

Conceptual Foundations of the Calculus and the Computer

David Tall

Mathematics Education Research Centre
University of Warwick
COVENTRY
CV4 7AL
UK

Introduction

In this paper I consider conceptual problems which students are known to encounter in learning the calculus and show how computer software based on “local straightness” can provide an environment for an approach which is cognitively satisfying, gives a meaning to the Leibniz notation which is in accord with the original conception, and provides a conceptual introduction to modern theories of analysis, both standard and non-standard.

Conceptual problems in the calculus

(i) problems with the limit concept

In recent years the difficulties of teaching the calculus have become a focus of attention at college level. In the late seventies, my colleague Rolph Schwarzenberger and I did some research with students which showed that the notion of limit is a singularly difficult concept for the beginning student. (Schwarzenberger & Tall 1978). We found that the students developed a conception of the processes arising from their experience of mathematics and the use of the language that caused them to formulate implicit ideas which were at variance with the mathematical theory. Subsequent research revealed that the quantifiers occurring in the limit definition caused enormous problems. Students at university have difficulty even remembering the definitions by rote, let alone using them in a mathematical way.

In Britain it is common practice in the sixth-form (senior high school, age 16 to 18) to approach the limit concept dynamically, “as $x \rightarrow a$, so $f(x) \rightarrow c$ ” to describe the limit of a function, and to build the calculus on (what was considered to be) an intuitive limit notion of derivative where the gradient of the curve $y=f(x)$ is defined to be the limit of $\frac{f(x+h)-f(x)}{h}$

as $h \rightarrow 0$. This was supported by a visualization of a secant “tending” to a limiting position as a tangent to help the students to “see” what is going on. Research performed in the late 1970s showed that these concepts were less intuitive than many teachers believed and, as soon as the first introduction of the formalities had passed by, the students simply concentrated on the mechanics of calculating the derivative symbolically. Here they faced further difficulties.

(ii) problems with the Leibniz notation

The derivative is often written using the Leibniz notation dy/dx , and it sometimes seemed to behave like a fraction, as in the chain rule

$$\frac{dy}{dx} = \frac{dy}{dt} / \frac{dx}{dt}$$

But students are usually told that they must not consider dy/dx to be a fraction, but some kind of “useful fiction” invented by Leibniz which happens to work by some kind of wizardry. For instance, in *Teaching the Calculus*, Shuard & Neill wrote:

The student ... has to learn that, in spite of all the evidence to the contrary, which seems to him to build up from statements such as

$$\frac{dy}{dx} \times \frac{dx}{dt} = \frac{dy}{dt}$$

dy/dx is not a symbol for a fraction, but for the limit of the gradient of a chord. (Shuard & Neill, p. 13).

A school textbook advised students:

‘ dy/dx ’ must, at least for some considerable time, be regarded as an inseparable whole ... It does not in any simple or straight-forward way mean anything like ‘ dy divided by dx ’, and a statement such as

$$dy/dx \times dx/dt = dy/dt, \text{ by cancelling } dx$$

is just so much gibberish. (SMP Advanced Mathematics, p. 221)

Students are told that the dx in $\int f(x) dx$ means “with respect to x ”, and should not be thought of as a separate symbol, although they need to be willing to make the substitution $du = \frac{du}{dx} dx$, to compute the integral by substitution.

Then the indivisible symbol $\frac{dy}{dx}$ in the differential equation

$$\frac{dy}{dx} = -\frac{x}{y}$$

is suddenly written as

$$y \, dy = -x \, dx$$

to “separate the variables” and an integral sign is introduced to give

$$\int y \, dy = -\int x \, dx$$

(where presumably dx suddenly changes its meaning to “with respect to x ”), to obtain the solution(s)

$$\frac{y^2}{2} = -\frac{x^2}{2} + c.$$

Under such circumstances, what are students to believe? Is it all a matter of bogus trickery that one is taught to indulge in to “get the right answer”?

(iii) problems with algebra

In Britain a much smaller proportion of the population studies the calculus than in America, and it begins earlier, at age 16 in school. But as access to higher education is being broadened, the chasm between what students are asked to believe and what they feel capable of understanding grows ever wider.

In the last two decades in Britain, the move to comprehensive education has resulted in algebra being taught to an increasing proportion of the population, with the result that it is perceived as being increasingly difficult to teach and learn. At the same time, investigations and problem-solving are on the increase and there is less attention given to drill and practice. So students were getting progressively less able to cope with the algebraic demands of the calculus. The deterioration has now reached such a level that the new curriculum designed by the School Mathematics Project for the 16 to 19 age group considers that the manipulation of the expression $\frac{(x+h)^3-x^3}{h}$ to calculate the derivative of x^3 is too difficult for most students. University students seem to have less facility with symbolic manipulation. With the decrease in drill and practice, they seem to have difficulties simplifying expressions which earlier generations would have regarded as routine.

In the early 1980s I therefore found myself faced with a number of seemingly inescapable conclusions:

- The modern limit concept is extremely difficult for students to understand and is not a natural starting point for their understanding of the calculus,
- There is a mismatch between the notation of Leibniz and the formalities of modern analysis, yet students are expected to perform operations using this notation in a mechanical way whilst overlooking any deficiencies in meaning.
- algebraic skills seemed to be on the decrease and there was a need for a new approach which demanded less manipulation yet more meaning.

Insights from non-standard analysis

For several years I taught non-standard analysis in which the quantifiers involved are somewhat easier to handle than in standard analysis. Dynamic limiting ideas in which quantities “tended to zero” or grew “as small as one pleases” tended to give students beliefs in “arbitrary small quantities” (Cornu 1983). In some senses, therefore, non-standard analysis offered certain advantages. In particular, it showed that under an infinite magnification (when the “standard part” of each number is taken) the graph of a differentiable function is straight. For instance, $y=x^2$ near $x=a$ has the form $(x+h)^2$ where h is small. The graph has gradient $\frac{(a+h)^2-a^2}{h} = 2a+h$, and if h is infinitesimal, the standard part is $2a$. Thus under infinite magnification, neglecting infinitesimals, the gradient of the graph near $x=a$ is $2a$ and the graph looks straight. However, I found that there are certain difficulties in the initial stages of non-standard analysis too. It is natural for students dealing with variable quantities that become small to imagine that these produced objects which are “arbitrarily small”. The intuitive belief in infinitesimals is therefore strong. However, the cognitive imagery is often at variance with the non-standard theory. (For instance, $0.9999\dots$ is believed by many students to be “the largest number less than one” whilst $0.999\dots$ to N places differs from 1 by $1/10^N$, so that even if N is a (non-standard) infinite integer, there is a number between this and 1.) The formal approach, involving a field \mathbb{R}^* which is an extension field of the real numbers also requires a greater sophistication than is present at the time that students begin the calculus. So, although non-standard analysis has some clues, as a formal theory there are still difficult entry points for the learner.

The computer and a locally straight approach to the calculus

Then the computer arrived and it gave an impetus for a new development in the theory of calculus, perhaps more appropriate for beginning students. Based on the non-standard idea of infinite magnification, I considered finite magnification for a suitably large scale factor. This turned out not to be very large at all. It was possible to magnify a graph by a factor, say 100, and many standard graphs would look so much less curved that they seemed almost straight (figure 1).

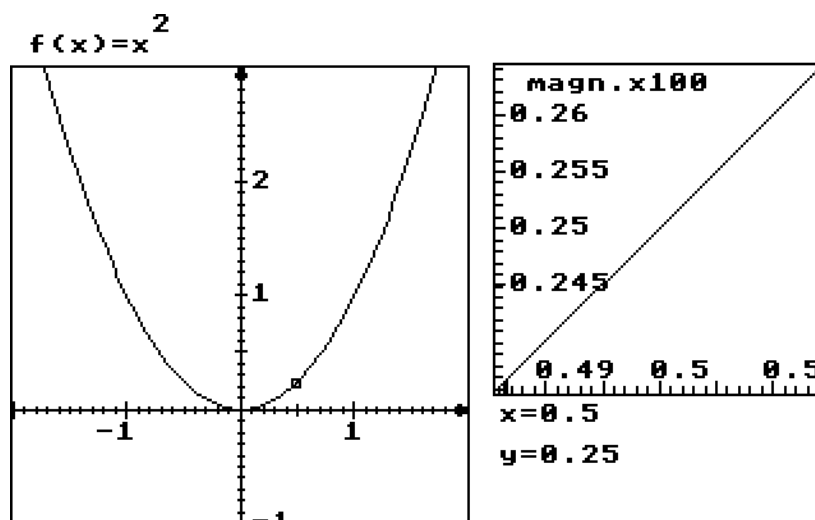


Figure 1: magnifying a small part of a graph

I therefore began to build an approach to the calculus based on this notion of “local straightness” which is now available for IBM compatibles (Tall et al, 1990). It transpires that the magnification process has the limiting notion *implicit* within it. Thus it is not the initial focus of attention. Instead one can use the computer to consider examples in which the graph is locally straight and non-examples in which the graph fails to magnify locally straight, say it has a “corner” with a different left and right gradient, or it oscillates so wildly near a point that the graph never magnifies to look straight, or it is so wrinkled that it is *nowhere* locally straight. By including such a function in the software I developed it became possible for students *in their very first encounter with the calculus* to meet the idea of a function which is intuitively *everywhere continuous but nowhere differentiable*. (figure 2).

$$f(x) = b1(x)$$

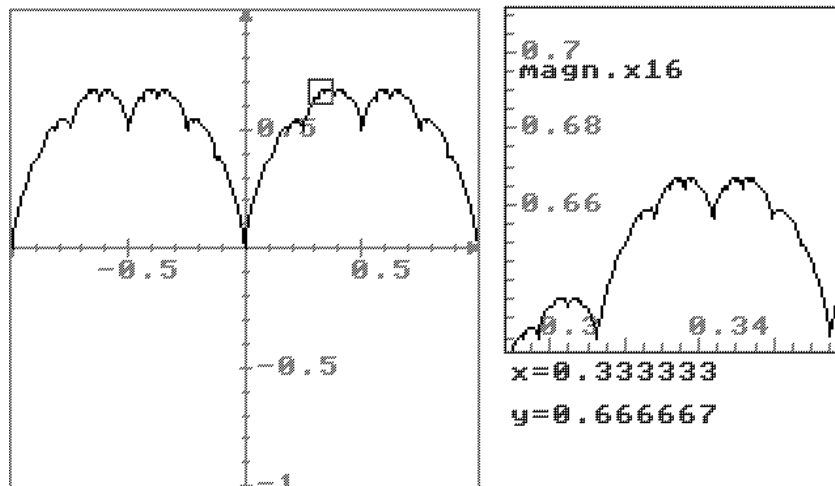


Figure 2: A graph which is so wrinkled it is nowhere locally straight

New meanings for old notations

This approach also gives a new meaning to the Leibniz notation, which proves to be equivalent to the old meaning formulated by Leibniz. Initially I could only see it applying to the notion of the derivative, but more recently I have looked at the pictures a little more sensitively and now I can see a coherent use of the Leibniz notation throughout the whole of differential and integral calculus, and in the meaning of differential equations.

In the first publication on the calculus in 1684 Leibniz referred to a diagram which is shown simplified in figure 3 by referring only to the standard variables x, y .

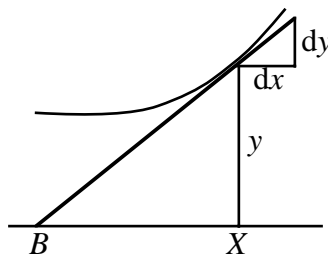


Figure 3 : The Leibniz definition of dx and dy

The curve represents a variable y depending on x , and B is the point where the tangent to the curve meets the x -axis.

Condensing what Leibniz said to concentrate on the variables x, y we get the statement:

Jam recta aliqua pro arbitrio assumpta vocetur dx , & recta quae sit ad dx ut y est ad XB vocetur dy .

which translates to

Now some straight line selected arbitrarily is called dx , and the line which is to dx as y is to XB is called dy .

Thus the length dx is arbitrary and the length dy is the corresponding increment in y such that the quotient dy/dx equals y/XB . Disentangling the definition, we see that dx is any increment and dy is the corresponding increment to the tangent (figure 4.)

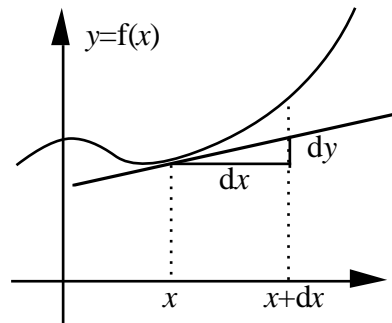


Figure 4 : The differentials of Leibniz as increments to the tangent

There is no mention of infinitesimals: they came later in the paper when Leibniz had to develop a method of calculating the direction of the tangent. Today we (usually) calculate the tangent direction by a limiting process, but there is no reason why we should not use the Leibniz notation in its original meaning. In modern terminology, the tangent is drawn to the curve and the components of the tangent vector are dx , dy . Provided that dx is suitably small (dependent on the curvature of the graph in a naive sense), the graph will approximate to the tangent and a small part of the tangent will look to the naked eye much like the locally straight graph.

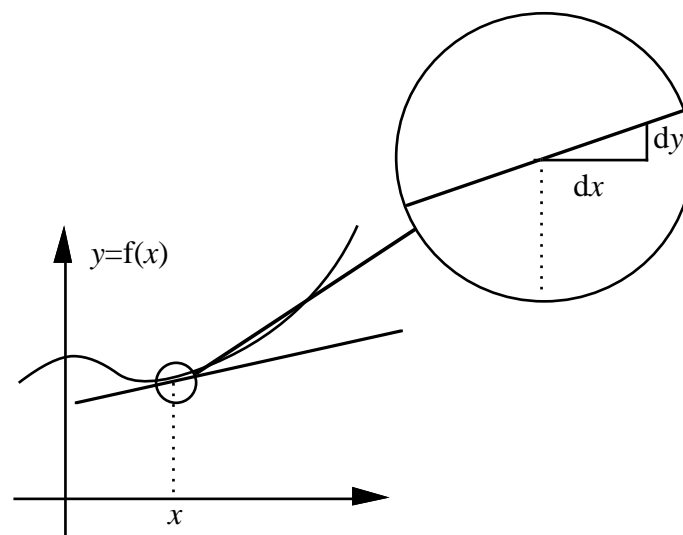


Figure 5 : magnifying a locally straight graph

Investigational approaches with new technology

In the new School Mathematics Project in Britain for 16-19 year old students in school, we therefore designed a numerical and graphical approach to the calculus in which the gradient of the graph was investigated both numerically and pictorially. Using investigational problem-solving techniques it became possible to conjecture (guess) the gradient functions of many standard graphs: x^2 , x^3 , x^n , $\sin x$, $\cos x$, e^x , $\ln|x|$, etc, sufficient to give a meaningful start to the subject without too much manipulation and without explicitly using the limit concept.

For instance, a graph such as $y=x^2$ will look almost straight when any small part is magnified (figure 6). This means that the rate of change of y with respect to x , which is found by measuring

$$\frac{\text{y-change}}{\text{x-change}}$$

will not vary much over the magnified part.

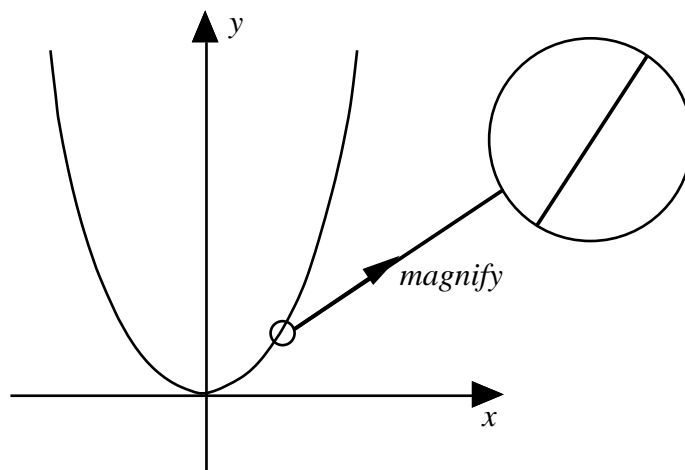


Figure 6 : Magnifying the graph of $y=x^2$

It is this idea that such a curve magnifies to look straight when a small portion is placed under a microscope that makes the calculus possible.

A student quickly learns to scan an eye along the graph and see various parts of it changing in gradient. Just by *looking*, the gradient can be seen decreasing from a large negative gradient, getting less and less steep until the gradient is zero at the origin, then increasing for positive values.

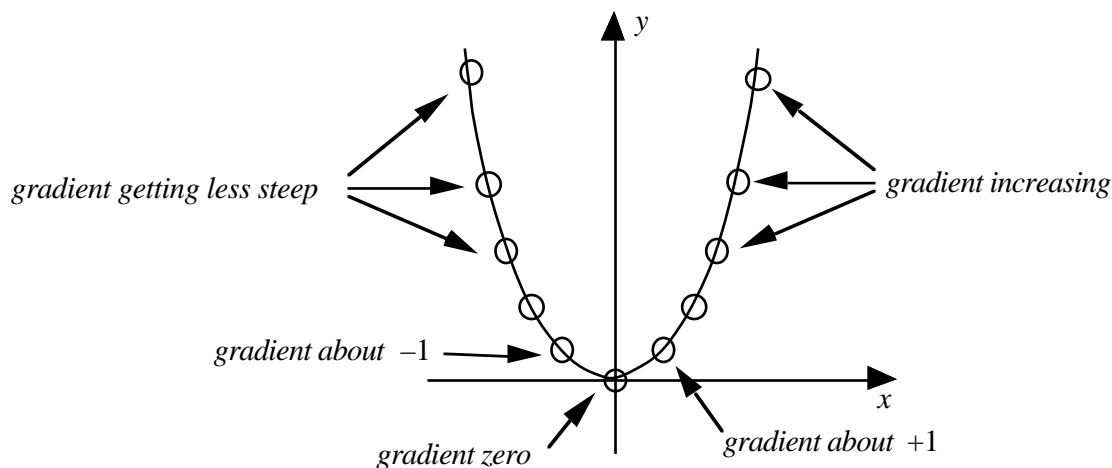


Figure 7: The changing gradient of $y=x^2$

By using software that sketches this gradient more accurately, the students may *see* that it stabilises on $y=2x$ and thus the gradient curve can be conjectured before any symbolic manipulation is performed. By such methods students in the SMP curriculum are now being encouraged to *see* the changing gradient of a graph and to use computer software to give them the environment in which they may conjecture the formula for the gradient of the graph. (It is my personal belief that this particular syllabus goes too far by relying mainly on visual imagery and conjectures. Visual imagery is only of lasting value when it can be turned into mathematics that can produce answers. It is my hope that future iterations of the syllabus will demand more linkage between pictures and symbols, and more expectation of manipulative ability, but only time will tell.)

Undoing differentiation

Traditionally integration is seen as the reverse process to differentiation. This is not the most appropriate way of viewing the problem. The reverse of differentiation is knowing the gradient $\frac{dy}{dx}$, to find the original function $y=f(x)$. This is the theory of *differential equations*, not the theory of integration. A solution of a first order differential equation must, by definition, be differentiable, In the small, a tiny portion of its graph must approximate to a straight line segment and, if we know the gradient of the segment, then we can draw it. For instance, if we have

$$\frac{dy}{dx} = \frac{1}{2} y.$$

which tells us that the gradient of the original function through (x,y) is $\frac{1}{2} y$, then we can draw a line segment of gradient $\frac{1}{2} y$. By visually sticking such line segments end to end we can build up an approximate

solution curve. Figure 8 shows the *Solution Sketcher* software written for the new 16-19 A-level with a line segment drawn through $x=1.5$, $y=2$, where the gradient is therefore $\frac{1}{2}y=1$.

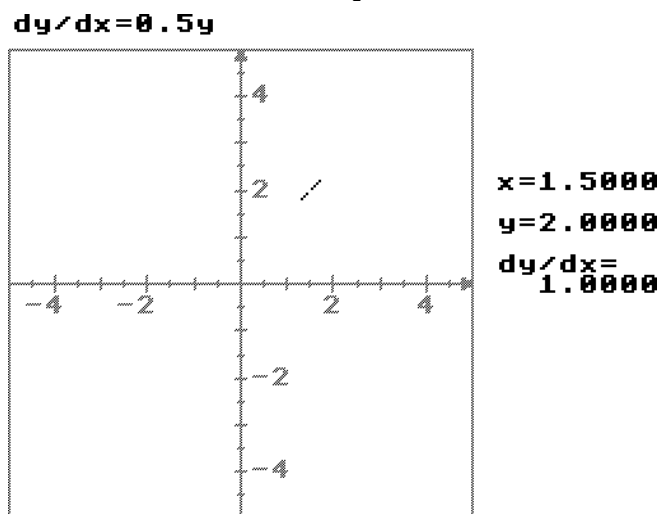


Figure 8 : part of a curve of gradient $\frac{1}{2}y$

By simply moving the segment around the computer screen under software control, and leaving a trace of such curves fitting end to end, an approximate solution curve can be built. A student performing this activity will sense physically that the resulting curve everywhere has gradient given by the differential equation and therefore will have a deeper cognitive understanding of the meaning of a solution.

Figure 9 shows such a curve and an array of line segments showing the directions of other possible solution curves. Through each point in the plane there is a unique solution of the differential equation

$$\frac{dy}{dx} = \frac{1}{2}y.$$

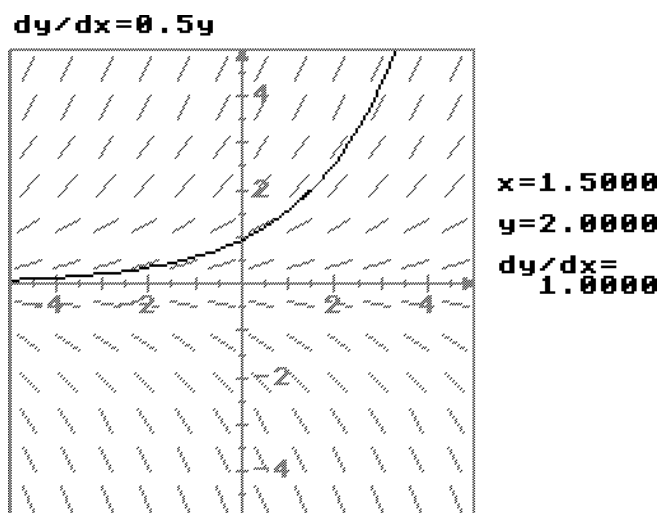


Figure 9 : A solution curve and an array of other segments with the appropriate gradient

This is the essential idea of a differential equation. It is an equation which specifies the gradient of a graph and, provided it does this unambiguously at every point, there is a unique solution of the differential equation through every point in the plane. This leads to the idea which is very new to the beginning student, that *there is a whole family of solutions to a given differential equation.*

Cumulative growth: integration

The final ingredient of the calculus is given by the *cumulative growth* of a function. The most straightforward example is to take a function and calculate the growing area under the graph. Say one might take a graph as in figure 10 and, by some method or other, work out the area from a fixed point a to a variable point x . The area will then be a function $A(x)$.

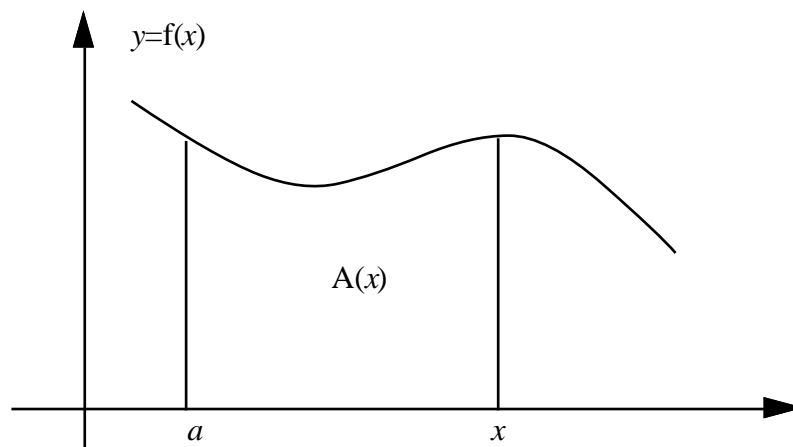


Figure 10 : The area $A(x)$ under the curve from a to x

The area from $x=a$ to $x=b$ is easily calculated approximately on a computer by simply simply chopping up the interval from a to b into small lengths which, in the Leibniz notation will be denoted by dx , and then adding together the rectangles of height $f(x)$ and width dx (figure 11).

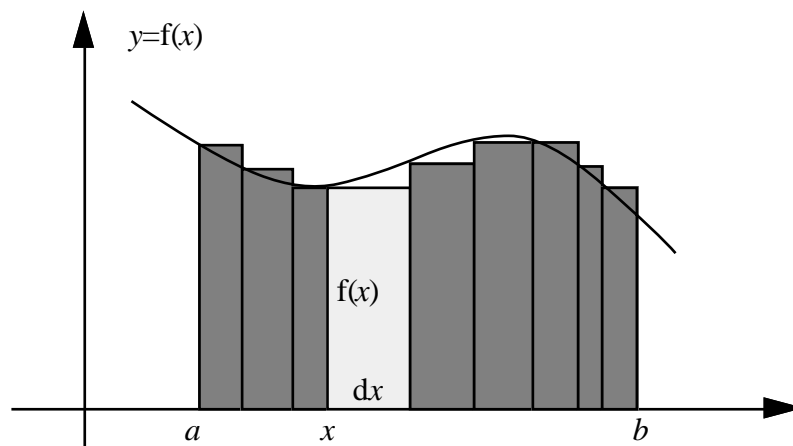


Figure 11 : The Leibniz sum

By taking thinner and thinner strips the Leibniz sum is likely to give better and better approximations to the area.

Leibniz denoted the area by

$$\int_a^b f(x) dx$$

where the elongated S denoted the first letter of the Latin word Summa, for sum. This Leibniz notation should therefore be read “the sum from a to b of $f(x) dx$ ”. It is also called “the *integral* from a to b of $f(x) dx$ ”.

It is this notation that has been the subject of most misunderstanding, compared with the dy/dx notation in differentiation. There is a reason for this. The limiting process for the derivative:

$$\lim_{h \rightarrow 0} \frac{f(x+h)-f(x)}{h}$$

occurs only at *one* point, with only *one* variable (h) tending to zero. The limiting process for the numerical approximation to the integral takes many strips and must occur over a whole interval. It is the extra complication of this process that seems to have prevented an adequate modern interpretation of the Leibniz notation in integration.

Let us use the notation $\sum_a^b f(x) dx$ to denote the sum of strips for a

partition of the interval $[a,b]$ in which dx denotes the (finite) strip width of a typical strip and x denotes the left-hand endpoint of the strip. (This is the essential definition of Cauchy at the beginning of the nineteenth century). Figure 10 shows this sum of strips seen as an area calculation.

However, a much more productive way to consider this notation is to look instead at the graph of a function $I(x)$ whose derivative is $I'(x)=f(x)$ and instead draw the corresponding graph of $I(x)$ (figure 12).

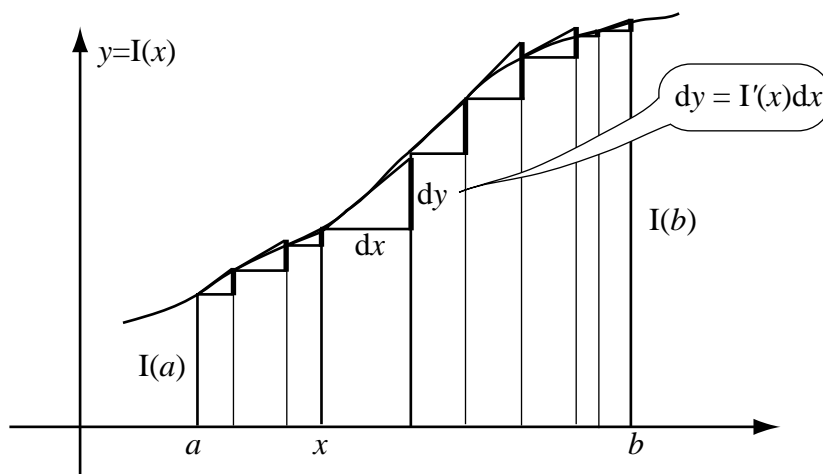


Figure 12 : The Leibniz sum as a sum of vertical line segments

This picture has the same subdivision of the interval from a to b into subintervals. As we saw earlier in figure 6, the value of $I'(x) dx$ is equal to the vertical distance to the tangent. Because $f(x)=I'(x)$, the Leibniz sum

$$\sum_a^b f(x) dx = \sum_a^b I'(x) dx$$

is the sum of these vertical segments.

The picture as it stands is not very helpful. But now we must see it through our new spectacles and set Leibniz's reputation straight. Figure 13 shows the same idea with a large number of very thin strips.

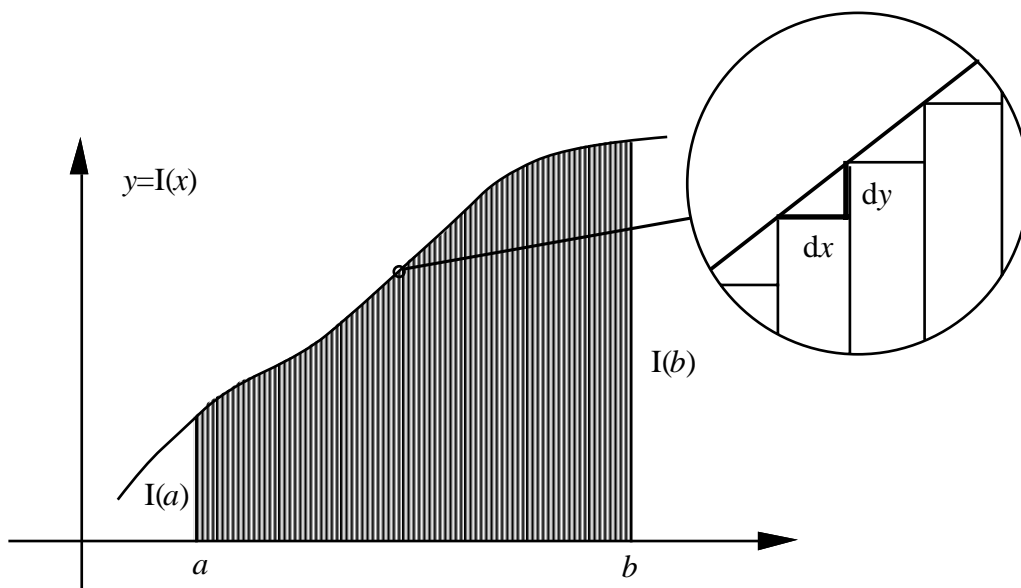


Figure 13: The Leibniz sum as the sum of the risers to the curve $y=I(x)$

Because the graph of $y=I(x)$ has a derivative, under a microscope it will look straight. So now the vertical steps are approximately equal to staircase risers. Adding together all these risers when the strips are very small gives the total sum to be the rise from $x=a$ to $x=b$, which is

$$I(b)-I(a).$$

Thus, when a very large number of strips are taken, the value of the Leibniz sum will stabilise on a value which we denote by $\int_a^b f(x) dx$, which satisfies

$$\int_a^b f(x) dx = I(b)-I(a)$$

which is simply the *Fundamental Theorem of the Calculus*.

Towards a formal proof of the Fundamental Theorem

In other papers (e.g. Tall 1986 and Tall 1991a) I have shown how a different method of visualising the area under a graph by pulling it horizontally can lead to a proof of the fundamental theorem for a continuous function. The idea is simple. interesting pictures might occur by maintaining a constant y -range whilst taking a much smaller x -range. For instance, figure 14 shows the graph of $y=\sin x$ with the same y -range in each (-3 to 3) but x -range being changed from -3 to 3 down to 1 to 1.01 . What happens is that the graph in the second case is pulled flat by the stretching of a thin x -range to fill the computer window.

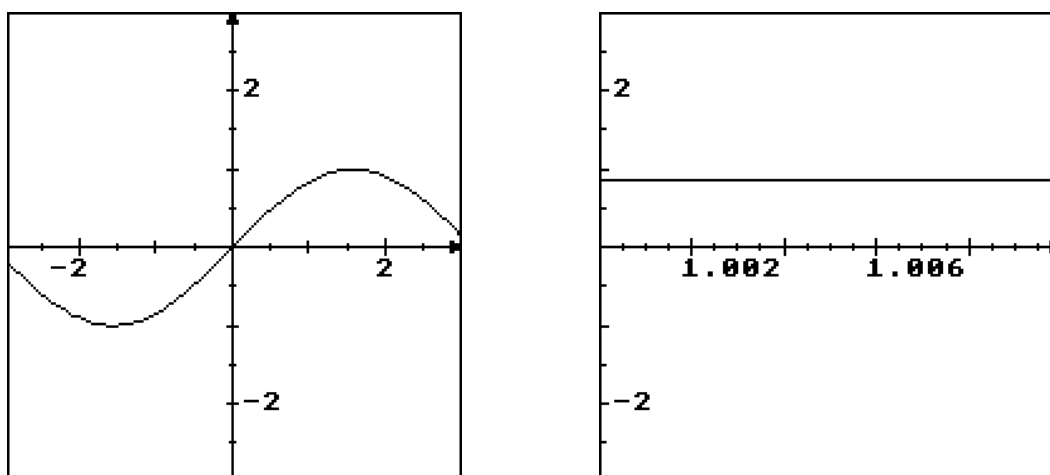


Figure 14 : The graph of $y=\sin x$, pulled out flat

If one calculates the area under a flat graph like this, the area from x to $x+h$ is approximately $f(x)h$. This represents a change in area from $A(x)$ to $A(x+h)$, so

$$A(x+h)-A(x) \approx f(x)h.$$

This suggests that we may have

$$\frac{A(x+h)-A(x)}{h} \approx f(x),$$

and perhaps, as $h \rightarrow 0$, we might get

$$A'(x)=f(x).$$

What kind of function, when stretched out horizontally near $x=x_0$, looks flat ? If we suppose this means that the graph lies in a pixel representing a height $f(x)\pm\epsilon$, then we need to know that, given such an $\epsilon > 0$, then we can find a small enough x interval, say $x\pm\delta$, so that when t lies between $x-\delta$ and $x+\delta$, then $f(t)$ lies between $f(x)-\epsilon$ and $f(x)+\epsilon$. In other words, a natural condition for the function to satisfy the fundamental theorem is that it be *continuous* in the formal sense:

Given any $\varepsilon > 0$, a $\delta > 0$ can be found such that whenever $x - \delta < t < x + \delta$ we know that $f(x) - \varepsilon < f(t) < f(x) + \varepsilon$.

For such a function if the width of a strip h is taken positive and less than δ then the value of $f(t)$ will lie between $f(x) - \varepsilon$ and $f(x) + \varepsilon$ throughout the strip and so

$$(f(x) - \varepsilon)h < A(x+h) - A(x) < (f(x) + \varepsilon)h.$$

Hence

$$\frac{A(x+h) - A(x)}{h}$$

is sandwiched between $f(x) - \varepsilon$ and $f(x) + \varepsilon$. A similar argument holds for h negative. As ε is arbitrary, this is the formal definition that

$$\lim_{\varepsilon \rightarrow 0} \frac{A(x+h) - A(x)}{h} = f(x),$$

i.e. $A'(x) = f(x)$.

In this way we see that the formal notion of continuity arises as a natural ingredient of the fundamental theorem of calculus, not as an esoteric definition introduced for purely theoretical purposes.

Note that the fundamental theorem requires only requires the *continuity* of $f(x)$, not that it be differentiable. Using another piece of software, *The Function Analyser* (Tall 1991b), I have programmed a numerical way of calculating the approximate area *as a function* $A(x)$. This uses the speed of a 32 bit RISC chip in the British Archimedes computer to calculate the area quickly from a fixed point a to a variable point x with specified strip-width using the mid-ordinate rule. It is therefore possible to numerically differentiate the area function and get back to the original function. Thus it is possible to integrate the nowhere differentiable blancmange function to get an area function differentiable everywhere once (whose derivative is the blancmange) and is therefore differentiable nowhere twice.

I have even amused myself by programming in a function that distinguishes between “pseudo-rationals” and “pseudo-irrationals”, giving the value TRUE to 1 , $22/7$, yet FALSE to π , e and $\sqrt{2}$. This uses the Greek method of continued fractions to approximate to a number to work out a close rational approximation and declares the number “pseudo-irrational” if the approximation has a large denominator. Using this approach I have been even able to simulate functions which take the value 1 on (pseudo)irrationals and 0 on (pseudo)rational to start the beginnings of a theory of integration which gives intuitions for

Lesbegue integration rather than Riemann integration. But that is another story! (Mills & Tall (to appear).)

Parametric functions

In other software developed for the School Mathematics Project (Tall 1991b) I programmed four views of a parametric curve $x=x(t)$, $y=y(t)$, to simultaneously show the curve in three-dimensional t - x - y space and its component pictures projected onto the three coordinate planes. The three-dimensional view can be rotated in space to coincide with any of the three projections. A numerical tangent can be calculated through two points with parameter t differing by a small value dt . This gives a tangent vector with coordinates that may be denoted by (dt, dx, dy) , whose projections give the standard Leibniz picture in any of the three coordinate planes. It is in this way that the equation:

$$\frac{dy}{dx} = \frac{dy}{dt} / \frac{dx}{dt}$$

can be seen to have meaning – as an equation involving the three sides of a box in three-dimensional space. Of course there will be a problem if the relationship between t and x is such that $dx/dt=0$, for then $dx=0$ for non-zero dt . But (provided that dy is not also zero), this will simply correspond to a vertical tangent (dx,dy) in the x - y plane.

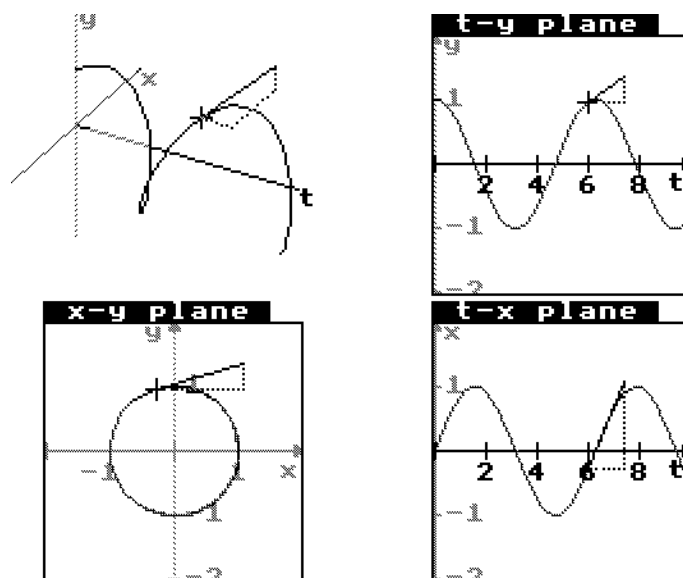


Figure 15: four views of a numerical tangent to a parametric curve

It is not far from here to consider higher dimensional pictures. The geometry gets more difficult to see, but the corresponding linear algebra works just as well and it is time for the symbolism to play a much more dominant role. Once more the notion of local straightness remains centre stage, now in terms of locally linear approximations in differentiable manifolds.

Thus it is that a locally straight approach to the calculus can be begun in a relatively unsophisticated and naive way, using the students natural intuitions, building implicitly, but not explicitly on the notion of limit in such a way that the theory is coherent with both standard and non-standard analysis (and, for those interested in computability, it has a close relationship with constructive analysis as well).

Our experience in the UK is that it works well with a wider cross-section of abilities than a traditional approach. The graphical ideas are readily absorbed by the those with mathematical minds who seek to refine the logic and may pursue a more formal course based on an intuition which can visualize non-differentiability, But it also gives good intuitions and a non-technical approach to those who need to know what a differential equation *is*, so that it can be used in applications. IBM versions of software based on this *Graphical Approach to the Calculus* are also available in the USA (Tall et al 1990).

References

- Cornu B.,1983, *Apprentissage de la notion de limite: Conceptions et Obstacles*, Doctoral thesis, Grenoble, France.
- Leibniz, G. W. 1684: Nova methodus pro maximis et minimis, itemque tangentibus, qua nec fractas, nec irrationales quantitates moratur, & sinulare pro illis calculi genus, *Acta Eruditorum* 467–473.
- Mills, J, & Tall, D. O., (to appear): "Modelling Irrational Numbers in Analysis using Elementary Programming", *The Mathematical Gazette*.
- School Mathematics Project, 1967: *Advanced Mathematics Book 1*, Cambridge University Press: Cambridge, UK.
- Schwarzenberger, R. L. E. & Tall D. O., 1978: "Conflicts in the learning of real numbers and limits", *Mathematics Teaching*, **82**, 44-49.
- Shuard H. & Neill, H. 1982, *Teaching Calculus*, Blackie: London.
- Tall, D. O., 1986: "A graphical approach to integration and the fundamental theorem", *Mathematics Teaching*, **113** 48-51.
- Tall, D. O., 1991a: "Recent developments in the use of the computer to visualize and symbolize calculus concepts", *The Laboratory Approach to Teaching Calculus*, M.A.A. Notes Vol. 20, 15–25.
- Tall, D. O., 1991b: *Real Functions and Graphs*, (software for BBC compatible computers), Cambridge University Press: Cambridge, UK.
- Tall, D. O., Blokland, P. & Kok, D. 1990, *A Graphical Approach to the Calculus*, (software for IBM compatible computers), Sunburst Communications Inc, Pleasantville, NY.