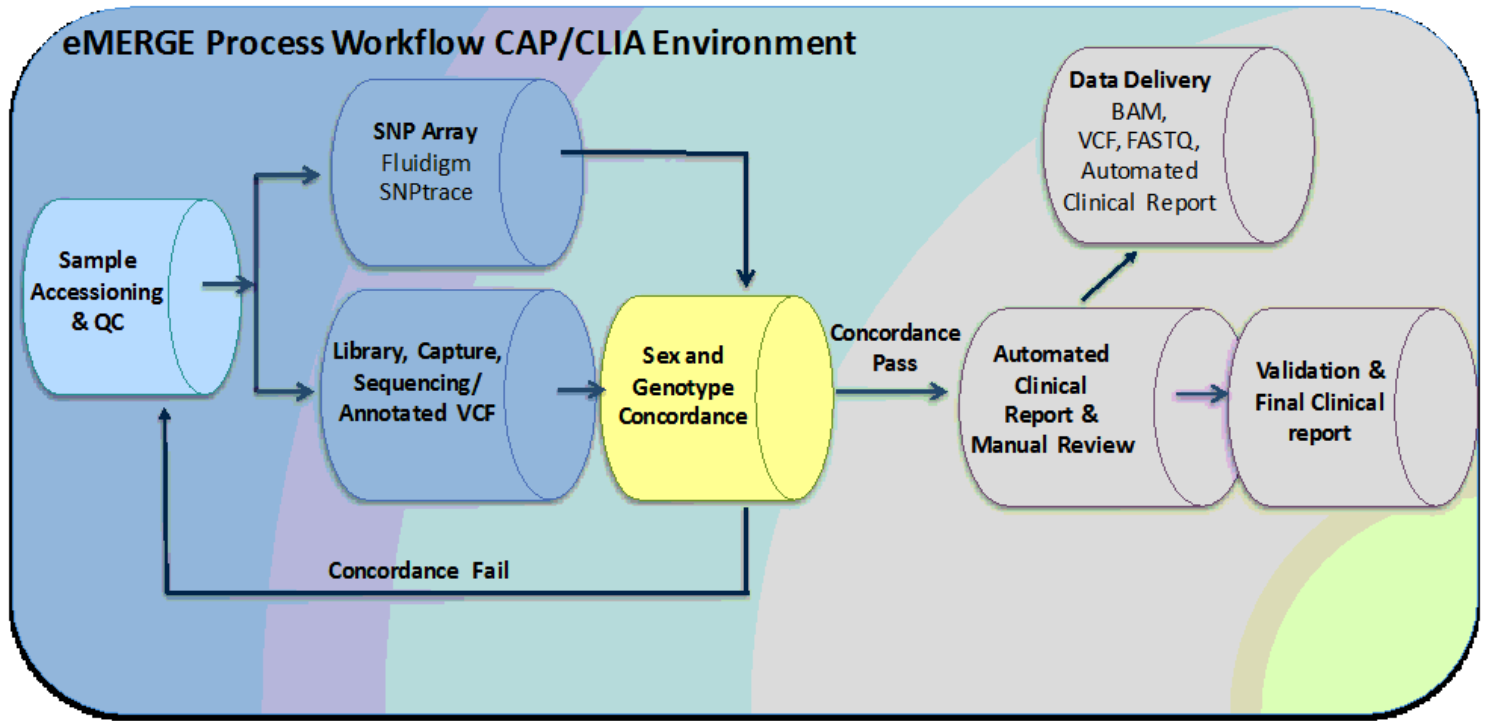


Figure 1

eMERGE sample processing workflow

Steps indicating where aliquots of DNA are taken from samples that are presented to the Clinical DNA Sequencing Laboratory for accession, to test via the Fluidigm 96-SNP panel assay. Data from the Fluidigm 96-SNP panel assay are compared with DNA sequence data from the DNA sequencing pipeline as a quality control step, ahead of the Automated Clinical Reporting step.

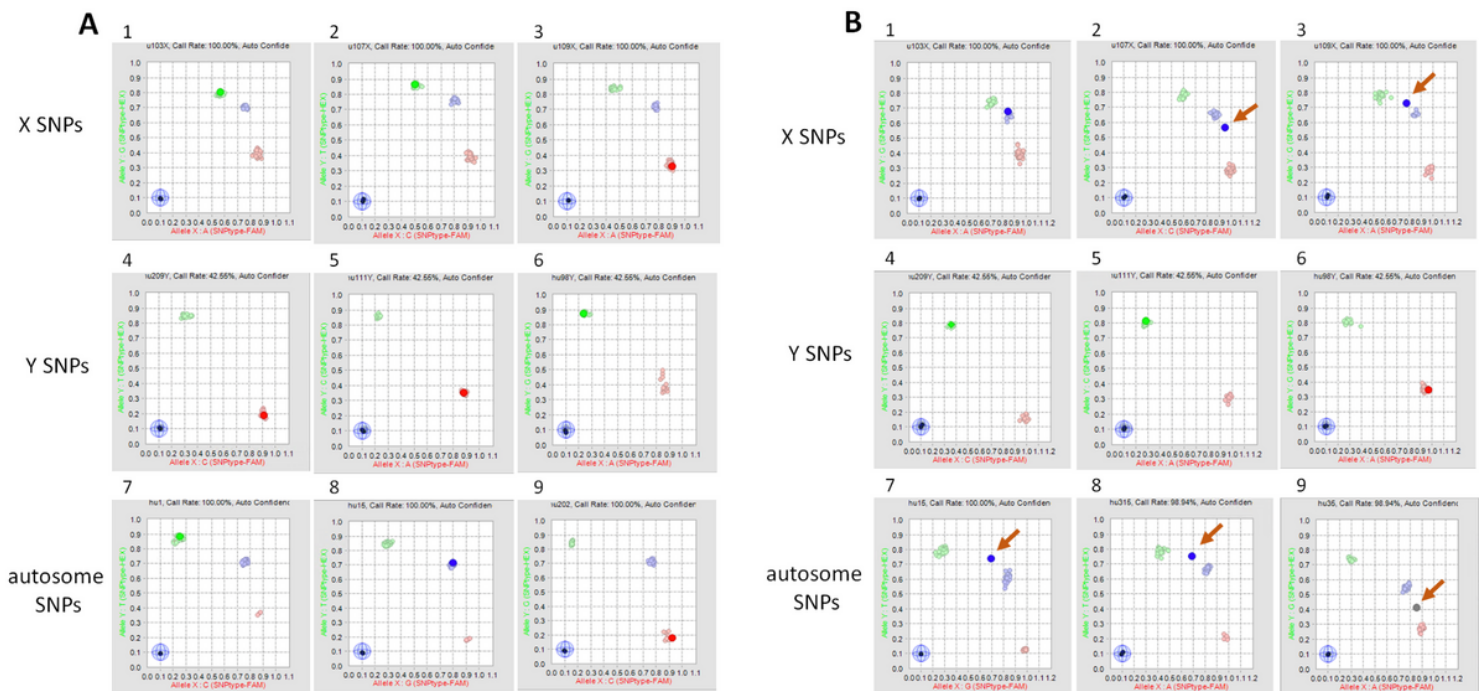


Figure 2

## Scatter Plot Analysis of 96-SNP Panel Reveals Sample Contamination

Scatter plot analysis from vendor software, showing a normal DNA male sample (A) or a contaminated sample containing a mixture of male and female DNAs (B). Panel 1-3: SNPs on X Chromosome; Panel 4-6: SNPs on Y Chromosome; Panel 7-9: autosomal SNPs. Each panel shows the data from a single SNP, as compared to clusters from all other SNPs. Clusters are shown as either homozygous (red or green), or heterozygous (blue) positions. In Panels B2, 3, 7-9, single SNPS are represented as outside the expected (arrows) resulting in erroneous or 'no-call' from the software.

## Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [TableS1.xlsx](#)
- [TableS2.xlsx](#)