

# On the word “because” in mathematics, and elsewhere

Barry Mazur

May 18, 2017

## Contents

<b>I</b>	<b>Introduction</b>	<b>3</b>
<b>1</b>	<b>What is an explanation?</b>	<b>3</b>
<b>2</b>	<b>Why is a mathematical truth true?</b>	<b>4</b>
<b>3</b>	<b>What is a question?</b>	<b>6</b>
<b>4</b>	<b>Teleology</b>	<b>7</b>
<b>5</b>	<b>Mathematical questions</b>	<b>7</b>
<b>6</b>	<b>Abstracting and/or Generalizing:</b>	<b>10</b>
<b>7</b>	<b>Are there truly different proofs of the <i>same</i> theorem?</b>	<b>11</b>
7.1	The Pythagorean Theorem . . . . .	11
7.2	Fermat’s Little Theorem . . . . .	12
7.2.1	Induction . . . . .	12
7.2.2	‘Automorphisms’ of the additive group of integers mod $p$ . . . . .	13

7.2.3	The order of an element in a finite group . . . . .	13
<b>8</b>	<b>Change of Setting: Desargue's Theorem</b>	<b>14</b>
<b>II</b>	<b>Ideas underlying, or analogous to, other ideas in mathematics</b>	<b>14</b>
<b>9</b>	<b>Viewing a mathematical issue as controlled by some <i>different</i> underlying conceptual structure</b>	<b>14</b>
<b>10</b>	<b>Depths and Heights</b>	<b>17</b>
<b>11</b>	<b>Horizontal, rather than Vertical</b>	<b>18</b>
<b>12</b>	<b>The personal guy-wires of explanation</b>	<b>19</b>

## Part I

# Introduction

I once had a conversation about mathematics with a friend and found myself making a certain claim. Now this 'claim' seemed reasonable at the time, but—on reflection—I came to feel that it actually deserves much more thought—more exploration—than I had given it. This is why I started to work on this essay. The 'claim' I made was about prime numbers 2, 3, 5, 7, 11, . . . . Specifically, I said that a certain structure—which I'll hint in Part II of this essay—was what *lay behind*—i.e., was the *explanation for*—the intricate complexity of this set of numbers. To phrase it as more of a puzzle: this basic most fundamental, mathematical notion *prime number* has its idiosyncratic behavior 'caused by'—or at least 'understandable only when one considers'—some underlying—and deeper—concept, or constellation of concepts.

On the face of it, there's nothing outlandish about this type of claim. After all, we're familiar with simple-to-describe macroscopic phenomena that are understood only as consequences of Quantum Mechanics. Nevertheless, the analogous thing in mathematics seemed to me to be a puzzle. One is very familiar, in mathematics, with  $A \implies B$  type arguments. But " $\implies$ " is a far cry from "is the underlying reason for." What can one say about the architecture of fundamental levels of explanation—perhaps: of causation—in the mathematical realm? I wrote this essay as an occasion to think about it, living for a while with—and paying attention to—the notion of explanation in mathematics.

## 1 What is an explanation?

This question has different answers, depending on what is being studied, what is meant to be explained. At a minimum, an explanation is a *satisfactory* answer to a question beginning with the word "why,"—*satisfactory* being the critical word. Satisfactory, then, to whom?

The task of explaining mathematical ideas has a variety of flavors: depending on whether you want it to be satisfactory to a class, to a certain student, to a certain colleague, to the mathematical community in general. . . . or to yourself.

In mathematics, we often depend on the proof of a statement to offer not only a justification of its truth, but also a way of understanding its implications, its connections to other established truths—a way, in short of explaining the statement. But sometimes even though a proof does its job of showing the truth of a result it still leaves us with the nagging question of why.' It may be elusive—given a specific proof—to describe in useful terms the type of explanation the proof actually offers. It would be good to have an adequate vocabulary to help us think about the explanatory features of mathematics (and, more generally, of science).

In preparing to write this essay I tried to remember what it felt like, in school, to be suddenly explained something that seemed somewhat opaque before that explanation. That is, to recall

some early moment of ‘eureka,’ as a way of revisiting the inner felt experience of ‘having been explained’ something—however elementary—the explanation having all the dramatic force and sense of new insight.

Such a moment, for me, was when someone taught me the simplest version of Fermat’s *principal of least action*.

Namely that if you look at a point source of light as reflected in a mirror, the light beam from the source to your eye will be reflected at a point in the mirror in a manner so that the total length of travel of the light beam is minimal—and *this is the point where the angle of incidence is equal to the angle of reflection*. The proof that minimal-length-of-travel implies angle-of-incidence-is-equal-to-angle-of-reflection indeed *was*—for me, with very little mathematical experience at the time—a revelation, an *explanation*: you simply reflect the diagram below through the looking-glass and note that the shortest distance between two points is a straight line.

*Put diagram here*

## 2 Why is a mathematical truth true?

At first it would seem to be a simple matter. Why is there a *unique* straight line between any two points? *Because* an axiom stipulated this. Why is every prime of the form  $4n + 1$  expressible uniquely as a sum of two squares? *Because* we can prove it, as follows ... Why is  $5 + 7 = 12$ ? *Because* of the following computation ...

Mathematics moves forward by the formulation of definitions, the adoption of axiomatic statements and procedures, computations following those procedures, and the application of logically ordered argument that relate to those definitions and axioms. But, of course, it isn’t that simple.

Ask a question about some analogy between mathematical structures and you’re fishing for a different kind of answer. Ask *Why are integers like functions?* as Weber must have done at the end of the nineteenth century, and you’re off on a pursuit that hasn’t ended yet.

Some of the more goading questions we ask, driven by the aim of understanding something mathematical, are skew to the architecture of what is already constructed in mathematics. Some of these questions search, in vain, for vocabulary so that they can actually be expressed. They may not really be questions; they are rather wonderings, spurts of curiosity. Call them ‘queries.’ These queries can precede definitions, axioms, and all the machinery needed to communicate and verify mathematical statements. Despite—or perhaps because of—their vagueness, they have force. A vital force that keeps mathematics enlarging.

These queries can occur at all stages of mathematical development. It could range from Richard Dedekind’s question: *Was sind und was sollen die Zahlen?* to a beginner’s puzzlement. A young student once came to me with a sequence of regular  $n$ -gons with  $n$  going to  $\infty$ . He wanted to know how to go about understanding such things, and said (a literal quote) “I have pain and hurt unless

I understand the answer.” There is something admirable in this type of intellectual pain<sup>1</sup>. But such queries can occur earlier too.

What types of pre-mathematical queries have very young children come up with, on their own, that lead them to do their own musing? My guess is that some pattern that has a puzzling element must be prominent as inspiration.

For example, when I was a child (I don’t remember what age) I kept wondering what to make of the fact that if I count the tips of my fingers on one hand I got 5, but if I count the wedges between them I got 4. I’m not sure what I ever made of this, but I know it was a recurrent issue in my musings.

Much more subtle is the following early ‘pre-mathematical quest’ a colleague of mine—call him Jim— had when he was a preschooler. If you want to indicate the number 5 by placing five stones in some pattern—e.g., a quincunx, or an L-shaped gnomon, a regular pentagon, or just in a straight line—what, asked the child, is the **true** way of representing that number? The striking assumption Jim had is that there *is* a true way, a canonical representation.

These two examples are what one might call *pattern-wonderings*.

From the examples mentioned above to the grand conjectures of the subject, one can see a unity to the gamut of question-asking related to mathematics. And there is often a delicate and surprising structure to any particular mathematical question, when viewed from a broad enough perspective.

For example, artfully-wrought mathematical questions—problems posed within some fairly well-developed field—can be intentionally slant: it is often not the bald answer to the question posed that has much importance. Rather, the question is a goad to develop ideas that shed light on it, and that are important for their own sake.

Fermat’s Last Theorem is a perfect example: it is a simple question about numbers. Its importance doesn’t stem from the the straight statement of the theorem. To be sure, the historical narrative that accompanies this problem already has great weight. But its extraordinary importance to the substance of mathematics comes from the enormous amount of deep mathematics that emerged from the challenge it posed.

An admirable, almost instinctive, practice of many mathematicians is to use the inspiration of new results—answers to long-asked questions—as springboard to the formulation of yet further, often more expansive, or more trenchant, new questions: an activity that one might label *questioning*

---

<sup>1</sup> This type of self-generated perplexity about mathematics is, perhaps, what we want to generate, even when we pose problems for students in what seems to be the traditional way: ‘Solve this problem:  $5 + 7 = ?$ . Happily, teachers sometimes modify this usual route by prompting student’s native curiosity, which gets the student to emerge with their own questions as a lead-in to mathematic thinking. For example, Bob and Ellen Kaplan begin their first session of a Math Circle class of four and five year olds by making some assertion like: “there is no number between 4 and 5” with the expectation—usually met—that some child will say: “but I’m 4 and a half!” Their response: ”oh yes, I forgot about 4 and a half, but of course there’s no number between 4 and 4 and a half!” A natural rebuttal occurs, and that sets the tone and mood of an exploration into the world of numbers where it is the children who are asking the questions, making constructions, and making assertions as claims to be considered, further questioned, for further discoveries.

*answers* rather than the more standard *answering questions*. Such a quality of thought would be good to instill in our students.

Although the question *what is a question?* has a venerable history and tradition of discourse, we probably don't have as detailed a discussion of the nature of (specifically) mathematical questions, and even less of a discussion about the possible range of pattern-wonderings. It is worth having such a discussion, and encouraging students to develop the habit of asking their own questions as a way of becoming intimate with what they are studying, and in that way fully experiencing the sense of wonder that our subject inspires.

### 3 What is a question?

The simplest type of questions are what one might call quotidian-questions, such as “How are you?”, or “When is the next bus coming?” or opinion-questions regarding judgment, such as “What do you think of the movie?” or regarding prediction: “Will the Red Sox win the pennant?” Some types of questions are less useful than others: when, for example, has any rhetorical question helped a discussion?

Other questions have subtler reasons for being. When my (then) two-year-old grand-daughter Naia wanted to deflect a certain parental order (not to jump from the dresser to the bed), she proclaimed from her perch on the dresser: “I have two questions... (pause)... about elephants.” (Those two questions never came.) The writer Jamaica Kincaid records in her book *Among Flowers: A Walk in the Himalaya* that as a member of a group in a trek in search of flowers, she would often need to stop to catch her breath. Her tactic was simply to look at the nearest bud and exclaim: “what *is* this?”

But, more often, we ask questions hoping that the answers are key to a better understanding of something, or that they will allow you to proceed better in the world. And our deepest questions are not exactly intended to be answered. Perhaps they don't even end with a question-mark, as in Gauguin's “D'où Venons Nous / Que Sommes Nous / Où Allons Nous.”

As has been long noted<sup>2</sup> our frequent simple query “*what is this?*” has (at least) four genres-of-asking enfolded in it:

- “Out of what” (material) is it made?
- “How is it organized?”
- “How did it come to be? What person, or type of person, or agency made it?”

and the tricky

- “What is it for?”

---

<sup>2</sup>as in Aristotle's *Metaphysics* Book V [1013].

What an *essential ambiguity* there lurks in any question that begins with the word “why.” Any such question is a ‘pre-question,’ an invitation to make a precise question. To answer the question *Why is the sky blue?*, for example, you can choose to respond in broadly different ways, the choice depending on some comprehension of the intent of the questioner. The answer can range from: “it’s a non-cloudy day” to a formula expressing the fact that air scatters short-wavelength light more than longer wavelengths.

The answer to a “Why” question begins explicitly or implicitly with “because.” This is not necessarily true for answers to “How” or “What” questions. We often view a satisfactory response to a “Why does B have a certain form?” question as finding some “A” that we feel is responsible. One then can say that A *causes* B; or perhaps we might be more loose about it and just say that A *is responsible* for B. Responsible, in some way, for the attributes of B under discussion; or is a way of understanding B; or: if we have A, then we have—or are soon to have—B. The fall of the *n*-th domino is responsible for the fall of the succeeding one.

## 4 Teleology

“*What is it for?*” questions are indeed tricky, as was already mentioned.

One is so tempted to create an intention of some incorporeal agent to account for the occurrence of certain phenomena. An agent such as, for example, *Nature*. Or, more modestly, to ‘explain’ some feature of anatomy—like the shape of the chin—by attributing to it some evolutionary advantage. Steven Jay Gould cautioned against doing this, labelling some of these accounts “Just-So Stories,’ i.e., lightly spun tales that offer confabulated “reasons for” biological occurrences.

[T]he standard misapplication of evolutionary theory assumes that biological explanation may be equated with devising accounts, often speculative and conjectural in practice, about the adaptive value of any given feature in its original environment (human aggression as good for hunting, music and religion as good for tribal cohesion, for example)<sup>3</sup>

## 5 Mathematical questions

A “Why is X true?” mathematical question can be thought of as a “How do I understand that X is true?” question.

Some answers simply connect the conclusion to previously agreed-upon assumptions, via an application of some straight grammar of logic, such as modus ponens:

---

<sup>3</sup> *THE EVOLUTION OF LIFE ON EARTH* Scientific American, (1994); <http://brembs.net/gould.html>

All men are mortal  
 Socrates is a man  
 -----  
 $\therefore$  Socrates is mortal

The curious ‘therefore sign’  $\therefore$ , a seventeenth century invention, is known to all math students. It is a business-like way of saying that all has been set in place to deduce whatever it is that comes after that symbol. There is, apparently, also a symbol meaning ‘because,’  $\because$ , also dating to the seventeenth century. It might be interesting to examine what subtle distinctions there are (bridging the analytic/synthetic divide) in the usage of these two symbols.

Some “How do I understand that X is true?” questions could be answered simply by a reference to a definition.

It could be calculation: Why is  $5 + 7 = 12$ ? Well, you can calculate according to the rules. But, of course, even this is sometimes not the most satisfying answer. Why is the sum of the all the numbers from 1 to 100 equal to 5050? Young Gauss knew why, and that wasn’t by a brute-force calculation<sup>4</sup>.

For other X’s the answer to why comes by reviewing the proof. Why is the sum of the angles of a triangle 180 degrees?<sup>5</sup>

Well, if you understand the proof, you see why. Here you might think of the steps of the proof as forming an epistemological time-line where you first get to know A and then get to know B and so on, until you learn X.

$$A \implies B \implies \dots \implies X.$$

Nevertheless, if that first step is construction of the line through the ‘apex’ parallel to the base,

---

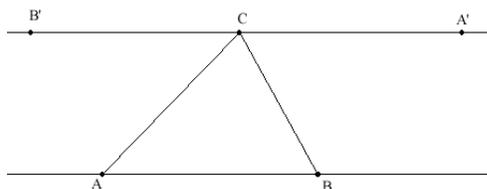
<sup>4</sup> It was, one might say, by a *reasoned calculation* making use of the symmetries in the problem: sum pairs of the summands in the sum

$$1 + 2 + 3 + \dots + 50 + 51 + \dots 98 + 99 + 100$$

symmetrically about the middle:  $1 + 100 = 101$ ,  $2 + 99 = 101$ ,  $\dots$ ,  $50 + 51 = 101$  and you get fifty copies of 101. That is,  $50 \times 101 = 5050$ .

<sup>5</sup> Somewhat baffling is Aristotle’s comment (in Part 9, Book II of *The Physics*) about the proof of this fact:

Necessity in mathematics is in a way similar to necessity in things which come to be through the operation of nature. Since a straight line is what it is, it is necessary that the angles of a triangle should equal two right angles. But not conversely; though if the angles are not equal to two right angles, then the straight line is not what it is either.



you have something pretty close to a “visual proof”<sup>6</sup>

But sometimes a “why” also begs for quite a different quality of an answer. That’s the sort of question that I want to think about. There are times when we ask ourselves *why is X true?* even after we have proven X. We can thoroughly understand the proof we have given. Yet we still hope to find some narrative, or possibly many narratives, offering us a more satisfactory explanation of ‘why’ X is true. Illuminating narratives might well come out of better proofs, or perhaps better ways of configuring and thinking about the structures involved. If we happen to hit upon such a satisfying narrative we can’t make any claim of its uniqueness—i.e., that it is *the* good explanation—there may be many competing ones. Nor can we be sure that anyone else will feel as we do—i.e., that it is anything more than a thoroughly subjective judgment<sup>7</sup>.

Consider two related ways of embedding a mathematical issue in another fuller, or perhaps deeper, context that has the power of explanation:

---

<sup>6</sup> See the Wikipedia entry ‘Proof without words,’ [https://en.wikipedia.org/wiki/Proof\\_without\\_words](https://en.wikipedia.org/wiki/Proof_without_words) if one knows the Euclidean theorem about alternate angles. After this construction one sees the triangle unfolding in three ways, corresponding to its three angles radiating from that constructed line.

<sup>7</sup> As an example of ‘many narratives,’ Curt McMullen suggested considering the question *Why are there infinitely many primes?* and comparing the various proofs of this in terms of their explanatory power. ‘Explanations’ range from the classical Euclidean proof—which, in effect, explains why, given any finite set of positive integers there is a prime not dividing any of them (any prime divisor of  $1 +$  their product will be such a prime) and hence there are infinitely many primes—to proofs that involve considering the values of the zeta-function  $\zeta(n)$  at positive integers  $n$  (if there were only finitely many primes,  $\zeta(n)$  would be a product of factors  $\frac{1}{1-p^{-n}}$  for each of those finitely many primes  $p$ , and hence would be a (finite) rational number for every positive integer  $n$ ; but  $\zeta(n)$  is either not finite ( $n = 1$ ) or not rational ( $n > 1$ )).

## 6 Abstracting and/or Generalizing:

Abstraction comes early in our mathematical education; it is inherent in the very definition of any mathematical concept that attempts to capture some informal intuition that we have. The number 4, a noun, is already the fruit of such an abstraction, as is any law or principle that will be applied in concrete instances (for example, the law of *commutativity of addition*:  $A + B = B + A$  which will be applied to special instances such as  $2 + 3 = 3 + 2$ ).

We demonstrate mathematical properties by reference to principles that are valid in a more general, and often more elementary, arena than that of the objects being studied. A typical such principle, amply celebrated, is the Dirichlet Box Principle, which is the simple assertion that if you put a finite number of objects into a smaller number of boxes, there will be at least one box that contains at least two of those objects. Usually the mathematical objects to which this is applied are rich in specifics, and not abstract ‘objects,’ and the ‘boxes’ similarly. But the simple act of ignoring these specifics clarifies.

The historical first use of that ‘principle’ was to show that any (say, irrational) real number  $\alpha$  can be very closely approximated by rational numbers. Here ‘very closely’ means that there is a sequence of rational numbers  $a_n$  (for  $n = 1, 2, 3, \dots$ ) that converges to  $\alpha$  and such that each of these rational numbers is at least as close to  $\alpha$  as the reciprocal of the square of its denominator. In the penultimate step in the proof of this theorem one sheds all concern about the specific features of the problem except to note that *any collection of  $N$  numbers between 0 and 1 has the property that at least two of them are of distance  $\leq 1/N$  apart*<sup>8</sup>.

A slightly different form of abstraction occurs when one changes the format of a problem to put the problem in a more natural light, or to allow a new source of intuitions to work on the problem. A classic example of this is recasting the theorem (that goes back to Fermat) that any prime  $p$  of the form  $4n + 1$  can be expressed *uniquely* as a sum of two squares  $p = a^2 + b^2$  (with  $0 < a \leq b$ ) in the context of the *Unique Factorization Theorem* for the ring of Gaussian integers.

Often abstraction of some mathematical issue constitutes generalizing it, to view it in a ‘larger’ framework (removing extraneous limitations that get in the way of seeing the truth of a general property). One then has reframed the problem in broader, perhaps simpler, terms, but in which its features and properties are retained, thereby viewing the mechanism that explains those features more naturally and clearly.

---

<sup>8</sup> For any  $N$ , consider the first  $N$  multiples of  $\alpha$  and subtract from each of them an appropriate integer to produce  $N$  positive real numbers  $\leq 1$ . (I.e., letting  $\lfloor x \rfloor$  denote the largest integer  $\leq x$ , we have  $N$  numbers

$$\alpha - \lfloor \alpha \rfloor, 2\alpha - \lfloor 2\alpha \rfloor, \dots, N\alpha - \lfloor N\alpha \rfloor$$

lying between 0 and 1.) So there is an  $i < j$  between 1 and  $N$  so that the difference between the  $j$ -th number in that sequence and the  $i$ -th is  $\leq 1/N$ . Therefore  $(j - i)\alpha$  is at least  $1/N$ -close to an integer. Or, dividing that inequality by  $d := j - i$  we get that  $\alpha$  is at least  $1/dN$ -close to a rational number with denominator  $d < N$ .

## 7 Are there truly different proofs of the *same* theorem?

Different proofs of what we happily call the ‘same theorem’ often are actually proofs of significantly different statements, that imply the theorem—or generalizations of it—and also often offer different *explanations* for why the theorem is true. Perhaps there is *never* actually more than one proof of a theorem, if the theorem,  $X$ , is formulated and understood specifically enough. Perhaps different proofs are inevitably proofs of different theorems related to  $X$ , or at least different shades of  $X$ . Nevertheless, a multiplicity of proofs of  $X$  puts a *dayenu* strain on the naive question: why is  $X$  true? giving an embarrassment of diverse responses<sup>9</sup>.

The classical Pythagorean theorem is a natural starting place if one wants to think about the structure created by multiple proofs of the ‘same’ theorem. This theorem has hundreds of different proofs (e.g., see [3]; see [2] for 118 proofs; also [4], [5]).

### 7.1 The Pythagorean Theorem

If  $\Delta$  is a right-angle triangle, and  $A$ ,  $B$  are the areas of the squares built with edges equal to (respectively) the lengths of the sides of the triangle flanking the ‘right-angle,’ then their sum,  $A + B$ , is the area of  $C$ , the square built with edges equal to the length of the hypotenuse.

To think about this in a concrete context, it is helpful to simply *stipulate* (as lawyers do) that the theorem is indeed true, and then formulate the specific ‘explanation’ each of the multitude of proofs of this theorem provides. As an example, we might take the first proof of this theorem, given as Proposition 47 in Euclid’s Book ([1]).

What this proof ‘explains’ is:

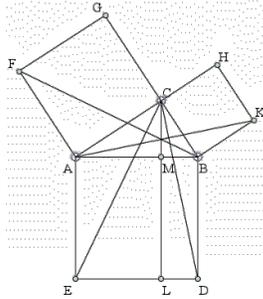
How to cut  $C$  into two rectangles, the first rectangle having area equal to  $A$  and the second having area equal to  $B$ .

**The answer** is neat: drop a perpendicular from the right-angle vertex to the hypotenuse and continue that straight line through the square  $C$  thus cutting it into the two rectangles, each equal,

---

<sup>9</sup> Why did World War I break out?—We are all familiar with the pat headline: the assassination of Archduke Franz Ferdinand at Sarajevo on 28 June 1914. But, of course, that is more the *occasion* of the onset of the war, rather than a satisfactory answer to the complex historical question *why did it occur at all?* Indeed, as for the onset of any war, historians offer us a thick description of motivating causes and crises— including the vacuum in the Balkans created by the weakened state of the Ottoman empire; the competition by major powers in Europe to fill that vacuum; the format of alliances of major European powers that created two armed camps (the Triple Alliance of Germany, Austria-Hungary and Italy *versus* the Triple Entente of Britain, Russia and France), etc. Historians provide such lists of factors—the totality of which represent a constellation of causes—with no claim, I imagine, that any one such cause is definitive, but it is rather the cluster of all of them that provides some explanation. The status of different proofs of a theorem in a mathematical context is qualitatively different: they offer a number of reasons why the theorem is true—each of them sufficient to guarantee that truth; and therefore, for some theorems, a plethora of surprisingly different explanations.

in area, to the square built on the appropriately adjacent side of the triangle  $\Delta$ . That is, we have made a simple construction expressing  $C$  as  $A + B$  as described in the statement of the theorem<sup>10</sup>



## 7.2 Fermat's Little Theorem

Consider the following result in elementary number theory:

If  $p$  is a prime number, and  $a$  is an integer then  $a^p - a$  is divisible by  $p$ .

This has at least three proofs.

### 7.2.1 Induction

One can prove this theorem by induction<sup>11</sup>, by noting first that one may assume that  $a$  is non-negative, since if it is true for  $a$  it is also true for  $-a$ . Then, noting that it is true for  $a = 0$ , and

<sup>10</sup> The proof is based on the observation that  $\Delta BAF = \Delta EAC$ : In fact a 90 rotation around  $A$  brings one onto the other. More formally: show SAS. But  $\Delta BAF$  is (in area) one half the square built on  $AF$  while  $\Delta EAC$  is one half the rectangle  $AMLE$ . So the square built on  $AF$  is equal in area to the rectangle  $AMLE$ . Exactly the same argument 'on the other side of the diagram' gives that the square built on  $BC$  is equal in area to the rectangle  $BMLD$ .

<sup>11</sup> Induction often gets a 'bad Press' as far as its ability to explain the reason behind what it proves. But in many instances, induction has striking explanatory force, as in the proof that sums of consecutive odd positive integers are squares—e.g.,  $1 + 3 + 5 + 7 + 9 = 5^2$ —the inductive proof neatly visualized as squares built out of successive layers of gnomons.

*put diagram here*

then showing that if it is true for  $a = n$  it is also true for  $a = n + 1$ , since

$$(n + 1)^p - (n + 1) = (n^p + \text{terms divisible by } p + 1) - (n + 1). \quad (1)$$

For people familiar with integers mod  $p$ , here's a very minor variant of the above proof<sup>12</sup>.

### 7.2.2 'Automorphisms' of the additive group of integers mod $p$

Taking up Equation (1) in slightly greater generality,

$$(n + m)^p = n^p + \text{terms divisible by } p + m^p, \quad (2)$$

one sees that the mapping  $x \mapsto x^p$  'preserves addition' mod  $p$ . i.e., is a homomorphism for the additive group of integers mod  $p$ :

$$(n + m)^p \equiv n^p + m^p \quad \text{modulo } p, \quad (2).$$

The mapping  $x \mapsto x^p$  sends the generator of the additive group—namely 1—to 1, so must be the identity, which gives the theorem:  $a^p \equiv a \quad \text{modulo } p$ .

For people familiar with group theory the theorem is a consequence of the following structural result:

### 7.2.3 The order of an element in a finite group

*Any* element of *any* finite group, when raised to the power equal to the order of the group is the unit element of the group.

One applies this result to the multiplicative group of nonzero elements in the field of integers mod  $p$ . This group is of order  $p - 1$  (i.e., has exactly  $p - 1$  elements) so

$$a^{p-1} \equiv 1 \quad \text{modulo } p$$

and therefore  $a^p \equiv a \quad \text{modulo } p$ .

What is striking is that these different proofs place the statement of Fermat's Little Theorem into three slightly different—and broader—contexts, the first being the context of induction (i.e., where the successor structure  $n \mapsto n + 1$  of the natural numbers plays a role), the second focuses more on the finite cyclic group<sup>13</sup>; and the third uses the multiplicative group of the ring of integers mod  $p$  and puts the problem squarely into the context of the structure of finite (abelian) groups.

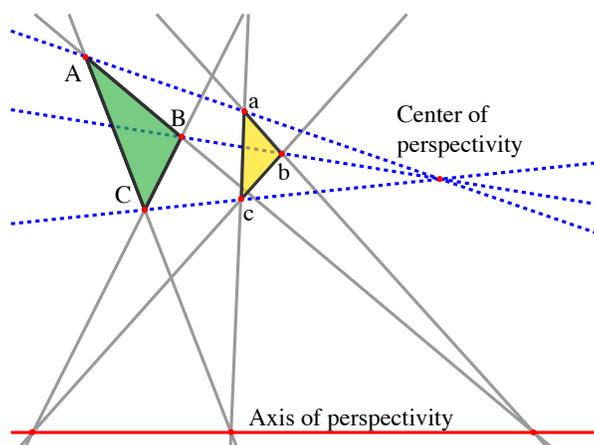
---

<sup>12</sup> which if you unravel it completely, has a hint of induction in it, as well. . .

<sup>13</sup>and the mapping  $x \mapsto x^p$  which plays an important role in a more general context; it is called the **Frobenius mapping** and defines an automorphism of any finite field containing the field of integers modulo  $p$ .

## 8 Change of Setting: Desargue's Theorem

A more complex type of abstract move is to change the context of a mathematical problem so that it is formulated with the least structure necessary to make sense of the problem and to have a uniform answer. This removes whatever complications are introduced by embedding the problem in a framework with irrelevant structure. Once framed in such a minimally structured setting, the problem can often be seen more clearly and the answer to it achieved more directly. A classical example is Desargue's Theorem. Although this theorem is often expressed as a fact in Euclidean Geometry (that if two triangles have a center of perspectivity, they have an axis of perspectivity) such a relationship is unchanged by projective transformations and therefore can be recast as a statement in Projective Geometry, where it is more cleanly dealt with.



### Part II

## Ideas underlying, or analogous to, other ideas in mathematics

### 9 Viewing a mathematical issue as controlled by some *different* underlying conceptual structure

There are times when we don't just embed "X"—the issue we are grappling with—into a larger, but essentially similar, context, as is the case of the genres of abstraction we have discussed. Sometimes

the 'underlying structure' meant to be the explanation of the initial mathematical issue  $X$  can only be formulated in vocabulary (concepts, intuitions) that are utterly different from those of  $X$ . And yet one feels that this structure is the *underpinning* of  $X$ ; it is what (we feel) lies . . . to use a curious word: . . . deeper<sup>14</sup>.

This sense that one constellation of intuitions, concepts, etc. is the underpinning of another one that has utterly different vocabulary and intuitions is harder to defend or even to think about. For me, a major example of such a relationship between two layers of theory—one layer sometimes thought of as the underpinning of the other—is the relationship between the zeroes of the Riemann zeta-function  $\zeta(s)$  and the placement of prime numbers. That is, one approach to the question

$X$ := *how are the prime numbers positioned on the number line?*

is to see  $X$  as controlled by some—if not deeper, then at least significantly different—structure: the zeroes of  $\zeta(s)$ , the *Riemann zeta-function*.

Here is the issue. For any positive number  $n$  let  $\pi(n)$  be the number of prime numbers less than or equal to  $n$ . So  $\pi(2) = 1$ ,  $\pi(3) = 2$  etc. The first type of question one might consider, after knowing that there are—in fact—infinately many prime numbers, is

how fast does  $\pi(n)$  grow as  $n$  goes to infinity?

Clearly  $\pi(n) < n$  since not all numbers less than or equal to  $n$  are primes. But is there a simple formulation of a good approximation to the rate of growth of  $\pi(n)$ ? And, if so, how good is the approximation?

The approach to this question that has, so far, yielded the most fruit is to rephrase it in entirely different vocabulary: if one knows the placement of the complex zeroes of  $\zeta(s)$  one knows  $\pi(n)$ . Moreover, what we already know about those zeroes tells us that there is significant 'regularity' in the growth of  $\pi(n)$ . Specifically:

**The Prime Number Theorem:**  $\pi(n)$  is asymptotic to  $n/\log n$  as  $n$  goes to infinity.

That is, the ratio  $\frac{\pi(n)}{n/\log n}$  tends to 1 as  $n$  goes to infinity.

But we expect more profound understanding of the statistics of the placement of prime numbers to come from more detailed knowledge of the zeroes of  $\zeta(s)$ ; namely, from a proof, when it comes, of the *Riemann Hypothesis*. The Riemann Hypothesis (still unproved) proposes a simple to state, and very clean, geometric pattern for the placement of the zeroes of the zeta-function: *the nontrivial zeroes lie on the vertical line  $\frac{1}{2} + it$  in the complex plane*. Such a pattern implies a significantly deeper *statistical* relationship for the placement of prime numbers on the number line than the Prime Number Theorem gives us. That a pattern of behavior as simple as "lying on a straight

---

<sup>14</sup>"Deep" is a word that becomes more curious the more you think about it. We'll return to it in Section 10 below.

line” in the world of zeroes of the zeta-function translates to a less immediately visible statistical behavior of prime numbers is striking.

The temptation is to think that the placement of the zeroes of  $\zeta(s)$  has a kind of precedence: it *underlies* the structure of the placement of prime numbers.

- *The arguments against this thought are strong:*

The relationship between the zeroes and the primes are, in a sense, two-way: there is little reason to imagine, or to impose, a hierarchy: each determine the other just as the Fourier transform of a function determines the function and vice versa. The epistemological time-line view would simply have it that we can actually *prove* something about the placement of zeroes of  $\zeta(s)$ , from which we deduce something about the placement of primes, and nothing more than that. Namely, we have the standard proof of the Prime Number Theorem, which does exactly that; the zeta-function is a convenient way-station in the epistemological route to the Prime Number Theorem. Moreover, as the Wikipedia page mentions (as quoted in the footnote above), there are other ways of deducing that theorem.

So, if we have this ‘feeling’ of the essentiality of comprehending the zeroes of  $\zeta(s)$  as a key to understanding primes, it is not, it seems to me, imposed on us by any purely logical structure in the situation.

- *The arguments in favor of this thought:*

First, a general issue: when we contemplate mathematics—and most certainly when we teach it—we naturally complement the merely logical structure of our demonstrations and deductions with a filigree of causal lines-of-thought that engage our intuitions and are meant to explain—meant to *ground*—the properties of one concept by means of its relationship to some other concept. These guy-wires of explanation constitute—for each of us—the way we view the subject. Without any such—perhaps subjective—structure of explanation, we are often lost.

Here is a reason for doing exactly this with the zeroes of the zeta-function. That the nontrivial zeroes might lie on a straight line itself connects with, and would itself be explained by, a classical (again, of course, only conjectured) scenario, credited to Hilbert and Pólya, that proposes the existence of a *natural* Hilbert space with unbounded self-adjoint operator that would have the property that the eigenvalues of this operator (appropriately normalized) are equal to the nontrivial zeroes of the zeta-function. The word *natural* is important. From the existence of this Hilbert space with self-adjoint operator the Riemann Hypothesis would follow, and the prime numbers themselves would thereby be embedded in an extremely structured world of concepts that might reveal far more than what we know, or even currently have conjectured, about prime numbers. Of course, this Hilbert space and operator may well be a chimera, but might also be the key to a deeper, perhaps more essential, way of thinking about primes<sup>15</sup>.

---

<sup>15</sup> Here is an analogous judgment. The structure of rational solutions of systems of polynomial equations modulo a prime  $p$ , and, more generally, over finite fields of characteristic  $p$ , is largely restrained by—explained by—the action of the Frobenius operator on the étale cohomology of the relevant algebraic variety. I, and I think many others, feel a palpable mathematical causal link here: *the cohomological structure underlies the epiphenomenon of ‘number of*

I use the phrase *deeper structure* but if I were challenged about this: *Is the reason you feel this just habit, or merely a result of the way you happened to have learned it?* the only defense I could offer is just that I “feel it” to be deeper; a very poor defense. Whether or not such ‘feelings of essentiality’ exist generally in the minds of people thinking about mathematics, and—if so—whether or not they correspond to anything other than mere subjective *feelings* connects with much standard philosophical discussion about the nature of mathematics. But whatever status these extra-logical feelings have, at least for me—and I suppose for many mathematicians—they are key to the vibrant explanatory structure of mathematics. What, then, is a satisfactory way of thinking about—or at least classifying, or naming—such extra-logical keys to understanding?

## 10 Depths and Heights

My friend Eva Brann told me that—in her view—the opposite of the notion *deep* is *complex*, a word that signifies ideas that, in reflection, spread to a broad diffusion of ramifications. In contrast to complex ideas, there are the modes of thought that get—more concentratedly—to the bottom of things. Those are the deep ones.

Deeper mathematics, higher mathematics—both phrases are used referring to perhaps multi-layered ‘strata of understanding’ where the lower strata are the underlying support of the middle-ground from where, it is assumed, we have started. The higher strata represent the over-arching goal from which, if attained, we can look down and see with utter clarity, all below. The vertical direction, though, is what seems to be insisted upon in these adjectives.

In Mathematics, the *Theory of Motives*, launched by Alexander Grothendieck already signals in its very name the flavor of its intent: to be the source, and the organizational underpinning, of certain more visible mathematical structures. The Theory of Motives was meant to focus on—to bring about—a primary theory, a common root—a precursor, so to speak—for the various (cohomology) theories that have been defined to understand the objects of algebraic geometry.

Depth and height—referred to as ‘strata’—make them metaphorical inert geological formations. More pertinent is when you and I enter this metaphor as spelunkers, or climbers: searchers in this geography. If we are doing our job, we *descend to the depths* and *ascend to the higher realms*. Both activities seem to be deemed epistemologically virtuous. We dig deep to uncover what is hidden under the surface terrain<sup>16</sup>. Various types of “bedrock” have been proposed. Aristotle reported that

---

*rational points over finite fields.*’ This structure—Frobenius, cohomology—was brought into play as the confirmation of the classical aperçu of André Weil, that rational points over finite fields might be thought of as entities of geometry: fixed points of the Frobenius mapping, therefore subject to those underlying structural constraints that count the number of fixed points of a mapping of a space as the trace of the linear transformation induced by the map on global (cohomological) invariants of the ambient space. It may be that the deeper structure of rational points over finite fields is only adequately understood if one understands this.

<sup>16</sup> Or, thinking of the medium as water, we descend as undersea explorers the activity being

“reminiscent of pearl-diving: one explores the depths of a perilous medium to the limit of one’s breath and life in order to bring back a rare and valuable treasure that could be found nowhere else.”

(*Hinges*— Grace Dane Mazur)

according to Thales, *water* is the nature of all things, and that the pythagoreans held that *number* is the fundamental structure of the cosmos. The early materialists, Democritus, Leucippus, and (more subtly) Lucretius would have it that indivisible particles (i.e., atoms) form the underlying structure. Contemporary Physics proposes, of course, a magnificent awe-inspiring underpinning for 'things.'

The Wikipedia entry [https://en.wikipedia.org/wiki/List\\_of\\_mathematical\\_jargon](https://en.wikipedia.org/wiki/List_of_mathematical_jargon) defines a *deep result* (specifically, in mathematics) to be one whose

proof requires concepts and methods that are advanced beyond the concepts needed to formulate the result. The prime number theorem, proved with techniques from complex analysis, was thought to be a deep result until elementary proofs were found.

As far as it goes, this definition points to an important feature of the way we use that word: a “deeper” viewpoint is often marked by a change of the constellation of concepts invoked. But as I have already mentioned, I don’t agree with the example cited. A more basic disagreement I have, though, is that “depth” is attributed, on this Wikipedia page, to the disembodied 'result;' to me it makes more sense to acknowledge that the judgment of depth is essentially subjective; it is often an extremely important judgment, but the thinker is the active agent here.

## 11 Horizontal, rather than Vertical

There are many 'hybrid subjects' in mathematics: two subjects joined, getting the benefit of two, quite disparate, intuitions. In many striking instances, a problem naturally posed in one of these joined subjects has a transparent solution, using the intuition of the other<sup>17</sup>. *Algebraic Geometry* is—as its name implies—such a hybrid subject. It makes use of algebraic intuition on the one hand and geometric intuition on the other, these being fused together forming a subject that has its own synthesis of those two intuitions. Algebraic geometry is very much like its precursor, René Descartes' analytic geometry which—as Descartes proclaimed—corrects the defects of both algebra and geometry. Such a merger of two subjects, with different vocabularies—each of which serving to explain the other—is not unique, in mathematics. It is not that one of these subjects is on a deeper, more explanatory, level than the other. They are partners, on the same level. They are joined 'horizontally', rather than 'vertically.'

In current mathematics such a proposed join of two subjects is being vigorously pursued. It is called the *Langlands Program* which establishes (still largely conjecturally) a bridge between *analysis* (specifically in the context of group representation theory) and *arithmetic* (more specifically, algebraic number theory)—where neither one of these subjects is viewed as underpinning for the other. The Langlands Program is a bridge—a Rosetta Stone—relating deep properties of one of these subjects to deep properties of the other.

---

<sup>17</sup> One simple example of this in *Combinatorial Geometry* is that the easiest way to understand subgroups of a free group is to think of the free group as the fundamental group of a graph (i.e., of a bouquet of circles).

## 12 The personal guy-wires of explanation

Might it be illuminating to think more about how, as we try to understand—and become at home with—mathematics, or any subject, we inevitably construct some constellation of explanatory hierarchies?

### References

- [1] <http://aleph0.clarku.edu/~djoyce/java/elements/toc.html>
- [2] <http://www.cut-the-knot.org/pythagoras/>
- [3] *Hidden Harmonies: The Lives and Times of the Pythagorean Theorem* by Ellen and Robert Kaplan, Bloomsbury (2011)
- [4] <https://www.youtube.com/watch?v=ItiF05y36kw>
- [5] <http://www.math.harvard.edu/~mazur/papers/moresymmetry.pdf>