

# DIOPHANTINE ANALYSIS APPLIED TO VIRUS STRUCTURE

WILHELM LJUNGGREN

*Presented by*

NICHOLAS G. WRIGLEY and VIGGO BRUN

## Introduction (by V. Brun).

After discussion with N. G. Wrigley in 1972 on some mathematical problems in virology, he asked me to explore the mathematical background to a diagram he had worked out (see fig. 2). By introducing new letters  $x, y, z$  (see equation 6) it was easy to transform the problem into diophantine form, and as I knew that my friend W. Ljunggren was much more capable of solving this problem, I asked him if he would try. After some weeks, shortly before his death in January 1973, he gave me the manuscript which we present below, with necessary virological explanations before and after Ljunggren's deduction.

## Virological background (by N. G. Wrigley).

Virus particles are invariably enclosed by shells of protein subunits, and these are packed geometrically according to strict symmetry rules, because of the chemical properties of the protein. Helical virus particles are found in nature, and also "spherical" or isometric particles. In 1962 Caspar & Klug [1] stated of the latter that: "... there has accumulated a large body of evidence that icosahedral symmetry is preferred in spherical virus structure. Indeed no well-established examples exist at present of isometric viruses which are not icosahedral." This remains true to-day, though other classes of polyhedra have been considered theoretically, for example by Goldberg [2] and by Brun [3], [4]. All known examples are close-packed with each subunit surrounded by six neighbours, except the twelve vertices which have five neighbours.

The total number of nearly identical subunits which may be regularly packed in this way on the closed icosahedral surface is given [2] by:

$$(1) \quad N = 10(a^2 + ab + b^2) + 2$$

where  $a$  and  $b$  may take any non-negative integral values. This gives the numbers:

$$N = 12, 32, 42, 72, 92, 122, 132, \dots$$

Of course, without upsetting the icosahedral 5, 3, 2 symmetry, the above subunits may be considered as spaces, and the spaces as subunits, and the subunits may themselves be further subdivided into symmetrical groups. However this does not alter the general applicability of equation (1).

Now, an icosahedron has 30 axes of twofold symmetry, 20 of threefold symmetry and 12 of fivefold symmetry. Therefore the subunits on the surface of an icosahedral virus may be thought of as divided into 30 identical groups each having twofold symmetry, 20 groups with threefold and 12 groups with fivefold symmetry. I have discussed this extensively [5] in relation to certain viruses which I found actually collapsed into these groups. Other examples have since been found by myself [6] and by Stoltz [7, 8]. The groups, shown in fig. 1 are called "symmetrons";

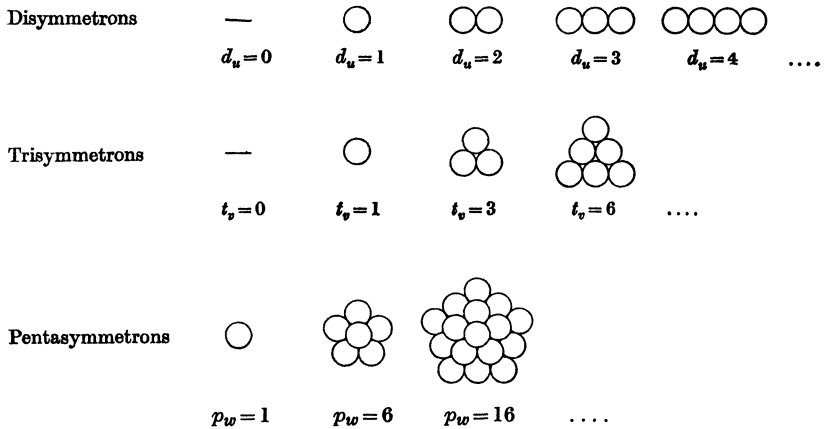


Fig. 1. Diagram showing the construction of linear, triangular and pentagonal symmetrons, with the numbers of subunits they contain. Di- and trisymmetrons may be absent altogether, leaving a "minimum" shell of 12 subunits.

the 30 Disymmetrons contain  $d_u$  subunits, the 20 Trisymmetrons contain  $t_v$  subunits, and the 12 Pentasympmetrons contain  $p_w$  subunits, so that:

$$(2) \quad N = 30d_u + 20t_v + 12p_w$$

where

$$(3) \quad d_u = u - 1$$

$$(4) \quad t_v = v(v-1)/2$$

$$(5) \quad p_w = 1 + 5w(w-1)/2$$

in which  $u, v$  and  $w$  can independently take the values  $1, 2, 3, 4, 5, \dots$

For each allowed value of  $N$  from equation (1) the number  $f(N)$  of solutions of equation (2) were calculated by computer, using (3), (4) and (5). This number  $f(N)$  corresponds to the number of theoretically possible ways of making a virus with  $N$  subunits, but with different combinations of symmetrons. For example,  $f(42)=1$ , that is:

$$42 = 30 \cdot 1 + 20 \cdot 0 + 12 \cdot 1,$$

which is a unique solution of (2). Another example is  $f(72)=3$ , which means that (2) has three solutions, viz.:

$$\begin{aligned} 72 = & 30 \cdot 2 + 20 \cdot 0 + 12 \cdot 1 \\ & \text{or } 30 \cdot 0 + 20 \cdot 3 + 12 \cdot 1 \\ & \text{or } 30 \cdot 0 + 20 \cdot 0 + 12 \cdot 6 \end{aligned}$$

As expected  $f(N)$  increases linearly with  $N$  (fig. 2), but it was surprising to find that (i) this increase is bi-modal, and that (ii) the points are

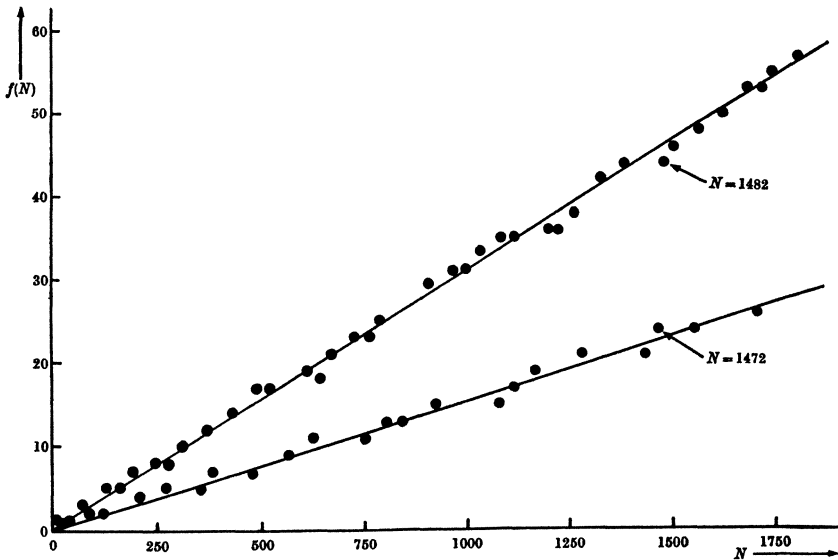


Fig. 2. Diagram showing the bi-modal distribution of the number  $f(N)$  of solutions of equation (2) for each  $N$ . The upper line has twice as many points as the lower. I am grateful to Dr. P. M. Bayley for the computer calculations in this work.

distributed in the proportion 2:1 between the two modes. It was for an explanation of observations (i) and (ii) that I consulted Viggo Brun, and here follows:

### W. Ljunggren's analysis.

Inserting in (2) the values of  $d_u$ ,  $t_v$  and  $p_w$  from (3), (4) and (5), and putting  $x = 2v - 1$ ,  $y = 2w - 1$ ,  $z = u - 1$  and  $N = 10M + 2$ , where  $M = a^2 + ab + b^2$  we get:

$$(6) \quad x^2 + 3y^2 + 12z = 4M.$$

The problem is now to find all odd, positive integers  $x$ ,  $y$  and all non-negative integers  $z$ , satisfying (6) for given values of  $a$  and  $b$ . However, there is no simple formula for the exact number of triples  $(x, y, z)$  in question. We may use the following known result [9]:

Let  $A$ ,  $B$ ,  $n$ ,  $x$ ,  $y$  denote natural numbers. Further, let  $T_n(A, B)$  denote the number of different pairs  $(x, y)$  satisfying the inequality

$$Ax^2 + By^2 \leq n.$$

Then:

$$(7) \quad T_n(A, B) = \frac{\pi n}{4\sqrt{AB}} - \theta \frac{\sqrt{A} + \sqrt{B}}{\sqrt{AB}} \sqrt{n}, \quad 0 < \theta < 1.$$

Equation (6) may be written:

$$(8) \quad \begin{aligned} x^2 + 3y^2 + 12z &= 4(a^2 + ab + b^2) \\ &= (a + 2b)^2 + 3a^2 \\ &= (2a + b)^2 + 3b^2, \end{aligned}$$

or

$$(8') \quad [x^2 - (2a + b)^2] + 3(y^2 - b^2) + 12z = 0.$$

By (8') it follows that  $x^2 - (2a + b)^2$  is divisible by 3, or

$$(9) \quad x = (2a + b)\varepsilon + 3A_1, \quad \varepsilon = \pm 1, A_1 \text{ integer}.$$

Inserting this value for  $x$  in (8') and dividing by 3, we obtain:

$$(10) \quad 2A_1(2a + b)\varepsilon + 3A_1^2 + y^2 - b^2 + 4z = 0,$$

or

$$(10') \quad y^2 - (b - A_1\varepsilon)^2 + 4A_1^2 + 4z + 4A_1a\varepsilon = 0,$$

hence

$$(11) \quad y = b - A_1\varepsilon + 2B_1, \quad B_1 \text{ integer}.$$

Inserting the values for  $x$  and  $y$  from (9) and (11) respectively in (8) we find by simple calculations that:

$$-z = A_1^2 - A_1 B_1 \varepsilon + B_1^2 + A_1 a \varepsilon + B_1 b, \quad z \geq 0.$$

The original problem is now transformed into the following: Determine all odd positive integers  $x, y$ , such that

$$x^2 + 3y^2 \leq 4M, \quad x = (2a + b)\varepsilon + 3A_1.$$

Now we can use formula (7).

We have to distinguish between two cases: firstly  $(b - a)$  divisible by 3, and secondly  $(b - a)$  indivisible by 3. In the first case  $x$  and  $M$  are both divisible by 3, and  $N = 10M + 2$  is of the form  $3C + 2$ ,  $C$  integer. In our counting we must delete the pairs  $(x, y)$  where  $x$  is not divisible by 3. In the second case we have to delete the pairs where  $x$  is divisible by 3. Our first conclusion is therefore: For  $N$  divisible by 3 the number of solutions is approximately double the number of solutions for the case  $N = 3C + 2$ . The case  $N = 3C + 1$  gives no solutions.

We now confine ourselves to discussion of the second case, since the first case can be dealt with similarly. We count the pairs  $(x, y)$ ,  $x, y$  both odd numbers and  $x$  indivisible by 3. For the number of solutions we get approximately:

$$\frac{\pi 4M}{4\sqrt{3}} \left(1 - \frac{1}{2} - \frac{1}{2} + \frac{1}{4} - \frac{1}{3} + \frac{1}{6} + \frac{1}{6} - \frac{1}{12}\right) = \frac{\pi M\sqrt{3}}{18}.$$

Here we have used the first term in (7). From the total number we must subtract the number of pairs where  $x$  is even and  $y$  is even. Then we have added the number of pairs where both  $x$  and  $y$  are even numbers  $(1 - \frac{1}{2} - \frac{1}{2} + \frac{1}{4})$ . Now we have subtracted the number of pairs  $(x, y)$ , where  $x$  is divisible by 3  $(-\frac{1}{3})$ . Then we must add the numbers of pairs  $(x, y)$ , where  $x$  is even and divisible by 3, and the number of pairs  $(x, y)$  where  $y$  is even and  $x$  divisible by 3  $(\frac{1}{3} + \frac{1}{6})$ . Finally we have to subtract the number of pairs where both  $x$  and  $y$  are even,  $x$  divisible by 3  $(\frac{1}{12})$ .

Introducing  $N$  for  $M$  we find the dominant term in our counting is  $\pi\sqrt{3}N/180$ . Denoting by  $f(N)$  the total number of solutions of (6), we have:

$$f(N) = \frac{\pi\sqrt{3}N}{180} + k_1\sqrt{N}.$$

Here  $k_1$  is bounded, independently of  $N$ . Furthermore

$$\lim_{N \rightarrow \infty} f(N)/N = \pi\sqrt{3}/180 = 0.03.$$

The curve  $y = kx + k_1\sqrt{x}$ ,  $k = 0.03$ , is a branch of the parabola

$$(y - kx)^2 - k_1^2 x = 0 .$$

The axis is parallel to the line  $y = kx$ , its vertex is situated in a bounded domain around the origin, and the parameter is also bounded.

Thus the conclusion is: *The points  $(x, y)$  all lie in the neighbourhood of the two lines  $y = 0.03x$  and  $y = 0.015x$  (in the first case  $k = \frac{1}{2} \cdot 0.03 = 0.015$ ).* (See fig. 2.)

Ljunggren concluded his manuscript with the words: "I have tried to give an 'elementary' explanation. It is possible to give better bounds for the difference

$$T_n(A, B) - \pi n / 4\sqrt{AB} ."$$

#### **Virological implications** (by N. G. Wrigley).

The relevance of this beautiful explanation of the numerical facts to the problem of virus structure is this: The protein subunits which form penta- tri- and disymmetrons have to assemble themselves inside a living cell into a closed shell of a particular volume determined by the contents of the virus. The statistical chance of errors in this assembly process is significant, particularly if  $N$  is large. A virus might have a choice, for example, of making an  $N = 1472$  shell or an  $N = 1482$  shell. These alternatives are nearly identical in volume, *but there are only half as many ways to make the  $N = 1472$  shell* (fig. 2), and therefore double the statistical chance of success in the assembly process. This would be a substantial reason for Evolution to choose this alternative.

The large viruses studied by Wrigley [5, 6] and by Stoltz [7, 8] very probably fall on the lower curve of figure 2. However, too little is known about the assembly process of viruses to say whether this result shows the hand of Evolution at work. Further evidence is required from the discovery of other large viruses. The only other viruses whose structure is known at present are much smaller ( $N < 300$ ), and they fall on both lines of figure 2. The virological relevance of this mathematical study is, therefore, extremely speculative. Moreover we have only considered the class of symmetrons which are pentagonal, triangular or linear in outline. It is quite possible to have symmetrons which retain five-, three- or twofold symmetry but have different shapes from those we have considered, and one example (adenovirus) of non-triangular trisymmetrons is known.

## REFERENCES

1. D. L. D. Caspar & A. Klug, *Physical principles in the construction of regular viruses*, Cold Spr. Harb. Symp. quant. Biol. 27 (1962), 1–24.
2. M. Goldberg, *A class of multi-symmetric polyhedra*, Tôhoku Math. J. 43 (1937), 104–108.
3. V. Brun, *Some theorems on the partitioning of the sphere, inspired by virus research*, Nordisk Mat. Tidskr. 20 (1972), 87–91. (Norwegian with English summary.)
4. V. Brun, *On some problems in solid geometry within virology*, (in preparation).
5. N. G. Wrigley, *An electron microscope study of the structure of Sericesthis iridescent virus*, J. gen. Virol. (1969), 123–134.
6. N. G. Wrigley, *An electron microscope study of the structure of Tipula iridescent virus*, J. gen. Virol. 6 (1969), 169–173.
7. D. B. Stoltz, *The structure of icosahedral cytoplasmic Deoxyriboviruses*, J. Ultrastruct. Res. 37 (1971), 219–239.
8. D. B. Stoltz, *The structure of icosahedral cytoplasmic Deoxyriboviruses II: An alternative model*, J. Ultrastruct. Res. 43 (1973), 58–74.
9. P. Bachmann, *Die analytische Zahlentheorie*, Zweiter Teil, p. 447, B. G. Teubner, Leipzig, 1894.

NATIONAL INSTITUTE FOR  
MEDICAL RESEARCH  
MILL HILL, LONDON NW 7 1AA  
ENGLAND

INSTITUTE OF MATHEMATICS  
OSLO UNIVERSITY  
BLINDERN, OSLO 3  
NORWAY