

Supplemental Material for “Individualized non-uniform quantization for vector search”

A FAST LOGARITHMS AND EXPONENTIATION

Even if the expressions of F_{KS} and F_{KS}^{-1} are mathematically straightforward, the use of fundamental functions makes their computation challenging. To speedup these computations, we use the identity

$$x^c = \exp(c \log x), \quad (25)$$

enabling the use of fast approximations of the common functions \exp and \log for $x \in \mathbb{R}_+$.

For the exponential and logarithm, we use an implementation based on the minimax polynomial approximation.¹ We have modified the Remez polynomials for an optimal tradeoff between speed and accuracy using <https://github.com/DKenefake/OptimalPoly/>. We reproduce the code for completeness.

```
1 #include <cmath>
2 #include <cstdint>
3 #include <cstring>
4
5 float int_as_float (int32_t a) {
6     float r;
7     memcpy (&r, &a, sizeof r);
8     return r;
9 }
10 int32_t float_as_int (float a) {
11     int32_t r;
12     memcpy (&r, &a, sizeof r);
13     return r;
14 }
15
16 float fast_logf (float a) {
17     float m, r, s, t, i, f;
18     int32_t e;
19
20     e = (float_as_int(a) - 0x3f2aaaab) & 0xff800000;
21     m = int_as_float(float_as_int(a) - e);
22     i = (float) e * 1.19209290e-7f;
23     f = m - 1.0f;
24     s = f * f;
25     r = fmaf(0.230836749f, f, -0.279208571f);
26     t = fmaf(0.331826031f, f, -0.498910338f);
```

¹<https://stackoverflow.com/a/39822314> and <https://stackoverflow.com/a/47025627>

```

27     r = fmaf(r, s, t);
28     r = fmaf(r, s, f);
29     r = fmaf(i, 0.693147182f, r);
30     return r;
31 }
32
33 float fast_exp(float x) {
34     float invlog2e = 1.442695041f; // 1 / log2(e)
35     float expCvt = 12582912.0f; // 1.5 * (1 << 23)
36
37     /* exp(x) = 2^i * 2^f; i = rint (log2(e) * x), -0.5 <= f <= 0.5 */
38     float t = x * invlog2e; // t = x / log2(e)
39     float r = (t + expCvt) - expCvt; // r = round(t)
40     float f = t - r; // f = t - round(t)
41     int i = (int) r; // i = (int) r
42
43     float temp = fmaf(f, 0.009651907610706037f, 0.05593479631997887f);
44     temp = fmaf(temp, f, 0.2402301551437674f);
45     temp = fmaf(temp, f, 0.6931186232012877f);
46     temp = fmaf(temp, f, 0.9999993887682104f);
47
48     temp = int_as_float(float_as_int(temp) + (i << 23)); // temp = temp * 2^i
49     return temp;
50 }
```

B FAST NQT LOGISTIC AND LOGIT FUNCTIONS

```

1 #include <cmath>
2 #include <cstdint>
3 #include <cstring>
4
5 float int_as_float (int32_t a) {
6     float r;
7     memcpy (&r, &a, sizeof r);
8     return r;
9 }
10 int32_t float_as_int (float a) {
11     int32_t r;
12     memcpy (&r, &a, sizeof r);
13     return r;
14 }
15
```

```
16 float logisticNQT(float value, float alpha, float x0) {
17     // The value -alpha * x0 can be precomputed
18     float temp = fmaf(value, alpha, -alpha * x0);
19     int p = round(temp + 0.5f);
20     int m = float_as_int(fmaf(temp - p, 0.5f, 1f));
21
22     temp = int_as_float(m + (p << 23)); // temp = m * 2^p
23     return temp / (temp + 1f);
24 }
25
26 float logitNQT(float value, float inverseAlpha, float x0) {
27     float z = value / (1f - value);
28
29     int temp = float_as_int(z);
30     int e = temp & 0x7f800000;
31     float p = (float) ((e >> 23) - 128);
32     float m = int_as_float((temp & 0x007fffff) + 0x3f800000);
33
34     return fmaf(m + p, inverseAlpha, x0);
35 }
```

C ADDITIONAL EXPERIMENTAL RESULTS

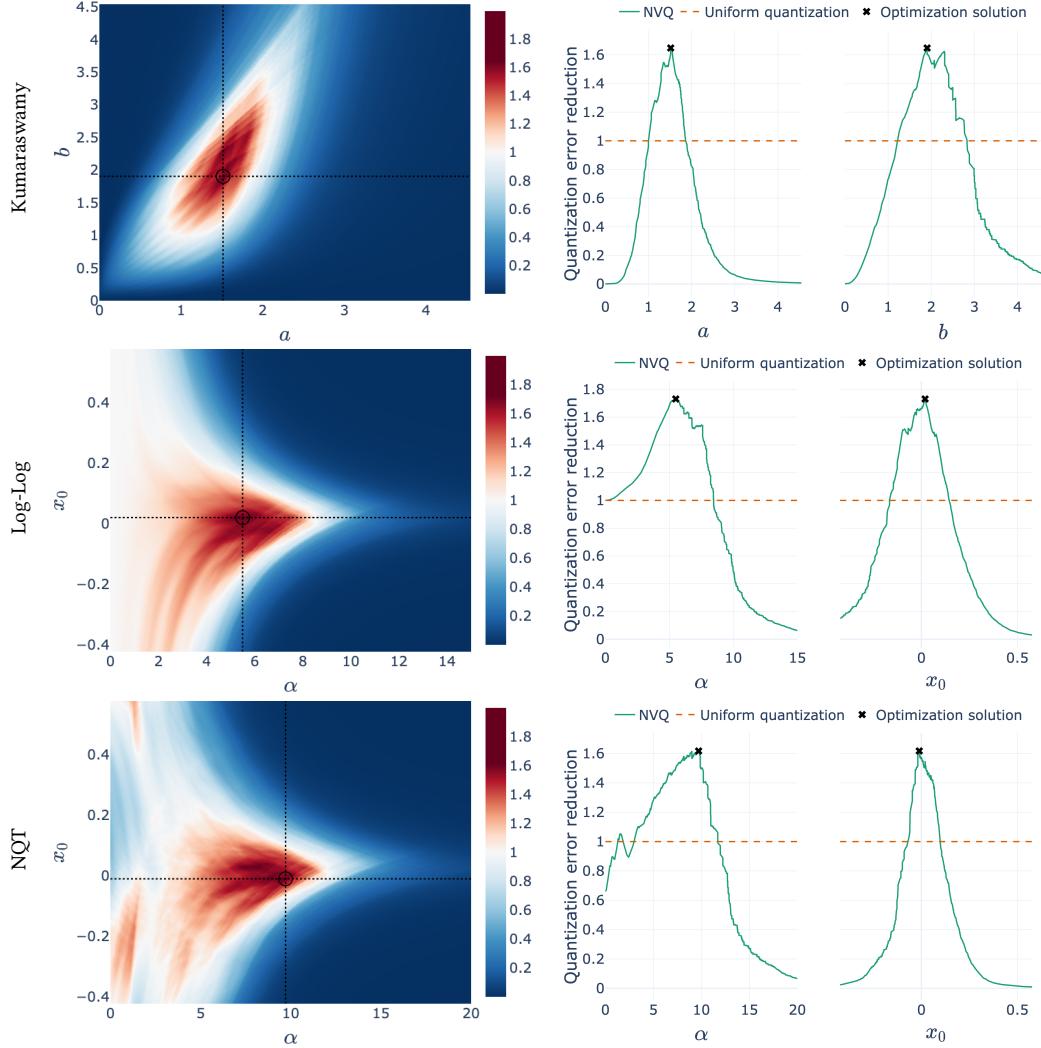


Fig. 12. Algorithm 1 finds a good maximum of Problem 6 using the nonlinearities presented in Section 4 for $\beta = 4$ bits. **Left:** the landscape of the objective function in Problem 6, which is the ratio of the MSE improvement over the uniform quantization (a value of 1 means parity, higher is better), as a function of the nonlinearity parameters for one vector from ada002-100k (see Table 2). The solution found by Algorithm1 is marked by a black circle. **Right:** two cross cuts taken at the values corresponding to the found solution.

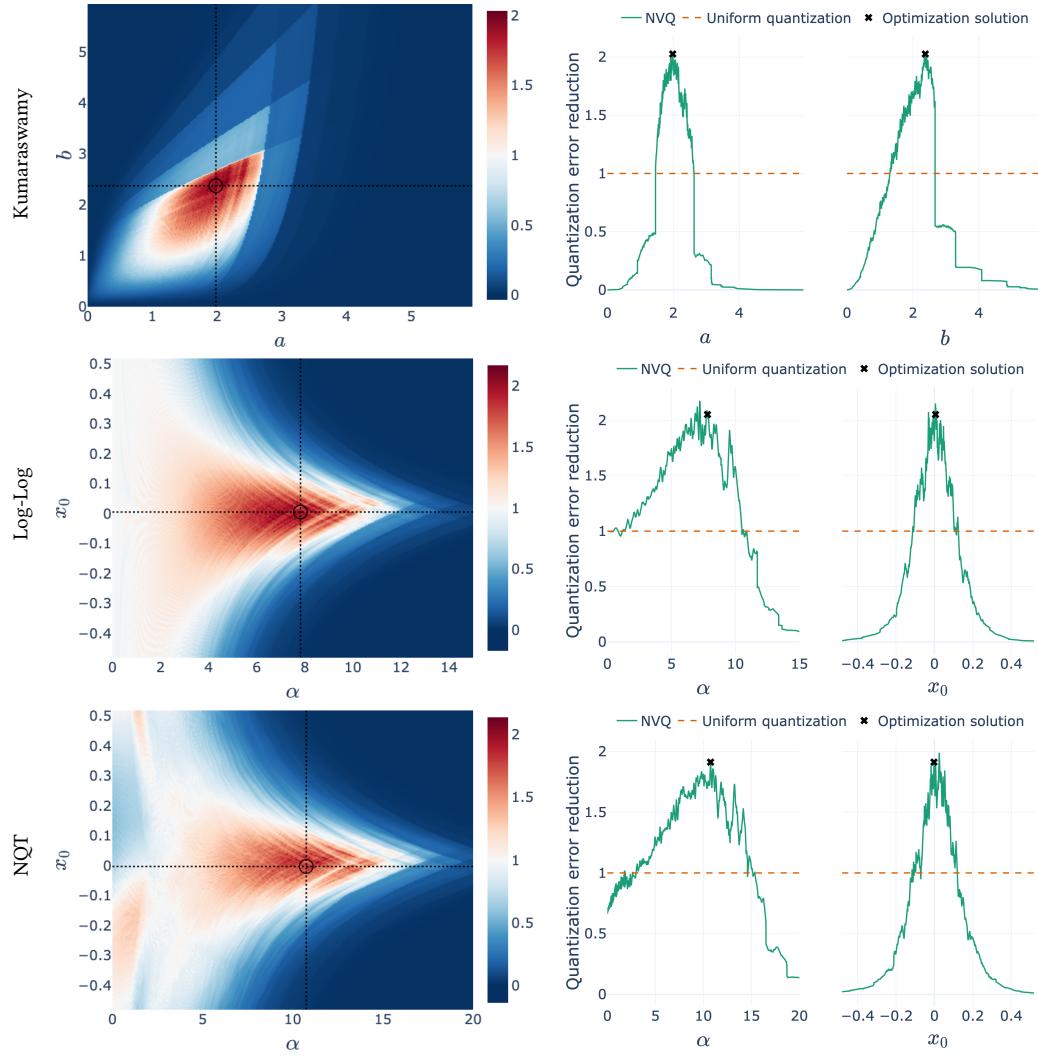


Fig. 13. Algorithm 1 finds a good maximum of Problem 6 using the nonlinearities presented in Section 4 for $\beta = 8$ bits. **Left:** the landscape of the objective function in Problem 6, which is the ratio of the MSE improvement over the uniform quantization (a value of 1 means parity, higher is better), as a function of the nonlinearity parameters for one vector from openai-v3-100k . The solution found by Algorithm 1 is marked by a black circle. **Right:** two cross cuts taken at the values corresponding to the found solution.

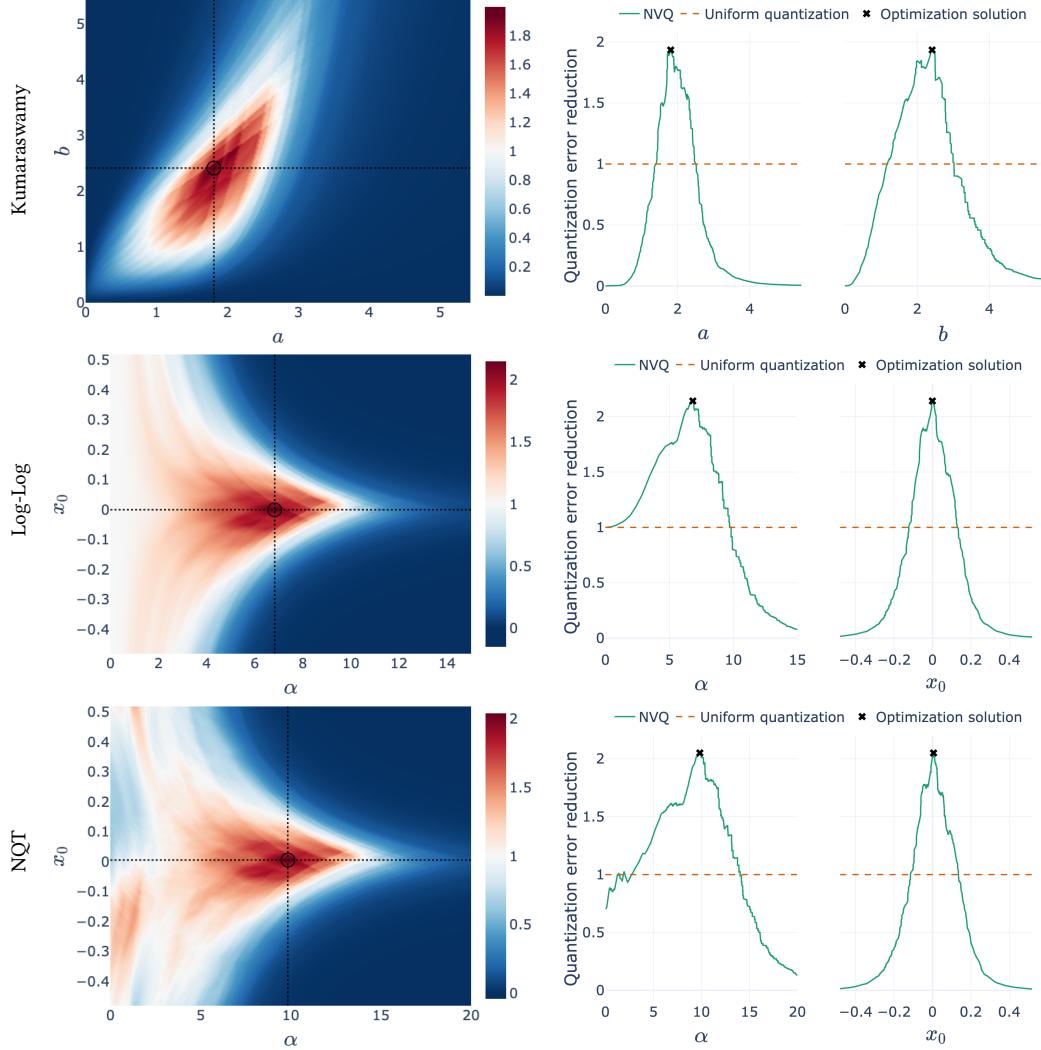


Fig. 14. Algorithm 1 finds a good maximum of Problem 6 using the nonlinearities presented in Section 4 for $\beta = 4$ bits. **Left:** the landscape of the objective function in Problem 6, which is the ratio of the MSE improvement over the uniform quantization (a value of 1 means parity, higher is better), as a function of the nonlinearity parameters for one vector from openai-v3-100k . The solution found by Algorithm 1 is marked by a black circle. **Right:** two cross cuts taken at the values corresponding to the found solution.

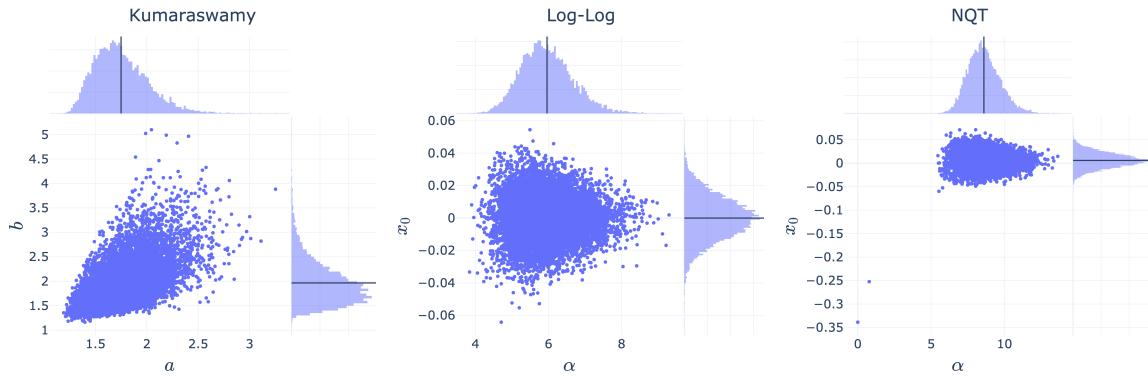


Fig. 15. The solutions for 10^4 vectors from gecko-100k. Different nonlinearity parameters are chosen for each vector.

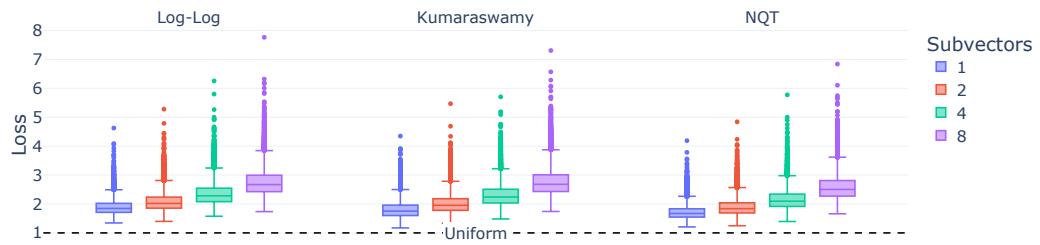


Fig. 16. Increasing the number of subvectors (depicted here for $\beta = 8$) increases the objective function (higher is better) in Problem 6 for 10^4 vectors from ada002-100k.

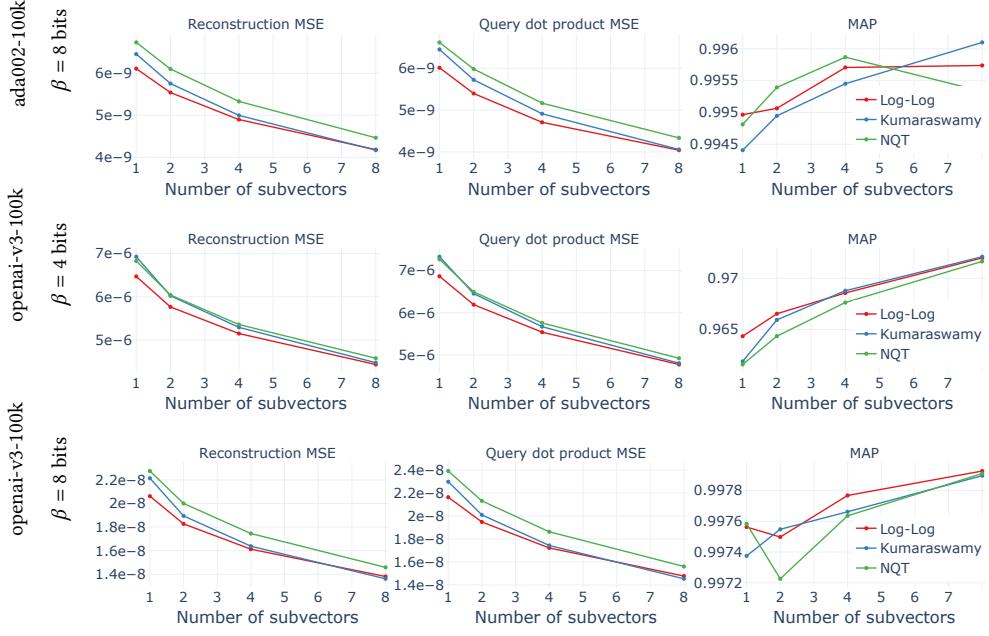


Fig. 17. Increasing the number of subvectors lowers the average reconstruction error and the query dot product error, defined as $\sum_{x \in \mathcal{X}} (\langle q, x \rangle - \langle q, \bar{x} \rangle)^2$, and increases the mean average precision (MAP) in openai-v3-100k.