倍精度浮動小数点指数関数計算回路の設計

河瀬 朋範[†] 高木 直史^{††} 高木 一義^{††}

概要 指数関数計算は、科学技術計算の分野でしばしば現れる計算である。本報告では、テー ブル参照と多項式近似に基づく、指数関数計算回路の設計について検討する。入出力は IEEE754 標 準の倍精度浮動小数点基本フォーマットとし、誤差は 1 ulp (unit in the last place) とする。本報告で は、近似多項式の次数とテーブルの分割数を変えて、テーブルサイズとクロック・サイクル数を評 価する。また、この回路に適した積和演算器の構成を提案する。積和演算器をパイプライン化しな い場合と、2 段パイプライン化する場合に対して、回路構成を示す。

Design of Double Precision Floating Point Exponential Function Computing Circuit

TOMONORI KAWASE,[†] NAOFUMI TAKAGI^{††} and KAZUYOSHI TAKAGI^{††}

Abstract Exponential function often appears in the field of scientific computing. In this report, we discuss the design of exponential function computing circuit based on table look-up and polynomial approximation. The inputs and the outputs of the circuit are in IEEE-754 double-precision floating-point format. The final error is within 1 ulp (unit in the last place). We evaluated the total table size and clock cycles for several combinations of the degree of approximation polynomial and the number of table division. We propose the structure of multiplier-accumulator suitable for this circuit and show the circuit design with either unfolded or pipelined multiplier-accumulator.

1. はじめに

近年の集積回路技術の発展により、高速な並 列乗算器や除算器、開平器などがプロセッサや ASICに搭載されるようになってきている。今後、 さらにLSIの集積度が上がれば、より複雑な演 算を行う回路も搭載されるようになると考えら れる[1]。そのような演算の一つとして指数関数 計算が挙げられる。指数関数計算は、科学技術 計算の分野でしばしば現れる計算である [2]。

指数関数計算のアルゴリズムとしては、STL (Seqential Table Look-up)法[3]、table-driven ア ルゴリズム[4][5]等が提案されている。STL法 はテーブルを逐次参照しながら指数関数計算を 行うアルゴリズムである。table-driven アルゴリ ズムはテーブル参照と多項式近似に基づいて指 数関数計算を行うアルゴリズムである。

本報告では、テーブル参照と多項式近似に基 づく、指数関数計算回路の設計について検討す る。入出力は IEEE754 標準の倍精度浮動小数点 基本フォーマットとし、誤差を 1 ulp (unit in the last place) 以内とする。本報告では、近似多項式

[†] 名古屋大学大学院工学研究科情報工学専攻 Department of Information Engineering, Nagoya University

^{††} 名古屋大学大学院情報科学研究科情報システム学専攻 Department of Information Engineering, Nagoya University

表 2 入力 *x* に対する exp(*x*)

表1 IEEE754 倍精度浮動小数点基本フォーマット

$E = 2047, F \neq 0$	x = NaN	非数
E = 2047, F = 0	$x = (-1)^S \cdot \infty$	±∞
0 < E < 2047	$x = (-1)^{S} \cdot 2^{E - 1023} \cdot (1 + F)$	正規化数
$E = 0, F \neq 0$	$x = (-1)^S \cdot 2^{-1022} \cdot (0+F)$	DEN 数
E = 0, F = 0	$x = (-1)^S \cdot 0$	± 0

x	exp(x)
NaN	NaN
$+\infty$	+∞
-∞	+0
正規化数 (E > 1032)	+∞
正規化数 (969 ≤ E ≤ 1032)	exp(x)
正規化数 (E < 969)	1
DEN 数	1
± 0	1

の次数とテーブルの分割数を変えて、テーブル サイズとクロック・サイクル数を評価する。ま た、この回路に適した積和演算器の構成を提案 する。積和演算器をパイプライン化しない場合 と、2段パイプライン化する場合に対して、回路 構成を示す。

2. 倍精度浮動小数点指数関数計算

本報告では、入力 x、出力 exp(x) はともに、 IEEE754 標準の倍精度浮動小数点基本フォーマッ トで表されているものとする。入力 x を IEEE754 標準の倍精度浮動小数点基本フォーマットで表 すと、 $x = (-1)^{S} \cdot 2^{E-1023} \cdot (1+F)$ (S:符号, 1 ビッ ト、E:指数部, 11 ビット、F: 仮数部, 52 ビット) となる。E、F の値により、入力 x は表 1 のよう に分類される。

まず、入力*x* に対する exp(*x*) の値について述べ る。IEEE754 標準の倍精度浮動小数点基本フォー マットの表現可能最大値は $(2 - 2^{-52}) \times 2^{1023}$ ($\approx 1.79763 \times 10^{308}$) である。従って、入力*x* が $x \ge 710$ の場合、exp(*x*) の値がこの表現可能最 大値よりも大きくなるので exp(*x*) の値は + ∞ に なる。また、Denomarlized 数 (DEN 数) の表現 可能最小値は 2^{-1074} ($\approx 4.94065 \times 10^{-324}$) であ る。従って、入力*x* が *x* ≤ -745 場合、exp(*x*) の 値がこの表現可能最小値よりも小さくなるので exp(*x*) の値は 0 になる。すなわち、exp(*x*) の値 が IEEE754 標準の倍精度浮動小数点基本フォー マットで表現可能とするような入力*x* の範囲は -745 < x < 710 となり、|x| < 745 であるから、 $E \leq 1032$ となる。

 $x \simeq 0$ の場合は、 $\exp(x) \simeq 1$ となる。IEEE754 標準で定められている丸めモードのうちの round to nearest を適用すると、 $1-2^{-54} \le \exp(x) \le 1+$ 2^{-53} となる入力xに対しては $\exp(x) = 1$ を出力す る。そのような入力xの範囲は $-2^{-54} \le x \le 2^{-54}$ である。従って、 $E \ge 969$ のとき、 $\exp(x) \ne 1$ と なる。

以上より、この回路で扱う*E*の範囲は、969 ≤ *E* ≤ 1032 となる。表 2 に入力 *x* に対する exp(*x*) の値を示す。

次に、969 ≤ *E* ≤ 1032 の場合について、指数 関数計算の手順を示す。まず、exp(*x*) を以下の ように変形する。

$$\exp(x) = \exp((-1)^{S} \cdot 2^{E-1023} \cdot (1+F))$$
$$= 2^{(-1)^{S} \cdot \frac{2^{E-1023}}{\ln 2} \cdot (1+F)}$$
(1)

ここで、 $D = (-1)^{S} \cdot \frac{2^{E-1023}}{\ln 2} \cdot (1+F)$ とし、Dの整数部を D_i 、小数部を D_f とすると、

$$D = D_i + D_f$$

$$\exp(x) = 2^D = 2^{D_i + D_f} = 2^{D_i} \cdot 2^{D_f}$$
(2)

$$= \pm + z$$

と表せる。

従って、 $D_i > -1023$ の場合、 $\exp(x)$ の符号ビットは 0、指数部は $D_i + 1023$ 、仮数部は 2^{D_f} となる。また、 $D_i \leq -1023$ の場合、 $\exp(x)$ の値はDEN数となる。このとき、 $\exp(x)$ の符号ビットは 0、指数部は 0、仮数部は 2^{D_f} を $(-1022 - D_i)$ ビット右シフトした値となる。

3. 計算法と回路構成

3.1 計 算 法

 $D = (-1)^{S} \cdot \frac{2^{E-1023}}{\ln 2} \cdot (1+F)$ を以下の手順で計算 して整数部 D_i と小数部 D_f に分離する。

- 1. $(1+F) \ge \frac{1}{\ln 2} \ge 0$ 乗算。
- 上記の乗算結果を (E-1023) ビットだけシ フト (左シフトは最大9ビット、右シフト は最大54 ビット)。
- 1. 上記のシフトした結果を整数部 D_i と小数 部 D_f に分離。(D_i = 12 ビット) 符号ビット S = 1 のときは、2 の補数をと る。ここでは、最終桁の誤差は許容範囲で あるので、実際には各ビットを反転し、1 の補数をとるだけでよい。

以上のようにして *D* を分離した結果、式 (2) が 得られる。

exp(x)の仮数部に該当する 2^{D_f} は、小数部 D_f の上位ビットを D_{f_u} 、下位ビットを D_{f_i} とすると、

 $2^{D_f} = 2^{D_{f_u} + D_{f_l}} = 2^{D_{f_u}} \cdot 2^{D_{f_l}}$ (3) と表せる。式 (3) の $2^{D_{f_l}}$ は多項式近似して積和 演算で求め、 $2^{D_{f_u}}$ はテーブル参照で求める。

また、 D_{f_u} の上位ビットを $D_{f_{u1}}$ 、下位ビットを $D_{f_{u2}}$ とすると、テーブル参照で求める $2^{D_{f_u}}$ は、以下のように分割することができる。

 $2^{D_{f_u}} = 2^{D_{f_{u1}}} + 2^{D_{f_{u2}}} = 2^{D_{f_{u1}}} \cdot 2^{D_{f_{u2}}}$ (4) 従って、分割する前の $2^{D_{f_u}}$ のテーブルの値は、 $2^{D_{f_{u1}}}$ のテーブルの値を $2^{D_{f_{u2}}}$ のテーブルの値を 乗算して求めることができる。

図1に指数関数計算の流れ図を示す。

3.2 回路構成

式(3)の2^{D_f}は多項式近似して積和演算で求め るので、積和演算器が必要となる。前節で示し た計算法では、複数の乗算を同時に行うことは ない。よって、全ての乗算は1つの積和演算器 で共有する。また、積和演算器を2段パイプラ ン化することにより、テーブル分割で生じる乗 算と近似多項式中での乗算がパイプライン処理 可能となる。

任意ビットのシフト操作にはバレルシフタが



必要となる。シフト操作を同時に行うことはないので、シフト操作は1つのバレルシフタで共有する。

4. 仮数部の計算

本章では、多項式近似にテーラー展開を用い て、仮数部の計算(3)における近似多項式の次数 とテーブルの分割数について考察する。

ここで、前章で示した小数部 D_f の上位ビット D_{f_u} のビット長をmビットとし、下位ビット D_{f_l} に対する値 $2^{D_{f_l}}$ をテーラー展開によりk次式で 近似する。このとき、 $D_{f_l} = y$ とすると最大誤差 ε_{taylor} は、

$$\varepsilon_{taylor} = 2^{2^{-m}} - \sum_{i=0}^{k} \frac{1}{k!} \cdot (\ln 2)^{i} \cdot y^{i}$$
(5)

となる。この ε_{taylor} が 2^{-52} 以下に収まるように、 k に対して m を求めると表 3 に示すようになる。

また、表3にはテーブルの分割数に対する、 テーブルサイズ(1エントリを56ビットとする) と、仮数部の計算に必要なクロック・サイクル数 も示した。表3のN1は積和演算器をパイプライ ン化しない場合のクロック・サイクル数、N2は 積和演算器を2段パイプライン化した場合のク ロック・サイクル数である。

表3から、テーブルの分割数を増やすと、テー ブルサイズは削減できるが、仮数部の計算にか かるクロック・サイクル数が増えることがわか

次数 (k)	分割数 (m)	D_{f_u}	テーブルサイズ	N_1	N_2
	1	16	$2^{16} \times 56$	3	6
2	2	8+8	$2^8 \times 56 + 2^8 \times 48$	4	6
	3	5+5+6	$2^5 \times 56 + 2^5 \times 51 + 2^6 \times 46$	5	7
3	1	12	$2^{12} \times 56$	4	8
	2	6+6	$2^6 \times 56 + 2^6 \times 50$	5	8
4	1	9	$2^9 \times 56$	5	10
	2	4+5	$2^4 \times 56 + 2^5 \times 52$	6	10

表3 仮数部計算に必要なサイクル数とテーブルサイズ

る。従って、テーブルサイズとクロック・サイ クル数を考慮して、最適な近似多項式の次数と テーブルの分割数を決める必要がある。表3か ら考察すると、積和演算器をパイプライン化し ない場合は、2次近似でテーブルを3分割にする か、または3次近似でテーブルを2分割にすれ ばよいと考えられる。一方、積和演算器を2段 パイプライン化する場合は、2次近似でテーブル を3分割にすればよいと考えられる。

5. 積和演算器

積和演算 $M_1 \times M_2 + M_3$ は加数 M_3 を部分積と して乗算に組み込むことで実現できる。 M_1 の ビット長を小数点以下 m_1 ビット、 M_2 及び M_3 の ビット長を小数点以下 m_2 ビットとする。

ここでの M_1 は小数部 D_f の下位ビット D_{f_i} なので、 M_1 の上位ビットは0である。従って、図2に示す部分積マトリックスの下の方は全て0となるので、図2の斜線部で、 M_3 と最下段の部分積との論理和をとることで、積和演算は実現できる。また、(1+F)の乗算と定数 $\frac{1}{\ln 2}$ との乗算は、図2に示すように、小数点以下 m_2 ビット目以下で打ち切れば、図2の破線部の部分積の生



図2 積和演算器の部分積マトリックス

成と累算の必要がない。

全ての乗算は1つの積和演算器を共有するの で、この誤差を考慮して、*m*₁及び*m*₂を決定す る必要がある。

部分積マトリックスの小数点以下 m_2 ビット目 以下を打ち切ることによって生じる誤差を E_{trunc} とすれば、

$$E_{trunc} = 2^{-m_1 - m_2} \cdot \sum_{i=0}^{m_1} (i+1) \cdot 2^i \tag{6}$$

と表せる。さらに、出力のビット長を $m_1 + k$ ビットとし、それ以降を切り捨てることによって生じる誤差を E_{round} とすれば、

$$E_{round} = 2^{-m_1-k} \cdot \sum_{i=0}^{m_2-m_1-k-1} 2^{-i-1}$$
$$= 2^{-m_1-k} - 2^{-m_2}$$
(7)

と表せる。従って、この積和演算器の出力の誤 差 *E_{total}* は、

 $E_{total} = E_{trunc} + E_{round}$ (8) と表せる。

6. 回路設計

表3より、小数部D_fの上位16ビットを5ビット、5ビット、6ビットに3分割して、テーブル 参照で求め、下位ビットを2次近似多項式で求 める方法で回路設計を行う。

6.1 誤差解析

誤差を 1 ulp 以内とする回路を設計するために は、2 次近似多項式による誤差を ε_{approx} 、仮数部 を計算するときに生じる誤差を ε_{comput} とすると

係数	
A_0	1.000000000000000000000000000000000000
A_1	0.101100010111001000010111111101111100100
A_2	0.0011110101111111100100010100111101001111

以下の式を満たす必要がある。

 $\epsilon_{approx} + \epsilon_{comput} \le 2^{-53}$ (9) ここで、積和演算器の乗数のビット長を小数 点以下 m_1 ビット、被乗数のビット長を小数点以 下 m_2 ビットとした場合の ϵ_{comput} について考え る。まず、(1+F) と $\frac{1}{\ln 2}$ とを乗算するときに生 じる誤差 ϵ_l は、

 $\varepsilon_l = 2^{-m_2} \cdot 53 + 2^{-m_1-9} - 2^{-m_2}$

 $< 2^{-m_2+6} - 2^{-m_2+3} + 2^{-m_1-9} - 2^{-m_2}$ (10)

と表せる。これが最大9ビット左シフトされる ので、シフト後の誤差 ε_y は、

 $\varepsilon_y = \varepsilon_l \cdot 2^9$

 $< 2^{-m_2+15} - 2^{-m_3+12} + 2^{-m_1} - 2^{-m_2+9}$ (11)

と表せる。次に2次近似多項式 $Y = A_0 + y(A_1 + yA_2)$ を計算するとき、 $A_1 + yA_2$ の計算で生じる 誤差を ε_{p_1} 、 $A_0 + y(A_1 + yA_2)$ の計算で生じる誤 差を ε_{p_2} とすると、

 $\varepsilon_{p_1} = A_2 \varepsilon_y + \varepsilon$

 $\varepsilon_{p_2} = y\varepsilon_{p_1} + (A_1 + yA_2)\varepsilon_y + \varepsilon$ (12) と表せる。ここで ε は前節の式 (6) で示した E_{trunc} である。

また、3 つテーブルから参照した値を上位の テーブルから各々 C_0 、 C_1 、 C_2 、各々の値の誤差 を $\varepsilon_t (\leq 2^{-m_1-1})$ とする。そして、 $Y \times C_2$ を計算 するときに生じる誤差を ε_{q_1} 、 $Y \times C_2 \times C_1$ を計算 するときに生じる誤差を ε_{q_2} 、 $Y \times C_2 \times C_1 \times C_0$ を 計算するときに生じる誤差を ε_{q_3} とすると、

$$\varepsilon_{q_1} = Y\varepsilon_t + C_2\varepsilon_{p_2} + \varepsilon$$

$$\varepsilon_{q_2} = YC_2\varepsilon_t + C_1\varepsilon_{q_1} + \varepsilon$$
(12)

 $\varepsilon_{q_3} = Y C_2 C_1 \varepsilon_t + C_0 \varepsilon_{q_2} + \varepsilon \tag{13}$

と表せる。従って、仮数部を計算するときに生 じる誤差 ϵ_{comput} は、

$$\varepsilon_{comput} = \varepsilon_{q_3} \tag{14}$$

となる。

以上より、式 (9) を満たすには、 $m_1 = 56$ 、 $m_2 = 72$ 、 $\varepsilon_{approx} \le 2^{-56}$ とすればよい。

6.2 近似多項式の係数

実際に使う近似多項式の係数を求めるときに は、MiniMax 近似法を用いる。MiniMax 近似法 とは最大誤差を最小とする近似多項式を求める 手法である。本報告では、MiniMax 近似多項式 を求める際に、Mathematica を使用した。

前節で示したように、近似多項式の係数は、近 似多項式による最大誤差が 2^{-56} 以下となるよう に決めなければならない。 $0 \le y \le 2^{-16}$ において、 2^{y} を2次MiniMax近似多項式 $A_0 + A_1y + A_2y^2$ で 近似したときの最大誤差は、 2^{-57} 以下となる。こ のとき、Mathematicaで求められる係数のうち、 小数点以下 56 ビット目までを使うと、近似多項 式による最大誤差が 2^{-56} 以下となる。従って、 回路設計に使う2次近似多項式 $A_0 + A_1y + A_2y^2$ の係数は表4に示す値になる。

6.3 積和演算器をパイプライン化しない場合

積和演算器をパイプライン化しないで回路設 計を行った場合、

1. $(1+F) \times \frac{1}{\ln 2} : 1 \forall f \neq h$

(E-1023) バレルシフト+2の補数:1サイクル

3. テーブル参照 + 多項式近似:5 サイクル

4. バレルシフト (DEN 数出力): 1 サイクル

となり指数関数計算は8サイクルで終了する。

6.4 積和演算器を2段パイプライン化する 場合

積和演算器を2段パイプライン化する場合、 テーブル分割で生じる乗算と近似多項式中での 乗算がパイプライン処理可能となる。図3に仮 数部の計算手順を示す。



図3 仮数部の計算

積和演算器を2段パイプライン化して回路設 計を行った場合、

- 1. $(1+F) \times \frac{1}{\ln 2} : 2$ サイクル
- (E-1023) バレルシフト+2の補数:1サ イクル
- 3. テーブル参照 + 多項式近似:7 サイクル

4. バレルシフト (DEN 数出力): 1 サイクル となり、指数関数計算は11 サイクルで終了する。

積和演算器を2段パイプライン化する場合、最 終的な指数関数計算回路のブロック図は図4の ようになる。

7.まとめ

本報告では、入出力は IEEE754 標準の倍精度 浮動小数点基本フォーマットとし、テーブル参 照と多項式近似に基づく指数関数計算回路の設 計について検討した。本報告では、近似多項式 の次数とテーブルの分割数を変えて、テーブル サイズとクロック・サイクル数を評価した。ま た、この回路に適した積和演算器の構成を提案 した。積和演算器をパイプライン化しない場合 と、2段パイプライン化する場合に対して、回路 構成を示した。本報告で示した手法を用いれば、 誤差を1 ulp (unit in the last place) 以内とする指 数関数計算回路が設計できる。



参考文献

- [1] 高木直史,"初等関数計算回路のアルゴリズ ム",情報処理学会,第37巻第4号,1996.
- [2] John Harrison, "The Computation of Transcendental Function on IA-64 Architecture", Intel Technology Journal Q4, 1999.
- [3] Chen, T.C., "Automatic Computation of Exponentials, Logarithms, Rations and Square Roots", IBM J.Research and Development, Vol.16, No. 4, pp. 380-388 (July 1972).
- [4] Douglas M.Priest, "Fast Table-Driven Algorithms for Interval Elementary Functions", Proc.13th IEEE Symposium on Computer Arithmetic 1997.
- [5] J.M.Muller, "Elementary Function-Algorithms and Implementations", Birkhauser, 1997.