

A compensated dot product for vectors of floating-point expansions

Nicolas Louvet, UCB Lyon 1
Aric, LIP

RAIM 2018, 10th edition,
November 13, 2018

Current **general purpose processors** implement two IEEE 754 binary floating point arithmetics: Binary32 ($p = 24$) and Binary64 ($p = 53$).

The IEEE 754 also specifies Binary128 ($p = 113$); implemented in **hardware** by IBM: available on the z13 and on the POWER9 processor.

When more than 53 bits of precision are required, different **software** solutions exist:

- libraries based on **integer arithmetic**:
 - ▶ `__float128` in GCC for binary128.
 - ▶ the MPFR library [Lefèvre, Zimmermann 2017].
- libraries based on **floating point expansions**:
 - ▶ `dd_real` and `qd_real` in the qd library [Hida, Li, Bailey 2001],
 - ▶ `quad_float` in the NTL library (code from K. Briggs),
 - ▶ more recently : CAMPARY [Popescu 17].
- **compensated algorithms**: summation, dot products ([Kahan 65], [Møller 65], [Pichat 72], [Neumanier 74], [Priest 92], [Ogita-Rump-Oishi 05, 08], . . .), polynomial evaluation ([Graillat-Langlois-L. 06], [Jiang-Barrio-Liao-Cheng 11]. . .)

How to trade off speed against a moderate loss of accuracy when computing with floating point expansions?

We do so by using **compensation techniques**, and by focusing on **dot products**.

Outline

- 1 Introduction
- 2 Compensation with expansions
- 3 Experiments
- 4 Conclusion

Notation and assumptions

\mathbb{F} denotes the set of binary floating-point numbers in precision p , with an unbounded exponent range (no over/underflow).

The floating point operations are rounded to the nearest even.

$\text{fl}()$ indicates that the expression inside the parentheses is evaluated in f. p. arithmetic.

$u = 2^{-p}$ is the unit roundoff. Given $\text{op} \in \{+, -, \times, /\}$, and $a, b \in \mathbb{F}$,

$$\text{fl}(a \text{ op } b) = (1 + \varepsilon)(a \text{ op } b), \quad \text{with } |\varepsilon| \leq u.$$

Given $x, y \in \mathbb{R}^n$, we note $\langle x, y \rangle = \sum_{i=0}^n x_i y_i$.

$e_n = [1, \dots, 1]^T \in \mathbb{R}^n$, and we will often write $\langle x, e_n \rangle$ for the sum of the elements in x .

Classical dot product algorithm

Let $x, y \in \mathbb{F}^n$ be two column vectors.

We assume a Fused Multiply-Add (FMA) instruction is available.

```

s = dotprod(x, y)
  s ← 0
  for i = 1 ... n - 1
    s ← fl(s + xiyi)
  return s

```

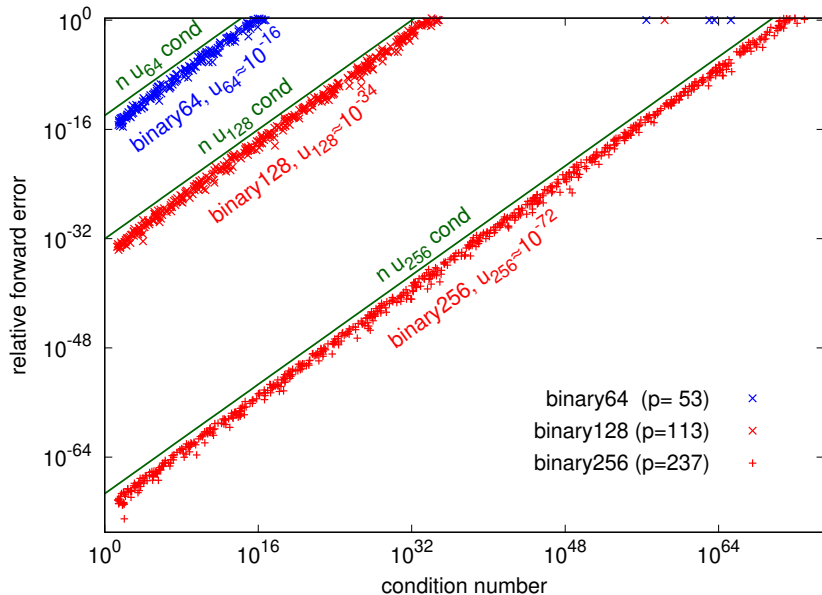
Relative error bound ($\langle x, y \rangle \neq 0$):

$$\frac{|s - \langle x, y \rangle|}{|\langle x, y \rangle|} \leq n u \text{cond}(x, y),$$

where the **condition number** is given by

$$\text{cond}(x, y) = \frac{\langle |x|, |y| \rangle}{|\langle x, y \rangle|} \quad \left(= \frac{\sum_{i=1}^n |x_i y_i|}{|\sum_{i=1}^n x_i y_i|} \right).$$

$$|s/\langle x, y \rangle - 1| \leq nu \text{ cond}(x, y)$$



Floating point expansions of length $K > 2$

Such an expansion a is a **sequence** (a_0, \dots, a_{K-1}) of elements in \mathbb{F} , representing $\sum_{i=0}^{K-1} a_i$.

A condition ensures that the elements are “sufficiently non-overlapping”: a **renormalization step** is required at the end of each arithmetic operation (see [Popescu 17]).

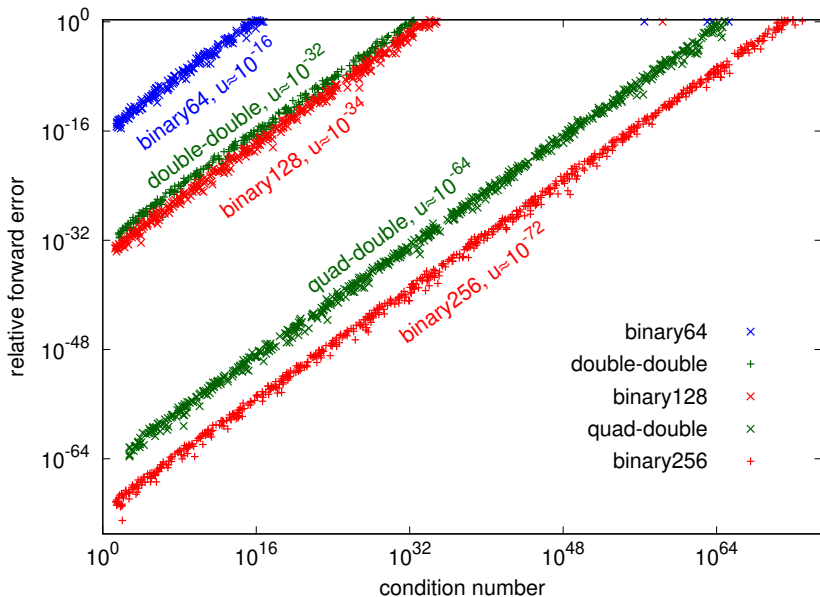
In this talk, we only assume that the floating point expansions we take as input satisfy

$$|a_i| \leq 2u|a_{i-1}|, \quad \text{for } 0 < i < K,$$

and we use a renormalization algorithm from [Popescu 17] for the output.

In the qd library [Hida, Li, Bailey 2001]:

- a **double-double** is made of **2** Binary64 numbers (goes back to [Dekker 1971]).
 \hookrightarrow 106 bits available, precision $\approx u^2$.
- a **quad-double** is the unevaluated sum of **4** Binary64 numbers.
 \hookrightarrow 212 bits available, precision $\approx u^4$.



Error-free transformations (EFT)

- Given $\text{op} = +$ or \times , and $x, y \in \mathbb{F}$, an EFT for op computes $r, \rho \in \mathbb{F}$ such that

$$x \text{ op } y = r + \rho \quad \text{with} \quad r = \text{fl}(x \text{ op } y).$$

+	$(r, \rho) = \text{TwoSum}(x, y)$	6 flops	[Knuth 74]
+	$(r, \rho) = \text{FastTwoSum}(x, y)$ if $ x \geq y $	3 flops	[Dekker 71]
\times	$(r, \rho) = \text{TwoProd}(x, y)$ with FMA	2 flops	[Karp, Markstein 1997]

- EFT for the dot product of two floating point vectors $x, y \in \mathbb{F}^n$:

$$(p, \pi) = \text{VecTwoProd}(x, y)$$

for $i = 1 \dots n$

$$(p_i, \pi_i) = \text{TwoProd}(x_i, y_i)$$

\rightsquigarrow computes $p, \pi \in \mathbb{F}^n$ s. t.

$$\langle x, y \rangle = \langle p, e_n \rangle + \langle \pi, e_n \rangle.$$

- EFT for the summation of the elements of a vector $v \in \mathbb{F}^n$ [Ogita, Rump, Oishi 2005]:

$$(s_n, v') = \text{VecSum}(v)$$

$$s_1 = v'_1 = 0$$

for $i = 2 \dots n$

$$(s_i, v'_i) = \text{TwoSum}(s_{i-1}, v_i)$$

\rightsquigarrow computes $s \in \mathbb{F}$ and $v' \in \mathbb{F}^n$ s. t.

$$\langle v, e_n \rangle = s_n + \langle v', e_n \rangle.$$

Outline

- 1 Introduction
- 2 Compensation with expansions
- 3 Experiments
- 4 Conclusion

Compensation with double-double ($K = 2$) vectors

$x = x^{(0)} + x^{(1)}$ and $y = y^{(0)} + y^{(1)}$ are two double-double vectors of size n :

$$\langle x, y \rangle = \langle x^{(0)}, y^{(0)} \rangle + \langle x^{(1)}, y^{(0)} \rangle + \langle x^{(0)}, y^{(1)} \rangle + \langle x^{(1)}, y^{(1)} \rangle.$$

We relate the magnitude of all the terms to $\langle |x|, |y| \rangle$ (numerator in $\text{cond}(x, y)$).

Using $|x^{(1)}| \leq 2u|x^{(0)}|$, $|y^{(1)}| \leq 2u|y^{(0)}|$:

$$\begin{aligned} |\langle x^{(0)}, y^{(0)} \rangle| &\leq \langle |x|, |y| \rangle + \mathcal{O}(u), \\ |\langle x^{(1)}, y^{(0)} \rangle| &\leq 2u \langle |x|, |y| \rangle + \mathcal{O}(u^2), \\ |\langle x^{(0)}, y^{(1)} \rangle| &\leq 2u \langle |x|, |y| \rangle + \mathcal{O}(u^2), \\ |\langle x^{(1)}, y^{(1)} \rangle| &\leq 4u^2 \langle |x|, |y| \rangle + \mathcal{O}(u^3). \end{aligned}$$

Compensation with double-double ($K = 2$) vectors

$x = x^{(0)} + x^{(1)}$ and $y = y^{(0)} + y^{(1)}$ are two double-double vectors of size n :

$$\langle x, y \rangle = \langle x^{(0)}, y^{(0)} \rangle + \langle x^{(1)}, y^{(0)} \rangle + \langle x^{(0)}, y^{(1)} \rangle + \langle x^{(1)}, y^{(1)} \rangle.$$

$$(p^{(0,0)}, \pi^{(0,0)}) = \text{VecTwoProd}(x^{(0)}, y^{(0)})$$

$$(s_0, q^{(0)}) = \text{VecSum}(p^{(0,0)})$$

$$\Rightarrow \langle x^{(0)}, y^{(0)} \rangle = s_0 + \sum_{\ell=1}^n \pi_{\ell}^{(0,0)} + \sum_{\ell=1}^n q_{\ell}^{(0)}$$

$$(p^{(1,0)}, \pi^{(1,0)}) = \text{VecTwoProd}(x^{(1)}, y^{(0)})$$

$$(p^{(0,1)}, \pi^{(0,1)}) = \text{VecTwoProd}(x^{(0)}, y^{(1)})$$

$$(s_1, q^{(1)}) = \text{VecSum}([p^{(1,0)}; p^{(0,1)}; \pi^{(0,0)}; q^{(0)}])$$

$$\Rightarrow \langle x^{(0)}, y^{(0)} \rangle + \langle x^{(1)}, y^{(0)} \rangle + \langle x^{(0)}, y^{(1)} \rangle = s_0 + s_1 + \sum_{\ell=1}^n (\pi_{\ell}^{(1,0)} + \pi_{\ell}^{(0,1)}) + \sum_{\ell=1}^{4n} q_{\ell}^{(1)}$$

Overall, we get

$$\langle x, y \rangle = s_0 + s_1 + \langle x^{(1)}, y^{(1)} \rangle + \sum_{\ell=1}^n (\pi_{\ell}^{(1,0)} + \pi_{\ell}^{(0,1)}) + \sum_{\ell=1}^{4n} q_{\ell}^{(1)}$$

In the algorithm, we only keep what is needed to compute s_0 and s_1 :

$$\begin{aligned} (r_0, r_1) &= \text{DotComp2}(x, y) \\ (p^{(0,0)}, \pi^{(0,0)}) &= \text{VecTwoProd}(x^{(0)}, y^{(0)}) \\ (s_0, q^{(0)}) &= \text{VecSum}(p^{(0,0)}) \\ p^{(1,0)} &= \text{VecProd}(x^{(1)}, y^{(0)}) \\ p^{(0,1)} &= \text{VecProd}(x^{(0)}, y^{(1)}) \\ s_1 &= \text{Sum}([p^{(1,0)}; p^{(0,1)}; \pi^{(0,0)}; q^{(0)}]) \\ (r_0, r_1) &= \text{Renorm}(s_0, s_1) \end{aligned}$$

Bounding the magnitude of the $\mathcal{O}(u^2)$ terms gives:

If $\langle x, y \rangle \neq 0$, under some assumptions on u and n ,

$$\frac{|\langle x, y \rangle - (r_0 + r_1)|}{|\langle x, y \rangle|} \leq (1 + 5u)(4 + 24n + 4n^2)u^2 \text{cond}(x, y).$$

In practice, we re-arrange the algorithm so that only one pass through the data is needed:

DotComp2 ($14n + \mathcal{O}(1)$ flops)

```
dd_real
DotComp2(dd_real *x, dd_real *y, int n) {
    double xh, xl, yh, yl, s, c, p1, e1, e2;

    s = c = 0.0;
    for(int i=0; i<n; i++) {
        xh = x[i]._hi(); xl = x[i]._lo();
        yh = y[i]._hi(); yl = y[i]._lo();

        TwoProd(p1, e1, xh, yh);
        c += xh*yl + xl*yh;

        TwoSumIn(s, e2, p1);
        c += e1 + e2;
    }
    TwoSumIn(s, c, c);
    return dd_real(s, c);
}
```

DotDD ($20n + \mathcal{O}(1)$ flops)

```
dd_real
DotDD(dd_real *x, dd_real *y, int n) {
    double xh, xl, yh, yl, sh, sl, ph, pl, e;

    sh = sl = 0.0;
    for(int i=0; i<n; i++) {
        xh = x[i]._hi(); xl = x[i]._lo();
        yh = y[i]._hi(); yl = y[i]._lo();

        TwoProd(ph, pl, xh, yh);
        pl += xh*yl + xl*yh;
        FastTwoSum(ph, pl, ph, pl);

        TwoSumIn(sh, e, ph);
        e += sl + pl;
        FastTwoSumIn(sh, sl, sh, e);
    }

    return dd_real(sh, sl);
}
```

DotComp avoids the renormalization steps used in **DotDD**: this exposes **more instruction-level parallelism**, which in general results in better performances (as in the case of CompHorner [L. 2007], and other compensated algorithms [Goossens *et al.* 2012]).

Compensated dot product with vectors of expansions of length $K > 2$

x and y are vectors of f. p. expansions: $x = \sum_{i=0}^{K-1} x_i$.

Same principles as in the case $K = 2$:

- all the renormalization steps are avoided,
- if K is fixed, “single pass” implementation.

Error bound of the form:

$$\frac{|r - \langle x, y \rangle|}{|\langle x, y \rangle|} \leq f_K(n) u^K \text{cond}(x, y).$$

```

( $r_0, \dots, r_{K-1}$ ) = DotComp( $x, y, K$ )
for  $k = 0 \dots K - 1$ 
   $s_k \leftarrow 0$ 
end do
for  $\ell = 1 \dots n$ 
  for  $i, j$  such that  $0 \leq i + j \leq K - 2$ 
    ( $p, \pi$ )  $\leftarrow$  TwoProd( $x_\ell^{(i)}, y_\ell^{(j)}$ )
    ( $s_{i+j}, p$ )  $\leftarrow$  TwoSum( $s_{i+j}, p$ )
    for  $k = i + j + 1 \dots K - 2$ 
      ( $s_k, p$ )  $\leftarrow$  TwoSum( $s_k, p$ )
      ( $s_k, \pi$ )  $\leftarrow$  TwoSum( $s_k, \pi$ )
    end do
     $s_{K-1} \leftarrow \text{fl}(s_{K-1} + p + \pi)$ 
  end do
  for  $i, j$  such that  $i + j = K - 1$ 
     $s_{K-1} \leftarrow \text{fl}(s_{K-1} + x_\ell^{(i)} y_\ell^{(j)})$ 
  end do
end do
( $r_0, \dots, r_{K-1}$ ) = Renorm( $s_0, \dots, s_{K-1}$ )

```

With **quad-double** inputs ($K=4$), **116n flops** are required, and

$$f_4(n) = (1 + 5u)(96 + 768n + 41472un^3 + 1296n^4).$$

Outline

- 1 Introduction
- 2 Compensation with expansions
- 3 Experiments**
- 4 Conclusion

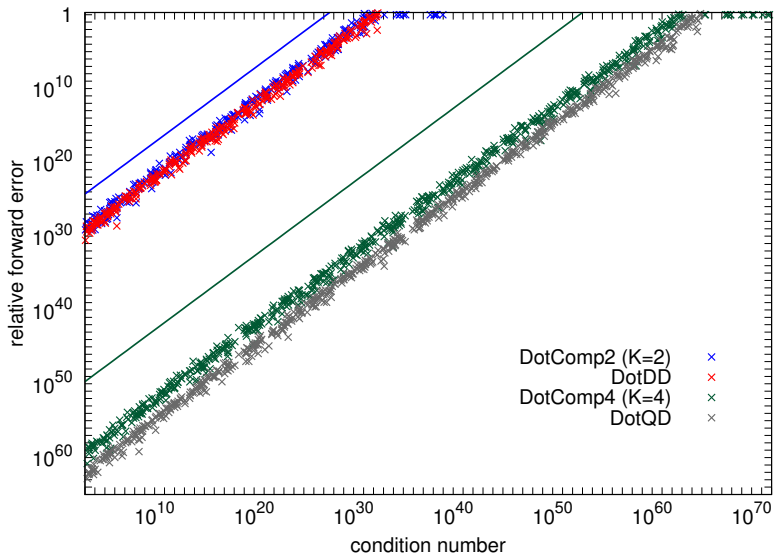
Numerical experiments

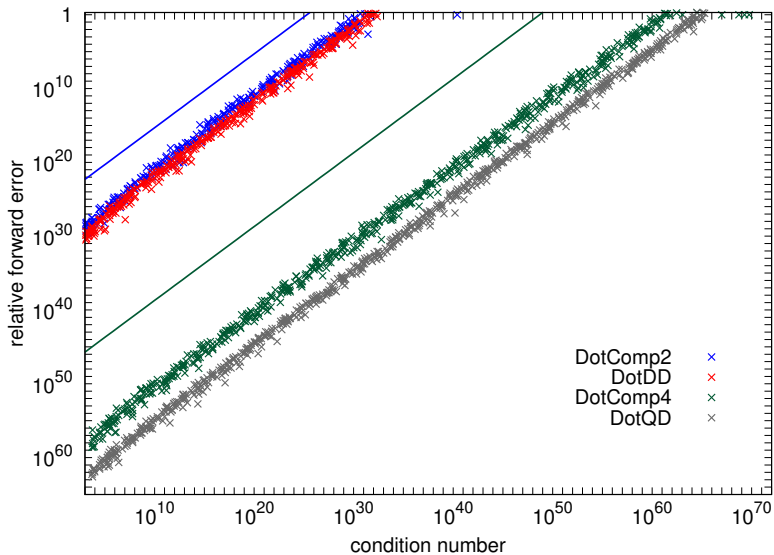
IEEE 754 Binary64 (double) is used as the working precision ($p = 53$, $u \approx 10^{-16}$).

In these experiments, we use:

- **DotDD** and **DotQD** implements the naive dot product using the qd library,
- **DotComp2** and **DotComp4** implements our algorithm with $K = 2$ and 4 respectively; the implementations are compatible with the qd library.

The algorithms have been implemented in C++, compiled with GCC 6.3.0 (with options `-fno-associative-math` and `-ffp-contract=off`).

$n = 100$ 

$n = 1000$ 

Practical performances

We report experiments in two environments:

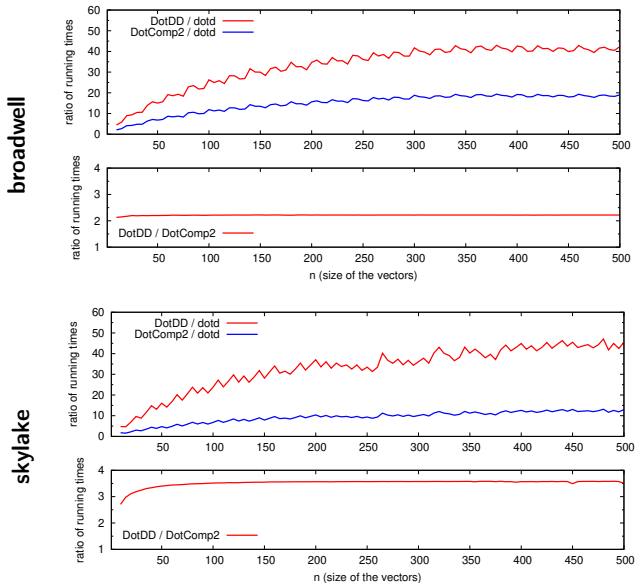
- **broadwell**: Intel CPU with micro-architecture Broadwell, GCC 6.3.0 (`-march=broadwell -mavx2` and `-mfma`).
- **skylake**: Intel CPU with micro-architecture Skylake, GCC 6.3.0 (`-march=skylake -mavx512f` and `-mfma`).

We measured the **computing times in clock-cycles** using the RDTSC and RDTSCP: we repeat each measure 400 times, ignore the 10% largest values, and take the average.

In all the case, the **caches are warm** (and the data fit within the L1 32 kB cache).

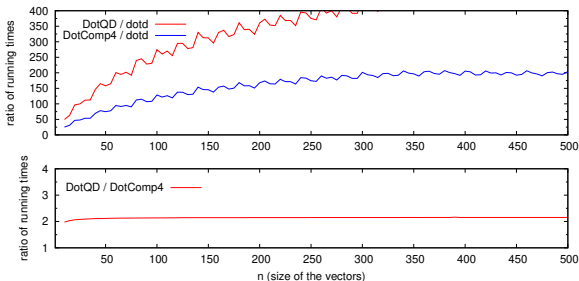
dotd is the naive algorithm in Binary64 arithmetic, vectorized by GCC.
DotDD, **DotQD**, **DotComp2**, **DotComp4** are non-vectorized.

DotComp2 runs between 2 and 3.5 times faster than DotDD

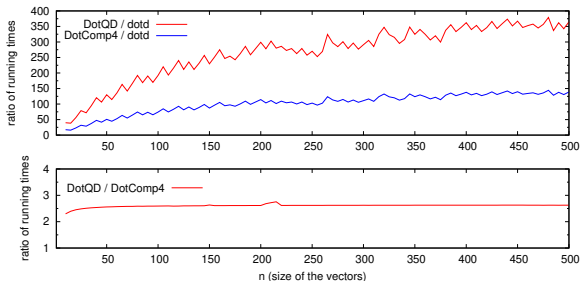


DotComp4 runs between 2 and 2.5 times faster than DotQD

broadwell



skylake



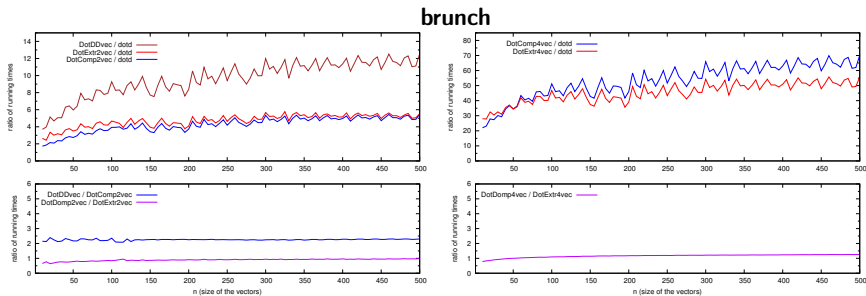
Conclusion

DotComp2 and DotComp4 are interesting alternatives to the use of the qd library: slightly less accurate algorithms, but a least two times faster.

We also experimented using PAPI (Performance Application Programming Interface) to confirm that the compensated algorithms provide more instruction level parallelism.

Other EFTs for summation of floating point numbers can be used in DotComp: “ExtractVector” from [Rump 2009] leads to better performances when $K = 4$.

Even better speedups can be obtained using SIMD vectorization:



Appendix

Compensated dot product [Neumanier 74, Pichat 76, Ogita-Rump-Oishi 05]

Given $x, y \in \mathbb{F}^n$, assume we start the evaluation of $\langle x, y \rangle$ as follows:

$$\begin{aligned}(p, \pi) &= \text{VecTwoProd}(x, y) & // \langle x, y \rangle &= \langle p, e_n \rangle + \langle \pi, e_n \rangle \\(s_n, \sigma) &= \text{VecSum}(p) & // \langle p, e_n \rangle &= s_n + \langle \sigma, e_n \rangle\end{aligned}$$

Since $\langle x, y \rangle = s_n + \langle \pi + \sigma, e_n \rangle$, we evaluate $\langle \pi + \sigma, e_n \rangle$ to compensate the errors in s_n .

```
r = dot2(x, y)
(p, pi) = VecTwoProd(x, y)
(s_n, sigma) = VecSum(p)
c = Sum(fl(pi + sigma))
r = fl(s_n + c)
```

Classical algorithm [Neumaier 1974, Pichat 1976].

Error bound from [Ogita, Rump, Oishi 2005]:

$$\frac{|r - x^T y|}{|x^T y|} \leq u + n^2 u^2 \text{cond}(x, y) + \mathcal{O}(u^3).$$

Error-free transformations (EFT)

- Given $\text{op} = +$ or \times , and $a, b \in \mathbb{F}$, an EFT for op computes $x, y \in \mathbb{F}$ such that

$$x + y = a \text{ op } b \quad \text{with} \quad x = \text{fl}(a \text{ op } b).$$

Known EFT using only floating point operations in \mathbb{F} :

+	$(x, y) = \text{TwoSum}(a, b)$	6 flops	[Knuth 74]
+	$(x, y) = \text{FastTwoSum}(a, b)$ if $ a \geq b $	3 flops	[Dekker 71]
\times	$(x, y) = \text{TwoProd}(a, b)$ with FMA	2 flops	[Karp, Markstein 1997]

$$\begin{aligned}(x, y) &= \text{FastTwoSum}(a, b) \\ x &= \text{fl}(a + b) \\ y &= \text{fl}(\text{fl}(a - x) + b)\end{aligned}$$

$$\begin{aligned}(x, y) &= \text{TwoProd}(a, b) \\ x &= \text{fl}(a \times b) \\ y &= \text{fl}(a \times b - x)\end{aligned}$$

Basic blocs for extending the working precision:

- expansions: [Dekker 71], [Priest 91], [Hida, Li, Bailey 01], \dots , [Popescu 17], [van der Hoeven 17]
- compensated summation and dot product: [Kahan 65], [Møller 65], [Pichat 72], [Neumanier 74], [Priest 92], [Ogita-Rump-Oishi 05, 08], [Rump 09]

- EFT for the dot product of two floating point vectors $x, y \in \mathbb{F}^n$:

$$\begin{aligned} (p, \pi) &= \mathbf{VecTwoProd}(x, y) \\ &\text{for } i = 1 \dots n \\ (p_i, \pi_i) &= \mathbf{TwoProd}(x_i, y_i) \end{aligned}$$

VecTwoProd computes $p, \pi \in \mathbb{F}^n$ s. t.

$$\langle x, y \rangle = \langle p, e_n \rangle + \langle \pi, e_n \rangle \quad \text{and} \quad |\langle \pi, e_n \rangle| \leq u \langle |x|, |y| \rangle.$$

- EFT for the summation of the elements in $v \in \mathbb{F}^n$:

$$\begin{aligned} (s_n, v') &= \mathbf{VecSum}(v) \\ s_1 &= v'_1 = 0 \\ &\text{for } i = 2 \dots n \\ (s_i, v'_i) &= \mathbf{TwoSum}(s_{i-1}, v_i) \end{aligned}$$

VecSum computes $s \in \mathbb{F}$ and $v' \in \mathbb{F}^n$ s. t.

$$\langle v, e_n \rangle = s_n + \langle v', e_n \rangle \quad \text{and} \quad |\langle v', e_n \rangle| \leq nu \langle |v|, e_n \rangle.$$

Compensated dot product : some details

Three pass version of the algorithm:

```

$$r = \mathbf{dot2}(x, y)$$

$$(p, \pi) = \mathbf{VecTwoProd}(x, y)$$

$$(s_n, \sigma) = \mathbf{VecSum}(p)$$

$$c = \mathbf{Sum}(\mathbf{fl}(\pi + \sigma))$$

$$r = \mathbf{fl}(s_n + c)$$

```

Single pass version of the algorithm:

```

$$r = \mathbf{dot2}(x, y)$$

$$s_0 = c_0 = 0;$$

$$\text{for } i = 1 \dots n$$

$$(p_i, \pi_i) = \mathbf{TwoProd}(x_i, y_i)$$

$$(s_i, \sigma_i) = \mathbf{TwoSum}(s_{i-1}, p_i)$$

$$c_i = \mathbf{fl}(c_{i-1} + \mathbf{fl}(\pi_i + \sigma_i))$$

$$r = \mathbf{RN}(s_n + c_n)$$

```

Further experiments with the PAPI library

- We used PAPI to measure the Instruction Per Cycle (IPC) rate of each implementation :

$$\text{IPC} = \frac{\#instr}{\#cycles}$$

- Example with vectors of length $n = 500$:

	broadwell			skylake		
	#instr	#cycles	IPC	#instr	#cycles	IPC
DotDD	11816	12690	0.9	11816	16329	0.7
DotCmp2	8455	5017	1.7	8455	3833	2.2
DotQD	140585	115581	1.2	125586	103076	1.2
DotCmp4	67689	53180	1.3	66687	40235	1.7

- When we move from **broadwell** to **skylake**:

- the IPC of DotDD decreases ☹
- the IPC of CompDot2 increases ☺

Renormalization steps in DotDD = 2×3 ADD/SUB in series in the loop :

	instruction	ports	throughput	latency
broadwell	ADD/SUB	p1	1	3
skylake	ADD/SUB	p0 or p1	2	4

- Cheat sheet concerning the floating point instructions (<http://www.agner.org/>):

	instruction	ports	throughput	latency
broadwell	ADD/SUB	p1	1	3
	MUL/FMA	p0 or p1	2	3
skylake	ADD/SUB	p0 or p1	2	4
	MUL/FMA	p0 or p1	2	4