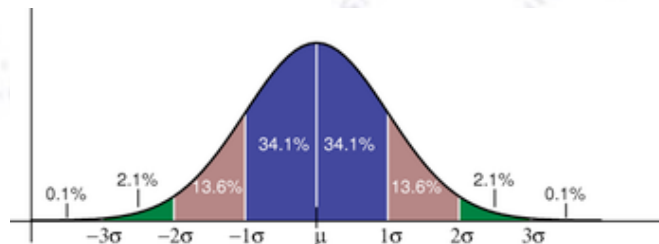


# Applied Statistics

## Testing random Number



Troels C. Petersen (NBI)

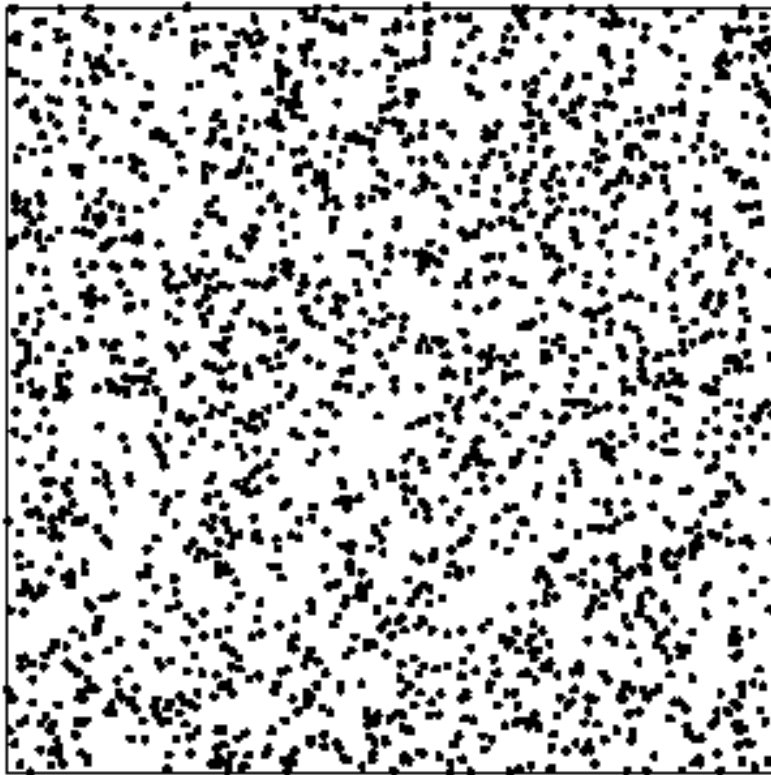


*"Statistics is merely a quantisation of common sense"*

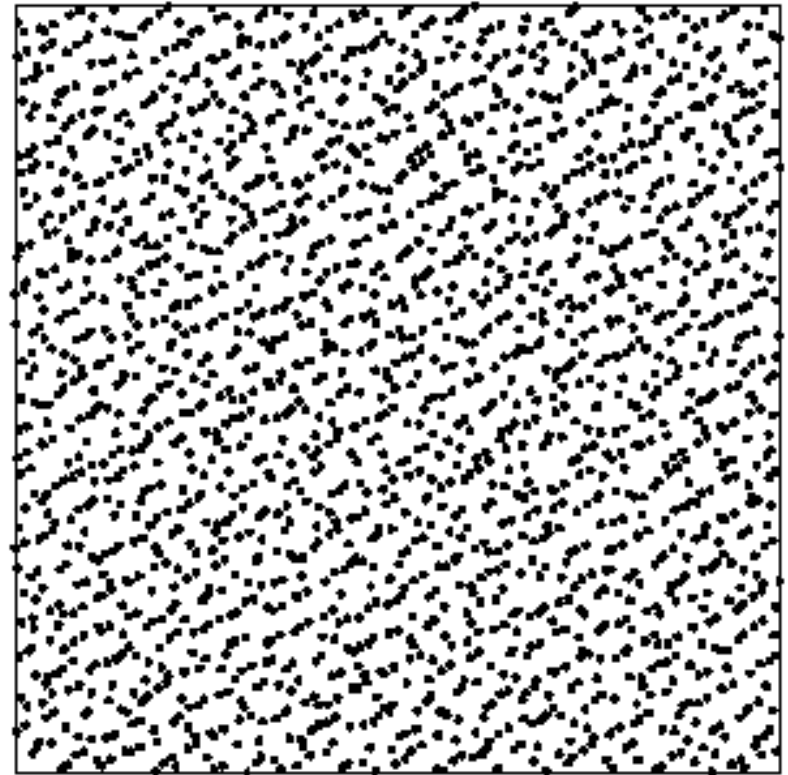
# Random numbers

“Anyone who considers arithmetical methods of producing random digits is, of course, in a state of sin.” [John Von Neumann]

Random



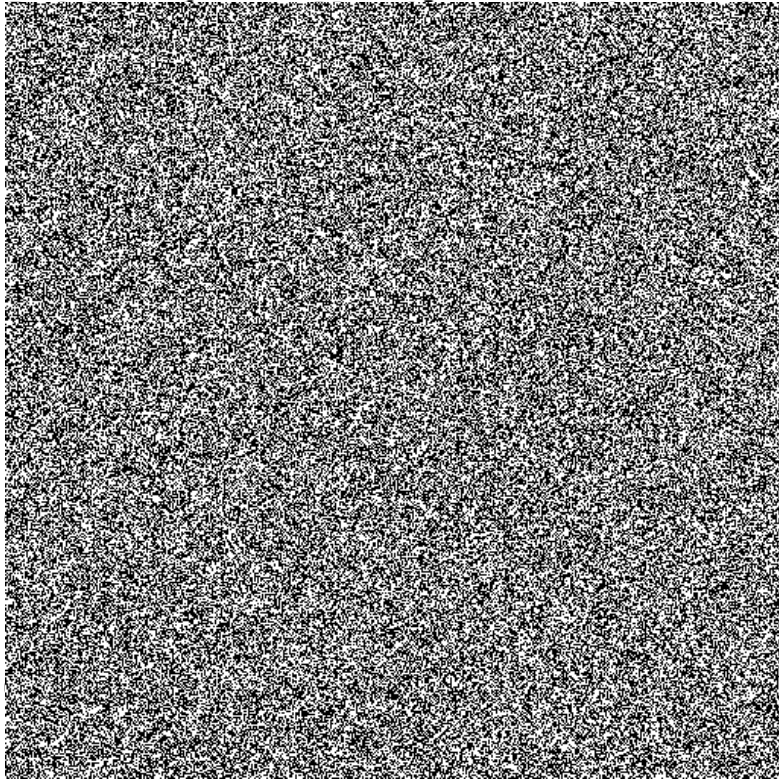
Quasi-Random



# Random numbers

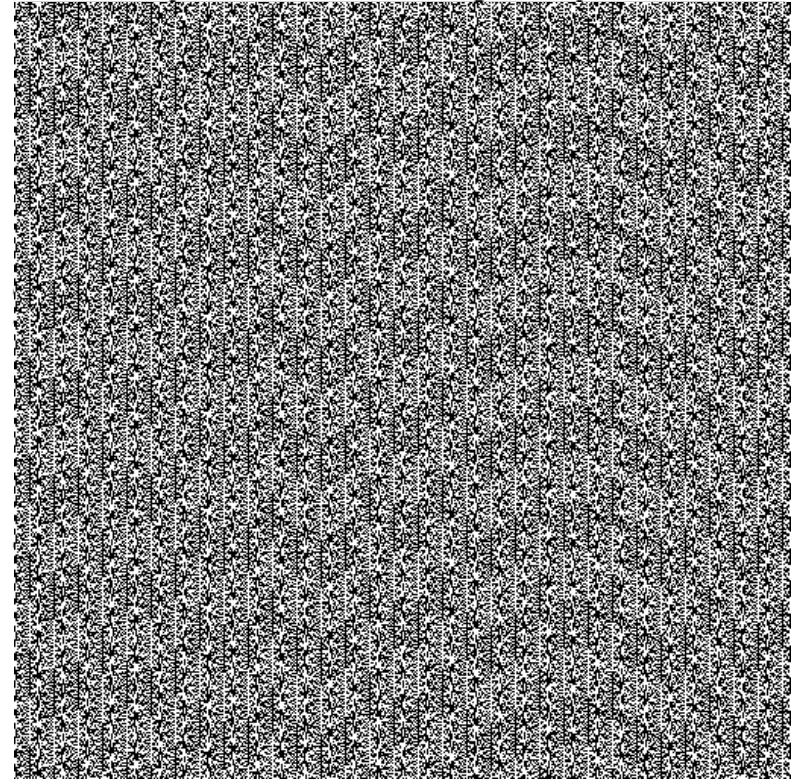
“Anyone who considers arithmetical methods of producing random digits is, of course, in a state of sin.” [John Von Neumann]

## Random



RANDOM.ORG

## Quasi-Random



PHP rand() on Microsoft Windows

# The build-in random number generators are sometimes not optimal!

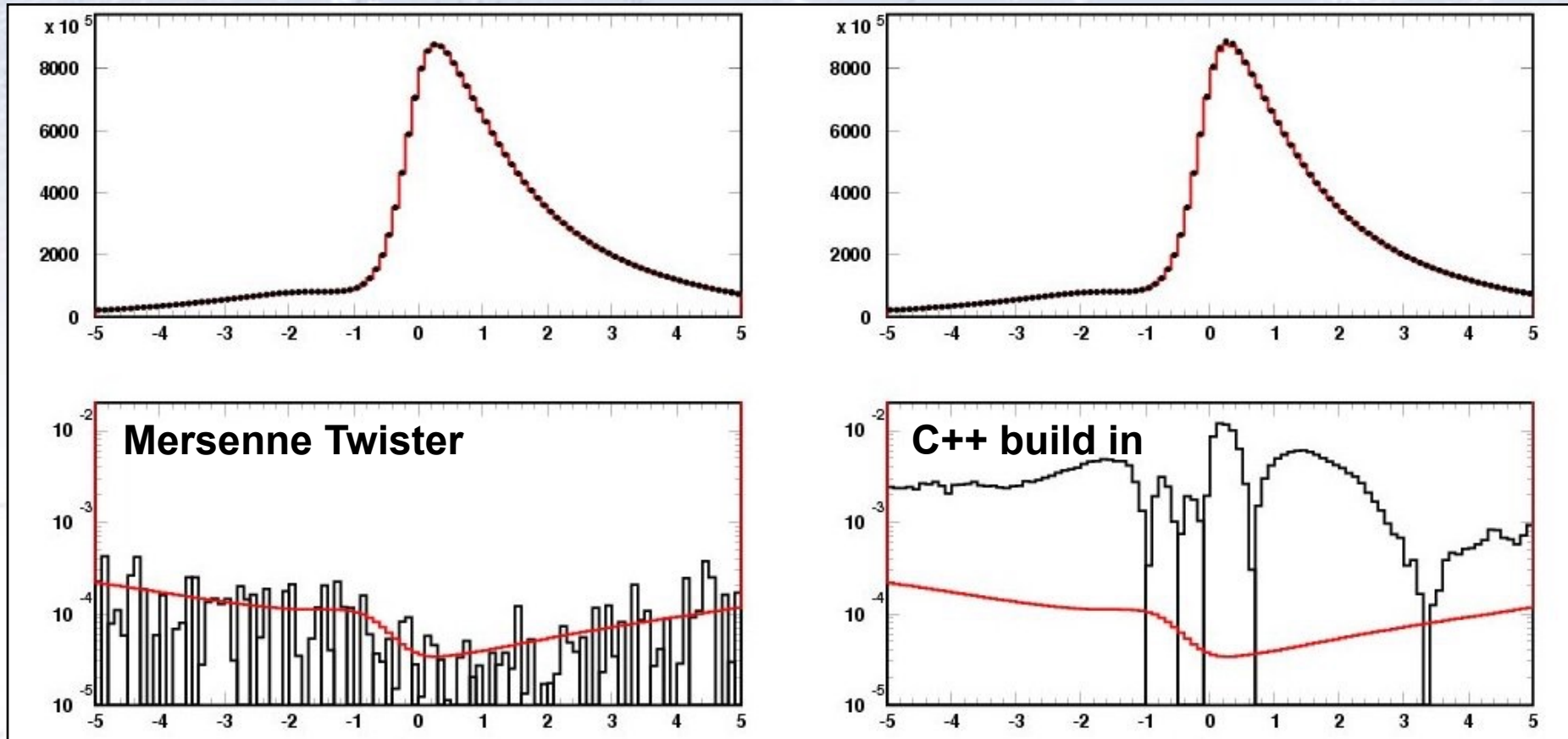


Figure 8.1: To test the resolution function convolutions and the Monte Carlo generation, a large Monte Carlo sample of 25 billion events is generated, binned, and compared to the predictions of the convolutions. The actual function used is  $e^{-2|t|} [\cosh(1.5t) + 0.5 \sinh(1.5t) + 0.1 \cos(t) + 0.9 \sin(t)]$  convoluted with a GExp function with  $(\sigma, s, \tau) = (0.2, 0.7, 0.5)$ . In the upper plots, the red histograms show the predicted values and the black dots the number of accumulated Monte Carlo events. In the lower plots the relative difference is shown in a logarithmic plot. The red curve shows the expected statistical deviation and the black histogram the actual relative difference. In (a) the Monte Carlo was generated using a decent random number generator while in (b), it was generated using the standard built-in C++ routine. The agreement in (a) seems to be perfect within the statistical precision of 4-5 significant digits. In (b) on the other hand, the Monte Carlo deviates already at 2-3 significant digits, which is unacceptable.

T. Kittelmann, master thesis (2002)

# Random numbers

# Random numbers

*“Most studies find that human subjects have some degree of non randomness when attempting to produce a random sequence of e.g. digits.”*

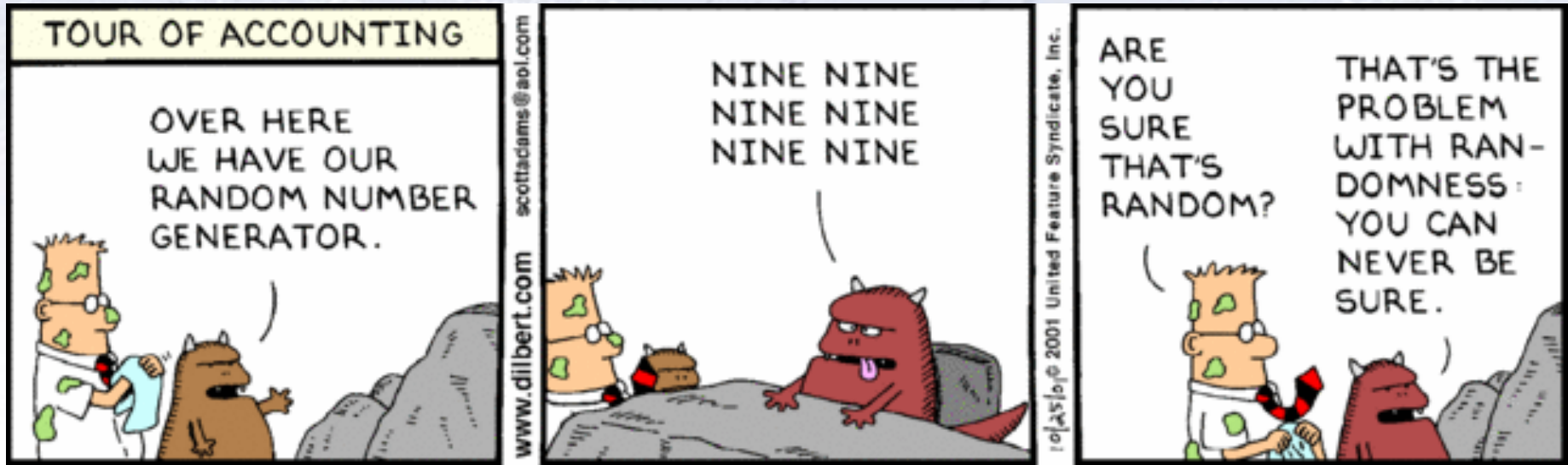
[Wikipedia, On random numbers produced by humans]

# Random numbers

*“Most studies find that human subjects have some degree of non randomness when attempting to produce a random sequence of e.g. digits.”*

[Wikipedia, On random numbers produced by humans]

# Testing random numbers



In contrast to popular believe, there are ways of testing random numbers. But it is not easy! A recommended list (Louis Foley, Random.org, 2001) is:

- A chi-square test
- A test of runs above and below the median
- A reverse arrangements test
- An overlapping sums test
- A binary rank test for  $32 \times 32$  matrices

# Binary Rank Test

A binary rank test for  $31 \times 31$  matrices:

The leftmost 31 bits of 31 random integers from the test sequence are used to form a  $31 \times 31$  binary matrix over the field  $\{0,1\}$ .

The rank is then determined. That rank can be from 0 to 31, but ranks  $< 28$  are rare, and their counts are pooled with those for rank 28.

Ranks are found for 40,000 such random matrices and a chi-square test is performed on counts for ranks 31,30,29 and  $\leq 28$ .

I.e. given that this distribution is known for truly random numbers, you can see if the distribution of ranks is the same!

DieHard tests

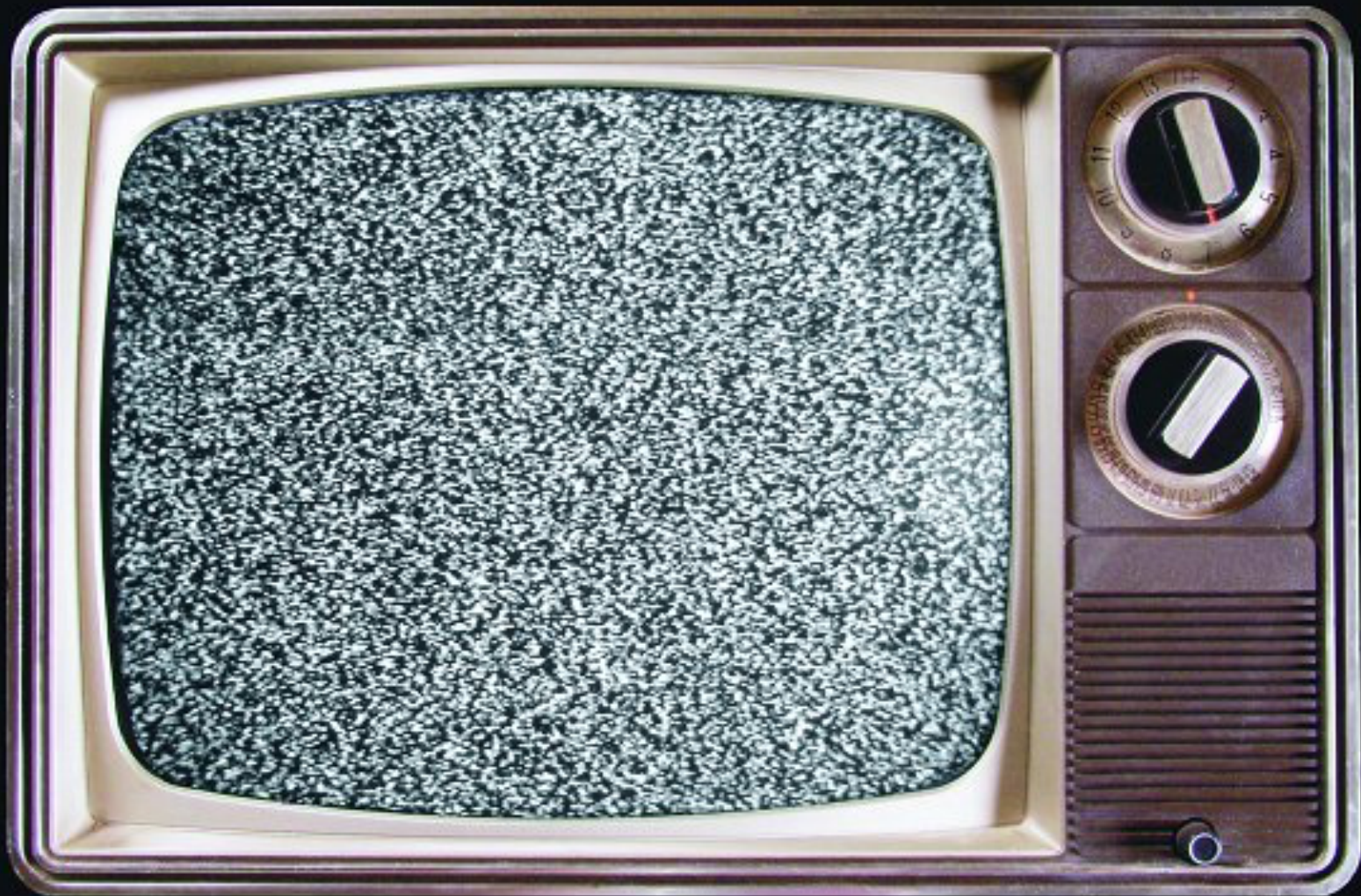


# More tests for randomness

NIST Statistical Test Suite		
Test	Defect Detected	Property
Frequency (monobit)	Too many zeroes or ones	Equally likely (global)
Frequency (block)	Too many zeroes or ones	Equally likely (local)
Runs test	Oscillation of zeroes and ones too fast or too slow	Sequential dependence (locally)
Longest run of ones in a block	Oscillation of zeroes and ones too fast or too slow	Sequential dependence (globally)
Binary matrix rank	Deviation from expected rank distribution	Linear dependence
Discrete fourier transform (spectral)	Repetitive patterns	Periodic dependence
Non-overlapping template matching	Irregular occurrences of a pre-specified template	Periodic dependence and equally likely
Overlapping template matching	Irregular occurrences of a pre-specified template	Periodic dependence and equally likely
Maurer's universal statistical	Sequence is compressible	Dependence and equally likely
Linear complexity	Linear feedback shift register (LFSR) too short	Dependence
Serial	Non-uniformity in the joint distribution for m-length sequences	Equally likely
Approximate entropy	Non-uniformity in the joint distribution for m-length sequences	Equally likely
Cumulative sums (cusum)	Too many zeroes or ones at either an early or late stage in the sequence	Sequential dependence
Random excursions	Deviation from the distribution of the number of visits of a random walk to a certain state	Sequential dependence
Random excursions variants	Deviation from the distribution of the number of visits (across many random walks) to a certain state	Sequential dependence

# What to use for random seed?

"Randomisation is too important to be left to chance." [ J. D. Petruccelli]



# NSA and Dual\_EC\_DRBG

Cryptography is based on random numbers. One of the four accepted random number generators (RNG) used for cryptography was Dual\_EC\_DRBG.

Several academics point out, that Dual\_EC\_DRBG was a very poor and possibly “back-doored” pseudorandom number generator! Nevertheless, one of the large cryptography companies in the US, RSA Security, continued using it.

In 2013, Edward Snowden published papers that showed, that the NSA had put a backdoor in the Dual\_EC\_DRBG algorithm! The lack of perfect randomness allowed NSA to break the encryption.

Outputs of multiple independent RNGs can be combined (for example, using a bit-wise XOR operation) to provide a combined RNG at least as good as the best RNG used. This is referred to as software whitening.

The background is a faded nautical chart. It features contour lines representing depth, with labels such as 210, 240, 270, 300, 330, 360, 390, 420, 450, 480, 510, 540, 570, 600, 630, 660, 690, 720, 750, 780, 810, 840, 870, 900, 930, 960, 990, 1020, 1050, 1080, 1110, 1140, 1170, 1200, 1230, 1260, 1290, 1320, 1350, 1380, 1410, 1440, 1470, 1500, 1530, 1560, 1590, 1620, 1650, 1680, 1710, 1740, 1770, 1800, 1830, 1860, 1890, 1920, 1950, 1980, 2010, 2040, 2070, 2100, 2130, 2160, 2190, 2220, 2250, 2280, 2310, 2340, 2370, 2400, 2430, 2460, 2490, 2520, 2550, 2580, 2610, 2640, 2670, 2700, 2730, 2760, 2790, 2820, 2850, 2880, 2910, 2940, 2970, 3000, 3030, 3060, 3090, 3120, 3150, 3180, 3210, 3240, 3270, 3300, 3330, 3360, 3390, 3420, 3450, 3480, 3510, 3540, 3570, 3600, 3630, 3660, 3690, 3720, 3750, 3780, 3810, 3840, 3870, 3900, 3930, 3960, 3990, 4020, 4050, 4080, 4110, 4140, 4170, 4200, 4230, 4260, 4290, 4320, 4350, 4380, 4410, 4440, 4470, 4500, 4530, 4560, 4590, 4620, 4650, 4680, 4710, 4740, 4770, 4800, 4830, 4860, 4890, 4920, 4950, 4980, 5010, 5040, 5070, 5100, 5130, 5160, 5190, 5220, 5250, 5280, 5310, 5340, 5370, 5400, 5430, 5460, 5490, 5520, 5550, 5580, 5610, 5640, 5670, 5700, 5730, 5760, 5790, 5820, 5850, 5880, 5910, 5940, 5970, 6000, 6030, 6060, 6090, 6120, 6150, 6180, 6210, 6240, 6270, 6300, 6330, 6360, 6390, 6420, 6450, 6480, 6510, 6540, 6570, 6600, 6630, 6660, 6690, 6720, 6750, 6780, 6810, 6840, 6870, 6900, 6930, 6960, 6990, 7020, 7050, 7080, 7110, 7140, 7170, 7200, 7230, 7260, 7290, 7320, 7350, 7380, 7410, 7440, 7470, 7500, 7530, 7560, 7590, 7620, 7650, 7680, 7710, 7740, 7770, 7800, 7830, 7860, 7890, 7920, 7950, 7980, 8010, 8040, 8070, 8100, 8130, 8160, 8190, 8220, 8250, 8280, 8310, 8340, 8370, 8400, 8430, 8460, 8490, 8520, 8550, 8580, 8610, 8640, 8670, 8700, 8730, 8760, 8790, 8820, 8850, 8880, 8910, 8940, 8970, 9000, 9030, 9060, 9090, 9120, 9150, 9180, 9210, 9240, 9270, 9300, 9330, 9360, 9390, 9420, 9450, 9480, 9510, 9540, 9570, 9600, 9630, 9660, 9690, 9720, 9750, 9780, 9810, 9840, 9870, 9900, 9930, 9960, 9990. There are also labels for magnetic variation: 'MAGNETIC' and 'VAR 10°15'W'. A small crosshair symbol is located near the 'VAR 10°15'W' label. In the upper right corner, there is a label '1.52 BITTER END YACHT CLUB'.

# Solutions for exercise

# Solution to exercise

**A: data\_RandomDigits2022\_ObviouslyNonRandom.txt**

Series of length 11 7 tables, shifted randomly (mod 10).

**B: data\_RandomDigits2022\_TrulyRandom.txt**

**C: data\_RandomDigits2022\_DigitsOfPi.txt**

**D: data\_RandomDigits2022\_100WithShift.txt**

The same 100 random digits with a random shift (mod 10).

**E: data\_RandomDigits2022\_Every11thNonRandom.txt**

Random digits with every 11th being the previous +2 (mod 10).

**F: data\_RandomDigits2022\_FromStudents.txt**

**G: data\_RandomDigits2022\_IceCubeData.txt**

The 3rd decimal of a (random) column of data from IceCube.