# Some Observations on Packability of Spheres

Jason Burns

August 24, 2003

## I  Introduction

The question of how many spheres of a given size can be placed tangent to a given central sphere has been around at least since 1694, when Isaac Newton and David Gregory argued to each other about whether it was possible to fit thirteen spheres of the same size so that all touched a center sphere of that same size [1]. (It is easy to put twelve around the center sphere, at the vertices of an icosahedron, and the fact that they fit quite loosely suggests that it might be possible to fit thirteen; but in fact, as was proven only in 1874, one cannot.) The corresponding question for higher dimensions—how many hyperspheres, all of the same (given) size, can fit around a central hypersphere in $N$ dimensions?—is very much unsolved, even if we restrict ourselves to the most interesting special case, where the central sphere is the same size as the other spheres; and even in three dimensions the answer is unknown for all but a handful of cases. In this paper I extend a result from [1], which shows that for three dimensions the answer to the "how many?" question is never 5, to $N$ dimensions, and then discuss some of what is known about the general problem. I close with a different direction in which to generalize this idea.

## II  A promising beginning: the first two weeks

The original problem sprang from a note in the *Monthly* [1] which proved that, if one can pack five spheres of radius $r$ so that all touch a central sphere of fixed radius, then one can pack six of that radius. What won me over about it was that the proof was both elementary and very simple, essentially relying

only on the Pigeonhole Principle; I was not expecting the proof to be nearly as short.

The first question that arose was the question of whether this phenomenon occurred in $N$ dimensions, or only (for instance) for odd $N$. The second question was whether there were similar implications for other numbers of spheres (11, for example). The third question was whether any of this still held true when either the central sphere or the surrounding spheres were replaced by ellipsoids, slightly elongated in one direction.

The first question answered itself almost immediately (see Theorem 1 below), and with some relatively easy companion theorems I worked on over the second week, gives a complete description of how many $(N-1)$-spheres of a given radius may be packed around a central $(N-1)$-sphere in $N$ dimensions for all radii greater than or equal to $1 + \sqrt{2}$.

**Theorem 1.** *If, in $N$-dimensional space, $N+2$ nonoverlapping $(N-1)$-spheres of a given radius all touch a central sphere of some other radius, then $2N$ nonoverlapping $(N-1)$-spheres of the given radius can be placed all tangent to the central $(N-1)$-sphere.*

**Theorem 2.** *There is a radius such that $2N$ $(N-1)$-spheres and no more can fit around the center $(N-1)$-sphere.*

**Theorem 3.** *There are radii such that each of 2, 3, ..., $N+1$ $(N-1)$-spheres and no more can fit around the center $(N-1)$-sphere.*

/smallskip There is quite a bit of material in the statement of the theorems above which is really incidental to the proofs. For one thing, we may take the center sphere's radius to be 1, with no loss of generality. For another, we really don't need to know the radii of the surrounding spheres; all we need is their points of tangency. See, if we're given the points at which the surrounding spheres hit the center sphere, we can just increase the radii of the spheres from the center $O$ until the first radius at which two spheres just touch. That radius $r$, in turn, can be expressed in terms of the angle $\theta$ between the points of tangency of each sphere with the center sphere.[1] It is clear that, the bigger $\theta$ is, the bigger $r$ can be. So we can then restate our

---

[1]How so? Well, consider a triangle with vertices $OC_1C_2$, the centers of the center sphere and the two touching spheres, respectively. Since all three spheres are tangent, we get $OC_1 = 1+r = OC_2$, and $C_1C_2 = 2r$. The angle $C_1OC_2$ is $\theta$. So, bisecting our isosceles triangle, we find that $\sin(\theta/2) = \frac{r}{1+r} = 1 - \frac{1}{1+r}$, or $r = \frac{1}{1-sin(\theta/2)} - 1$. We don't really need a formula for $r$, though; all we need to know is that, as the minimum $\theta$ increases, so can $r$. (But formulas are nice.)

theorems as we will prove them, in terms of the angles between points on a sphere.

**Theorem 1′.** *Given any $N + 2$ points on an $(N - 1)$-sphere in $N$-dimensional space, some pair is separated by an angle of $\leq \frac{\pi}{2}$.*

**Theorem 2′.** *We can fit $2N$ points on an $(N-1)$-sphere with every pair separated by an angle $\geq \frac{\pi}{2}$, but of any $2N + 1$ points, some pair is less than $\frac{\pi}{2}$ apart.*

**Theorem 3′.** *The best possible (largest minimum angle between points) arrangement of $k$ points $(2 \leq k \leq N + 1)$ has $\cos\theta = -\frac{1}{k-1}$, where $\theta$ is the minimum angle.*

*Proof of Theorem 1′.* Our strategy is to choose our coordinate axes well, and then count the points per "octant". The total turns out to be more than the number of octants, so we apply the Pigeonhole Principle (which just states, that if there are more pigeons than pigeonholes, some pair of pigeons are in the same hole) to conclude that some octant has two points in it. The result follows.

Let us take a set of $N + 2$ points, which we'll call $\vec{P}_1$, $\vec{P}_2$, ..., $\vec{P}_{N+2}$, on an $(N - 1)$-sphere for some given $N$. As noted before, we may assume for convenience that the sphere has radius 1. We may also choose the center of the coordinate system to be the center of the circle, which we'll call $\vec{O}$. That leaves $N$ orthogonal axes to choose; call them $x_1$, $x_2$, ..., $x_N$. Here's how we'll do it:

We pick $x_1$ so that $\vec{P}_1$ has first coordinate as large as possible (that is, 1). In other words, we choose $\vec{OP}_1$ to be our $x_1$-axis, with $\vec{P}_1$ on the positive side. So the $x_2$- to $x_N$-coordinates of $\vec{P}_1$ are zero, and the $x_1$-coordinates of all the points are determined.

We pick $x_2$ so that $\vec{P}_2$ has second coordinate as large as possible (that is, $\sqrt{1 - x_1^2}$). In other words, we choose $\vec{OP}_1\vec{P}_2$ to be our $x_1$-$x_2$-plane, with $x_2$ perpendicular to $x_1$ in that plane, with $\vec{P}_2$ on the positive side. So the $x_3$ to $x_N$-coordinates of $\vec{P}_2$ are zero, and the $x_2$-coordinates of all the points are determined.

In general, we pick $x_k$, $1 \leq k \leq N$, to be the axis perpendicular to all of $x_1$, $x_2$, ..., $x_{k-1}$ in the hyperplane determined by the $x_1$-, $x_2$-, ..., $x_{k-1}$-axes together with point $\vec{P}_k$, with orientation such that $P_k$ has nonnegative $x_k$-coordinate. (The orientation, to be truthful, doesn't really matter, but it's nice to have something specific to refer to.) This means that all coordinates

of $\vec{P}_k$ from $x_{k+1}$ on down to $x_N$ are zero.

The only way this can fail at any step is if the $(k-1)$ different axes and $\vec{P}_k$ don't determine a hyperplane of dimension k. (Or, said differently, if we cannot pick a positive $x_k$-coordinate.) Since the axes are orthogonal, the culprit is $P_k$, which is in the $(k-1)$-hyperplane with the other axes. But this is no problem, since we can then choose *any* $x_k$-axis and orientation which is perpendicular to the other axes. All the coordinates of $\vec{P}_k$ from $k+1$ on will still be zero, which is the only fact we really need (as I mentioned before, the orientation is just for convenience).

There are $2^N$ "octants" of the $(N-1)$-sphere: each coordinate can be either $\geq 0$ or $\leq 0$. For instance, for $N = 3$, $x_1 \leq 0$, $x_2 \geq 0$, $x_3 \leq 0$ is an octant of the sphere in 3-space. But if we take a point with, for example, $x_1 = 0$, $x_2 > 0$, $x_3 < 0$, it is not only in that octant but also in the octant $x_1 \geq 0$, $x_2 \geq 0$, $x_3 \leq 0$; it is on the line separating the two octants. By choosing our axes wisely, we have landed as many points as possible onto as many hyperplanes separating the "octants" of the $(N-1)$-sphere as we possibly could; this is what makes everything work.

If we count the points in each "octant" and add up the total, we count point $\vec{P}_k$ at least $2^{n-k}$ times, since its coordinates $x_{k+1}$ through $x_N$ are guaranteed to be zero. That leaves $\vec{P}_{N+1}$ and $\vec{P}_{N+2}$, each of which must get counted at least once. The total is

$$2^{N-1} + 2^{N-2} + \cdots + 2^1 + 2^0 + 1 + 1$$

or (summing the series)
$$2^N + 1.$$

We conclude that, since there are only $2^N$ total "octants", some one of them has two points in it.

It now remains only to show that any two points in an "octant" are within $\frac{\pi}{2}$ of each other. We recall that the angle between two vectors $\vec{P}$, $\vec{Q}$ is given by $|\vec{P}||\vec{Q}|\cos(\theta) = \vec{P} \cdot \vec{Q}$, or (since $|\vec{P}| = |\vec{Q}| = 1$) $\cos(\theta) = \vec{P} \cdot \vec{Q}$. The dot product just multiplies the corresponding components of $\vec{P}$ and $\vec{Q}$; since corresponding components have the same sign, the dot product—and hence $\cos(\theta)$—is nonnegative; and finally, if $\cos(\theta) \geq 0$, then $\theta \leq \frac{\pi}{2}$.

*Proof of Theorem 2′.* First of all, we can put $2N$ points on an $(N-1)$-sphere easily: just put one at each end of $N$ orthogonal axes. (More explicitly, if $(x_1, x_2, \ldots, x_n)$ is an orthonormal coordinate system centered at the sphere's

center, then we let $\vec{P}_i = (x_1 = 0, x_2 = 0, \ldots, x_i = 1, \ldots, x_N = 0)$ and $\vec{P}_{N+i} = (x_1 = 0, x_2 = 0, \ldots, x_i = -1, \ldots, x_n = 0)$ for all $i$ from 1 to $N$. Clearly the dot product between any pair of these is either 0 or $-1$, and the angle between is either $\frac{\pi}{2}$ or $\pi$, as desired.)

To prove that we cannot do the same for $2N+1$ points, we use induction. An easy case to start with is $N = 1$: a 0-sphere consists of just two points, being defined by $x^2 = 1$ or $x = \pm 1$, so of any $2 \cdot 1 + 1 = 3$ points, some two are identical and hence at an angle of $0 < \frac{\pi}{2}$ apart.

So, consider a counterexample of smallest dimension N. Take a point and assign it coordinates $(x_1 = 0, \ldots, x_{N-1} = 0, x_N = -1)$. Then, since there is no point closer than $\frac{\pi}{2}$ to this point, the $2N$ remaining points all have $x_N \geq 0$. Our plan is to project all the points downward into the $x_1 x_2 \cdots x_{N-1}$-plane, and from there outward into the $(N-2)$-sphere $x_1^2 + x_2^2 + \ldots + x_{N-1}^2 = 1$; unfortunately, the point given by $x_N = 1$ maps to the origin, and so we specifically exclude it. This is no problem, as we still have $2N - 1$ points left to project, which is just enough for the induction to work. We need only prove that, if the angle between two points was $\geq \frac{\pi}{2}$ before the projections, it still will be after them. (Hence every counterexample in $N$ dimensions leads to one in $N - 1$ dimensions, contradicting our inductive hypothesis.)

Consider two points $\vec{P} = (p_1, p_2, \ldots, p_N)$ and $\vec{Q} = (q_1, q_2, \ldots, q_N)$. We have that $\vec{P} \cdot \vec{Q} = |\vec{P}||\vec{Q}| \cos(\theta) \leq 0$ (since $\theta \geq \frac{\pi}{2}$), so $p_1 q_1 + p_2 q_2 + \cdots + p_N q_N \leq 0$. We project downward to get $\vec{P'} = (p_1, p_2, \ldots, p_{N-1}, 0)$ and $\vec{Q'} = (q_1, q_2, \ldots, q_{N-1}, 0)$; the dot product of $\vec{P'}$ and $\vec{Q'}$ is $|\vec{P'}||\vec{Q'}| \cos(\theta') = p_1 q_1 + p_2 q_2 + \ldots + p_{N-1} q_{N-1}$; this is less than or equal to $\vec{P} \cdot \vec{Q}$, which we know is $\leq 0$. We conclude that $\theta' \geq \frac{\pi}{2}$ and, since the outward projection preserves such angles, the image of all $2N - 1$ points after both projections is a configuration on the $(N-2)$-sphere with all angles $\geq \frac{\pi}{2}$, a contradiction.

*Remark.* I originally tried to prove that the mapping used here (squash one dimension and project out) does not decrease angles, which is very false. (For instance, take two points on the same longitude of a sphere, one on the equator, one just short of the pole; the mapping sends both to the same place, though the angle is virtually $\frac{\pi}{2}$.) Had I been thinking more clearly, I would never have attempted to prove that—and never would have found even this much.

Also note that the argument for the second part of Theorem 2 shows that, for $2N$ points, each point must come with an antipodal point, and all other points are equidistant from the two, which repeated $N$ times shows that our

configuration for the first part is essentially unique.

*Proof of Theorem 3′.* Let the center of the $(N-1)$-sphere be our center of coordinates, as usual, and let $\vec{P}_1$ through $\vec{P}_k$ be our points for some $k \leq N+1$. Note again that, if the angle $\theta$ between $\vec{P}$ and $\vec{Q}$ is greater than (or equal to) $x$, $\vec{P} \cdot \vec{Q}$ is less than (or equal to) $\cos x$. We have

$$0 \leq (\vec{P}_1 + \vec{P}_2 + \cdots + \vec{P}_k)^2 = \vec{P}_1^2 + \vec{P}_2^2 + \ldots + \vec{P}_k^2 + 2 \sum_{1 \leq i < j \leq k} \vec{P}_i \cdot \vec{P}_j.$$

Let $\theta$ be our minimum angle; then the above is less than or equal to $k + \binom{k}{2} \cos(\theta)$, by what we have said before. So $0 \leq k + k(k-1)\cos(\theta)$, or $-\frac{1}{k-1} = \cos(\theta)$. This is achievable; just inscribe a regular $(k-1)$-simplex in the $(N-1)$-sphere, and equality is achieved in both $\leq$ signs.[2]

# III  A series of dead ends: the next three weeks

Near the end of the second week, I had shifted my attention to the second problem, the search for other gaps; I was making very little progress, as the tools I had been using were generally useless for smaller $\theta$. (The proof of theorem 2 is a typical example.) I began looking for useful references in the library, and found *Sphere packings, lattices and groups* by J. H. Conway and N. J. A. Sloane [**2**], which told me more than I wanted to know in just four pages (25–28).

The question, according to Conway and Sloane, is that of "generalized kissing number" $A(N, \theta)$, where $A(N, \theta)$ is the maximum number of $(N-1)$-spheres that can be packed around a central $(N-1)$-sphere in $N$ dimensions so that no two spheres are less than $\theta$ apart on the center sphere. (As we have seen, this translates readily into finding $A$ in terms of the radius, or into the more convenient problem of placing points on a sphere so that all points

---

[2]We can construct such a simplex inductively. Pick an axis $x_N$, choose as one point $\vec{P}$ of our simplex the point $x_N = 0$, and as the other points the points of any regular $(N-1)$-simplex in the $(N-2)$-sphere with $x_N = -\frac{1}{N}$. The cosine of the angle between $\vec{P}$ and any other point is clearly $-\frac{1}{N}$, and the cosine of the angle between two points in the $(N-1)$-simplex is the angle between the points in the $(n-2)$-sphere divided by the radius of the $(N-2)$-sphere, plus $\left(-\frac{1}{N}\right)^2$, or $-\frac{1}{N-1}\left(1 - \frac{1}{N^2}\right) + \frac{1}{N^2}$, which simplifies to $-\frac{1}{N}$. So all the angles are $\cos(\frac{-1}{k-1})$, as desired.

are at least $\theta$ apart.) It is a generalization (hence the name) of the "kissing number", which deals only with the case of the packing spheres having the same radius as the center sphere, that is, with $\theta = \frac{\pi}{3}$. (This is by far the most important case, because it relates to the problem of filling $N$-space with identical spheres, which has been extensively studied.) Since the answer to the kissing-number problem is known only for dimensions 1, 2, 3, 8, and 24, according to page 23, it is not surprising that I had some difficulty with any sort of general approach to question 2, which is roughly 'Find some more values for the generalized kissing number.'

*Known results about $A(N, \theta)$.* For $N = 1$ the problem is trivial; for $N = 2$ we need only put our points equally spaced about the circumference of the circle, so for $A(N, \theta) = k$ the best value of $\theta$ we can get is $2\pi/k$. Once $N$ hits 3, the problem becomes very hard; the only values known for general $N$ are the ones described above in section III. (These values, as I learned from Conway and Sloane for the first time, had already been found by R. A. Rankin [**3**] in 1955; I was unable to locate this article, however. (There is a book [**4**] which apparently contains Rankin's proof, and seems like a good one on the subject; alas, it was checked out.) For three dimensions, some good configurations are known for most numbers of points less than 100, but for most of them, there is no proof that they are in fact the best. There are proofs for up to 12 points, and for 24, but no others I know of; more on this later. A few results are also known for four dimensions, and there are the kissing numbers for 8 and 24 dimensions, as noted above. (For the curious, they are 240 and $196, 560$, respectively.) Otherwise, all that are known are bounds.

To compare the different bounds known, it will be convenient to take $\theta = \pi/3$, that is, to compare the results the bounds give for the (non-generalized) kissing number in $N$ dimensions. The current upper bound is due to Kabatiansky and Levenshtein and is on the order of $2^{0.401N}$, and the current lower bound is on the order of $2^{0.2075N}$, due to A. D. Wyner. The catch is that Wyner's result is nonconstructive. Conway and Sloane outline a simple and apparently new way to get a constructive lower bound: inscribe a hypercube in our $(N - 1)$-sphere of interest, and use an error-correcting code to select the vertices to actually use in the point set. One cannot just use all the vertices, since the angle between two adjacent vertices goes to zero as $N \to \infty$; so we choose only the vertices corresponding to codewords, thus guaranteeing that the vertices differ in at least $d$ coordinates. (The

correspondence between the hypercube vertices and codewords is given as follows: Take a binary codeword of length $N$; it's just a list of $N$ zeroes and ones. Change all the zeros to $+1$ and all the ones to $-1$, and we have the coordinates of a vertex of a hypercube of side 2. If we divide all the coordinates by $\sqrt{N}$, our hypercube is now inscribed in a sphere of radius 1, and we can calculate the minimum angle $\theta$ in terms of the minimum distance $d$ (number of places in which two codewords differ); the answer turns out to be $\theta = \cos^{-1}(1 - 2d/N)$.) What makes this a superior approach is that there are codes for which the number of codewords grows exponentially with $N$ *and* the minimum distance remains at least a constant fraction of $N$ (so $\theta$ doesn't decrease.) So we can, in fact, get an exponentially growing constructive lower bound, on the order of $2^{0.003N}$. But as one can see, there is quite a gap between this result and Wyner's. Since the construction seemed so simple (apart from the fact that codes with the required properties are still rare and hard to construct, though well-known in coding theory), a natural thing to do seemed to be to try to improve the construction to obtain a better bound. Enter two-and-a-half weeks of dead ends. I won't give a detailed account of them, since there isn't much to say. But first, let me give some idea of how these special codes are made, aided by pages 82–83 of the ever-helpful [**2**].

*The Justesen codes.* First, we recall some basic facts about finite fields, which are the basic tool in developing most codes in use today. We know that the integers (mod p), where $p$ is a prime, are a field, so we can add, subtract, multiply, and divide integers (mod p), and all the usual rules hold. (For example, consider $Z_5$, the integers (mod 5). Addition, subtraction, and multiplication work just as in the integers, except that when we're done we divide by 5 and take the remainder. Division is somewhat less obvious: but if $2 \cdot 3 \equiv 1 \pmod 5$, then the multiplicative inverse of 2 must be 3, or in other words $2^{-1} \equiv 3 \pmod 5$. Similarly, $\frac{1}{3}$ is 2, and $\frac{1}{4}$ is 4.) What you may not know is that we can make larger fields of any prime power order by extending $Z_p$. Specifically, to make a field with $p^n$ elements, we take an $n$th degree polynomial $P(x)$ in $Z_p$ which doesn't factor, and consider the set of all polynomials in $Z_p$ modulo this polynomial—that is, the set of all remainders we get after dividing by $P(x)$. Just as always, in other words, "modulo" means to divide and take the remainder. (Formally, our field is $Z_p[x]$ mod $P(x)$.) Let's take an example: if we want to make a field of 25 elements, we would take a quadratic polynomial irreducible over $Z_5$. As

you can check for yourself, $x^2 + x + 1$ is such a polynomial. So from now on two polynomials in $x$ are "the same" if they differ only by a multiple of $x^2 + x + 1$; and so every time we see an $x^2$, we can replace it by $-x - 1$, since $x^2 + x + 1 \equiv 0$. (And every time we see an $x^3$, we can replace it by $-x^2 - x$, which then becomes $x + 1 - x$, or just 1.) Hence all the $x^2$ (and higher) terms disappear, and we do indeed get 25 elements as promised. Addition, subtraction, and multiplication work just as they do for polynomials; it isn't obvious that division always works now, but it does. One more fact about finite fields: any two of the same size are isomorphic. The reason for this is that the multiplicative group of the field is always cyclic (of order $p^n - 1$, of course), so by picking the right element to generate the whole group (called a "primitive element") we can establish an isomorphism.

The simplest to construct of these so-called "asymptotically good codes" are the Justesen codes. Take a field with $2^m$ elements, and let $\alpha$ be a primitive element. Then take a linear code over this field, $2^m - 1$ elements long. Our new codewords will be twice as long as the old, because for every element $c_k$ of the old codeword we will put two elements, $c_k$ and $\alpha^k \cdot c_k$. So if our old codeword was $(c_0, c_1, \ldots, c_{2^m-2})$, our new codeword will be $(c_0, c_0, c_1, \alpha \cdot c_1, \ldots, c_{2^m-2}, \alpha^{2^m-2} \cdot c_{2^m-2})$. So we multiply by all the $2^m - 1$ different powers of $\alpha$, once per codeword element. Now, since each codeword element is really an element of our field of order $2^m$, each element really represents a sequence of $m$ bits; so we took a set of $m(2^m-1)$-bit codewords and turned them into a set of $2m(2^m-1)$-bit codewords, one apiece. Moreover, because we multiplied each element of the code by a different nonzero element of the group, it turns out that a significant fraction of the bits in the new codeword will be nonzero, unless the old codeword was zero to begin with. (In fact, about ten percent, if we don't use too many codewords.) Because the code is linear, which means the sums and differences of codewords are also codewords, if any two codewords were closer than this, their difference would also be a codeword closer than this to the zero codeword; so we are assured that every pair of codewords differ in about ten percent of the bits.

What does this mean to us? Well, the actual formula says that, for large code lengths $N$, $\frac{d}{N}$ (the proportion of bits in the new codewords that differ) will be at least $0.110(1 - R)$ where $R$ is the rate of the old code we used. The rate of a code is just the proportion of bits that carry information; so the total number of bits carrying information in the old code was $m(2^m - 1) \cdot R$. Now if $k$ bits actually carry information, then we can distinguish $2^k$ different outcomes, so we have $2^k$ different codewords. Let's work out the math: if we

start with a code with rate 0.01, we get a total of $m(2^m - 1) \cdot 0.01$ bits—wait! That factor of $m(2^m-1)$ is just the length of the old code, or half the length of the new code, or $N/2$. So we have a total of $0.005N$ bits carrying information, which means that we have $2^{.005N}$ codewords. So our list of codewords grows exponentially. And $d/N$ is $(0.110)(0.99)$ or about $0.109$. This means that, if we have $2^{.005N}$ points to place on a sphere for a given dimension $N$, we can place them at the vertices of a hypercube that correspond to the codewords of one of the codes we just made up, with a rate of 0.99. And if we do that, then we are guaranteed that their coordinates will be different in at least $0.109N$ places, or (in other words) that one list of coordinates has $-1/\sqrt{(N)}$'s and the other $1/\sqrt{(N)}$'s or vice versa at least 10.9 percent of the time. We can compute the minimum angle for this configuration rather easily now: as I mentioned before, it comes out to $\theta = \cos^{-1}(1 - 2d/N)$, which in this case is about 38.5 degrees (or .67 radians). So $A(N, 38.5°) \geq 2^{0.005N}$, and we have a constructive, exponential lower bound. And that's how the Justesen codes work for us. (There are other codes that generalize this method, which allow us to get a larger maximum $d/N$ (hence a larger $\theta$); Conway and Sloane use one of these rather than the one above for the case $\theta = \pi/3$.)

*Ideas that didn't work.* At first, I thought I could improve the construction pretty easily. After all, selecting vertices from a hypercube is clearly poor compared to selecting them over the whole sphere. In fact, the hypercube isn't even the best configuration—in three dimensions, we can get a larger distance between the top points and the bottom points by twisting the top half of the cube 45° to get a "skew cube". (Which means that we get a larger distance between all points, because the next thing we do is to move the top and bottom planes closer together. This trades top-to-bottom distance for top-to-top and bottom-to-bottom distance, making all points more evenly spaced.) But, as I soon discovered, the problem with trying to use something like that in the above result, instead of a hypercube, is that the crucial part of the result is not the geometry, but the code. By twisting, we almost certainly move some pair of codewords closer, so that the minimum distance is $(d-1)$ edges now, not $d$; and the more twists we put in, the more the distance can drop. Even though the individual edges lengthen, the net distance drops. We might still be able to make this work, but only by designing a (necessarily nonlinear, hence difficult) code specifically for the skew hypercube.

Trying to use geometry without the code is useless; I don't know of any good geometric ways of putting points evenly on a sphere in $N$ dimensions so

that the number of points grows exponentially as N increases (unlike various configurations of points on the equators of different axes, which grow about as $N^2$), but slower than the known upper bounds. For example, the hypercube doubles its points every time we add a dimension, and most configurations do worse; so we end up with $2^N$ points. This is more than the upper bound for a given $\theta$, so the angles must go to zero as the number of dimensions increases. (I have in mind here the upper bound mentioned above for $\theta = \pi/3$, which is asymptotic to $2^{0.4N}$.) Without the code to filter out the gross excess of points, we have a problem. And even if we do find a construction that gives us the points, we need to keep them reasonably well distributed; if they're all on one or a few circles or spheres, we still won't get the angle we need.

I managed to convince myself that substituting a ternary code for a binary code would give worse bounds. With a ternary code, we can use what looks like a hypercube, but with points not only at the corners, but at the midpoints of the edges, faces, etc., for a total of $3^N - 1$ points rather than $2^N$. (The center of the hypercube is useless for these purposes, since when we project outward to the hypersphere's surface, we can't project the center anywhere useful; hence the $-1$.) On the other hand, codes that differ in $d$ places used to have $d$ edges between them of necessity; but now one point could be on the midpoint of an edge and one at one end of an edge (corresponding to a 1 and a 0 or 2 in that place); the distance (as measured by the number of different digits) would only be half as large (as measured by the number of edges), since one digit's difference now equates to either an edge or half an edge. If we do the calculations, using some bounds for the rate of the code which I looked up and don't remember, we find that the best case for a ternary code is worse than the corresponding binary best case.

Finally, on page 116 of [5], the author proves by a counting argument that, if we take a code on some finite field, we can pick an element $\alpha$ so that, just as in the Justesen codes, the new codeword $=$ (old codeword, $\alpha \cdot$ old codeword) has a certain proportion of nonzero digits, but better than the Justesen codes. Unfortunately, no method for choosing such an $\alpha$ is known. I played with this, using the sequence of fields $F_2[x]$ mod $x^{2 \cdot 3^m} + x^{3^m} + 1$ for $m \geq 0$, and wasted a good three days. (The only interesting thing I remember finding out about it was that $(x+1)$ is a generator for those fields, not that I doubt it was known.)

*Remark on section III.* As I promised earlier, here is what is known and proven about the best way to lay points on a sphere; it seems like a fitting

conclusion to this section. [6], as its title suggests, contains proofs for seven, eight, and nine points. The eight-point configuration is a "skew cube", as mentioned before, which has $\theta \approx 74.86°$; the seven- and nine-point configurations are harder to describe, but have $\theta$ equal to $80°$ and $\cos^{-1}(1/4)$. [7] contains the 24-point configuration, which takes the form of a "snub cube", which is a cube with the corners sliced off to yield triangular faces at the old corners and make the old faces octagonal. Finally, [8] gives and proves the best arrangements of up to 12 points, so it contains not only the relevant results in [6] and the case of 12 points (the vertices of an icosahedron—this was well-known well before any of these), but the optimal configurations for ten and eleven points as well. The ten-point configuration is hard to describe, but the eleven-point configuration is merely an icosahedron with one point removed. Which means that, after all, the second question has been answered, though in a limited sense. There *are* other gaps, at least in three dimensions: if 11 spheres can be packed around a center sphere, then so can 12.

# IV    The last three weeks

I hesitated to move on to the third problem for a long time, because the problem seemed even more intractable than the second. The difficulty is that, unlike for spheres, the problem for ellipsoids of finding the distance between two points cannot be replaced by that of finding the central angle. (Since measuring the angle two points make from the center essentially ignores distance, we can expect that such a trick would only work when all distances from the center are the same—that is, on a sphere.) This means that, in some fashion or another, we have to calculate the equation of a shortest path between two points of the ellipsoid, which is something I was unable to calculate successfully.

There are well-known techniques for calculating "geodesics", as they are called. One is to use the Euler-Lagrange equations $\frac{\partial f}{\partial y} = \frac{d}{dx_i} \frac{\partial f}{\partial y'}$, where $f$ is the function we want to minimize (here, the arc length) and the $y$-coordinate is a function of the $x_i$-coordinates (on the top half of the ellipsoid, for example). But I had so many problems even setting up the equations, to say nothing of solving them, that I finally gave up with this approach.

I also tried working with the differential equation for a geodesic in general coordinates, which [9] gives as $\ddot{u}^\gamma + \Gamma^\gamma_{\alpha\beta}\dot{u}^\alpha\dot{u}^\beta = 0$, but could not get so far

as calculating the Christoffel symbols without getting contradictory answers. (I should explain the symbolism in that mysterious-looking equation. First of all, it is standard practice in tensor notation to omit the "For every $\gamma$ from 1 to $n$" explanations, and to assume the $\Sigma$ before repeated indices in a product. Second of all, $u^\gamma$ represents $u^1$ and $u^2$, the coordinates on the surface; since a geodesic is just a curve, which we can express parametrically by writing both $u^1$ and $u^2$ in terms of $t$ for some $t$, and so $\dot{u}^1$ is really $\frac{du^1}{dt}$ and $\ddot{u}^2$ is $\frac{d^2(u^2)}{dt^2}$. So our equation is really a set of differential equations for $u^1(t)$ and $u^2(t)$, namely the two equations (for $\gamma = 1$ and $\gamma = 2$ respectively) $\frac{d^2(u^1)}{dt^2} + \sum_{1 \leq \alpha, \beta \leq 2} \Gamma^1_{\alpha\beta} \frac{du^\alpha}{dt} \frac{du^\beta}{dt} = 0$ and $\frac{d^2(u^2)}{dt^2} + \sum_{1 \leq \alpha, \beta \leq 2} \Gamma^2_{\alpha\beta} \frac{du^\alpha}{dt} \frac{du^\beta}{dt} = 0$. I won't go into the meaning of the Christoffel symbols ($\Gamma^\gamma_{\alpha\beta}$) here, except to say that they represent in some sense the curvature embedded in our choice of coordinates.)

Finally, I tried expressing the cross-section of the ellipsoid through two points as an ellipse, so that I could express the length along that ellipse as an elliptic integral. In fact, I spent the bulk of my time on this approach, but I was not able to overcome my chief obstacle. For there are many different cross-sections which run through two points, and it is not sufficiently obvious which one to take (the one through the center? the one perpendicular to the ellipsoid in both places?) that we can go without proof. And to do this, we must either find the elliptic integrals for all possible cross-sections and find the smallest of them, or else we're back to solving the differential equation.

I was stuck on this point when my advisor pointed out that I had yet to look at the case for spheres (that is, how many spheres of a given radius can fit around a given ellipsoid?) Well, in retrospect this seems obvious, but unlike for circles, this is not the same as the best-distribution-of-points problem. In fact, though the problem of finding $k$ points on an ellipse such that the minimum distance is maximized is simple (just pick the points equally spaced about the circumference), finding $k$ circles of identical, maximal radius so that all touch the ellipse is surprisingly difficult.

Let us first prove that the two questions are not the same. If they are, then for every set of $k$ equally spaced points on an ellipse with major axis $A$ and minor axis $a \neq A$, we can draw circles tangent both to the ellipse and to their neighbors on either side (or else there would be a gap between some pair of circles, and we could shift the other circles slightly around the ellipse to spread the gap between every pair of neighboring circles, allowing us to increase all the circle radii; hence our radius was not maximal). So pick $k$ to

be very large, and consider two pairs of points, one pair each near the major and minor axes. Both pairs of points are some distance $d$ apart, but near the minor axis the curvature is within some small $\epsilon$ of $a$, and near the major axis the curvature is within $\epsilon$ of $A$; if we pick $k$ large enough we can make $\epsilon$ as small as we like. So between the two points of each pair the curve is as close as we like to a circle of radius $a$ and $A$, respectively; but then, since the radius must change continuously as we smoothly change the circle, the radii near each pair are as close as we like to those for the respective circles. This means, since as mentioned before we may assume the two circles in each pair are tangent, that we may use the formula mentioned in a previous footnote to calculate what the radius would have been had the two pairs of points actually been on circles of radius $a$ and $A$, and the actual radii will be as close as we like to those numbers. The formulas tell us $a(\frac{1}{1-\sin(d/4\pi a)}-1)$ in the first case and $A(\frac{1}{1-\sin(d/4\pi A)}-1)$ in the second case, which, as it may readily be seen (try using the fact that, for small enough $d$, $\sin(d/4\pi a) \approx d/4\pi a$ and $\sin(d/4\pi A) \approx d/4\pi A$, to within $\epsilon$, say) are not equal unless $A = a$. So we have a circle, or a contradiction.

Unfortunately, as might be expected, the equations are simple to set up, but messy to solve. I had finally reduced myself to considering the simplest special case: maximizing the radius of three identical circles tangent to a given ellipse. I still couldn't find a solution, but at least the equations were easy to find; but I noticed after drawing a few pictures that good solutions seemed to have a certain symmetry. Try it yourself: draw three mutually tangent circles, and try to draw an ellipse in the middle that touches all three. It seems that the positions in which the ellipse can fall all have the major axis of the ellipse on the line of tangency of some pair of the circles. This reduces the complexity of the problem considerably; we now know where the ellipse must touch the third of the circles (at the other end of the major axis), and the other two points of tangency must be symmetrically placed above and below this axis. We get a set of equations; all that is necessary is to set the derivative of the radius equal to zero, find the critical points and test them, just like a first-year calculus problem. But a hard calculus problem nevertheless; I suspect I might have been able to solve them if I had had two or three weeks left to find a solution by some means or another, but it was already the next-to-last of our eight weeks and I had to begin learning TeX and writing up what results I had.

*Suggestions for further research.* This last problem (that of determining the

best way of placing a given number of circles around an ellipse so that the circles can have radius as large as possible) seems to be a tractable problem, although with the limited time available to me at this late point in the program, I did not pursue it as much as I would have liked to. In fact, the case for four may be even easier than the case for three, because one obvious arrangement for four points presents itself (namely, the one with a point of tangency at each end of the major and minor axes). The truly ambitious reader could even generalize this problem to that of surrounding other types of convex bodies with three identical circles as large as possible (after solving the problem for ellipses, of course).

# V    Discussion

We have generalized the result in [1] to show that, for any given number $N$ of dimensions, if for a given size of $(N-1)$-sphere we can fit more than $N+1$ hyperspheres around the center hypersphere, we can certainly fit $2N$. In addition, we have shown that there is a radius (namely $1 + \sqrt{2}$) for which we can fit $2N$ hyperspheres but not $2N+1$, and that there exist radii (given by a more complicated formula) for which we can fit any number from 3 to $N+1$ hyperspheres around the central sphere and no more. (The case 1 is impossible, and the case 2 holds for all sufficiently large radii.) We have also discussed some of the problems associated with the general problem of finding the relationship between the maximum number of $(N-1)$-spheres that can be packed about the central $(N-1)$-sphere and the radius of the packing $(N-1)$-spheres, and suggested a different direction for generalizing this problem, namely that of packing circles about ellipses (the two-dimensional case, which is no longer trivial).

# VI    References

[**1**]  J. Angelos, G. Grossman, Yu. Ionin, E. Kaufman, T. Lenker and L. Rakesh, "Packability of five spheres on a sphere implies packability of six", *American Mathematical Monthly* **103** (1996), pp. 894–896.

[**2**]  J. H. Conway and N. J. A. Sloane, *Sphere Packings, Lattices and Groups*, Springer-Verlag, 1988.

[3] R. A. Rankin, "The closest packing of spherical caps in n dimensions", Proc. Glasgow Math. Assoc., **2** (1955), pp. 139–144.

[4] L. Fejes Tóth, *Lagerungen in der Ebene, auf der Kugel und in Raum*, 2nd ed., Springer-Verlag, 1972.

[5] J. H. van Lint, *Introduction to Coding Theory*, Springer-Verlag, 1972.

[6] K. Schütte and B. L. van der Waerden, "Auf welcher Kugel haben 5, 6, 7, 8, oder 9 Punkte mit Mindesabstand Eins Platz?", *Mathematische Annalen* **123** (1951), 96-124.

[7] R. M. Robinson, "Arrangement of 24 points on a sphere", *Mathematische Annalen* **144** (1961), 17-48.

[8] L. Danzer, "Finite point sets on $S^2$ with minimum distance as large as possible", *Discrete Mathematics* **60** (1986), 3-66.

[9] J. G. Simmonds, *A Brief on Tensor Analysis*, 2nd ed., Springer-Verlag, 1994.