



## **VIRTUAL GENETICIST™**

**Improving Transparency,  
Collaborative Variant Interpretation,  
and High-throughput Diagnostics:  
A Clinician's Review of the VIRTUAL GENETICIST Platform**

**Talk by Dr. Adrienne Elbert, MD  
British Columbia Children's Hospital  
American Society of Human Genetics (ASHG)  
2023 Annual Conference**



## Introduction

I'm Adrian Albert and I'm a clinical geneticist in Vancouver, Canada. I don't have any industry funding or compensation of any sort for this study.

And I am here today to talk about Virtual Geneticist, which is an AI-based variant prioritization platform that we use in the clinic to help make additional diagnoses in those individuals who had undergone whole exome sequencing and had not received any diagnosis from the lab.

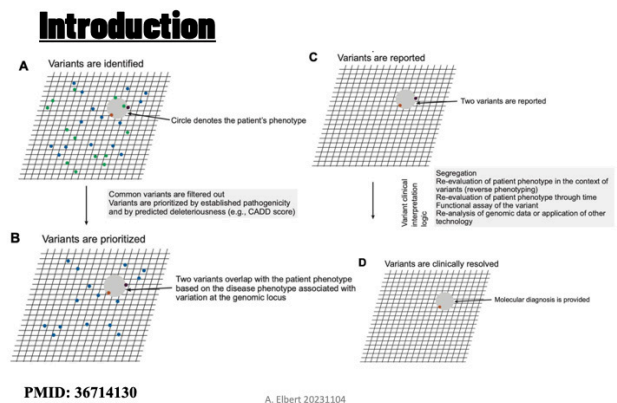
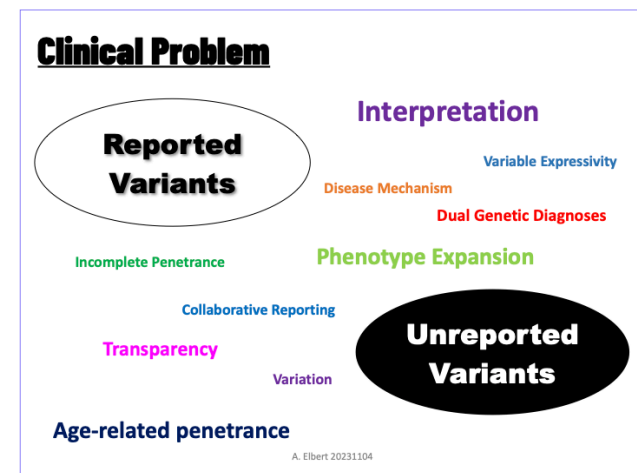
So as a clinical geneticist, the bottom line is you see children and their families. That's certainly true, but we also see people of all ages too, so our clinical cohort included everything from prenatal to pediatric to teens to adults.

It's an evolving situation because in the past we would just get the report back from the lab. And if there's nothing there that makes sense with what we're seeing in front of us in the context of the patient, then that's the end of the line.

As things have evolved, there is now more access to the big genomic data that the clinical laboratories produce and they are increasingly able to provide that data back to us, so it makes sense for us to get the data and actually look at it ourselves, especially when we are highly suspicious of an underlying genetic cause for what we are seeing with our patients.

I think the challenges are really with creating the internal structures of having the right tools. When you see that there is actually a higher diagnostic yield from doing that internal re-analysis, that really pushes the boundaries of what should be standard clinical practice. Whereas before, there was this review of what was found on the clinical report and that was it. Now we know there's a deeper ocean of variants that could be causal and we just don't know what they are because we're not going in there to look for them ourselves.

And this new challenge really promotes the idea of collaboration with the lab and with each other as clinicians, working together to improve what we report out and also the information that we provide to our patients and their families. And so both sides (the clinicians and the lab) will end up getting better at this process of analysis and the inputs required. And overall, the main goal is, of course, to improve the number of individuals we can diagnose, which is getting more and more important as there are more therapies emerging for treating rare disease.



I've had so many times where I've had a patient and I'm so suspicious for a genetic disorder and have gotten a negative report back from the lab and just been absolutely baffled about how it's possible we didn't get any answer out of the testing.

And so, if you've been in that position, this talk will be really helpful. So I will be presenting our unpublished results from a collaboration with Breakthrough Genomics and our evaluation of Virtual Geneticist, which is an AI-driven variant prioritization platform.

So the clinical problem that we faced is really in that scenario, what leads to some variants being reported and not others?

## Complexity of Variant Interpretation

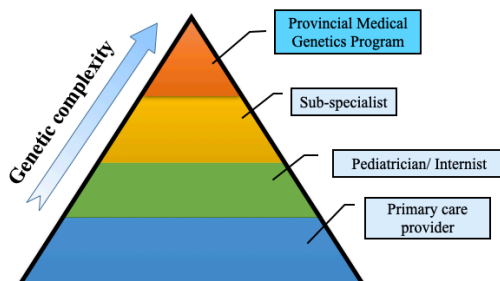
Variant interpretation is becoming increasingly more and more complex. And it's not just that, it's also that there's no transparency around the pipelines that the clinical labs are using to report out the variants. So we don't understand what the limitations and strengths of their pipeline are.

In current clinical genetics practice, the thousands of detected variants go through filtering and prioritization, and this could be based on population alone frequencies as well as predicted and demonstrated deleteriousness. And then of course patient disease phenotype matching with the gene disease phenotype. And then we hope that this leaves us with just a handful of variants that we can manually evaluate less than ten or so and then determine whether in the context of the patient we actually have found the disease-contributing variant or not.

And of course this is the ideal scenario where we're able to resolve based on the variants reported. But often it is a very important balance between reporting out variants that are helpful to the clinician and the patient and avoiding reporting of variants that are not helpful or there's associated burden with interpreting additional variants.

And I should start also by mentioning that I work in the top tier of the subspecialty genetics. So in British Columbia, the way the genetic care is organized is that in the provincial medical genetics program, we are often seeing the more complex individuals that the primary care providers, the pediatrician, internists and subspecialists have not been able to diagnose. And so that does introduce additional complexity in interpreting their genomic and phenotypic data.

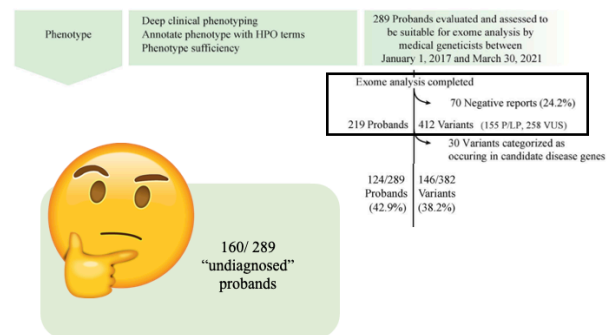
## Introduction



A. Elbert 20231104

And just to give you an idea, after 289 pro-bands that were evaluated by whole exome sequencing and 219 reports with 412 variants, 80% about of the pathogenic and likely pathogenic variants supported a diagnosis, and about 20% of the variance of uncertain significance reported contributed to a diagnosis after further segregation or phenotypic evaluation. This led us to wonder about that remaining 160 undiagnosed pro-bands. Is it possible that their variant is just below the reporting threshold? Maybe it's considered a variant of uncertain significance, or maybe it wasn't prioritized by a pipeline that we don't really understand that well, But there must be something more we can do for those patients.

## Introduction



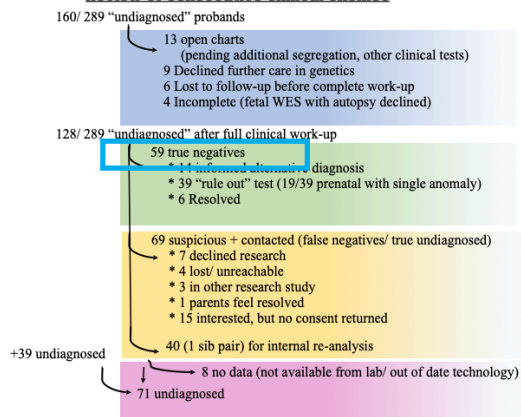
PMID: 35442193

A. Elbert 20231104

I reviewed those 160 undiagnosed pro-bands and after removing those that were not appropriate for further analysis at the time because they had declined follow up or were still being evaluated. I could say that about half of them were probably true negatives, so their negative zone represented a patient that was not suspicious for a genetic disease because either it was a rule out test, for example, prenatal with a single anomaly that then after birth didn't have any additional complexity or the exam informed an alternate diagnosis.

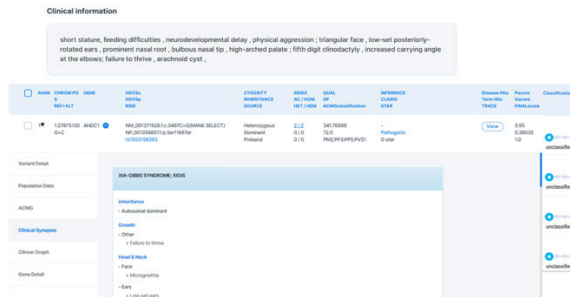
For example, an inflammatory or autoimmune condition, or the patients were in a state where the phenotype was evolving and it was no longer concerning or had resolved. But the other half were still suspicious for genetic disease and therefore we contacted those individuals and reached out to see if they would like to participate in having their raw data re-analyzed internally. And we also added to this another additional group of undiagnosed patients from the department through other clinicians to have a clinical cohort of 71 patients.

### Review of consecutive clinical exomes



VIRTUAL GENETICIST™

- Natural language processing for phenotypic terms
- AI-based variant prioritization

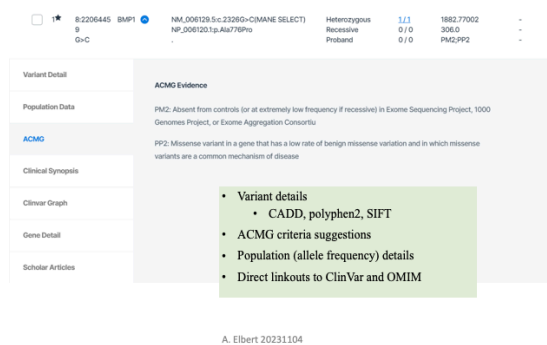


For this re analysis pipeline, we also wanted to include some positive controls, so we collaborated with a local research endeavor called the CAUSES study that had set out and previously published their experience with over 500 patients with rare disease. And this allowed us to use their diagnosed patients as a training cohort or positive controls. And the question really was, could we contribute to more diagnoses both in their undiagnosed group and our uninformative exome clinical group, but also can an AI-based tool like Virtual Geneticist help to identify the misdiagnoses? And what were the underlying factors that may have led to the diagnoses not being on the report in the first place?

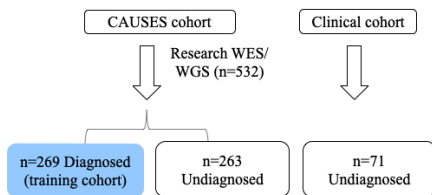
You can expand on the information that are contained in all of these categories. So for instance, this first variant here, you can expand it to see what that condition is, phenotypically in OMIM as well as link out to that specific database. And then here, for example, you can also see the suggested ACMG criteria and link out directly to some of those databases that we use so frequently in this field.



VIRTUAL GENETICIST™



### Designing a Re-analysis Pipeline



**Purpose:**  
End the diagnostic odyssey through re-analysis of genomic data

- Questions:**
- 1) Can internal re-analysis of unsolved cases improve diagnostic rates?
  - 2) Can an AI-based tool (Virtual Geneticist) help to identify missed diagnoses?
  - 3) What factors lead to missed diagnoses?

### Introduction to Virtual Geneticist

Virtual Geneticists is a natural language processing tool that incorporates phenotype as well as deep learning in both interpreting genetic variants and also prioritizing them. And I'll just give you a little bit of an idea of what that looks like. So in this case, we've just inputted terms, but it can also input paragraphs of clinical information and draw out phenotypic terms from that as well. And then you will have your variance in descending rows ranked with annotations in the columns. And those annotations can be various database information as well as ClinVar and HGMD criteria suggestions.

### High Throughput Tool

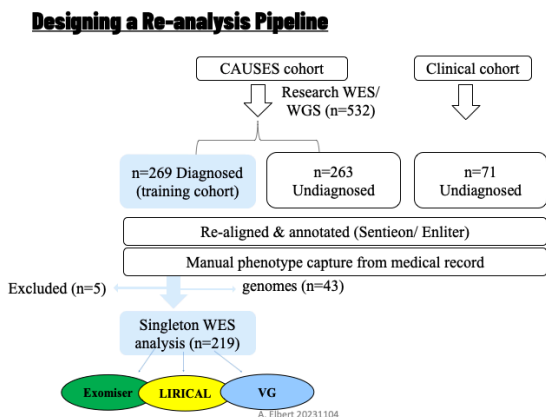
This is a really high throughput tool. So when we submitted our training cohort as singletons and trios, so a total of about 800 cases, this was analyzed within 5 hours. And the running time for each case is really rapid. So when I get negative data now, I can upload a VCF file with phenotypic information and have variants ranked within 10-15 minutes, depending on the size of the files.



VIRTUAL GENETICIST™

- Configured with a parallel processing capacity of 20 cases
- Submission of each case is stream-lined, with a time interval of 45 seconds
- The running time for each case was 10-15 minutes
- The 800 cases were analyzed within 5 hours.
- The platform can be further configured to accommodate a higher maximum parallel processing capacity

## Study Overview

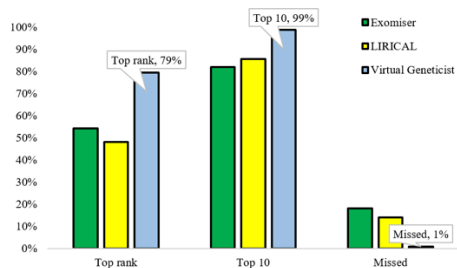


In order to compare Virtual Geneticist, we used other tools that are frequently used for varying prioritization in the field - Exomizer and Lyrical - for our singleton exams, of which there were 219 in this training cohort. So in blue you see Virtual Geneticists compared to in green Exomizer and Lyrical in yellow. And in that training cohort, virtual geneticists ranked the as the #1 top variant, the known diagnosis, 79% of the time. And within the top ten, 99% of the time, against the known diagnostic variant from all of these cases.

## Results of Analysis in Training Cohort

### Virtual Geneticist:

Ranks Diagnosis in top 10 for 99% of cases



Two "Missed" diagnoses were ranked 16 and 20 (re-ranked to 15 and 15)  
- Explored further: both suspected twosies, with part of phenotype unexplained

A. Elbert 20231104

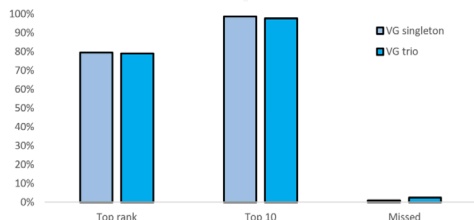
What was interesting for me, because I had always as a clinician felt that trio was better, is that actually Virtual Geneticist did not require a trio to rank the deleterious variants, even in cases where we would have not known it to be the one without determining that to be de novo. So it was still accurately predicting the deleterious variants with just singletons.

So we see Singleton and the pale blue versus the trio analysis and the darker blue and the denominator changes a bit here because we didn't have trios for all of the pro bands.

But it really was important for us to know because in our province singleton is much more fundable by the province than doing a trio. You can see compared to Exomizer which does have a modest benefit from incorporating a trio, that there is a difference there with Virtual Geneticist compared with other tools.

I think the other interesting thing for me was I didn't really appreciate that there are variants lost during the trio process that are filtered out because of things like incomplete penetrance. With those variants, some of those hard filtering barriers used in other pipelines will get your variant removed from the actual data set. And we had quite a few mosaic parents as well in the CAUSES study. And so some of those variants would have been filtered out in trio analysis.

### No advantage of "trio" in identifying diagnostic variant



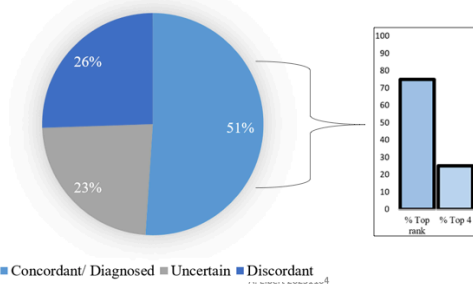
## Results of Analysis in Undiagnosed Cases

Moving on to the undiagnosed and what we found there. So again, we have two cohorts of undiagnosed patients. We've completed the exams. So with both the CAUSES study which is a research cohort as well as the clinical cohort, I wanted to first start with the uninformative results where there was a suspicious variant of uncertain significance (VUS) reported out. And what I wanted to know was whether Virtual Geneticist agreed with the team of clinicians and researchers that were suspicious of that variant.

And what we found is that 51% of those reported VUS variants were concordant with what VG thought was the best candidate diagnosis. And this is important because in determining whether a variant could be the diagnostic variant, it is really helpful, as an independent line of evidence, that a separate tool or pipeline also ranked that diagnosis as highly suspicious or most likely in the dark blue. I should also say that of those variants, three quarters of them were actually the top ranked variant and the other quarter were in the top four.

## Diagnoses: contextualizing VUS

CAUSES research: reported VUS  
 • 51% concordance with VG analysis



So there were also some variants which remained concordant in the gray with what we thought was the most likely. But there was still insufficient evidence to know for sure if they were disease contributing. And then there were an additional quarter which were non concordant, which either suggested a new diagnosis or the exome was negative and didn't have any good candidates.

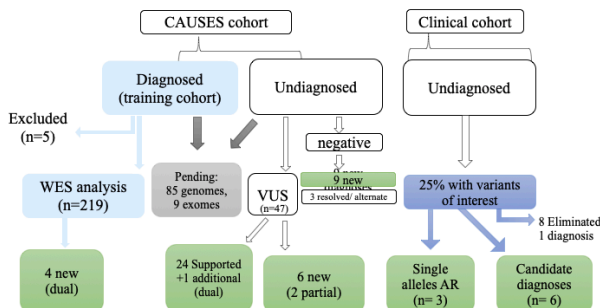
So to summarize our new diagnoses in these groups, we actually had several. So in addition to the 24 reported variance of uncertain significance from a of study that was supported by Virtual Geneticist as probable diagnoses, there was one additional diagnosis in one of those cases, which was a dual genetic diagnosis, and six in the in the discordant group of which two are partial explanations suggesting more complexity or potentially a dual genetic diagnosis. And then in the negative group, there were nine new diagnoses. And upon reviewing their clinical chart, some of those cases had actually resolved or had alternate diagnoses identified in the interim. And three is likely an underestimate.

In the training cohort, we also identified some new diagnoses, which was a bit surprising. So we found additional dual genetic diagnoses and then in our clinical cohort we found 25% of them had a variant that was not reported that was clinically suspicious. And after further evaluation, either by segregation or biochemical studies or reverse phenotyping, many of those were eliminated, but some of them led to additional candidate diagnoses and still are currently undergoing workup or new diagnoses.

But essentially what we believe is there's probably 10% to 15% false positive or false negative rate in these exomes where the diagnostic variant is not reported. And just to be conservative, I'll say 10%. And what's really important is that in the clinical cohort, these were not in genes that were new disease genes. These were variants that were in known disease genes and just had not been reported because of differences in the integration of the levels of evidence and the thresholds for reporting different genes by the different labs.

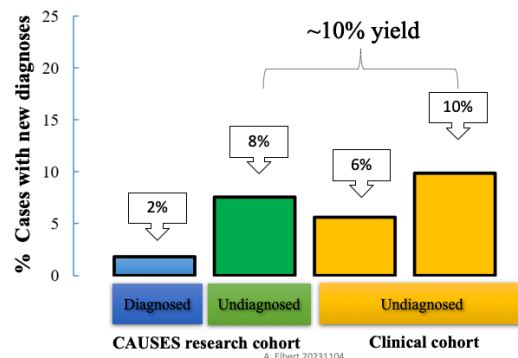
To reflect more on that clinical cohort. I had wondered whether the highly suspicious ones were more likely to have new diagnoses, and to some degree that's true. But there were also new diagnoses made in the cases that I was less suspicious of. Which is important as we try to triage which cases should go for internal, immediate reanalysis. And so why were some of those diagnoses missed?

## New Diagnoses



A. Elbert 20231104

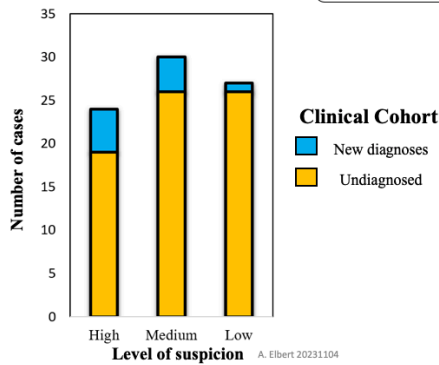
## New Diagnoses



One of the reasons is related to disease mechanism and these genes that have multiple diseases and gene phenotype dyads associated with them. And often the mechanistic information as to how a variant causes one phenotype and not the other is not clear and not easily accessible and requires a lot of reading and manual integration of the different lines of evidence.

## Reflection on Results

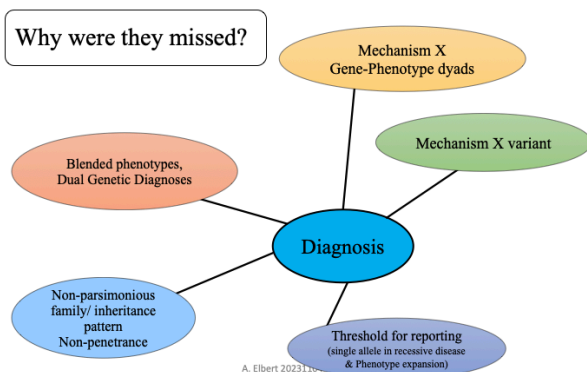
### Reflection on results:



So for example, it could be falsely concluded that a splicing variant causes loss of function when, in fact, it may cause one exon to be skipped, which is a regulatory element and have a different mechanism of disease, or it can also cause an atypical presentation of the disease.

And what we also observed was that there were more dual genetic diagnoses than we previously expected. So definitely the complexity of blended phenotypes results in more challenges with regards to reporting. And this might be due to the threshold for a phenotype gene match because it's only partial, or it may be below what a laboratory thinks is a good enough match to report out.

### Reflection on results:



And then the other types of problems include non-parsimonious families. So when we decide that a parent and child have the same phenotype or the same cause, take for example neuropathy, when in fact, they do not and the inheritance pattern suggests that there might be a parsimonious answer, but it is actually two different diagnoses.

And then the other types of problems include non-parsimonious families. So when we decide that a parent and child have the same phenotype or the same cause, take for example neuropathy, when in fact, they do not and the inheritance pattern suggests that there might be a parsimonious answer, but it is actually two different diagnoses. So when there is non-penetrance and we expect there to be a de novo variant and then the threshold for reporting both from the perspective of atypical phenotypes, blended phenotypes as well as single yields and recessive genes where we may have missed the alternate explanation.

### Blended Phenotypes

- Increased complexity of cases (more dual genetic diagnoses) appreciated on re-analysis
  - Clinical cohort has 20% (n=14 /71) expected
- Virtual Geneticist detects dual genetic diagnoses; both variants still rank within the top ten for VG

I wanted to stress that we had an increased complexity in our cases as more dual genetic diagnoses were appreciated on re-analysis, and that was 20% in our clinical cohort, which is a huge number. And we think this is partially because of the fact that we are on that top of that triangle referral system in Canada and the fact that our group is receiving more complex cases, as well as we think the negative cohort enriches for complex cases, making these inherently more difficult to solve and more likely to be negative because of the way many laboratory pipelines currently rely on parsimony for analysis and solving.

What was really interesting was that Virtual Geneticist does detect dual diagnoses, so both variants would still rank within the Top 10 in these cases. And of course, we don't have the denominator. We don't know what we're missing there in terms of dual genetic diagnoses. But because we have such a high number, it's very interesting to see that Virtual Geneticist definitely performed better than we did as clinicians.

And then I think there's a lot of things that are challenging to describe, except by case example, of why Virtual Geneticist was really helpful in the clinic.

## Case Examples

So, for example, this was an atypical case of a patient with Mirage syndrome (see slide below) and a pathogenic variant, and in particular SAMD9 was recurrent and a novel variant. And the patient was atypical in the sense that they didn't have cytokine and what you can see is a comparison between a set of phenotypic terms, HPO Set 1 and HPO Set 2, where HPO Set 1 is actually quite rich, but actually results in the diagnostic variant being ranked below that top ten and maybe leads to it being missed in reporting.

Another thing is that when we get back these views reported on the clinical exams, we were able to see how Virtual Geneticist does with regards to its ranking in comparison to these other tools and that was also helpful evidence in contextualizing that variant of uncertainty significance (VUS).

And the ability as a clinician to have a tool like Virtual Geneticist where we can play around with and manipulate those phenotypic terms and see how it influences that ranking of variants is very important.

### Choice of phenotypic terms influence ranking

<b>SAMD9-related MIRAGE syndrome</b>	Ranks phenotype 1	Ranks phenotype 2
<b>HPO Set 1</b> proportionate short stature, microcephaly, delayed dental eruption, intrauterine growth retardation, bronchiectasis, hypoplasia of dental enamel, ambiguous genitalia, stage 5 chronic kidney disease, tubulointerstitial nephritis, glomerulosclerosis, pain insensitivity, intellectual disability	ATM	MAGEL2
	GDF6	DALRD3
	GLI2	MYO18B
	MAGEL2	GLI2
	RELN	TOP2B
	CBS	ATM
	FOXJ1	UPB1
	WASF1	CBS
	MYOCD	<b>SAMD9</b>
	SLC12A6	SAMD9
	TOP2B	CACNA2D2
	CACNA2D2	CACNA2D2
	CACNA2D2	WFS1
	UPB1	SYCP2*
	DOCK7	GDF6
WFS1	XYLT1	
<b>HPO Set 2</b> global developmental delay, short stature, hypospadias, cryptorchidism	LMOD1	NF1
	<b>SAMD9</b>	WASF1
	SAMD9	GFM1
	MYO18B	SLC1A4

For example, this variant on the next slide was on a clinical report we received. It was a possible splicing variant and turned out to be maternally inherited. And what you can see is Virtual Geneticists ranked that at 61. Exomizer had it at 25 and Lyrical, had no posterior probability score for this variant. And so overall these were really helpful in establishing my level of suspicion for this variant.

### Contextualization of reported VUS

Some of the reported variants do not rank well by alternate variant prioritization pipelines, and they are thus less compelling candidates

Clinical report:		
<i>KDM6B</i> c.4469-4G>C, (maternally-inherited)		
	<b>KDM6B rank</b>	<b>KDM6B score</b>
VG	61	0.3660
Exomiser	25	0.4587
LIRICAL	N/A	N/A

A. Elbert 20231104

And so overall, these were really helpful in establishing my level of suspicion for this variant and that's critical when I discuss it with the family. And what I'm really getting at is that there are different flavors of negative or uninformative exome reports, and it really depends on what your prior suspicion was going into the testing. So where we very suspicious for a genetic disorder to begin with, in which case it's more likely a false negative, or whether there was a candidate present in the data when we looked at the raw data or not.

### There are different types of negative

- Genetic counseling considerations for non-diagnostic genomes/ exomes:
  - Prior suspicion for genetic disease
  - Phenotype: age, complexity, specificity
  - Candidates present in data or not
  - False negative rate of 10% vs. true negative

And I think that this has really important genetic counseling considerations for the negative exomes. And even if we know that this is not necessarily reassuring, we can color our counseling more with this additional information.

The takeaway is that immediate re-analysis, or re-analysis that's not reliant on knowledge evolution, but immediate re-analysis of non-diagnostic data through alternate platforms like Virtual Geneticist is certainly beneficial. On top of that, even looking at the data for a negative report allows the clinician to value that negative result for what it is. If there's no candidate variants in there, then it's a truly negative result and might mean the patient needs additional testing through a different modality or technology.



## Clinical Lessons

So the questions we had set out to ask at the beginning were:

Can internal re-analysis of unsolved cases improve diagnostic rates? And yes, we found additional diagnoses in about 10%, really 10 to 15%.

And can an AI-based tool like Virtual Geneticist help to identify misdiagnoses? Yes.

And really it is multimodal because it's not just allowing you to explore unreported variants, but it is allowing you to play with the phenotype and the data and see what the variant ranking relationships are, as well as contextualize the reported variants of uncertain significance and how they rank in an alternative pipeline. This exposes some of those weaknesses and strengths of the other labs' pipelines which are unknown to us, and it also allows us to contextualize a negative result when there's no real candidates in that report, in that data.

### **Questions answered:**

- 1) **Can internal re-analysis of unsolved cases improve diagnostic rates?**  
YES! False negative rate of ~10%
- 2) **Can an AI-based tool (Virtual Geneticist) help to identify missed diagnoses?**  
YES!  
Virtual Geneticist is powerful in detecting disease-causing variants, even in complex cases.  
With Virtual Geneticist, the clinician can:
  - Explore unreported variants/ new diagnoses
  - Explore phenotype/ variant ranking relationships
  - Contextualize reported VUS
  - Contextualize non-diagnostic results
- 3) **What factors lead to missed diagnoses?**  
Differences in integration of data about the variant, phenotype, family, disease, and its mechanism, and thresholds for reporting

A. Elbert 20231104

So, ultimately, there's lots of benefit to be gained from more collaboration between the clinic and the lab, and this is what Virtual Geneticist enables and facilitates us to do.

So when we do see a variant in Virtual Geneticist that was not reported in the data, we are able to contact the laboratory and discuss it with them, and that will allow us to have more diagnoses for our patients.



VIRTUAL GENETICIST™

## Acknowledgements

Breakthrough Genomics (Virtual Geneticist)

L. Li, N. Zhang, C. Bui, S. Padhi

CAUSES study

I. S. Rajan Babu, S. Adam, J. M. Friedman

PMGP (British Columbia)

C. F. Boerkoel, S. Huynh, H. Lee, L. Warnock

Clinicians of cases: L. Armstrong, K. Seath, L. Clarke, C. Chang, M. Van Allen, M.S. Patel, E. Digby, S. Langlois, A. Swenerton

Trainees: K. Al-Mamaar, S. J. Norte-Tangkpanya, H. Al-Khafaji

The patients and their families

A. Elbert 20231104



BREAKTHROUGH  
GENOMICS

Proudly Provided by