

Article

A Modification of the Fast Inverse Square Root Algorithm

Cezary J. Walczyk¹, Leonid V. Moroz² and Jan L. Cieśliński^{1,*}

¹ Uniwersytet w Białymstoku, Wydział Fizyki, ul. Ciołkowskiego 1L, 15-245 Białystok, Poland; c.walczyk@uwb.edu.pl, j.cieslinski@uwb.edu.pl

² Lviv Polytechnic National University, Department of Security Information and Technology, st. Kn. Romana 1/3, 79000 Lviv, Ukraine; moroz_lv@polynet.lviv.ua

Abstract: We present an improved algorithm for fast calculation of the inverse square root for single-precision floating-point numbers. The algorithm is much more accurate than the famous fast inverse square root algorithm and has a similar computational cost. The presented modification concerns Newton-Raphson corrections and can be applied when the distribution of these corrections is not symmetric (for instance, in our case they are always negative).

Keywords: floating-point arithmetic; inverse square root; magic constant; Newton-Raphson method

1. Introduction

Floating-point arithmetic has become widely used in many applications such as 3D graphics, scientific computing and signal processing [1–5], implemented both in hardware and software [6–10]. Many algorithms can be used to approximate elementary functions [1,2,10–18]. The inverse square root function is of particular importance because it is widely used in 3D computer graphics, especially in lightning reflections [19–21], and has many other applications, see [22–36]. All of these algorithms require an initial seed to start the approximation. The more accurate is the initial seed, the fewer iterations are needed. Usually, the initial seed is obtained from a look-up table (LUT) which is memory consuming.

In this paper we consider an algorithm for computing the inverse square root using the so called magic constant instead of a LUT [37–40]. The following code realizes the fast inverse square root algorithm in the case of single-precision IEEE Standard 754 floating-point numbers (type `float`).

```
1. float InvSqrt(float x){
2.     float halfnumber = 0.5f*x;
3.     int i = *(int*) &x;
4.     i = R - (i >> 1);
5.     y = *(float*) &i;
6.     y = y*(1.5f - halfnumber*y*y);
7.     y = y*(1.5f - halfnumber*y*y);
8.     return y;
9. }
```

The code *InvSqrt* consists of two main parts. Lines 4 and 5 produce in a very cheap way a quite good zeroth approximation of the inverse square root of a given positive floating-point number x . Lines 6 and 7 apply the Newton-Raphson corrections twice (often a version with just one iteration is used, as well). Originally R was proposed as $0x5F3759DF$, see [37,38].

22 *InvSqrt* is characterized by a high speed, more that 3 times higher than in computing the inverse
 23 square root using library functions. This property is discussed in detail in [41]. The errors of the fast
 24 inverse square root algorithm depend on the choice of the “magic constant” R . In several theoretical
 25 papers [38,41–44] (see also Eberly’s monograph [19]) attempts were made to determine analytically the
 26 optimal value of the magic constant (i.e., to minimize errors). In general, this optimal value can depend
 27 on the number of iterations, which is a general phenomenon [45]. The derivation and comprehensive
 28 mathematical description of all steps of the fast inverse square root algorithm is given in our recent paper
 29 [46]. We found the optimum value of the magic constant by minimizing the final maximum relative error.

30 In the present paper we develop our analytical approach to construct an improved algorithm
 31 (*InvSqrt1*) for fast computing of the inverse square root, see section 4. In both codes, *InvSqrt* and *InvSqrt1*,
 32 magic constants serve as a low-cost way of generating a reasonably accurate first approximation of the
 33 inverse square root. These magic constants turn out to be the same. The main novelty of the new algorithm
 34 is in the second part of the code which is changed significantly. In fact, we propose a modification of the
 35 Newton-Raphson formulae which has a similar computational cost but improve the accuracy even by
 36 several times.

37 2. Analytical approach to algorithm *InvSqrt*

38 In this paper we confine ourselves to positive floating-point numbers

$$x = (1 + m_x)2^{e_x} \quad (2.1)$$

39 where $m_x \in [0, 1)$ and e_x is an integer (note that this formula does not hold for subnormal numbers). In
 40 the case of the IEEE-754 standard, a floating-point number is encoded by 32 bits. The first bit corresponds
 41 to a sign (in our case this bit is simply equal to zero), the next 8 bits correspond to an exponent e_x and the
 42 last 23 bits encodes a mantissa m_x . The integer encoded by these 32 bits, denoted by I_x , is given by

$$I_x = N_m(B + e_x + m_x) \quad (2.2)$$

43 where $N_m = 2^{23}$ and $B = 127$ (thus $B + e_x = 1, 2, \dots, 254$). The lines 3 and 5 of the *InvSqrt* code interpret
 44 a number as an integer (2.2) or float (2.1), respectively. The lines 4, 6 and 7 of the code can be written as

$$I_{y_0} = R - \lfloor I_x/2 \rfloor, \quad y_1 = \frac{1}{2}y_0(3 - y_0^2x), \quad y_2 = \frac{1}{2}y_1(3 - y_1^2x). \quad (2.3)$$

45 The first equation produces, in a surprisingly simple way, a good zeroth approximation y_0 of the inverse
 46 square root $y = 1/\sqrt{x}$. Of course, this needs a very special form of R . In particular, in the single precision
 47 case we have $e_R = 63$, see [46]. The next equations can be easily recognized as the Newton-Raphson
 48 corrections. We point out that the code *InvSqrt* is invariant with respect to the scaling

$$x \rightarrow \tilde{x} = 2^{-2n}x, \quad y_k \rightarrow \tilde{y}_k = 2^n y_k \quad (k = 0, 1, 2), \quad (2.4)$$

49 like the equality $y = 1/\sqrt{x}$ itself. Therefore, without loss of the generality, we can confine our analysis to
 50 the interval

$$\tilde{A} := [1, 4). \quad (2.5)$$

51 The tilde will denote quantities defined on this interval. In [46] we have shown that the function \tilde{y}_0
 52 defined by the first equation of (2.3) can be approximated with a very good accuracy by the piece-wise
 53 linear function \tilde{y}_{00} given by

$$\tilde{y}_{00}(\tilde{x}, t) = \begin{cases} -\frac{1}{4}\tilde{x} + \frac{3}{4} + \frac{1}{8}t & \text{for } \tilde{x} \in [1, 2) \\ -\frac{1}{8}\tilde{x} + \frac{1}{2} + \frac{1}{8}t & \text{for } \tilde{x} \in [2, t) \\ -\frac{1}{16}\tilde{x} + \frac{1}{2} + \frac{1}{16}t & \text{for } \tilde{x} \in [t, 4) \end{cases} \quad (2.6)$$

where

$$t = 2 + 4m_R + 2N_m^{-1}, \quad (2.7)$$

54 and $m_R := N_m^{-1}R - \lfloor N_m^{-1}R \rfloor$ (m_R is the mantissa of the floating-point number corresponding to R). Note
55 that the parameter t , defined by (2.7), is uniquely determined by R .

56 The only difference between y_0 produced by the code *InvSqrt* and y_{00} given by (2.6) is the definition
57 of t , because t related to the code depends (although in a negligible way) on x . Namely,

$$|\tilde{y}_{00} - \tilde{y}_0| \leq \frac{1}{4}N_m^{-1} = 2^{-25} \approx 2.98 \cdot 10^{-8}. \quad (2.8)$$

58 Taking into account the invariance (2.4), we obtain

$$\left| \frac{y_{00} - y_0}{y_0} \right| \leq 2^{-24} \approx 5.96 \cdot 10^{-8}. \quad (2.9)$$

59 These estimates do not depend on t (in other words, they do not depend on R). The relative error of the
60 zeroth approximation (2.6) is given by

$$\tilde{\delta}_0(\tilde{x}, t) = \sqrt{\tilde{x}} \tilde{y}_{00}(\tilde{x}, t) - 1 \quad (2.10)$$

61 This is a continuous function with local maxima at

$$\tilde{x}_0^I = (6 + t)/6, \quad \tilde{x}_0^{II} = (4 + t)/3, \quad \tilde{x}_0^{III} = (8 + t)/3, \quad (2.11)$$

62 given respectively by

$$\begin{aligned} \tilde{\delta}_0(\tilde{x}_0^I, t) &= -1 + \frac{1}{2} \left(1 + \frac{t}{6}\right)^{3/2}, \\ \tilde{\delta}_0(\tilde{x}_0^{II}, t) &= -1 + 2 \left(\frac{1}{3} \left(1 + \frac{t}{4}\right)\right)^{3/2}, \\ \tilde{\delta}_0(\tilde{x}_0^{III}, t) &= -1 + \left(\frac{2}{3} \left(1 + \frac{t}{8}\right)\right)^{3/2}. \end{aligned} \quad (2.12)$$

63 In order to study global extrema of $\tilde{\delta}_0(\tilde{x}, t)$ we need also boundary values:

$$\tilde{\delta}_0(1, t) = \tilde{\delta}_0(4, t) = \frac{1}{8}(t - 4), \quad \tilde{\delta}_0(2, t) = \frac{\sqrt{2}}{4} \left(1 + \frac{t}{2}\right) - 1, \quad \tilde{\delta}_0(t, t) = \frac{\sqrt{t}}{2} - 1, \quad (2.13)$$

64 which are, in fact, local minima. Taking into account

$$\tilde{\delta}_0(1, t) - \tilde{\delta}_0(t, t) = \frac{1}{8}(\sqrt{t} - 2)^2 \geq 0, \quad \tilde{\delta}_0(2, t) - \tilde{\delta}_0(t, t) = \frac{\sqrt{2}}{8}(\sqrt{t} - \sqrt{2})^2 \geq 0, \quad (2.14)$$

65 we conclude that

$$\min_{\tilde{x} \in \tilde{A}} \tilde{\delta}_0(\tilde{x}, t) = \tilde{\delta}_0(t, t) < 0. \quad (2.15)$$

66 Because $\tilde{\delta}_0(\tilde{x}_0^{III}, t) < 0$ for $t \in (2, 4)$, the global maximum is one of the remaining local maxima:

$$\max_{\tilde{x} \in \tilde{A}} \tilde{\delta}_0(\tilde{x}, t) = \max\{\tilde{\delta}_0(\tilde{x}_0^I, t), \tilde{\delta}_0(\tilde{x}_0^{II}, t)\}. \quad (2.16)$$

67 Therefore,

$$\max_{\tilde{x} \in \tilde{A}} |\tilde{\delta}_0(\tilde{x}, t)| = \max\{|\tilde{\delta}_0(t, t)|, \tilde{\delta}_0(\tilde{x}_0^I, t), \tilde{\delta}_0(\tilde{x}_0^{II}, t)\}. \quad (2.17)$$

68 In order to minimize this value with respect to t , i.e., to find t_0^r such that

$$\max_{\tilde{x} \in \tilde{A}} |\tilde{\delta}_0(\tilde{x}, t_0^r)| < \max_{\tilde{x} \in \tilde{A}} |\tilde{\delta}_0(\tilde{x}, t)| \quad \text{for } t \neq t_0^r, \quad (2.18)$$

69 we observe that $|\tilde{\delta}_0(t, t)|$ is a decreasing function of t , while both maxima $(\tilde{\delta}_0(\tilde{x}_0^I, t))$ and $\tilde{\delta}_0(\tilde{x}_0^{II}, t)$ are
70 increasing functions. Therefore, it is sufficient to find $t = t_0^I$ and $t = t_0^{II}$ such that

$$|\tilde{\delta}_0(t_0^I, t_0^I)| = \tilde{\delta}_0(\tilde{x}_0^I, t_0^I), \quad |\tilde{\delta}_0(t_0^{II}, t_0^{II})| = \tilde{\delta}_0(\tilde{x}_0^{II}, t_0^{II}), \quad (2.19)$$

and to choose the greater of these two values. In [46] we have shown that

$$|\tilde{\delta}_0(t_0^I, t_0^I)| < |\tilde{\delta}_0(t_0^{II}, t_0^{II})|. \quad (2.20)$$

71 Therefore $t_0^r = t_0^{II}$ and

$$\tilde{\delta}_{0\max} := \min_{t \in (2, 4)} \left(\max_{\tilde{x} \in \tilde{A}} |\tilde{\delta}_0(\tilde{x}, t)| \right) = |\tilde{\delta}_0(t_0^r, t_0^r)|. \quad (2.21)$$

72 The following numerical values result from these calculations [46]:

$$t_0^r \approx 3.7309796, \quad R_0 = 0x5F37642F, \quad \tilde{\delta}_{0\max} \approx 0.03421281. \quad (2.22)$$

73 Newton-Raphson corrections for the zeroth approximation (\tilde{y}_{00}) will be denoted by \tilde{y}_{0k} ($k = 1, 2, \dots$). In
74 particular, we have:

$$\begin{aligned} \tilde{y}_{01}(\tilde{x}, t) &= \frac{1}{2} \tilde{y}_{00}(\tilde{x}, t) (3 - \tilde{y}_{00}^2(\tilde{x}, t) \tilde{x}), \\ \tilde{y}_{02}(\tilde{x}, t) &= \frac{1}{2} \tilde{y}_{01}(\tilde{x}, t) (3 - \tilde{y}_{01}^2(\tilde{x}, t) \tilde{x}). \end{aligned} \quad (2.23)$$

75 and the corresponding relative error functions will be denoted by $\tilde{\delta}_k(\tilde{x}, t)$:

$$\tilde{\delta}_k(\tilde{x}, t) := \frac{\tilde{y}_{0k}(\tilde{x}, t) - \tilde{y}}{\tilde{y}} = \sqrt{\tilde{x}} \tilde{y}_{0k}(\tilde{x}, t) - 1, \quad (k = 0, 1, 2, \dots), \quad (2.24)$$

76 where we included also the case $k = 0$, see (2.10). The obtained approximations of the inverse square
77 root depend on the parameter t directly related to the magic constant R . The value of this parameter can
78 be estimated by analysing the relative error of $\tilde{y}_{0k}(\tilde{x}, t)$ with respect to $1/\sqrt{\tilde{x}}$. As the best estimation we
79 consider $t = t_k^{(r)}$ minimizing the relative error $\tilde{\delta}_k(\tilde{x}, t)$:

$$\forall_{t \neq t_k^{(r)}} \left(\tilde{\delta}_{k \max} \equiv \max_{\tilde{x} \in \tilde{A}} |\tilde{\delta}_k(\tilde{x}, t_k^{(r)})| < \max_{\tilde{x} \in \tilde{A}} |\tilde{\delta}_k(\tilde{x}, t)| \right). \quad (2.25)$$

80 We point out that in general the optimum value of the magic constant can depend on the number of
81 Newton-Raphson corrections. Calculations carried out in [46] gave the following results:

$$\begin{aligned} t_1^r = t_2^r = 3.7298003, \quad R_1^r = R_2^r = 0x5F375A86, \\ \tilde{\delta}_{1 \max} \approx 1.75118 \cdot 10^{-3}, \quad \tilde{\delta}_{2 \max} \approx 4.60 \cdot 10^{-6}. \end{aligned} \quad (2.26)$$

82 We omit details of the computations except one important point. Using (2.24) for expressing \tilde{y}_{0k} by $\tilde{\delta}_k$ and
83 $\sqrt{\tilde{x}}$ we can rewrite (2.23) as

$$\tilde{\delta}_k(\tilde{x}, t) = -\frac{1}{2} \tilde{\delta}_{k-1}^2(\tilde{x}, t) (3 + \tilde{\delta}_{k-1}(\tilde{x}, t)), \quad (k = 1, 2, \dots). \quad (2.27)$$

84 The quadratic dependence on $\tilde{\delta}_{k-1}$ means that every Newton-Raphson correction improves the accuracy
85 by several orders of magnitude (until the machine precision is reached), compare (2.26).

86 The formula (2.27) suggests another way of improving the accuracy because the functions $\tilde{\delta}_k$ are
87 always non-positive for any $k \geq 1$. Roughly saying, we are going to shift the graph of $\tilde{\delta}_k$ upwards by an
88 appropriate modification of the Newton-Raphson formula. In the next section we describe the general
89 idea of this modification.

90 3. Modified Newton-Raphson formulas

91 The formula (2.27) shows that errors introduced by Newton-Raphson corrections are nonpositive,
92 i.e., they take values in intervals $[-\tilde{\delta}_{k \max}, 0]$, where $k = 1, 2, \dots$. Therefore, it is natural to introduce a
93 correction term into the Newton-Raphson formulas (2.23). We expect that the corrections will be roughly
94 half of the maximal relative error. Instead of the maximal error we introduce two parameters, d_1 and d_2 .
95 Thus we get modified Newton-Raphson formulas:

$$\begin{aligned} \tilde{y}_{11}(\tilde{x}, t, d_1) &= 2^{-1} \tilde{y}_{00}(\tilde{x}, t) (3 - \tilde{y}_{00}^2(\tilde{x}, t) \tilde{x}) + \frac{d_1}{2\sqrt{\tilde{x}}}, \\ \tilde{y}_{12}(\tilde{x}, t, d_1, d_2) &= 2^{-1} \tilde{y}_{11}(\tilde{x}, t, d_1) (3 - \tilde{y}_{11}^2(\tilde{x}, t, d_1) \tilde{x}) + \frac{d_2}{2\sqrt{\tilde{x}}}, \end{aligned} \quad (3.1)$$

96 where zeroth approximation is assumed in the form (2.6). In the following section the term $1/\sqrt{\tilde{x}}$ will be
97 replaced by some approximations of \tilde{y} , transforming (3.1) into a computer code. In order to estimate a
98 possible gain in accuracy, in this section we temporarily assume that \tilde{y} is the exact value of the inverse
99 square root. The corresponding error functions,

$$\tilde{\delta}_k''(\tilde{x}, t, d_1, \dots, d_k) = \sqrt{\tilde{x}} \tilde{y}_{1k}(\tilde{x}, t, d_1, \dots, d_k) - 1, \quad k \in \{0, 1, 2, \dots\}, \quad (3.2)$$

100 (where $\tilde{y}_{10}(\tilde{x}, t) := \tilde{y}_{00}(\tilde{x}, t)$), satisfy

$$\tilde{\delta}_k'' = -\frac{1}{2} \tilde{\delta}_{k-1}''^2 (3 + \tilde{\delta}_{k-1}'') + \frac{d_k}{2}, \quad (3.3)$$

101 where: $\tilde{\delta}_0''(\tilde{x}, t) = \tilde{\delta}_0(\tilde{x}, t)$. Note that

$$\tilde{\delta}_1''(\tilde{x}, t, d_1) = \tilde{\delta}_1(\tilde{x}, t) + \frac{1}{2} d_1. \quad (3.4)$$

102 In order to simplify notation we usually will suppress the explicit dependence on d_j . We will write, for
 103 instance, $\tilde{\delta}_2''(\tilde{x}, t)$ instead of $\tilde{\delta}_2''(\tilde{x}, t, d_1, d_2)$.

104 The corrections of the form (3.1) will decrease relative errors in comparison with the results of earlier
 105 papers [38,46]. We have 3 free parameters (d_1, d_2 and t) to be determined by minimizing the maximal
 106 error (in principle the new parameterization can give a new estimation of the parameter t). By analogy to
 107 (2.25), we are going to find $t = t^{(0)}$ minimizing the error of the first correction (2.25):

$$\forall_{t \neq t^{(0)}} \max_{\tilde{x} \in \tilde{A}} |\tilde{\delta}_1''(\tilde{x}, t^{(0)})| < \max_{\tilde{x} \in \tilde{A}} |\tilde{\delta}_1''(\tilde{x}, t)|, \quad (3.5)$$

108 where, as usual, $\tilde{A} = [1, 4]$.

109 The first of Eqs. (3.3) implies that for any t the maximal value of $\tilde{\delta}_1''(\tilde{x}, t)$ equals $\frac{1}{2}d_1$ and is attained at
 110 zeros of $\tilde{\delta}_0''(\tilde{x}, t)$. Using results of section 2, including (2.15), (2.16), (2.20) and (2.21), we conclude that the
 111 minimum value of $\tilde{\delta}_1''(\tilde{x}, t)$ is attained either for $\tilde{x} = t$ or for $\tilde{x} = x_0^{II}$ (where there is the second maximum
 112 of $\tilde{\delta}_0''(\tilde{x}, t)$), i.e.,

$$\min_{\tilde{x} \in \tilde{A}} \tilde{\delta}_1''(\tilde{x}, t) = \min \left\{ \tilde{\delta}_1''(t, t), \tilde{\delta}_1''(x_0^{II}, t) \right\} \quad (3.6)$$

113 Minimization of $|\tilde{\delta}_1''(\tilde{x}, t)|$ can be done with respect to t and with respect to d_1 (these operations obviously
 114 commute), which corresponds to

$$\underbrace{\max_{\tilde{x} \in \tilde{A}} \tilde{\delta}_1''(\tilde{x}, t^{(0)})}_{\tilde{\delta}_{1 \max}''} = - \min_{\tilde{x} \in \tilde{A}} \tilde{\delta}_1''(\tilde{x}, t^{(0)}). \quad (3.7)$$

115 Taking into account

$$\max_{\tilde{x} \in \tilde{A}} \tilde{\delta}_1''(\tilde{x}, t^{(0)}) = \frac{d_1}{2}, \quad \min_{\tilde{x} \in \tilde{A}} \tilde{\delta}_1''(\tilde{x}, t^{(0)}) = \tilde{\delta}_1''(t^{(0)}, t^{(0)}) = -\tilde{\delta}_{1 \max} + \frac{d_1}{2}, \quad (3.8)$$

116 we get from (3.7):

$$\tilde{\delta}_{1 \max}'' = \frac{1}{2}d_1 = \frac{1}{2}\tilde{\delta}_{1 \max} \simeq 8.7559 \cdot 10^{-4}, \quad (3.9)$$

117 where

$$\tilde{\delta}_{1 \max} := \min_{t \in (2, 4)} \left(\max_{x \in \tilde{A}} |\tilde{\delta}_1(\tilde{x}, t)| \right). \quad (3.10)$$

118 and the numerical value of $\tilde{\delta}_{1 \max}$ is given by (2.26). These conditions are satisfied for

$$t^{(0)} = t_1^{(r)} \simeq 3.7298003. \quad (3.11)$$

119 In order to minimize the relative error of the second correction we use equation analogous to (3.7):

$$\underbrace{\max_{\tilde{x} \in \tilde{A}} \tilde{\delta}_2''(\tilde{x}, t^{(0)})}_{\tilde{\delta}_{2 \max}''} = - \min_{\tilde{x} \in \tilde{A}} \tilde{\delta}_2''(\tilde{x}, t^{(0)}), \quad (3.12)$$

120 where from (3.3) we have

$$\max_{\tilde{x} \in \tilde{A}} \delta_2''(\tilde{x}, t^{(0)}) = \frac{d_2}{2}, \quad \min_{\tilde{x} \in \tilde{A}} \delta_2''(\tilde{x}, t^{(0)}) = -\frac{1}{2} \delta_{1 \max}''^2 \left(3 + \delta_{1 \max}''\right) + \frac{d_2}{2}. \quad (3.13)$$

121 Hence

$$\delta_{2 \max}'' = \frac{1}{4} \delta_{1 \max}''^2 \left(3 + \delta_{1 \max}''\right). \quad (3.14)$$

122 Expressing this result in terms of formerly computed $\tilde{\delta}_{1 \max}$ and $\tilde{\delta}_{2 \max}$, we obtain

$$\delta_{2 \max}'' = \frac{1}{8} \tilde{\delta}_{2 \max} + \frac{3}{32} \tilde{\delta}_{1 \max}^3 \simeq 5.75164 \cdot 10^{-7} \simeq \frac{\tilde{\delta}_{2 \max}}{7.99}, \quad (3.15)$$

where

$$\tilde{\delta}_{2 \max} = \frac{1}{2} \tilde{\delta}_{1 \max}^2 (3 - \tilde{\delta}_{1 \max}).$$

123 Therefore, the above modification of Newton-Raphson formulas decreases the relative error 2 times after
124 one iteration and almost 8 times after two iterations as compared to the standard *InvSqrt* algorithm.

125 In order to implement this idea in the form of a computer code, we have to replace the unknown
126 $1/\sqrt{\tilde{x}}$ (i.e., \tilde{y}) on the right-hand sides of (3.1) by some numerical approximations.

127 4. New algorithm of higher accuracy

128 Approximating $1/\sqrt{\tilde{x}}$ in formulas (3.1) by values at left hand sides, we transform (3.1) into

$$\begin{aligned} \tilde{y}_{21} &= \frac{1}{2} \tilde{y}_{20} (3 - \tilde{y}_{20}^2 \tilde{x}) + \frac{1}{2} d_1 \tilde{y}_{21}, \\ \tilde{y}_{22} &= \frac{1}{2} \tilde{y}_{21} (3 - \tilde{y}_{21}^2 \tilde{x}) + \frac{1}{2} d_2 \tilde{y}_{22}, \end{aligned} \quad (4.1)$$

129 where \tilde{y}_{2k} ($k = 1, 2, \dots$) depend on \tilde{x}, t and d_j (for $1 \leq j \leq k$). We assume $\tilde{y}_{20} \equiv \tilde{y}_{00}$, i.e., the zeroth
130 approximation is still given by (2.6). We can see that \tilde{y}_{21} and \tilde{y}_{22} can be explicitly expressed by \tilde{y}_{20} and
131 \tilde{y}_{21} , respectively.

132 Parameters d_1 and d_2 have to be determined by minimization of the maximum error. We define error
133 functions in the usual way:

$$\Delta_k^{(1)} = \frac{\tilde{y}_{2k} - \tilde{y}}{\tilde{y}} = \sqrt{\tilde{x}} \tilde{y}_{2k} - 1. \quad (4.2)$$

134 Substituting (4.2) into (4.1) we get:

$$\Delta_1^{(1)}(\tilde{x}, t, d_1) = \frac{d_1}{2 - d_1} - \frac{1}{2 - d_1} \tilde{\delta}_0^2(\tilde{x}, t) (3 + \tilde{\delta}_0(\tilde{x}, t)) = \frac{d_1 + 2\tilde{\delta}_1(\tilde{x}, t)}{2 - d_1}, \quad (4.3)$$

$$\Delta_2^{(1)}(\tilde{x}, t, d_2) = \frac{d_2}{2 - d_2} - \frac{1}{2 - d_2} \left(\Delta_1^{(1)}(\tilde{x}, t, d_1)\right)^2 \left(3 + \Delta_1^{(1)}(\tilde{x}, t, d_1)\right). \quad (4.4)$$

135 The equation (4.3) expresses $\Delta_1^{(1)}(\tilde{x}, t, d_1)$ as a linear function of the nonpositive function $\tilde{\delta}_1(\tilde{x}, t)$ with
136 coefficients depending on the parameter d_1 . The optimum parameters t and d_1 will be estimated by the
137 procedure described in section 3. First, we minimize the amplitude of the relative error function, i.e., we
138 find $t^{(1)}$ such that

$$\max_{\tilde{x} \in \tilde{A}} \Delta_1^{(1)}(\tilde{x}, t^{(1)}) - \min_{\tilde{x} \in \tilde{A}} \Delta_1^{(1)}(\tilde{x}, t^{(1)}) \leq \max_{\tilde{x} \in \tilde{A}} \Delta_1^{(1)}(\tilde{x}, t) - \min_{\tilde{x} \in \tilde{A}} \Delta_1^{(1)}(\tilde{x}, t) \quad (4.5)$$

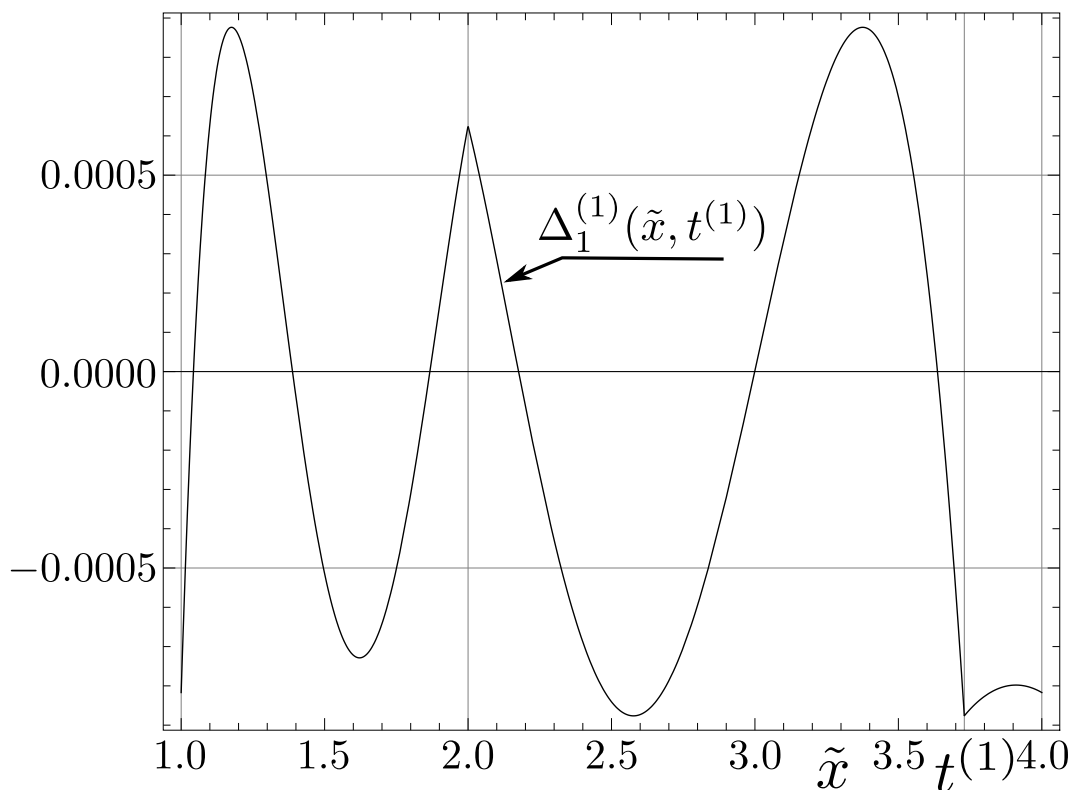


Figure 1. Graph of the function $\Delta_1^{(1)}(\tilde{x}, t^{(1)})$.

139 for all $t \neq t^{(1)}$. Second, we determine $d_1^{(1)}$ such that

$$\max_{\tilde{x} \in \tilde{A}} \Delta_1^{(1)}(\tilde{x}, t^{(1)}, d_1^{(1)}) = -\min_{\tilde{x} \in \tilde{A}} \Delta_1^{(1)}(\tilde{x}, t^{(1)}, d_1^{(1)}). \quad (4.6)$$

140 Thus we have

$$\max_{\tilde{x} \in \tilde{A}} |\Delta_1^{(1)}(\tilde{x}, t^{(1)}, d_1^{(1)})| \leq \max_{\tilde{x} \in \tilde{A}} |\Delta_1^{(1)}(\tilde{x}, t, d_1)| \quad (4.7)$$

141 for all real d_1 and $t \in (2, 4)$. $\Delta_1^{(1)}(\tilde{x}, t)$ is an increasing function of $\tilde{\delta}_1(\tilde{x}, t)$, hence

$$-\frac{d_1^{(1)} - 2 \max_{\tilde{x} \in \tilde{A}} |\tilde{\delta}_1(\tilde{x}, t_1^{(1)})|}{2 - d_1^{(1)}} = \frac{d_1^{(1)}}{2 - d_1^{(1)}}, \quad (4.8)$$

142 which is satisfied for

$$d_1^{(1)} = \max_{\tilde{x} \in \tilde{A}} |\tilde{\delta}_1(\tilde{x}, t_1^{(1)})| = \tilde{\delta}_1^{\max}. \quad (4.9)$$

143 Thus we can find the maximum error of the first correction $\Delta_1^{(1)}(\tilde{x}, t_1^{(1)})$ (presented at Fig. 1):

$$\max_{\tilde{x} \in \tilde{A}} |\Delta_1^{(1)}(\tilde{x}, t^{(1)})| = \frac{\max_{\tilde{x} \in \tilde{A}} |\tilde{\delta}_1(\tilde{x}, t^{(1)})|}{2 - \max_{\tilde{x} \in \tilde{A}} |\tilde{\delta}_1(\tilde{x}, t^{(1)})|}, \quad (4.10)$$

144 which assumes the minimum value for $t^{(1)} = t_1^{(r)}$:

$$\Delta_{1\max}^{(1)} = \frac{\max_{\tilde{x} \in \tilde{A}} |\tilde{\delta}_1(\tilde{x}, t_1^{(r)})|}{2 - \max_{\tilde{x} \in \tilde{A}} |\tilde{\delta}_1(\tilde{x}, t_1^{(r)})|} = \frac{\tilde{\delta}_{1\max}}{2 - \tilde{\delta}_{1\max}} \simeq 8.7636 \cdot 10^{-4} \simeq \frac{\tilde{\delta}_{1\max}}{2.00}. \quad (4.11)$$

145 This result practically coincides with $\tilde{\delta}_{1\max}''$ given by (3.9).

146 Analogously we can determine the value of $d_2^{(1)}$ (assuming that $t = t^{(1)}$ is fixed):

$$d_2^{(1)} = \frac{d_2^{(1)} - \max_{\tilde{x} \in \tilde{A}} |\Delta_1^{(1)2}(\tilde{x}, t^{(1)})(3 + \Delta_1^{(1)}(\tilde{x}, t^{(1)}))|}{2 - d_2^{(1)}} = \frac{d_2^{(1)}}{2 - d_2^{(1)}}. \quad (4.12)$$

147 Now, the deepest minimum comes from the global maximum

$$\max_{\tilde{x} \in \tilde{A}} |\Delta_1^{(1)2}(\tilde{x}, t^{(1)})(3 + \Delta_1^{(1)}(\tilde{x}, t^{(1)}))| = \frac{2\tilde{\delta}_{1\max}^2(3 - \tilde{\delta}_{1\max})}{(2 - \tilde{\delta}_{1\max})^3}. \quad (4.13)$$

148 Therefore we get

$$d_2^{(1)} = \frac{\tilde{\delta}_{1\max}^2(3 - \tilde{\delta}_{1\max})}{(2 - \tilde{\delta}_{1\max})^3} \simeq 1.15234 \cdot 10^{-6}, \quad (4.14)$$

149 and the maximum error of the second correction is given by

$$\Delta_{2\max}^{(1)} = \frac{d_2^{(1)}}{2 - d_2^{(1)}} \simeq 5.76173 \cdot 10^{-7} \simeq \frac{\tilde{\delta}_{2\max}}{7.98}, \quad (4.15)$$

150 which is very close to the value of $\tilde{\delta}_{2\max}''$ given by (3.15).

151 Thus we have obtained the algorithm *InvSqrt1* which looks like *InvSqrt* with modified values of
152 numerical coefficients.

```

1. float InvSqrt1(float x){
2.     float simhalfnumber = 0.500438180f*x;
3.     int i = *(int*) &x;
4.     i = 0x5F375A86 - (i >> 1);
5.     y = *(float*) &i;
6.     y = y*(1.50131454f - simhalfnumber*y*y);
7.     y = y*(1.50000086f - 0.999124984f*simhalfnumber*y*y);
8.     return y;
9. }
```

153 Comparing *InvSqrt1* with *InvSqrt* we easily see that the number of algebraic operations in *InvSqrt1*
154 is greater by 1 (an additional multiplication in line 7, corresponding to the second iteration of the modified
155 Newton-Raphson procedure). We point out that magic constants for *InvSqrt* and *InvSqrt1* are the same.

156 5. Numerical experiments

157 The new algorithms were tested on the processor Intel Core i5-3470 using the compiler TDM-GCC
158 4.9.2 32-bit (then, in the case of *InvSqrt*, the values of errors are practically the same as those obtained by
159 Lomont [38]). The same results were obtained also on Intel i7-5700. In this section we analyze rounding
160 errors for the code *InvSqrt1*.

161 Applying algorithm *InvSqrt1* we obtain relative errors characterized by “oscillations” with a center
162 slightly shifted with respect to the analytical solution, see Fig. 2. Calculations were carried out for all

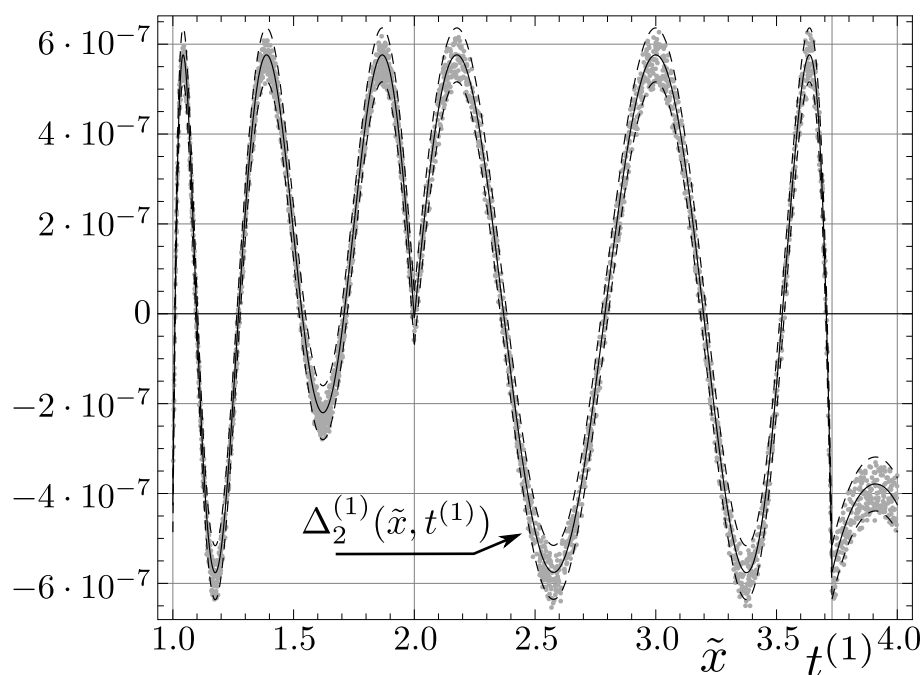


Figure 2. Solid lines represent function $\Delta_2^{(1)}(\tilde{x}, t^{(1)})$. Its vertical shifts by $\pm 6 \cdot 10^{-8}$ are denoted by dashed lines. Finally, dots represent relative errors for 4000 random values $x \in (2^{-126}, 2^{128})$ produced by algorithms *InvSqrt1*.

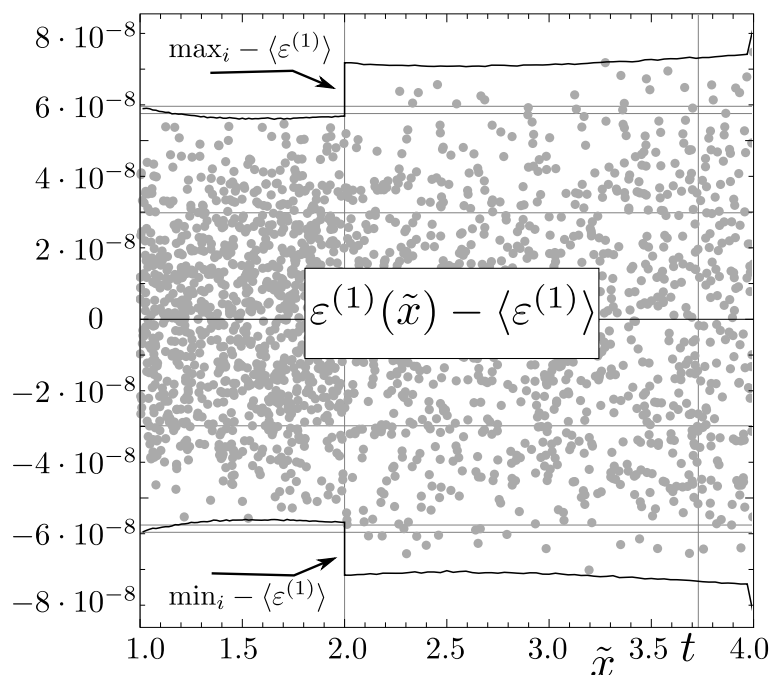


Figure 3. Relative error $\varepsilon^{(1)}$ arising during the **float** approximation of corrections $\tilde{y}_{22}(\tilde{x}, t)$. Dots represent errors determined for 2000 random values $\tilde{x} \in [1, 4)$. Solid lines represent maximum (\max_i) and minimum (\min_i) values of relative errors (intervals $[1, 2)$ and $[2, 4)$ were divided into 64 equal intervals, and then extremum values were determined in all these intervals).

163 numbers x of the type **float** such that $e_x \in [-126, 128)$. The range of errors is the same for all these
 164 intervals (except $e_x = -126$):

$$\Delta_{2,N}^{(1)}(x) = \text{sqrt}(x) * \text{InvSqrt1}(x) - 1. \in (\Delta_{2,N \min}^{(1)}, \Delta_{2,N \max}^{(1)}), \quad (5.1)$$

where

$$\Delta_{2,N \min}^{(1)} = -6.62 \cdot 10^{-7}, \quad \Delta_{2,N \max}^{(1)} = 6.35 \cdot 10^{-7}.$$

For $e_x = -126$ the interval of errors is slightly wider:

$$[-6.72 \cdot 10^{-7}, 6.49 \cdot 10^{-7}].$$

165 This can be explained by the fact that the analysis presented in this paper is not applicable to subnormals
 166 numbers, see (2.1). The observed blur can be noticed already for the approximation error of the correction
 167 $\tilde{y}_{22}(\tilde{x})$:

$$\varepsilon^{(1)}(\tilde{x}) = \frac{\text{InvSqrt1}(x) - \tilde{y}_{22}(\tilde{x}, t^{(1)})}{\tilde{y}_{22}(\tilde{x}, t^{(1)})}. \quad (5.2)$$

168 The values of this error are distributed symmetrically around the mean value $\langle \varepsilon^{(1)} \rangle$:

$$\langle \varepsilon^{(1)} \rangle = 2^{-1} N_m^{-1} \sum_{x \in [1,4)} \varepsilon^{(1)}(\tilde{x}) = -1.398 \cdot 10^{-8} \quad (5.3)$$

169 enclosing the range:

$$\varepsilon^{(1)}(\tilde{x}) \in [-9.676 \cdot 10^{-8}, 6.805 \cdot 10^{-8}], \quad (5.4)$$

170 see Fig. 3. The blur parameters of the function $\varepsilon^{(1)}(\tilde{x}, t)$ show that the main source of the difference
 171 between analytical and numerical results is the use of precision **float** and, in particular, rounding of
 172 constant parameters of the function *InvSqrt1*. It is worthwhile to point out that in this case the amplitude
 173 of the error oscillations is about 40% greater than the amplitude of oscillations of $(\tilde{y}_{00} - \tilde{y}_0)/\tilde{y}_0$ (i.e., in
 174 the case of *InvSqrt*), see the right part of Fig. 2 in [46].

175 6. Conclusions

176 In this paper we have presented a modification of the famous code *InvSqrt* for fast computation of
 177 the inverse square root. The new code has the same magic constant but the second part (which consists
 178 of Newton-Raphson iterations) is modified. In the case of one Newton-Raphson iteration the new code
 179 *InvSqrt1* has the same computational cost as *InvSqrt* and is 2 times more accurate. In the case of two
 180 iterations the computational cost of the new code is slightly higher but its accuracy is higher by 8 times.

181 The main idea of our work consists in modifying coefficients in the Newton-Raphson method and
 182 demanding that the maximal error is as small as possible. Such modifications can be constructed if the
 183 distribution of errors for Newton-Raphson corrections is not symmetric (like in the case of the inverse
 184 square root, when they are non-positive functions).

185 **Author Contributions:** Conceptualization, Leonid V. Moroz; Formal analysis, Cezary J. Walczyk; Investigation,
 186 Cezary J. Walczyk, Leonid V. Moroz and Jan L. Cieśliński; Methodology, Cezary J. Walczyk and Leonid V. Moroz;
 187 Visualization, Cezary J. Walczyk; Writing—original draft, Jan L. Cieśliński; Writing—review & editing, Jan L. Cieśliński

188 **Funding:** This research received no external funding.

189 **Conflicts of Interest:** The authors declare no conflict of interest.

190 **References**

- 191 1. M.D. Ercegovac, T. Lang: *Digital Arithmetic*, Morgan Kaufmann 2003.
- 192 2. B. Parhami: *Computer Arithmetic: Algorithms and Hardware Designs*, 2nd edition, Oxford Univ. Press, New York,
193 2010
- 194 3. K. Diefendorff, P. K. Dubey, R. Hochsprung, H. Scales: AltiVec extension to PowerPC accelerates media
195 processing, *IEEE Micro* 20 (2) (2000) 85-95.
- 196 4. D. Harris: A Powering Unit for an OpenGL Lighting Engine, *Proc. 35th Asilomar Conf. Singals, Systems, and*
197 *Computers* (2001), pp. 1641-1645.
- 198 5. M. Sadeghian, J. Stine: Optimized Low-Power Elementary Function Approximation for Chybyshhev series
199 Approximation, *46th Asilomar Conf. on Signal Systems and Computers*, 2012.
- 200 6. D. M. Russinoff: A Mechanically Checked Proof of Correctness of the AMD K5 Floating Point Square Root
201 Microcode, *Formal Methods in System Design* 14 (1) (1999) 75-125.
- 202 7. M.Cornea, C. Anderson, C. Tsen: Software Implementation of the IEEE 754R Decimal Floating-Point Arithmetic,
203 *Software and Data Technologies (Communications in Computer and Information Science*, vol. 10), Springer 2008, pp.
204 97-109.
- 205 8. J.-M. Muller, N. Brisebarre, F. Dinechin, C.-P. Jeannerod, V. Lefèvre, G. Melquiond, N. Revol, D.Stehlé, S.
206 Torres: Hardware Implementation of Floating-Point Arithmetic, *Handbook of Floating-Point Arithmetic* (2009), pp.
207 269-320.
- 208 9. J.-M. Muller, N. Brisebarre, F. Dinechin, C.-P. Jeannerod, V. Lefèvre, G. Melquiond, N. Revol, D.Stehlé, S. Torres:
209 Software Implementation of Floating-Point Arithmetic, *Handbook of Floating-Point Arithmetic* (2009), pp. 321-372.
- 210 10. T. Viitanen, P. Jääskeläinen, O. Esko, J. Takala: Simplified floating-point division and square root, *Proc. IEEE Int.*
211 *Conf. Acoustics Speech and Signal Process.*, pp. 2707-2711, May 26-31 2013.
- 212 11. M.D. Ercegovac, T. Lang: *Division and Square Root: Digit Recurrence Algorithms and Implementations*, Boston:
213 Kluwer Academic Publishers, 1994.
- 214 12. W. Liu, A. Nannarelli: Power Efficient Division and Square root Unit, *IEEE Trans. Comp.* 61 (8) (2012) 1059-1070.
- 215 13. L.X. Deng, J.S. An: A low latency High-throughput Elementary Function Generator based on Enhanced double
216 rotation CORDIC, *IEEE Symposium on Computer Applications and Communications (SCAC)*, 2014.
- 217 14. M. X. Nguyen, A. Dinh-Duc: Hardware-Based Algorithm for Sine and Cosine Computations using Fixed Point
218 Processor, *11th International Conf. on Electrical Engineering/Electronics Computer, Telecommuncations and Information*
219 *Technology*, IEEE 2014.
- 220 15. M. Cornea, Intel[®] AVX-512 Instructions and Their Use in the Implementation of Math Functions, Intel
221 Corporation 2015.
- 222 16. H. Jiang, S. Graillat, R. Barrio, C. Yang: Accurate, validated and fast evaluation of elementary symmetric
223 functions and its application, *Appl. Math. Computation* 273 (2016) 1160-1178.
- 224 17. A. Fog: Software optimization resources, Instruction tables: Lists of instruction latencies, throughputs and
225 micro-operation breakdowns for Intel, AMD and VIA CPUs, <http://www.agner.org/optimize/>
- 226 18. L. Moroz, W. Samotyy: Efficient floating-point division for digital signal processing application, *IEEE Signal*
227 *Processing Magazine* 36 (1) (2019) 159-163.
- 228 19. D.H. Eberly: *GPGPU Programming for Games and Science*, CRC Press 2015.
- 229 20. N. Ide, M. Hirano, Y. Endo, S. Yoshioka, H. Murakami, A. Kunimatsu, T. Sato, T. Kamei, T. Okada, M. Suzuoki:
230 2.44-GFLOPS 300-MHz Floating-Point Vector-Processing Unit for High-Performance 3D Graphics Computing,
231 *IEEE J. Solid-State Circuits* 35 (7) (2000) 1025-1033.
- 232 21. S. Oberman, G. Favor, F. Weber: AMD 3DNow! technology: architecture and implementations, *IEEE Micro* 19
233 (2) (1999) 37-48.
- 234 22. T.J. Kwon, J. Draper: Floating-point Division and Square root Implementation using a Taylor-Series Expansion
235 Algorithm with Reduced Look-Up Table, 51st Midwest Symposium on Circuits and Systems, 2008.
- 236 23. T.O. Hands, I. Griffiths, D.A. Marshall, G. Douglas: The fast inverse square root in scientific computing, *Journal*
237 *of Physics Special Topics* 10 (1) (2011) A2-1.

- 238 24. J. Blinn: Floating-point tricks, *IEEE Comput. Graphics Appl.* 17 (4) (1997) 80-84.
- 239 25. J. Janhunen: Programmable MIMO detectors, PhD thesis, University of Oulu, Tampere 2011.
- 240 26. J.L.V.M. Stanislaus, T. Mohsenin: High Performance Compressive Sensing Reconstruction Hardware with QRD
241 Process, IEEE International Symposium on Circuits and Systems (ISCAS'12), May 2012.
- 242 27. Q. Avril, V. Gouranton, B. Arnaldi: Fast Collision Culling in Large-Scale Environments Using GPU Mapping
243 Function, ACM Eurographics Parallel Graphics and Visualization, Cagliari, Italy (2012).
- 244 28. R. Schattschneider: Accurate high-resolution 3D surface reconstruction and localisation using a wide-angle flat
245 port underwater stereo camera, PhD thesis, University of Canterbury, Christchurch, New Zealand, 2014.
- 246 29. S. Zafar, R. Adapa: Hardware architecture design and mapping of "Fast Inverse Square Root's algorithm",
247 International Conference on Advances in Electrical Engineering (ICAEE), 2014, pp. 1-4.
- 248 30. T. Hänninen, J. Janhunen, M. Juntti: Novel detector implementations for 3G LTE downlink and uplink, *Analog.*
249 *Integr. Circ. Sig. Process.* 78 (2014) 645–655.
- 250 31. Z.Q. Li, Y. Chen, X.Y. Zeng: OFDM Synchronization implementation based on Chisel platform for 5G research,
251 *2015 IEEE 11th International Conference on ASIC (ASICON)*.
- 252 32. C.J. Hsu, J.L. Chen, L.G. Chen: An Efficient Hardware Implementation of HON4D Feature Extraction for
253 Real-time Action Recognition, 2015 IEEE International Symposium on Consumer Electronics (ISCE).
- 254 33. C.H. Hsieh, Y.F. Chiu, Y.H. Shen, T.S. Chu, Y.H. Huang: A UWB Radar Signal Processing Platform for Real-Time
255 Human Respiratory Feature Extraction Based on Four-Segment Linear Waveform Model, *IEEE Trans. Biomed.*
256 *Circ. Syst.* 10 (1) (2016) 219–230.
- 257 34. J.D. Lv, F. Wang, Z.H. Ma: Peach Fruit Recognition Method under Natural Environment, Eighth International
258 Conference on Digital Image Processing (ICDIP 2016), Proc. of SPIE Vol. 10033, edited by C.M.Falco, X.D.Jiang,
259 1003317 (29 August 2016).
- 260 35. D. Sangeetha, P. Deepa: Efficient Scale Invariant Human Detection using Histogram of Oriented Gradients for
261 IoT Services, 2017 30th International Conference on VLSI Design and 2017 16th International Conference on
262 Embedded Systems, p. 61–66, IEEE 2016.
- 263 36. J. Lin, Z.G. Xu, A. Nukada, N. Maruyama, S. Matsuoka: Optimizations of Two Compute-bound Scientific
264 Kernels on the SW26010 Many-core Processor, 46th International Conference on Parallel Processing, p. 432–441,
265 IEEE 2017.
- 266 37. id software, quake3-1.32b/code/game/q_math.c, Quake III Arena, 1999.
- 267 38. C. Lomont, Fast inverse square root, Purdue University, Tech. Rep., 2003. Available online:
268 <http://www.lomont.org/Math/Papers/2003/InvSqrt.pdf>.
- 269 39. H.S. Warren: *Hacker's delight*, second edition, Pearson Education 2013.
- 270 40. P. Martin: Eight Rooty Pieces, *Overload Journal* 135 (2016) 8–12.
- 271 41. M. Robertson: A Brief History of InvSqrt, Bachelor Thesis, Univ. of New Brunswick 2012.
- 272 42. B. Self: Efficiently Computing the Inverse Square Root Using Integer Operations. May 31, 2012.
- 273 43. C. McEniry: The Mathematics Behind the Fast Inverse Square Root Function Code, Tech. rep. 2007.
- 274 44. D. Eberly: An approximation for the Inverse Square Root Function, 2015, Available online:
275 <http://www.geometrictools.com/Documentation/ApproxInvSqrt.pdf>.
- 276 45. P. Kornerup, J.-M. Muller: Choosing starting values for certain Newton-Raphson iterations, *Theor. Comp. Sci.*
277 351 (2006) 101–110.
- 278 46. L. Moroz, C.J. Walczyk, A. Hrynchyshyn, V. Holimath, J.L. Cieśliński: Fast calculation of inverse square root
279 with the use of magic constant – analytical approach, *Appl. Math. Computation* 316 (2018) 245–255.