

How scientists are learning to predict your future with your genes

But what are the limits?

By Brian Resnick | @B_resnick | brian@vox.com | Aug 23, 2018, 9:10am EDT

Graphics and illustrations by Javier Zarracina

When the **Human Genome Project** — the massive endeavor to map all the genes that make humans human — was completed in 2003, scientists were elated.

Finally, they thought, it would be possible to find genes that cause or contribute to devastating diseases like diabetes. And, they thought, it would be relatively simple. In 2004, Francis Collins, a key collaborator on the project and now the head of the National Institutes of Health, **said** he expected there would be 12 genes for Type 2 diabetes, “and that all of them will be discovered in the next two years.”

It didn’t work out that way.

There’s a simple understanding of genetics that we learn early on in school. Depending on what version of the gene we inherit — dominant or recessive — certain traits simply get

turned on and off, like a light switch; think **Gregor Mendel and the pea plants**. But often, when it comes to humans, with our massive genome of 3 billion base pairs, the light-switch analogy pretty much falls apart.

In the years since the Human Genome Project, geneticists have learned that many of the traits that make you *you* arise from a stunningly complex constellation of genes — numbering in the hundreds, if not thousands — that interact with each other, the body, and the environment in incredibly complex ways.

“AT EVERY FORK IN THE ROAD... YOUR GENOME HAS ITS THUMB ON THE SCALES.”

The same goes for the genetics that put us at risk for diseases like diabetes. Collins’s prediction was way off: There are **hundreds** of locations in the genome that influence diabetes. The risk for the disease, in research speak, is polygenic: Each of the numerous genes **changes the odds of developing diabetes minutely**. There’s not one light switch but hundreds, each with the ability to slightly increase or decrease the chances of developing the illness.

The journal ***Nature Genetics*** recently published an enormous study demonstrating yet again how multiple sites on the genome can play a role in determining our fates. The meta-analysis assessed the genomes of 1.1 million people of white, European ancestry, looking at the 10 million spots where DNA sequences vary — where you might have the letter A (for the nucleotide adenine) but I have the letter C.

The question these scientists were asking: How many of these genetic variants correlate with our likelihood to stay in school? And can differences in our genetics predict who completes high school, and who completes college?

To a small degree, their analysis showed that some places on the genome are associated with educational attainment. In all, it found 1,271 spots in the genome that were significantly correlated with a greater number of years in school. What’s more, the researchers also showed that they could use these 1,271 spots in the genome to compute a score that predicts — mildly, and on average across a group — their likelihood of completing college. That’s all from a cheek swab.

How is this possible? The research technique used here is called a genome-wide association study, or GWAS for short.

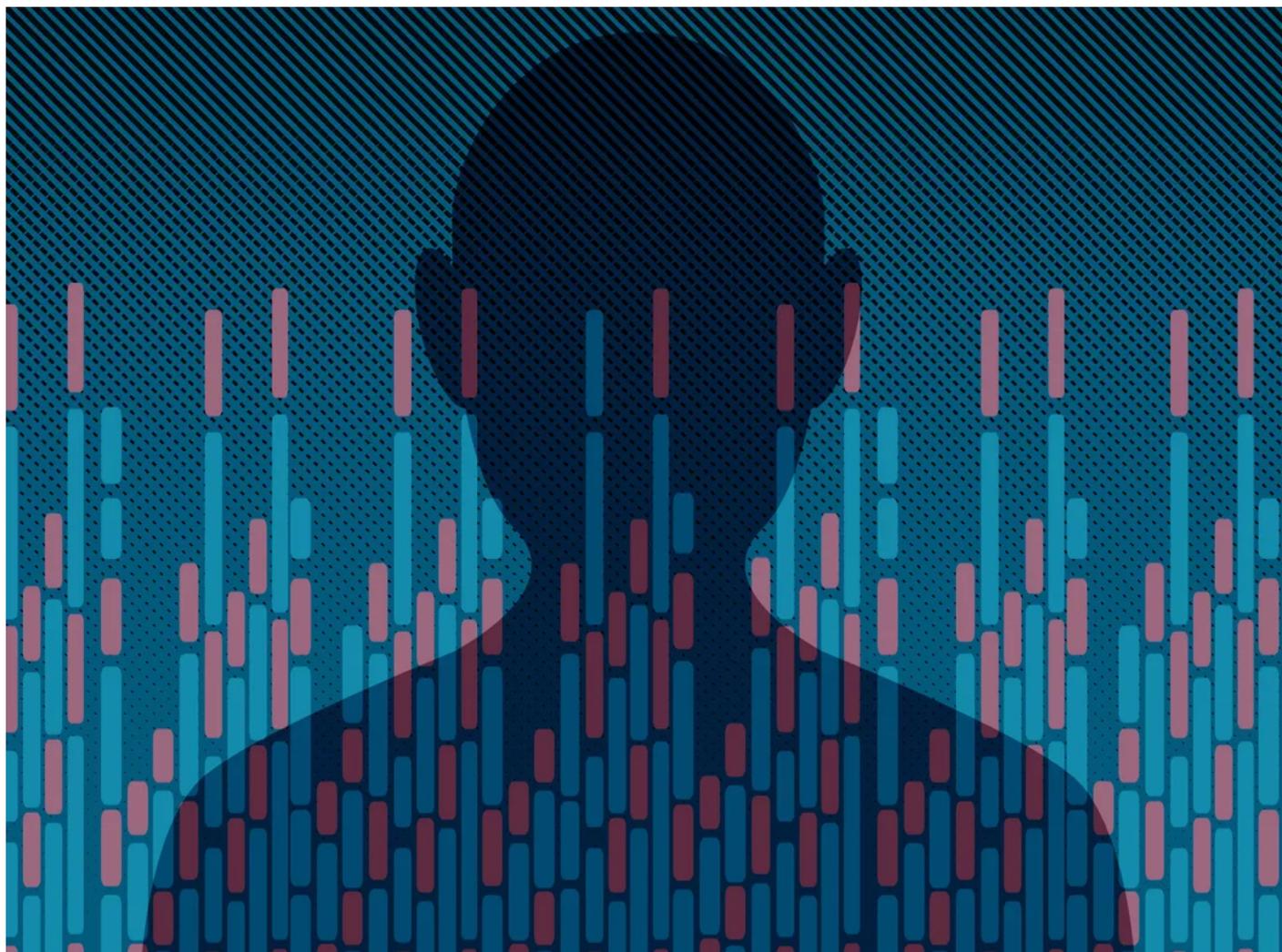
Learn the term, because its influence — in science, and in consumer genetic tests and medical research — is growing.

GWAS has two main applications in science. One is using the studies to better understand the underlying biological architecture of disease and human variation. The other: Scientists are also scanning genomes in the hope of predicting who is most at risk of developing conditions like heart disease and diabetes, potentially from a very early age — perhaps even before birth.

In 2007, there were 240 papers that mentioned GWAS, according to the **PubMed database**. In 2017, more than 3,800 were published.

And the size of these studies keeps increasing, thanks to the falling cost of genome sequencing. Ten years ago, GWAS studies used a few hundred participants. Today, they routinely use hundreds of thousands, and a new one topped 1 million. The more people in these studies, the more power researchers have to probe the genome.

But with any new scientific tool, there are some serious limits to what GWAS, and the genetic risk profiles generated from these studies, can tell us. And critics argue that the predictive ability of these tests is quickly hitting a ceiling and will not necessarily be useful for most individuals seeking to understand their genetic fate.



These predictions can also be wildly misinterpreted, even abused, in the wrong hands. There are fears that they'll give rise to an industry that feels more like genetic astrology than genetic prediction.

So it's worth getting a solid understanding of the questions GWAS can and cannot answer. And, at the very least, GWAS is revealing this fact: Our differences are awe-inspiringly complicated.

"The bottom line is basically, genomics is just wickedly complex," says Adam Rutherford, a geneticist and author of ***A Brief History of Everyone Who Ever Lived***. "And we're only just beginning to really understand."

How GWAS works

There are about 3 billion base pairs in the human genome. These are the literal letters of life, the very building blocks of our DNA, and they spell out a code for our bodies to follow.

If you're a human and you're reading this, know: Our DNA is overwhelmingly identical. Indeed, all the beautiful permutations of the human form — the differences between the tallest and shortest, the brown-eyed and the green-eyed — are explained by just **a tiny fraction** of those base pairs. Finding the genetic differences that make one person taller or shorter than another is like looking for needles in a haystack.

In the case of GWAS, these needles are called single nucleotide polymorphisms, or SNPs (pronounced “snips”).

You may recall from high school biology that a strand of DNA is made up of a string of nucleic acids: adenine, guanine, cytosine, and thymine. A SNP is literally a spot in the genome where one of those nucleic acids is swapped out for another. For the most part, SNPs are meaningless, harmless, and random. (Though the random ones **are useful in tracing a person's ancestry**, as these random SNPs tend to be shared in people with a common geographical origin.)

In the case of the recent *Nature Genetics* paper, the question is: Which SNPs are commonly found in people who have completed many years of school, and not found in those who have not? Other GWASs have asked questions like: What are the genetic correlates of diabetes risk? Or schizophrenia? One recent GWAS study even **identified** 250 genetic sites that correlate with male pattern baldness.

If GWAS finds a SNP that is correlated with a trait, it could mean that SNP is a part of a gene that influences that trait directly. Or it could mean that SNP is just correlated with another SNP that does influence that trait.

Or it could be a bit of a red herring: For instance, the authors of the *Nature Genetics* paper **point out**, a GWAS could identify a SNP correlated with lung cancer. But does that gene cause cancer? Or does it make someone more likely to be addicted to nicotine?

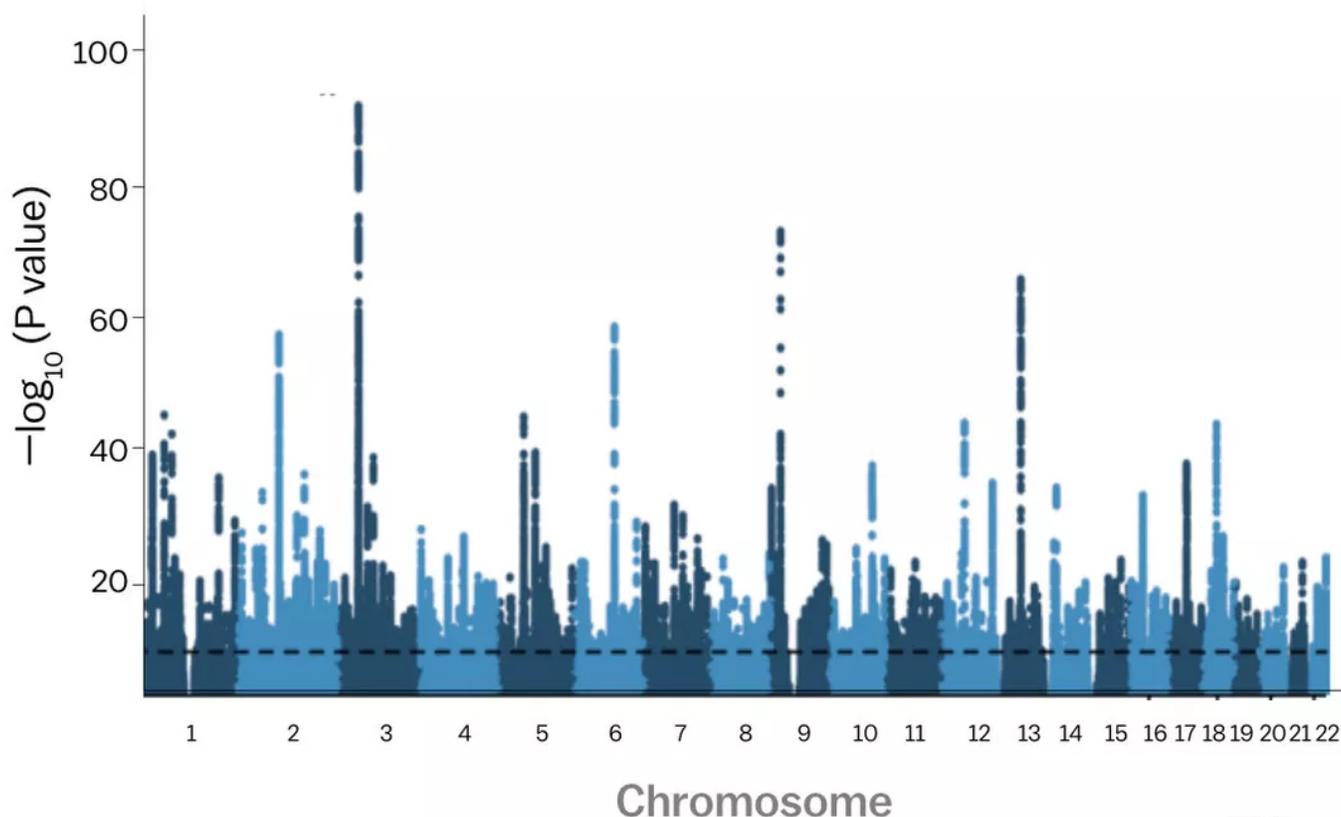
“The A in GWAS is association,” Rutherford says. This is crucial to remember. GWAS doesn't find the genes that cause a trait to occur. It finds the genetic changes that are *correlated* with a trait.

The output of a GWAS is sometimes called a Manhattan plot. The plot shows the locations on the genome where the greatest statistically significant correlations between the SNPs and the target trait lie. It's called “Manhattan” plot because the most statistically significant correlations stick out like a skyscraper. Here's an example of one from the

Nature Genetics paper. The x-axis shows location in the genome. The y-axis shows how significant the correlation is.

“Manhattan plots” correlate genes to a trait or illness

Genome-wide associations for educational attainment



Source: Nature Genetics

Vox

“What GWASs do is plant a flag in the landscape and say, ‘There’s something interesting here,’” Rutherford says. It takes a lot of painstaking follow-up work to determine what that interesting thing is.

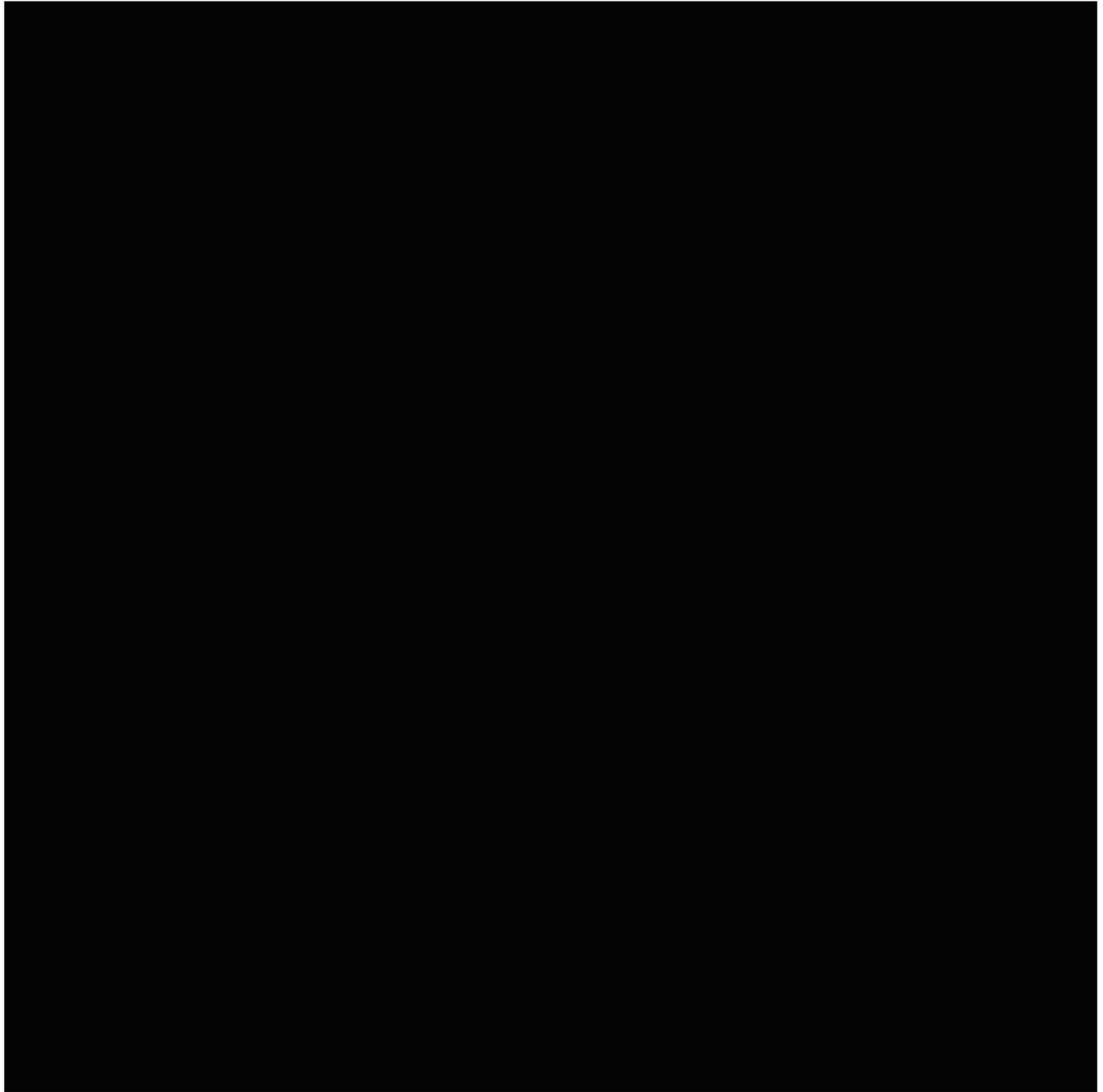
Often, though, the placement of these flags is not random; they point to biologically meaningful regions of the genome.

The genes flagged by GWAS on height tend to relate to the skeletal system: They involve genes that are particularly active in **the growth plate regions of bones**, and point to genes that are **involved** in the manufacture of connective tissues like collagen. The ones

correlated with educational attainment were clustered in regions having to do with the central nervous system. These are regions that have to do with the development of our minds.

A big hope behind GWAS: If scientists can identify spots on the genome associated with a disease or a behavior, they can begin to trace the pathways from genetics to organ tissues to symptoms. And along that pathway, they can possibly find places to intervene and discover new cures.

In the more immediate future, some researchers say, doctors will use the tests to predict who is most likely to develop diseases. There was a recent GWAS effort to help **make significant predictions about** who is most likely to develop coronary artery disease in their lifetime. If these polygenetic tests are used at an early age, **some doctors hope**, patients can be given early preventive measures and medications to avoid heart attacks.



Javier Zarracina/Vox

Each SNP is correlated with a tiny, tiny effect

Let's back up for a second. Why are personal life outcomes like educational attainment even influenced by genetics? It's because many different aspects of ourselves that go into success at school have a genetic basis.

Studies with twins have **shown us virtually every** human trait — height, intelligence, irritability, personality, etc. — is, in part, heritable, meaning the differences seen between

people can in part be traced back to differences in genetics. Of course, the environment often plays a huge role too and interacts with genetics in complicated ways. So it's not surprising that GWAS can yield predictions for educational attainment.

Or, as behavioral geneticist Eric Turkheimer has explained, even predicting **marital status** from GWAS is not out of bounds of the science.

“At every fork in the road,” Turkheimer **writes**, “from when you are a pinhead-sized bunch of cells until you are reacting to your first date, at every single choice point, your genome has its thumb on the scales, making it slightly more likely that you will develop in one direction than another.”

That's why scientists are able to find genetic correlations for numbers of years spent in school.

Think about it: Intelligence is highly heritable. We also inherit aspects of our personality that help us make it through school, like **conscientiousness**. We also inherit health. Being a healthy kid helps one make it through school. Same goes for sleep. We know genes help decide whether we're a natural early or late riser, and kids who have trouble waking up early also tend to have more trouble in school.

Educational attainment is an intriguing variable to study because it shows how our genetics influence our lives through many subtle channels. But here's the key: The effect of each SNP is tiny. In the *Nature Genetics* paper, each of the 1,200-plus SNPs identified in the GWAS seems to contribute, on average, to predicting just two more weeks of schooling.

This pattern is true for most of the traits geneticists study: Each genetic variation accounts for a tiny, tiny difference.

Take height, for example. It's one of the simplest attributes we use to describe a person. It's written on millions of driver's licenses with just a few digits. GWAS has revealed, in the genetic code, that height is likely written **across hundreds**, if not thousands, of locations in the genome.

And at each of these locations, if you inherit one version of SNP, it could raise your height by a mere millimeter. (Of course, there are also environmental factors that help determine height, like access to nutrition.)

Similarly, there might be 1,000 genes that influence intelligence. “If you have the bad variant of *one* gene for IQ, maybe your IQ score ... is 0.001 percent lower than it would have been,” Danielle Posthuma, a statistical geneticist at Vrije Universiteit in Amsterdam, told me in an interview last year.

Genetics have a tiny (and sometimes big) thumb on the scale for just about every human trait. So that means GWAS can find SNPs correlated for just about any trait measured — even ones that aren’t obviously the direct result of biology. For instance, there was **recently a GWAS study** on loneliness and social isolation. Sure, loneliness is partially the result of genes, as our personalities and our propensity for mood disorders like depression are linked to loneliness, and our genes. But loneliness is also often the result of our environments, and of our cultural understandings of friendship and companionship.

Again, because all human traits are in part heritable, scientists and the media need to be careful when they claim research has found the “gene for” anything.

As GWAS sample sizes grow, so does the ability to find correlations

Let me recap: The genetics of our individual differences reside in a relatively small area of our genome. At the same time, many small genetic changes are involved in the expression of a single trait, and each change is correlated with a tiny tweak to the human form.

To find the tiny effects that individual letters of the genome have on traits, disease, and behavior, you need enormous data sets to separate signal from noise.

When GWAS was first deployed on genetics data sets in the early 2000s, it didn’t discover many correlations. The problem: There were just too few people with publicly available genetic data. Now, hundreds of thousands of people have joined public genetics data sets, like the UK Biobank, which has data on 500,000 people. And more huge data sets are being assembled right now. Commercial DNA testing companies like 23&Me also **sometimes collaborate on GWAS research.**

And with more people’s data available, researchers can peer deeper into the genetic code.

Just two years ago, a GWAS trying to parse education attainment with around 300,000 participants only yielded 74 significantly correlated locations on the genome. But the recent *Nature Genetics* paper on educational attainment used data from 1.1 million people — and got 1,271 significant hits on the genome.

That shows that the bigger these sample sizes get, the more genetic markers will be identified, and the more we'll learn about the stunningly complicated ways genetics influence our lives.

Genetic predictions about health are getting better and better. But will they be useful?

It is possible, using insights derived from GWAS, to make some predictions. In the *Nature Genetics* article, the researchers found they could predict — to a mild degree, and on average, across a large sample — a person's likelihood of finishing college.

In their study, around 50 percent of the people who had a high number of the genetic variants associated with educational attainment had completed college. For the people with the lowest number of variants, that figure was 10 percent. And research has shown scientists can make similar predictions based on genetics for traits like height.

The researchers make these predictions by calculating a tally called a polygenic risk score. Remember that GWAS is about figuring out which genetic variants — which DNA letter swaps — are associated with a trait, like being taller. In an individual, scientists can count up all the spots in the genome that suggest a person is going to be tall, and then assign them a score.

Polygenic risk scores are usually conveyed as a percentile (as in: you're in the 90th percentile of risk for X disorder or acquiring a trait).

Again, these risk scores make weak predictions for most individuals. Currently, you can make a stronger prediction about how many years of schooling a person will receive by just asking how many years their parents got. The accuracy of risk scores can also change over time as our environments change.

Sure, we can predict educational attainment from genes. But what if our school systems were radically different, favoring different traits other than studiousness?

"The genetic variants we have identified as associated with educational attainment in the current environments in which they were studied would play a lesser role still in a zombie apocalypse where many schools have been overrun by walkers," the authors of the *Nature Genetics* paper **explain in an FAQ**. In that case, genes correlated with running faster would yield a stronger prediction for success.

That said, doctors are now wondering if genetic risk scores, in the future, could be useful in medical decision-making.

Genetic risk scores may guide medical decision making

Here's an example of what's coming.

In a recent study, also published in *Nature Genetics*, people who scored the highest in risk on a polygenic test for coronary artery disease (about 8 percent of participants) were three times more likely to have contracted the disease. Coronary artery disease occurs when plaque builds up in the blood vessels around the heart, and it is a scary precursor to heart attacks. Wouldn't you want to know, from a young age, if you were in the category of highest risk?

In the following chart, you can see that as polygenic risk score goes up (as measured in percentile), so does the incidence of coronary artery disease in the study population.

"We propose that it is time to contemplate the inclusion of polygenic risk prediction in clinical care, and discuss relevant issues," the authors of the paper write.

"This is going to take hold in common medical practice," Eric Topol, a cardiologist and geneticist with Scripps Research who is not involved in this research, says. "It's a matter of when, not if."

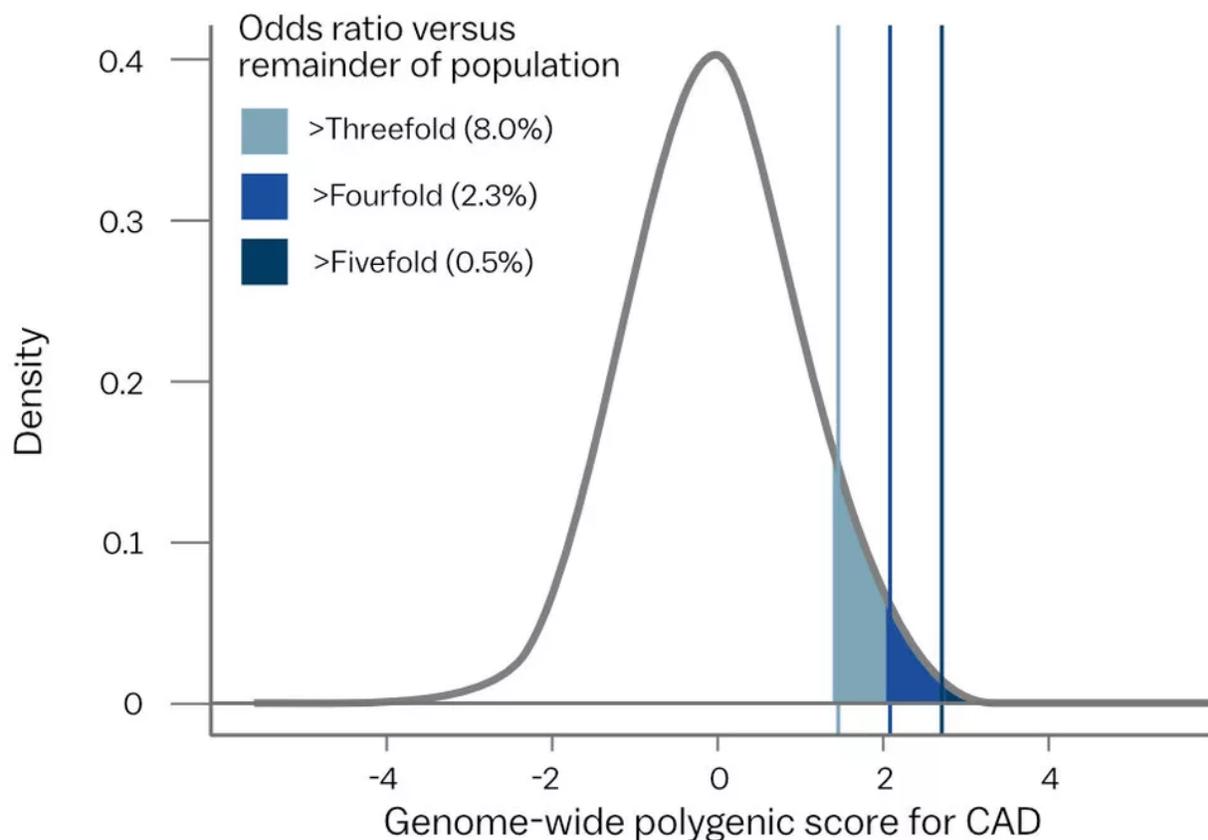
There are already polygenic risk tests available for **breast cancer**; it factors in not just the BRCA1 and BRCA2 genes but also other variations in the genome that appear more often among women who get breast cancer. In the coming years, more tests like this for more diseases are likely to trickle down to consumers to help them make decisions.

For instance, a woman with a higher polygenic risk score for breast cancer may want to go for more frequent mammograms. The new paper also identified that polygenic risk scores could be useful in finding people most at risk for developing diabetes, breast cancer, and inflammatory bowel disease.

But there are skeptics of whether these scores will truly be useful for patients. For one, they don't necessarily surpass the risk assessments of asking patients simple questions about their lifestyles and family histories (more research is needed to see if they do). And the polygenic risk scores may only be meaningful for people who score on the extreme

ends. Sure, you may find out you score in the 60th percentile for heart disease risk, but it's only at the 90th percentile where your odds of acquiring heart disease go up considerably.

Polygenic scores signal increased risk for coronary artery disease



Source: Nature Genetics

Vox

The hope is that the predictive power of these scores may continue to improve as GWAS studies grow larger.

But “we’re going to hit a ceiling really quickly,” Cecile Janssens, an epidemiologist at Emory, says of the predictive power of polygenic scores. “All the SNPs that we are discovering now have such a tiny effect,” she says, and the ones that we will continue to discover are likely to be less and less powerful in terms of the predictions you can make from them. “We see this trend already for years: Every new SNP that we discover has a smaller effect than we knew already.”

In the coronary artery disease study, she points out, when the researchers increased the number of SNPs in their risk model from 74 to 6 million, the predictive power of the test only increased by a smidgen.

It's also possible the predictive power of the scores will decrease when replicated in new samples of people. Also, we shouldn't take it for granted that intervening early, based on polygenic risk scores, will necessarily save lives.

So we should be cautiously optimistic that polygenic scores may be useful for predicting an individual's heart attack risk.

However, many of the researchers I spoke to agreed that polygenic scores might also prove useful in areas where it is hard to guess a person's risk for a disease just by asking them simple questions, like in psychiatry.

Gerome Breen, a psychiatric geneticist, explains that when someone is admitted to a hospital for psychosis, "it can often take quite a long time to determine what disorder they actually have."

Schizophrenia can sometimes look like bipolar disorder, but these two illnesses tend to be treated with different drugs. If a doctor could obtain a polygenic risk score for the patient, she could possibly see which diagnosis is more likely, given a person's genetics. And that could help doctors get the patient on the right treatment faster.

The bottom line is that these risk scores — for schizophrenia, for heart disease, for educational attainment — come with a lot of caveats, and they need to be assessed on a case-by-case basis. But overall, risk scores aren't very useful at the individual level, unless a person scores on the extreme high or low end of the scale. "For a lot of disorders, roughly speaking, a polygenic risk score is approximately [as useful a predictor as knowing] your grandparents had that disorder," Breen says.

There's another huge limitation to most GWAS studies — they're only done on white Europeans

One big limiting problem with GWAS studies is that the **vast majority of them have been carried** out using white participants from European backgrounds. And the polygenic predictions don't often translate from one population to another. Over time, different populations of people — i.e., people of European or Asian descent — acquire random genetic variants that are meaningless and have nothing to do with biology.

And if you included both Europeans and Asians in a GWAS study, you couldn't be sure if the results reflect the random changes that happen over time or are actually revealing underlying biology.

This all means that if researchers make progress on developing polygenic risk scores for heart attacks, those scores will be less meaningful if applied to any population outside of white Europeans. If we truly want to make GWAS **an equitable and useful tool** to predict disease risk, scientists must repeat their studies in more diverse populations.

GWAS can help guide researchers dive deeper into biological mechanisms — and help understand when the environment matters

The extent of what GWAS reveals is often explained in terms of prediction life outcomes. But it's also useful in two other ways: investigating biology and, potentially, helping us better understand the power of a good environment to influence our lives.

“A lot of people are interested in GWAS because then you find specific genes and then you take a deep dive into the biology,” says Paige Harden, a psychologist who studies the genetics of behavior at the University of Texas Austin.

Indeed, there's been progress from GWAS in identifying genes that **signal risk for diabetes, schizophrenia, and depression**. Once scientists discover the genetic architecture of conditions, they can better trace the biology of how they manifest through the body, which could then help researchers develop better pharmaceuticals to treat these conditions.

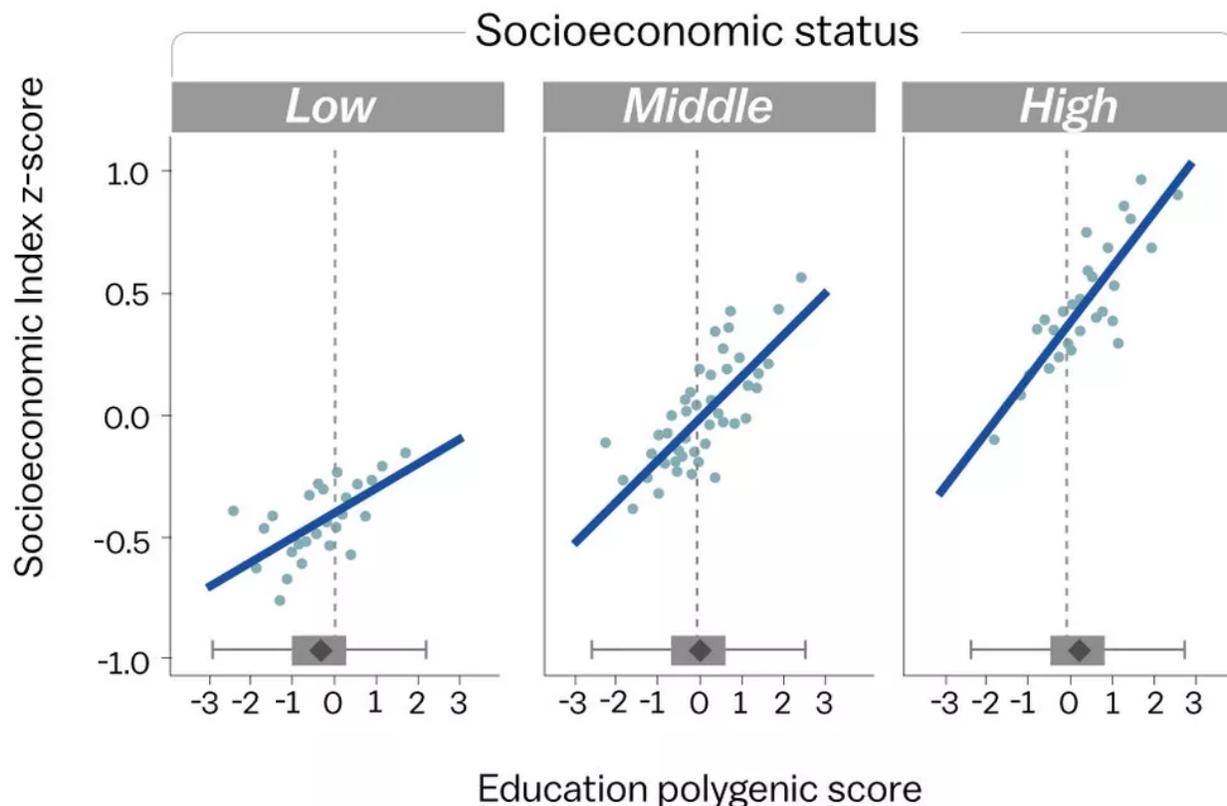
But here's what Harden is more excited about: It is possible to use polygenic risk scores as a control in studies about the power that environmental factors — like family income, education, or home life — can have over our well-being.

So often, studies in psychology can't control for genetics. For instance, take a study on parenting styles and child outcomes. If a study finds that parents being aggressive with their kids leads to negative life outcomes for those kids, is it because the parents are aggressive, or is it because the kids and parents share similar DNA?

Harden pointed me toward a graph from a recent study in **PNAS** on whether polygenic risk scores for education can predict upward social mobility (that is, can genetics explain why some people who start off poorer in life end up richer later on?). The x-axis of the chart is polygenic scores for educational attainment. The y-axis is a measure of socioeconomic

success as an adult. And the chart is broken down into what types of homes the participants were born into (lower class, middle class, or upper class).

What polygenic scores can teach us about the power of wealth



Source: PNA

Vox

What you can see: Yes, there is a correlation between genetics and economic attainment. No matter what socioeconomic class a person is born into, those with the higher polygenic scores are more likely to be richer later in life.

"But I think another thing you can look at is, 'Okay, let's take people who are genetically identical with respect to their polygenic score,'" Harden says. When you do that, you can see that a person with a polygenic score of 3 has a much harder time getting rich as an adult when they are born into a lower-class family than the 3s who were born into a higher social class family.

By controlling for genetics, scientists can more clearly see that our fate in life is largely determined by how much money our parents have. “People who are genetically identical [on their polygenic scores] really differ in their social attainments depending on their class background,” she says.

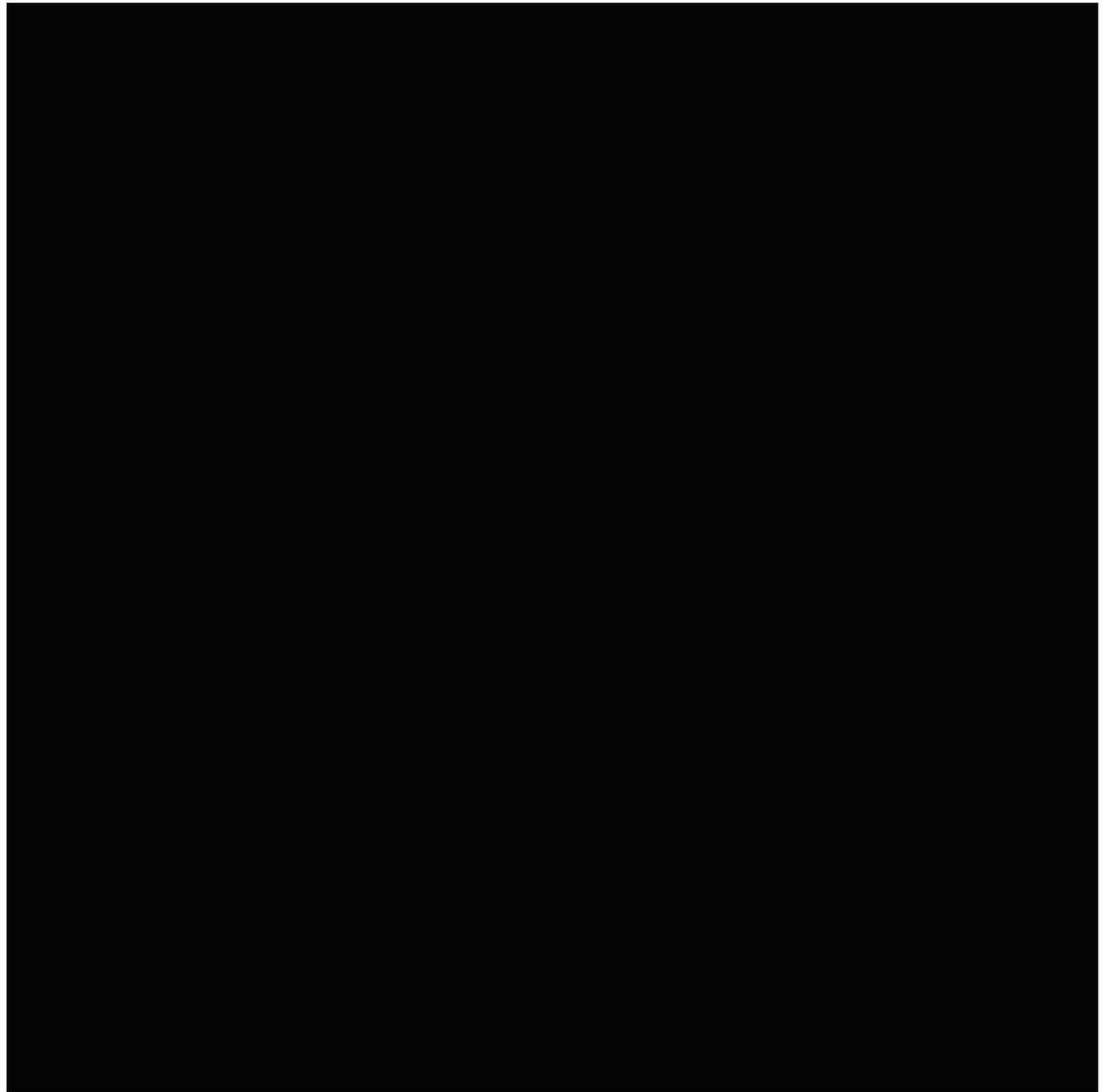
We can’t change our genetics. But we can develop social programs to help lessen the burdens of growing up poor. If we can better control for genetics, we can better understand if those social programs are working. Yes, with risk scores, we cannot completely control for genetics. But even slightly controlling for genetics can allow scientists to draw more meaningful conclusions on the power of the environment.

There’s a high potential for misuse in GWAS and genetic prediction

Right now, it’s possible for you to purchase **a genetic test for intelligence**. These genetic tests are derived from the results of GWAS studies, and, as science writer Carl Zimmer has explained in the Atlantic, they **are near useless**. There are about “538 genes with clear-cut influence on intelligence test scores,” Zimmer writes. “Those 538 genes have a predictive power of about 7 percent.”

Harden calls these types of predictive tests “genetic astrology,” and they’ll almost always be a waste of money. (If you want to find out how intelligent you are, take an **IQ test!**)

It’s not hard to imagine a world where genetic astrology becomes an ever-growing commercial product (heck, actual astrology is everywhere). Imagine a situation where, during in vitro fertilization, parents could assess the polygenic risk scores of their fertilized embryos and choose ones that seem the most likely to be tall, smart, and attractive. Sure, the prediction might help them guide their decisions, but it would be based on a very weak scientific understanding of how SNPs affect traits and won’t guarantee anything.



That's because GWASs can't even be used to explain all the genetic differences between individuals. From twin studies, we know that around 60 percent of the differences between people's educational achievement can be explained by differences in their genes (though other papers argue it's closer to 20 percent). The other 40 percent of the variation, then, can be explained by differences in the environment.

Yet the *Nature Genetics* study on educational attainment didn't find genetic sites that can explain 60 percent of the difference between people's achievement. They didn't come close. Even with a huge GWAS of 1 million people, the genes identified in the *Nature*

Genetics study can only account for 11 percent of the differences seen between people. Scientists call this problem “missing heritability.”

Similarly, twin studies suggest around 60 percent of the height differences seen in the human population can be explained by differences in genetics (the rest being explained by environmental factors, like socioeconomic status and nutrition).

But even the hundreds of genetic sites that have been linked to height via GWAS can only explain around 25 percent of the differences seen in people. It's possible complex traits like height are omnigenetic, meaning that almost **the entire genome plays** a role in determining how tall a person will be. If genes are a recipe for our bodies, omnigenetics suggests they're drawn like an M.C. Escher drawing: Each part of the scaffolding connects to many others in mind-bendy ways.

The point is that GWASs, at best, provide an incomplete picture from which to draw predictions.

There's an age-old problem to watch out for: Predictions derived from science can be used for really silly ends. “Many people are concerned that insurance companies will use [polygenic risk score],” Posthuma said. “That they will look into people's DNA and say, ‘Well, you have a very high risk of being a nicotine addict, so we want you to pay more.’ Or, ‘You have a high risk of dying early from cancer, so you have to pay more early in life.’ And of course, that's all nonsense.”

Similarly, it would be nonsense to be sold vitamins or diets specially tailored to your polygenic risk scores. But considering how poorly the wellness industry is regulated, it's not out of the realm of possibility to see such a thing marketed one day.

And remember: These polygenic scores can be generated without anyone understanding what the underlying biology of the genes involved actually does. “We think that is mindful biology? No. It's an approximation,” Janssens says.

If anything, the results of GWAS studies point toward how hard it is to play God with the genetic code. It reveals how beautifully complicated the architecture of our genetics is.

“Evolution has equipped us with a genome which had no intent of us understanding how it worked,” Rutherford says. And we're now just barely scratching at the surface.

SCIENCE & HEALTH

Heart disease risk is hidden in your genes. Scientists are getting better at finding it.

SCIENCE & HEALTH

Americans should be more afraid of HPV

SCIENCE & HEALTH

How to get a good night's sleep

[View all stories in Science & Health](#)