# De Bruijn Sequences

Last year I wrote an article about 4-digit PIN codes. It became quite popular, and I even got asked to give a TED[X] talk about it which you can watch here (shameless plug).

The basis of the talk was that people are very predictable in the selection of their codes. Mathematically, whilst there are 10,000 ways that the digits 0-9 can be arranged into a four digit PIN, people's selections are not random. Only 426 distinct codes are needed to guess over half the PINs in use!

## Perfect World

Imagine, instead of being predictable, that people selected their codes entirely at random. If you wanted to guess the PIN of a four digit combination lock, you might have to walk through all 10,000 four digit combinations. Since there are four digits in each number, the worst case is that you'd have to type 40,000 key presses, and on average it would take you half this number.

That's a lot of key presses.

Can you do better? Can you do it with less key presses?

Well, if the lock is smart, no, you can't do better. But some systems that accept PINs are less sophisticated. These less secure mechanisms don't quantize the inputs in batches of four (or however long the code is), but instead simply look at just the *last four keys pressed*.

An example of this in use is shown below. Imagine you entered the six digits **123479** into one of these systems. As the system examines just the last four digits entered, this string would give you three distinct (overlapping) unlock attempts.



Attempt #1          Attempt #2          Attempt #3

This simpler system has the advantage of not having to worry about what state you are in when you start, or having to code in time-outs into the entry system. Without this simplification, if your code was **1234**, but just before you tried to enter it, unbeknownst to you, someone else had entered **99** on the keypad. When you entered your PIN, the lock would interpret this as **9912** *(fail)*, and then **34** *(still a fail and waiting for more digits before telling you it failed)*. Examining just the last four digits solves this problem.

If our system supports overlapping numbers we can be more efficient in guessing. We can create an input stream that goes through all permutations, but requires less key presses. The question is, how much more efficient? What is the shortest sequence of numbers that we can go through in order to ensure that all possible combinations of the digits are seen? Is it possible to create a sequence that does not repeat any sub-sequence of codes?

## Nicolaas Govert de Bruijn

The mathematics, number theory, combinatronics and logic of these types of problems were studied extensively by a Dutch Professor called *Nicolaas Govert de Bruijn* (9 July 1918 – 17 February 2012).

Sequences of these numbers are named after him as *De Bruijn Sequences*.

The quick answer is that, yes, it is possible to make a non-repeating sequence of numbers that covers ever sub-sequence internally, just once. However, before we look at the PIN number solution, let's look at some simpler versions of the problem …

Image Credit: Konrad Jacobs

## Simpler Sequences

De Bruijn sequences can be described by two parameters:

- $k$ the number of entities in the alphabet *e.g.* {0,1,2,3,4,5,6,7,8,9} for $k$=10

- $n$ the order (length of sub-sequence required) *e.g.* $n$=4 for a four digit long PIN.

These are typically describe by the representation $B(k,n)$.

For our PIN example, the notation would be $B(10,4)$

## B(2,2)

Let's start with a really simple example: $B(2,2)$

We'll use the dictionary {0,1} for the possible values. We want to generate a string that contains substrings each possible combination of two digits. Here is a solution:

## 0011

The first two digits give us **00**, the next two **01**, then **11**. To get to **10**, we need to *'Wrap Around'* taking the last digit from the string, and the first digit. (If this is not appropriate to do, like the key-press example, then we can simply append the first character from the string to the end to make **00110**).

 **NOTE:**  There can be multiple De Bruijn solutions to any problem. You can easily see this using even simple rotations of the string. Since every adjacent pair of digits in this string is unique, it does not matter what the starting position is. As *k* and *n* increase, the number of possible solutions grows rapidly. The hairy looking equation on the left shows the number of distinct solutions.

## B(2,3)

Here's a solution for $B(2,3)$:

## 00010111

Starting from the front, we have **000**, **001**, **010**, **101**, **011**, **111**, then staring to wrap around, we have **110** and finally **100**. All eight possible combinations are present in this string.

## B(2,4)

Here's a solution for $B(2,4)$:

## 0000100110101111

You can see we have **0000**, **0001**, **0010**, **0100**, **1001** ...

And here it is as a pretty ribbon. We can now hint at some of the awesome potential uses for these sequences. Imagine this *De Bruijn* sequence is written onto a looped tape, which passes past a reader head.

Each increment of the tape one position along gives a unique output. We have found an efficient mechanism of encoding position as there is a distinct code for any contiguous set of four digits (in this case, as *n*=4). How cool is that?

This sequence also walks through all permutations of combinations from **0000-1111**. If you are a software engineer or tester, I'm sure this is giving you interesting ideas for how you can implement this as test cases to walk through all binary inputs to your functions. Each shift of one bit gives a unique binary word (of length of the substring size *n*), which does not repeat, goes through each possible number, then returns to the start again.

## B(2,6)

Skipping ahead a couple, here is a solution for $B(2,6)$:

0000001000011000101000111001001011001101001111010101110110111111

## B(6,2)

Of course, we can also change the size of the dictionary. In this example rather than simple binary *k*=2, I've increased this to *k*=6 to use possible values {0,1,2,3,4,5}. Here is the output for $B(6,2)$:

0010203040511213141522324253343544555

Reading this from left to right we can see that we have: **00**, **01**, **10**, **02**, **20** ... **41**, **15** ... **33**, **34** ... **55**, **50**

One last example, here's a sequence for $B(5,3)$:

**000100200300400110120130140210220230240310320330340410420430441112113114
221231241321331341421431442223224233234243244334334344444**

# Back to PIN

Returning to our PIN cracking, the $B(10,4)$ is 10,000 digits long, and contains every single substring for each combination of digits. As we have to type in this string, the concept of *Wrap-Around* is meaningless, so rather then wrap-around, we simply add the first three numbers again to the end of the strong. So, the maximum number of keypress events is 10,003. Far fewer than the 40,000 that would be required if typing them in completely!

For no useful reason at all, here are all 10,003 digits of this string:

0001000200030004000500060007000800090001100120013001400150016001700180019002100220023002400250026002700280029003100320033003400350036003700380039
0041004200430044004500460047004800490051005200530054005500560057005800590061006200630064006500660067006800690071007200730074007500760077007800790
0081008200830084008500860087008800890091009200930094009500960097009800990101020103010401050106010701080109011101120120130140150160170180191190
1210122012301240125012601270128012901310132013301340135013601370138013901410142014301440145014601470148014901510152015301540155015601570158015901590
1610162016301640165016601670168016901701720173017401750176017701780179018101820183018401850186018701880189019101920193019401950196019701980199020
2020302040205020602070208020902110212021302140215021602170218021902210222022302240225022602270228022902310232023302340235023602370238023902410241024
2024302440245024602470248024902510252025302540255025602570258025902610262026302640265026602670268026902710272027302740275027602770278027902810282028
2028302840285028602870288028902910292029302940295029602970298029903030403050306030703080309031103120313031403150316031703180319031303240320320323032403240
3250326032703280329033103320333033403350336033703380339034103420343034403450346034703480349035103520353035403550356035703580359036103620363036403640
3650366036703680369037103720373037403750376037703780379038103820383038403850386038703880389039103920393039403950396039703980399040040500406040704040
8040904110412041304140415041604170418041904210422042304240425042604270428042904310432043304340435043604370438043904410442044304440445044604470447044
8044904510452045304540455045604570458045904610462046304640465046604670468046904710472047304740475047604770478047904810482048304840485048604870480480
8048904910492049304940495049604970498049905010502050305040505050605070508050905110512051305140515051605170518051905210522052305240525052605270528052905305300
5340535053605370538053905410542054305440545054605470548054905510552055305540555055605570558055905610562056305640565056605670568056905710572057305730
5740575057605770578057905810582058305840585058605870588058905910592059305940595059605970598059905990606070608060901106120613061406150616061706180601
9062106220623062406250626062706280629063106320633063406350636063706380639064106420643064406450646064706480649065106520653065406550656065706580650580
9066106620663066406650666066706680669067106720673067406750676067706780679068106820683068406850686068706880689069106920693069406950696069706980690
9070708070907110712071307140715071607170718071907210722072307240725072607270728072907310732073307340735073607370738073907410742074307440745074607460
7470748074907507510752075307540755075607570758075907610762076307640765076607670768076907710772077307740775077607770778077907810782078307840785070860
7870788078907910792079307940795079607970798079908110812081308140815081608170818081908210822082308240825082608270828082908310832083308340830830834087087
5083608370838083908410842084308440845084608470848084908510852085308540855085608570858085908610862086308640865086608670868086908710872087308740740870
5087608770878087908810882088308840885088608870888088908910892089308940895089608970898089909110912091309140915091609170918091909210922092309240920
92509260927092809290931093209330934093509360937093809390941094209430944094509460947094809490951095209530954095509560957095809590961096209630960640
9650966096709680969097109720973097409750976097709780979098109820983098409850986098709880989099109920993099409950996099709980999111211131114111115
111611171118111911221123112411251126112711281129113211331134113511361137113811391142114311441145114611471148114911521153115411551156115711581159
1162116311641165116611671168116911721173117411751176117711781179118211831184118511861187118811891192119311941195119611971198119911221123112411251512
1612171218121912221223122412251226122712281229123212331234123512361237123812391242124312441245124612471248124912521253125412551256125712581259129120
6212631264126512661267126812691272127312741275127612771278127912811282128312841285128612871288128912921293129412951296129712981299131314131513161317
131813191322132313241325132613271328132913321333133413351336133713381339134213431344134513461347134813491352135313541355135613571358135913621363
1364136513661367136813691372137313741375137613771378137913821383138413851386138713881389139213931394139513961397139813991422142314241425142614271460
2214231424142514261427142814291432143314341435143614371438143914421443144414451446144714481449145214531454145514561457145814591462146314641465146514
6614671468146914721473147414751476147714781479148214831484148514861487148814891492149314941495149614971498149915211522152315241525152615271528152915525
15261527152815291532153315341535153615371538153915421543154415451546154715481549155215531554155515561557155815591562156315641565156615671568156815691569
1572157315741575157615771578157915821583158415851586158715881589159215931594159515961597159815991612161316141615161616171618161916221623162416280162916
3216331634163516361637163816391642164316441645164616471648164916521653165416551656165716581659166216631664166516661667166816691672167316741675167516
7617671768176917821783178417851786178717881789179217931794179517961797179817991812181318141815181618171818181918221823182418251826182718281829182937
1738173917421743174417451746174717481749175217531754175517561757175817591762176317641765176617671768176917721773177417751776177717781779179179178177137
1784178517861787178817891792179317941795179617971798179918121813181418151816181718181819182218231824182518261827182818291831183218331834183518361836183518
4618471848184918521853185418551856185718581859186218631864186518661867186818691872187318741875187618771878187918821883188418851886188718881889188918918
9218931894189518961897189818991912191319141915191619171918191919221923192419251926192719281929193219331934193519361937193819391942194319441945194619471948194919491955
1956195719581959196219631964196519661967196819691972197319741975197619771978197919821983198419851986198719881989199219931994199519961997199819991999
22223222422252226222722282229223222332234223522362237223822392242224322442245224622472248224922522253225422552256225722582259226222632264226522662267226822692227
327422752276227722782279228322842285228622872288228922922293229422952296229722982299232324232523262327232823292332233323342335233623372338233923423442
345234623472348234923522353235423552356235723582359236223632364236523662367236823692372237323742375237623772378237923822383238423852386238723892393239423952
396239723982399242324252426242724282429243324342435243624372438243924422443244424452446244724482449245324542455245624572458245924622463246424662467246746
824692472247324742475247624772478247924822483248424852486248724882489249224932494249524962497249824992522252325242525252625272528252925322533253425352535253925352592
545254625472548254925522553255425552556255725582559256225632564256525662567256825692572257325742575257625772578257925822583258425852586258725882589259325942959
59625972598259926222623262426252626262726282629263226332634263526362637263826392642264326442645264626472648264926522653265426552656265726582659266226632666266726
326742675267626772678267926822683268426852686268726882689269226932694269526962697269826992722272327242725272627272728272927322733273427352736273727382739274274274082
7492753275427552756275727582759276227632764276527662767276826922693269426952696269727282729273227332734273527362737273827392742274327442745274627472748274927474
82829283328342835283628372838283928422843284428452846284728482849285228532854285528562857285828592862286328642865286628672868286928722873287428752876287728782879287287992
9632964296529662967296829692972297329742975297629772978297929822983298429852986298728982899292329242925292629272928292929322933293429352936293729382939293923943952
3346334733483349335433553356335733583359336233633364336533663367336833693372337333743375337633773378337933823383338433853386338733883389339233933394339533963397398399343434534
3643374337443734377343834393344335433553356335733583359336233633364336533663367336833693372337334373544353545354635475348354935533554355535563557355835593563564356535663567356835693597435753576357735783579358435853863587
358835893594359535963597359835993636373638363936443645364636473648364936533654365536563657365836593663366436653666366736683669363743674367536763637367437437867438673867386
84368536863687368836893694369536963697369836993723373337433753376337733783379338234837783873738373838583937843785378637877388378937923793379437953796379737983793883867487387538
763877378837938843885388638873888388938943895389636373638363936443645364636473648364936533554355535563557358835893593359435953596359737397399379439597346374753948395378463894539
397639773978397939843985398639873988398938943895389636373638363936443645364636473648364936533554355535563558359335943595359637464746744747545746744847837487435849395739863987489
948854486448744884489449544964497449844994544965446744884489449454496449744984499454496544674488448944954496449744984499444854486448744884489449544964497449844994499454496544674488448944954496449744984499449948154818518518518518518518518
5964597459845994646474647464849454654664746454665466654665466654654664746454665466654654664746454665466654665466654664746454665466654665466654664746454665466654665466654664746454665466654665466654665466654665466654665
7564757475847594765476647674768476947754776477747784779478547864787478847894795479647974798479944854498544854486448744884489449544964497449844984499454496544674488448944954496449744984489449544864489449454496448744884489944804
5487648774878478947944795479647974798479944854498544854486448744884489449544854489449454496448744884489944854486448744884489944804489454496544674488448944954489448744884489944804
996499749984999555565575585595665675685695575765785795855865875885985965975985996665677566765676765667665676665667665676665
865687568856956965957557558558655656655656656656556655656656656556655656656656556655656656656556655656656656556655656656656
865887588858895896589758985899566596759865996657697697697698696697697697698696697697697698696697697698696697697697698696
96767686769677767786779678767886789679767986799686969687688689689697698697698697698697698697698697698698
77987799787879788878897898789879979788789978998888898898989999000

# How do these things work?

I glossed over how to calculate a De Bruijn sequence earlier. I'll make up for it now by explaining the principles for how these things work. The concept is not complicated, just tedious; things that computers can do easily. We're going to look at this problem using *Directed Graphs*.

**State$_{n-1}$** | **State$_n$** | **State$_{n+1}$**

For this example, I'm going to use $B(2,3)$. This sequence is simple enough to get the concept without being too complex and busy.

Every sub-string in the sequence can be described as a *node* on a directed graph. From each node, it is possible to add either a **0** or a **1** to make the next number in the sequence.

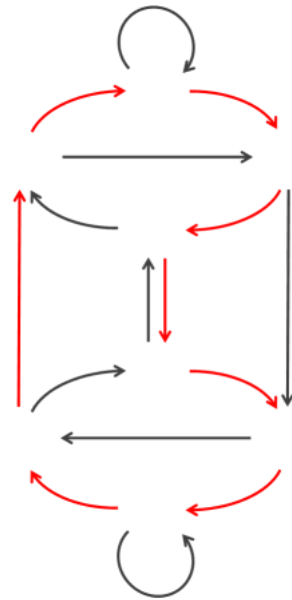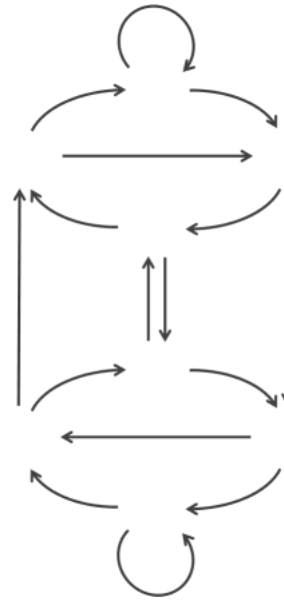Similarly, each state/node on the graph could have been formed from adding a digit from one of two earlier states.

First we write down all the nodes that are possible using our data dictionary.

In this case these are **000 - 1111**.

Then we connect the nodes with lines showing which possible next states are achievable from each location. Each node has two outbound links corresponding to the addition of a **1** or **0** from that state.

These edges are directed. They have a direction (depicted by the arrows)

Note - For States **000** and **111**, the addition of a **0** or **1**, respectively, takes you back to the same node.

In order to make sure every substring is present in the solution, we need to make sure each node is passed through once (and only once). We need to trace a path through the graph (following the arrows) to connect the nodes.

A path that traverses a graph and visits each node exactly once is called a Hamiltonian Path.

One such Hamiltonian Path on the $B(2,3)$ graph is shown in red on the left.

This highlights why there are many different solutions (not just rotations). There is more than one way to walk through this graph.

The way to create a De Bruijn sequence is to find a Hamiltonian Path through a graph their nodes.

As the values of $k$ and $n$ increase, so does the complexity of the graph, but the principle is the same.

Here's a slightly more complex one using a dictionary {1,2,3} for a sequence $B(3,2)$

Find a path through the graph the passes through each node just once, and you have your solution.

# Other uses

We already hinted that De Bruijn sequences can be used for encode/decode positions, and that they can be used by savvy computer programmers. What are some other applications?

I've seen magicians use a similar principle in various card tricks. By pre-arranging (loading) a deck with a known sequence of red and blacks cards in a De Bruijn sequence, it allows him/her to know what the position is, and thus, what the next card will be. Using a binary encoding of the red/black cards to generate a unique number, which can then, be used to encode the value of what the next card will be.

Because of the wrap-around nature of De Bruijn sequences, a loaded deck can be cut as many times as desired without upsetting or disturbing the encoding. I don't want to link directly examples, as I don't want to ruin the tricks for others that might be performing them, but I'm sure if you know how to use a web search tool, you can find some strategies.

## DNA

The concept of Hamilton cycles and De Bruijn are used extensively in modern DNA sequencing techniques.

A large DNA chain can be broken up into smaller pieces (The smaller pieces being easier to process and sequence). Then, the results of these smaller pieces can be 'glued' back together like some kind of giant jigsaw because the individual sub-pieces contain overlaps with other partial strings.

This technique is called **Shotgun Sequencing**.

Below you can see a representation of a long chain of DNA. This is smashed into smaller sub-pieces (of different sizes), which are easier to classify. The classified pieces can then be joined back together to form the complete sequence by looking at the overlaps between the substrands.

GCATTGCATTAGCAATAT
AATATGTTAGCAATATCCGC
CCGCCDCGCTATGCGAAATGGCTT

In some cases, the shotgun sequencing technique generates multiple canonical sequences. Of course, only one of these can be the real sequence. There are a variety of tests that can narrow down which of these sequences is the original one. Shotgun sequencing, although it does not *necessarily* define the exact sequence, greatly speeds up the process by reducing the number of possible candidates.
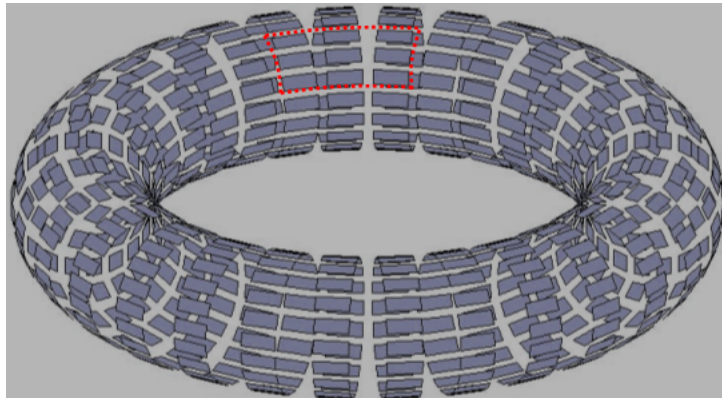
The advent of shotgun sequencing techniques advanced the initial mapping of the human genome by years, and it has provided biologists, geneticists, and doctors with a powerful new tool. The ability to sequence genetic data rapidly has potential benefits not only for life and health science professionals, but also for the public at large. It's pretty cool.

## Chess

Computer programs that play chess make use of De Bruijn sequences. A chess board is conveniently *8 x 8* squares, and these can be represented by the numbers 0-63, which is a nice fit for a six bit long De Bruijn sequence.
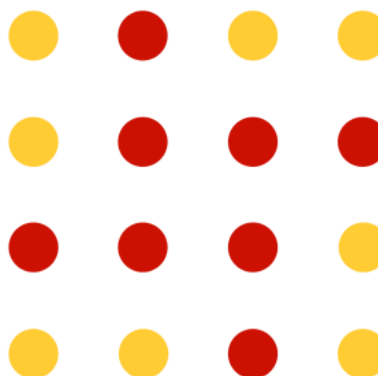
# De Bruijn Toroids

OK, prepare your mind for something cooler still. We can expand a De Bruijn one dimensional sequence into two dimensions! We call these Toroids because every row and column wraps around onto itself.



Sliding a window over this surface creates a unique matrix for a well-formed De Bruin array.

Things get much more complicated because not only do we need to specify the number of items in the dictionary, but also the two dimensions of the window.
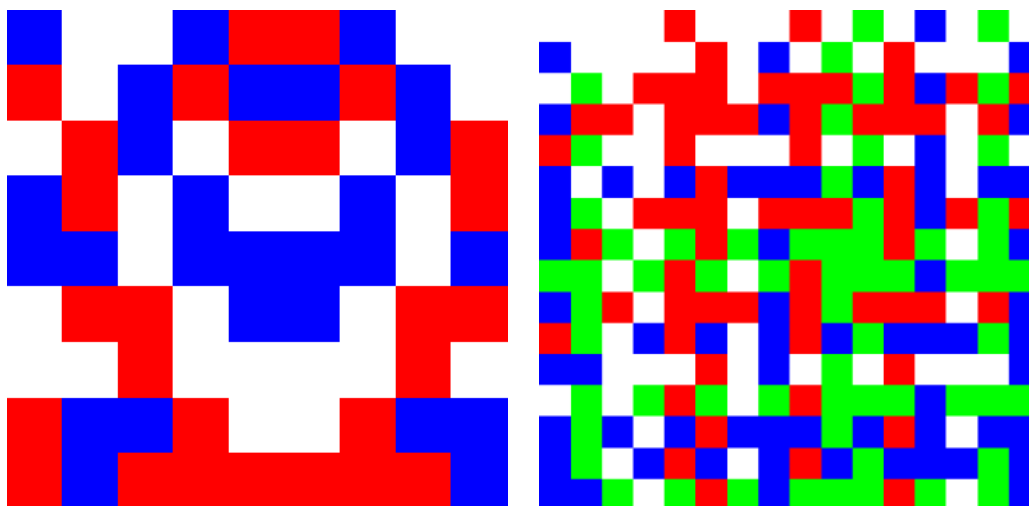
Below is a simple example of an array that has a dictionary of two {red, yellow}, and a (2 x 2) window:



If you look carefully, you will see that every combination of red/yellow dots appears all possible combinations of (2x2) sub-matrices. Remember, you need to wrap-around (for instance to get a matrix containing four yellow dots you need the matrix that is the four corners wrapped around the outside!)

I'm sure you can see instant applications for something like this in determining (x,y) position. Each window is unique and determines coordinate position.

Of course, we're not restricted to using just a dictionary of two items. Below are two larger De Bruijn toroids. The one on the left has a dictionary size of three, and matrix window of (2 x 2). The one on the right has the same matrix window and this time has a dictionary of size four.



The grid on the left is (9 x 9) and that on the right is (16 x 16), allowing representation of all 256 possible combinations that a dictionary of order four into a (2 x 2) matrix. De Bruijn Toroids also don't have to be square; for instance there is a solution the four dictionary solution that fits in a matrix (8 x 32) instead of (16 x 16).

Of course there are multiple equivalent translations of these solutions as any combination of row or column shift is a valid solution. It's also possible to tesselate these tiles perfectly edge to make a map the repeats with the same order frequency.

# Digital Paper



Anoto, a Swedish company, has invented, and has numerous patents, concerning the use of digital pens and paper. The Anoto concept is based around a distinct pattern of dots that can be printed onto paper stock. Using a special pen, which contains a small camera, the dot pattern can be read, and exact position on the paper determined.

Whilst not quite the same concept, the Anoto system works by small changes in the displacement in the position of dots from a nominal mean position.

## Want more? Here's four more random articles I wrote (Click here to refresh)

Posted on March 10, 2021    Posted on November 27, 2019    Posted on December 30, 2018    Posted on April 11, 2021

**Painted Cube Puzzle**
Yet another painted cube puzzle.

**Sampling from non-symmetric discrete outcomes**
A super clever algorithm.

**Sudoku Solver**
Online Sudoku Solver.

**Coffee Time Challenge**
Two unknowns in arithmetic and geometric progressions.

SHOW    SHOW    SHOW    SHOW

You can find a complete list of all the articles here.    Click here to receive email alerts on new articles.

☕ Buy me a coffee