

[Thesis text....] ...As one can easily notice, we have got exactly the same curve as Kanerva. Both his and our model expect that, after reading, say, from 550 bits of distance from a written bitstring, we should obtain the $n/2$ equator distance. We have not, however. This question has intrigued us, and here we look for a more analytic explanation than merely interference from the other written attractors. Let us turn back to mathematics to study this anomaly.

1 A deviation from the equator distance?

Kanerva writes¹:

You have done an incredibly thorough analysis of SDM. I like the puzzle in your message and believe that your simulations are correct and to be learned from. So what to make of the difference compared to my Figure 7.3 (and your Figure ??)? I think the difference comes from my not having accounted fully for the effect of the other 9,999 vectors that are stored in the memory. You say in it

“Our results show that the theoretical prediction is not accurate. There are interaction effects from one or more of the attractors created by the 10,000 writes, and these attractors seem to raise the distance beyond 500 bits (Figure ??).”

I think that is correct. It also brings to mind a comment Louis Jaeckel made when we worked at NASA Ames. He pointed out that autoassociative storage (each vector is stored with itself as the address) introduces autocorrelation that my formula for Figure 7.2 did not take into account. When we read from memory, each stored vector exerts a pull toward itself, which also means that each bit of a retrieved vector is slightly biased toward the same bit of the read address, regardless of the read address. We never worked out the math.

¹Email thread ‘SDM: A puzzling issue and an invitation’, started March 16th 2018, in which we discussed the aforementioned discrepancy. To think that some centuries ago, all scientific publishing was the exchange of such letters.

This is an important observation. A hard location is activated because it shares many dimensions with the items read from or written onto it. Imagine the ‘counter’s eye view’: each individual counter ‘likes’ to write on its own corresponding bit-address value more than it likes the opposite; as each hard-location has a say in its own area — and nowhere else.

Let x and y be random bitstrings and n be the number of dimensions in the memory; let x_i and y_i be the i -th bit of x and y , respectively; and $d(x, y)$ be the Hamming distance. Whilst the probability of a shared bit-value between same dimension-bits in two random addresses is $1/2$, an address only activates hard-locations close to it. Let us call these shared bitvalues a *bitmatch in dimension i* .

So, what is the probability of bitmatches given that we know the access radius r between the address and a hard-location?

Theorem 1.1. Each dimension has a small pull bias, which can be measured by $P(x_i = y_i | d(x, y) \leq r) = \frac{\sum_{k=0}^r \binom{n-1}{k}}{\sum_{k=0}^r \binom{n}{k}}$.

Proof. The left-hand expression $P(x_i = y_i | d(x, y) \leq r)$ computes the probability of a bitmatch in i , given that we know that x and y are in the access radius defined by r , i.e., $d(x, y) \leq r$.

Applying the law of total probability to the left-hand expression we obtain

$$\sum_{k=0}^r P(x_i = y_i | d(x, y) = k \leq r) P(d(x, y) = k | d(x, y) \leq r) \quad (1)$$

We also know that

$$P(x_i = y_i | d(x, y) = k) = \frac{n-k}{n} \quad (2)$$

$$P(d(x, y) = k | d(x, y) \leq r) = \frac{\binom{n}{k}}{\sum_{j=0}^r \binom{n}{j}} \quad (3)$$

Hence,

$$P(x_i = y_i | d(x, y) \leq r) = \frac{\sum_{k=0}^r \frac{n-k}{n} \binom{n}{k}}{\sum_{j=0}^r \binom{n}{j}} \quad (4)$$

Finally, the combinatorial identity

$$\frac{n-k}{n} \binom{n}{k} = \frac{(n-k)}{n} \frac{n!}{(n-k)!k!} = \frac{(n-1)!}{k!(n-1-k)!} = \binom{n-1}{k} \quad (5)$$

closes the theorem. \square

Theorem 1.1 is valid for both “x written at x” (autoassociative memory) and “random written at x” (heteroassociative memory). When $n = 1,000$ and $r = 451$, $P(x_i = y_i | d(x, y) \leq r) = p = 0.552905498137$. Each bit of a hard location does indeed have a small pull bias. What is meant by this is that each particular dimension has a small preference toward positive values if its address bit is set to 1, and negative values if set to 0.

So far we have looked only at a single pair of bitstrings, the probability of a single bitmatch between bitstrings within the access radius distance. Now let us consider the number of activated hard locations exhibiting this bitmatch.

Let h be the number of activated hard locations. As the probability of activating a specific hard location is constant, $h \sim \text{Binomial}(H, p_1)$. Thus, $\mathbf{E}[h] = \mu_h = Hp_1$ and $\mathbf{V}[h] = \sigma_h^2 = Hp_1(1 - p_1)$, where $p_1 = 2^{-n} \sum_{k=0}^r \binom{n}{k}$.

Let Z be the number of activated hard locations with the same bit as the reading address. Then, $Z = \sum_{i=1}^h X_i$, where $X_i \sim \text{Bernoulli}(p)$, where $p = P(x_i = y_i | d(x, y) \leq r)$.

Theorem 1.2. *Given a reading address x and a dimension i , the number of activated hard-locations with bitmatches at i follows a normal distribution with $\mathbf{E}[Z] = \mu_Z = p\mu_h$ and $\mathbf{V}[Z] = \sigma_Z^2 = p(1 - p)\mu_h + p^2\sigma_h^2$.*

Proof. As $P(973 < h < 1170) = 0.997$, by the central limit theorem, Z may be approximated by a normal distribution.

By the central limit theorem, Z is normally distributed.

Applying the law of total averages and the law of total variance, $\mathbf{E}[Z] = \mathbf{E}[\mathbf{E}[Z|h]] = \mathbf{E}[ph] = p\mathbf{E}[h] = p\mu_h$, and $\mathbf{V}[Z] = \mathbf{E}[\mathbf{V}[Z|h]] + \mathbf{V}[\mathbf{E}[Z|h]] = \mathbf{E}[hp(1 - p)] + \mathbf{V}[ph] = p(1 - p)\mathbf{E}[h] + p^2\mathbf{V}[h] = hp(1 - p) + p^2Hp_1(1 - p_1)$.

Applying the law of total variance, $\mathbf{V}[Z] = \mathbf{E}[\mathbf{V}[Z|h]] + \mathbf{V}[\mathbf{E}[Z|h]] = \mathbf{E}[hp(1 - p)] + \mathbf{V}[ph] = p(1 - p)\mathbf{E}[h] + p^2\mathbf{V}[h] = p(1 - p)\mu_h + p^2\sigma_h^2$. \square

See Figure 1 for a comparison between the theoretical model and a simulation.

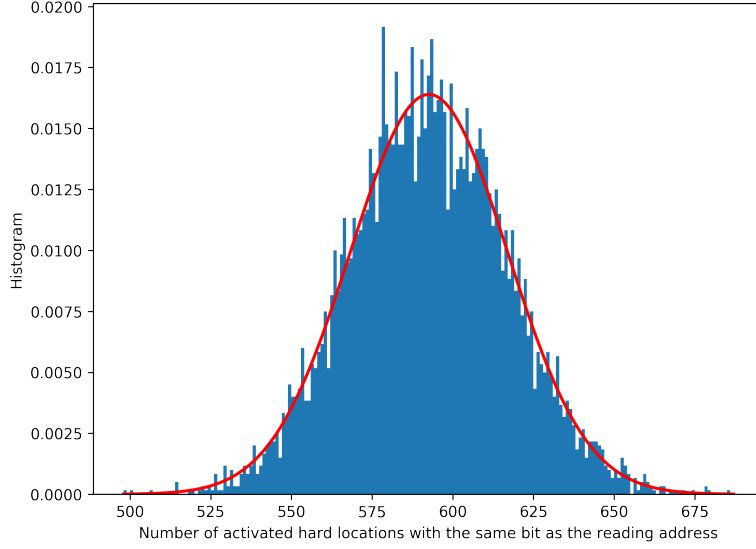


Figure 1: Given an address x and a dimension i , how many hard locations with bitmatches in i are activated by reading at x ? The histogram was obtained through numerical simulation. The red curve is the theoretical normal distribution found in Theorem 1.2.

2 Counter bias

The bias begins in the counters. Let's analyze the i th counter of a hard location.

Let s be the number of bitstrings written into memory (in our case, $s = 10,000$) and addr_i be the i th bit of the hard location's address.

Let θ be the average number of bitstrings written in each hard location. As there are s bitstrings written into the memory, and the probability of activating a specific hard location is constant, $\theta \sim \text{Binomial}(s, p_1)$. Thus, $\mathbf{E}[\theta] = \mu_\theta = sp_1$ and $\mathbf{V}[\theta] = \sigma_\theta^2 = sp_1(1 - p_1)$.

Let Y_i be the number of bitmatches in the i bit of a hard location's address after s written bitstrings. Then, $Y_i = \sum_{k=1}^{\theta} X_k$.

Theorem 2.1. *Giving the number of written bitstrings s , $\mathbf{E}[Y_i] = \mu_Y = p\mu_\theta$ and $\mathbf{V}[Y_i] = \sigma_Y^2 = p(1 - p)\mu_\theta + p^2\sigma_\theta^2$.*

Proof. Applying the law of total expectation, $\mathbf{E}[Y] = \mathbf{E}[\mathbf{E}[Y|\theta]] = \mathbf{E}[p\theta] = p\mathbf{E}[\theta] = p\mu_\theta$.

Applying the law of total variance, $\mathbf{V}[Y] = \mathbf{E}[\mathbf{V}[Y|\theta]] + \mathbf{V}[\mathbf{E}[Y|\theta]] = \mathbf{E}[\theta p(1-p)] + \mathbf{V}[p\theta] = p(1-p)\mathbf{E}[\theta] + p^2\mathbf{V}[\theta] = p(1-p)\mu_\theta + p^2\sigma_\theta^2$. \square

During a write operation, the counters are incremented for every bit 1 and decremented for every bit 0. So, after s writes, there will be θ bitstrings written in each hard location with Y_i bitmatches and $\theta - Y_i$ non-bitmatches. Thus, $[\text{cnt}_i | \text{addr}_i = 1] = (Y_i) - (\theta - Y_i) = 2Y_i - \theta$ and $[\text{cnt}_i | \text{addr}_i = 0] = \theta - 2Y_i$.

Theorem 2.2. $\mathbf{E}[\text{cnt}_i | \text{addr}_i = 1] = \mu_{\text{cnt}} = (2p - 1)\mu_\theta$ and $\mathbf{V}[\text{cnt}_i | \text{addr}_i = 1] = \sigma_{\text{cnt}}^2 = 4p(1-p)\mu_\theta + (2p - 1)^2\sigma_\theta^2$.

Proof. $\mathbf{E}[\text{cnt}_i | \text{addr}_i = 1] = \mathbf{E}[2Y_i - \theta] = \mathbf{E}[2Y_i] - \mathbf{E}[\theta] = 2\mathbf{E}[Y_i] - \mu_\theta = 2p\mu_\theta - \mu_\theta = (2p - 1)\mu_\theta$.

Applying the law of total variance, $\mathbf{V}[\text{cnt}_i | \text{addr}_i = 1] = \mathbf{V}[2Y_i - \theta] = \mathbf{E}[\mathbf{V}[2Y_i - \theta | \theta]] + \mathbf{V}[\mathbf{E}[2Y_i - \theta | \theta]]$.

Let us solve each part independently. Thus,

$$\mathbf{V}[2Y_i - \theta | \theta] = \mathbf{V}[2Y_i | \theta] = 4\mathbf{V}[Y_i | \theta] = 4\mathbf{V}[\sum_{k=1}^{\theta} X_k] = 4\theta p(1-p).$$

$$\mathbf{E}[\mathbf{V}[2Y_i - \theta | \theta]] = \mathbf{E}[4\theta p(1-p)] = 4p(1-p)\mathbf{E}[\theta] = 4p(1-p)\mu_\theta.$$

Finally,

$$\mathbf{E}[2Y_i - \theta | \theta] = 2\mathbf{E}[Y_i | \theta] - \mathbf{E}[\theta | \theta] = 2p\theta - \theta = (2p - 1)\theta.$$

$$\mathbf{V}[\mathbf{E}[2Y_i - \theta | \theta]] = \mathbf{V}[(2p - 1)\theta] = (2p - 1)^2\mathbf{V}[\theta] = (2p - 1)^2\sigma_\theta^2. \quad \square$$

Theorem 2.3. $\mathbf{E}[\text{cnt}_i | \text{addr}_i = 0] = -\mu_{\text{cnt}}$ and $\mathbf{V}[\text{cnt}_i | \text{addr}_i = 1] = \sigma_{\text{cnt}}^2$.

Proof. Notice that $[\text{cnt}_i | \text{addr}_i = 0] = -[\text{cnt}_i | \text{addr}_i = 1]$. Thus, $\mathbf{E}[\text{cnt}_i | \text{addr}_i = 0] = -\mathbf{E}[\text{cnt}_i | \text{addr}_i = 1]$ and $\mathbf{V}[\text{cnt}_i | \text{addr}_i = 0] = \mathbf{V}[\text{cnt}_i | \text{addr}_i = 1]$. \square

In summary,

$$[\text{cnt}_i | \text{addr}_i = 1] \sim \mathcal{N}(\mu_{\text{cnt}}, \sigma_{\text{cnt}}^2) \quad (6)$$

$$[\text{cnt}_i | \text{addr}_i = 0] \sim \mathcal{N}(-\mu_{\text{cnt}}, \sigma_{\text{cnt}}^2) \quad (7)$$

In our case, $p = 0.5529$, $s = 10,000$, and $H = 1,000,000$, so $[\text{cnt}_i | \text{addr}_i = 1] \sim \mathcal{N}(\mu = 1.1341, \sigma^2 = 10.7184)$. For “random at x”, $p = 0.5$, so $\mu = 0$ and $\sigma^2 = 10.7185$. See Figure 2.

Finally,

$$P(\text{cnt}_i > 0 | \text{addr}_i = 1) = P(\text{cnt}_i < 0 | \text{addr}_i = 0) = 1 - \mathcal{N}.\text{cdf}(0) \quad (8)$$

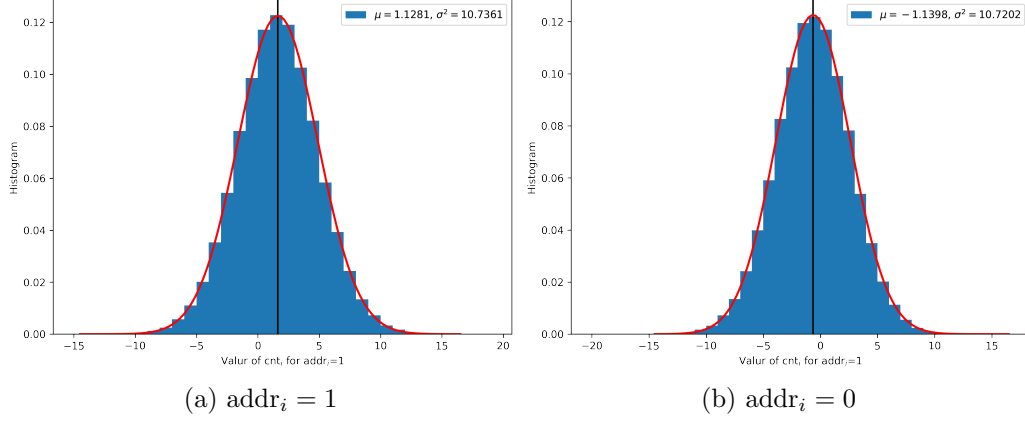


Figure 2: The value of the counters after $s = 10,000$ writes shows the autocorrelation in the counters in autoassociative memories (“x at x”). The histogram was obtained through simulation. The red curve is the theoretical normal distribution found in equations (6) and (7).

For “random written at x”, $p = 0.5$ implies $\mu_{\text{cnt}} = 0$, which implies $P(\text{cnt}_i > 0 | \text{addr}_i = 1) = P(\text{cnt}_i < 0 | \text{addr}_i = 0) = 0.5$, independently of the parameters because they will only affect the variance and the normal distribution is symmetrical around the average.

However, for “x written at x”, $p = 0.5529$ and the probabilities depend on s . For $s = 10,000$, they are equal to 0.6354. For $s = 20,000$, they are equal to 0.6867. For $s = 30,000$, they are equal to 0.7232. The more random bitstrings are written into the memory, the more the hard locations point to themselves.

See Figure 3 and notice that I still have to figure out why the mean is correct, but the standard deviation is not. As each of the n counters of a hard location may be equal or not with the same probability, I assumed it would follow a Binomial distribution (and it worked for “random at x”).

3 Read bias

Now that we know the distribution of $\text{cnt}_i | \text{addr}_i$, we may go to the read operation. During the read operation, on average, h hard locations are activated and their counters are summed up. So, for the i th bit,

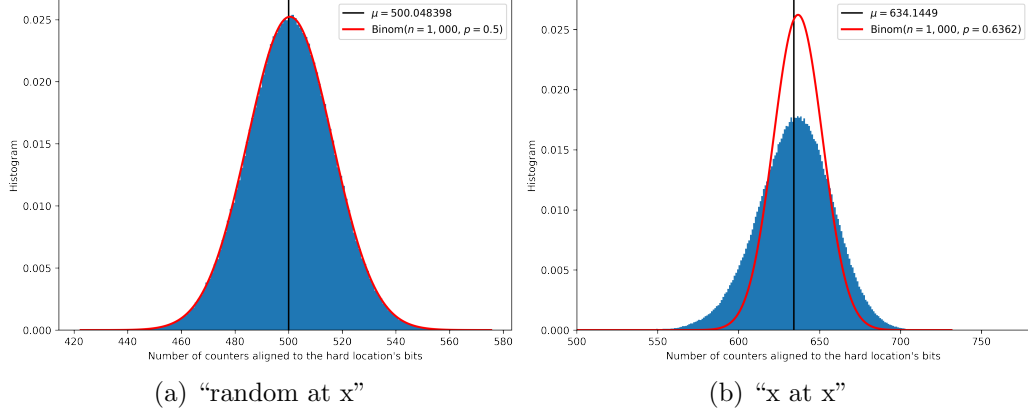


Figure 3: Autocorrelation in the counters in autoassociative memories (“x written at x”). The histogram was obtained through simulation. The red curve is the theoretical distribution.

$$\text{acc}_i = \sum_{k=1}^h \text{cnt}_k \quad (9)$$

Let η be the reading address and η_i the i th bit of it. Then, let’s split the h activated hard locations into two groups: (i) the ones with the same bit as η_i with ph hard locations, and (ii) the ones with the opposite bit as η_i with $(1-p)h$ hard locations.

$$[\text{acc}_i | \eta_i] = \sum_{k=1}^{ph} [\text{cnt}_k | \text{addr}_k = \eta_i] + \sum_{k=1}^{(1-p)h} [\text{cnt}_k | \text{addr}_k \neq \eta_i] \quad (10)$$

Each sum is a sum of normally distributed random variables, so

$$\sum_{k=1}^{ph} [\text{cnt}_k | \text{addr}_k = \eta_1] \sim \mathcal{N}(\mu_3 = \mu_2 ph, \sigma_3^2 = \sigma_2^2 ph + \mu_2^2 hp(1-p)) \quad (11)$$

$$\sum_{k=1}^{(1-p)h} [\text{cnt}_k | \text{addr}_k \neq \eta_1] \sim \mathcal{N}(\mu_3 = -\mu_2(1-p)h, \sigma_3^2 = \sigma_2^2(1-p)h + \mu_2^2 hp(1-p)) \quad (12)$$

In our case, $\sum_{k=1}^{ph} [\text{cnt}_k | \text{addr}_k = 1] \sim \mathcal{N}(\mu = 672.12, \sigma^2 = 6281.00)$, and $\sum_{k=1}^{ph} [\text{cnt}_k | \text{addr}_k = 1] \sim \mathcal{N}(\mu = -543.49, \sigma^2 = 5078.99)$. See Figure 4 — we can notice that the average is correct but the variance is too small.

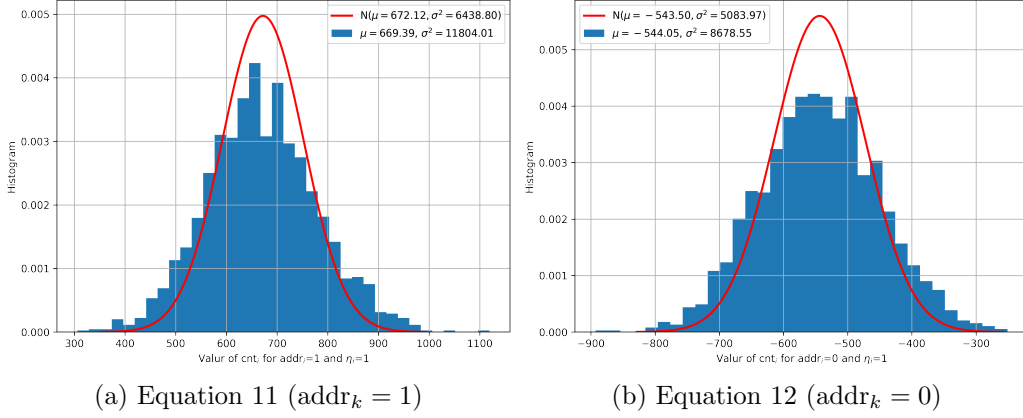


Figure 4: The histogram was obtained through simulation. The red curve is the theoretical normal distribution.

Hence,

$$[\text{acc}_i | \eta_i = 1] \sim \mathcal{N}(\mu = (2p - 1)^2 \theta h, \sigma^2 = \sigma_2^2 h + 2\mu_2^2 h p(1 - p)) \quad (13)$$

$$[\text{acc}_i | \eta_i = 0] \sim \mathcal{N}(\mu = -(2p - 1)^2 \theta h, \sigma^2 = \sigma_2^2 h + 2\mu_2^2 h p(1 - p)) \quad (14)$$

In our case, $[\text{acc}_i | \eta_i = 1] \sim \mathcal{N}(\mu = 128.62, \sigma^2 = 12181.95)$, and $[\text{acc}_i | \eta_i = 0] \sim \mathcal{N}(\mu = -128.62, \sigma^2 = 12181.95)$. See Figure 5 — we can notice that the variance issue from Figure 4 has propagated to these images.

Finally,

$$P(\text{wrong}) = P(\text{acc}_i < 0 | \eta_i = 1) \cdot P(\eta_i = 1) + P(\text{acc}_i > 0 | \eta_i = 0) \cdot P(\eta_i = 0) \quad (15)$$

$$= \frac{\mathcal{N}_{\eta_i=1}.\text{cdf}(0)}{2} + \frac{1 - \mathcal{N}_{\eta_i=0}.\text{cdf}(0)}{2} \quad (16)$$

$$= \frac{\mathcal{N}_{\eta_i=1}.\text{cdf}(0)}{2} + \frac{\mathcal{N}_{\eta_i=1}.\text{cdf}(0)}{2} \quad (17)$$

$$= \mathcal{N}_{\eta_i=1}.\text{cdf}(0) \quad (18)$$

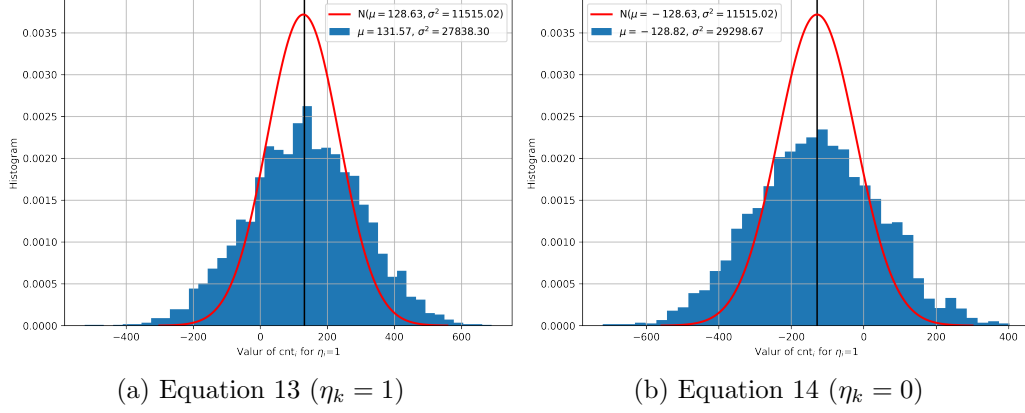


Figure 5: The histogram was obtained through simulation. The red curve is the theoretical normal distribution.

Using the empirical variance of $\sigma^2 = 27838.3029124$, we calculate $P(wrong) = 0.22037771219874325$.

In order to check this probability, I have run a simulation reading from 1,000 random bitstrings (which have never been written into memory) and calculate the distance from the result of a single read. As the $P(wrong) = 0.22037$, I expected to get an average distance of 220.37 with a standard deviation of 13.10. See Figure 6 for the comparison between the simulated and the theoretical outcomes.

Figure 7 shows the new distance between η_d and $\text{read}(\eta_d)$, where η_d is d bits away from η . As for $d \geq 520$ there is no intersection between η and η_d , our models applies and explains the horizontal line around distance 220.

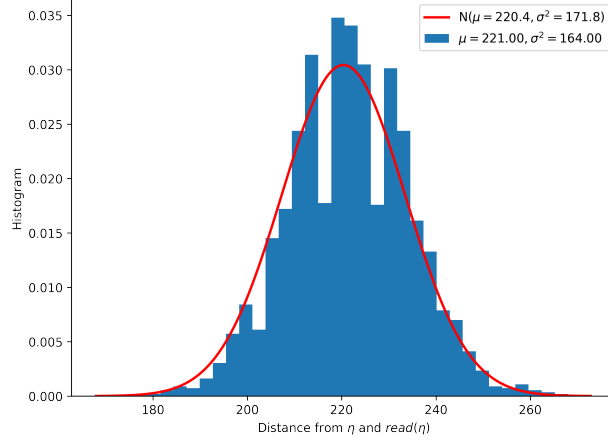


Figure 6: The histogram was obtained through simulation. The red curve is the theoretical normal distribution.

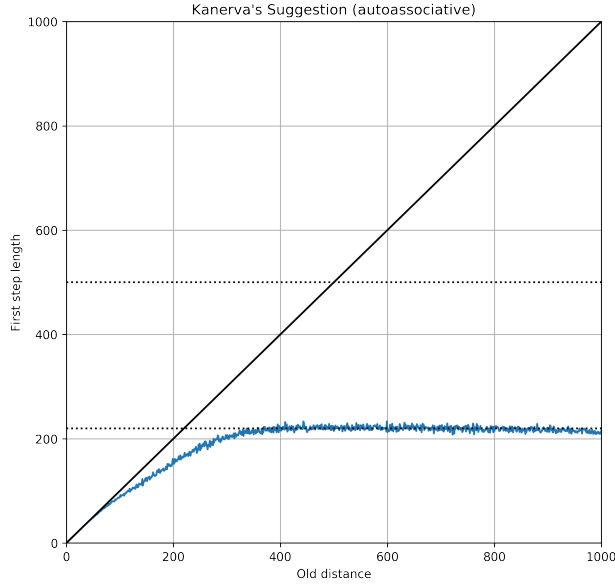


Figure 7: New distance after a single read operation in a bitstring η_d , which is d bits away from η . The new distance was calculated between η_d and $\text{read}(\eta_d)$. Notice that when $d \geq 520$, the intersection between η and η_d is zero, which means there is only random bitstrings written into the activated hard locations. The distance 220 equals $1000 \cdot 0.220$ which is the probability find in Figure 6.