

## Abstract

In computer aided geometric design a polynomial is usually represented in Bernstein form. This paper presents a family of compensated algorithms to accurately evaluate a polynomial in Bernstein form with floating point coefficients. The principle is to apply error-free transformations to improve the traditional de Casteljau algorithm. At each stage of computation, round-off error is passed on to first order errors, then to second order errors, and so on. After the computation has been “filtered”  $(K - 1)$  times via this process, the resulting output is as accurate as the de Casteljau algorithm performed in  $K$  times the working precision. Forward error analysis and numerical experiments illustrate the accuracy of this family of algorithms.

## Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>Basic notation and results</b>	<b>3</b>
<b>3</b>	<b>Compensated de Casteljau</b>	<b>5</b>
<b>4</b>	<b><math>K</math> Compensated de Casteljau</b>	<b>6</b>
<b>5</b>	<b>Numerical experiments</b>	<b>9</b>
<b>6</b>	<b>de Casteljau’s Method</b>	<b>9</b>
6.1	Condition Number . . . . .	11
6.2	Example of Compensated de Casteljau Failing . . . . .	11
6.3	Selection of Test Cases . . . . .	12
6.4	Operation Counts . . . . .	12
<b>7</b>	<b>Bogus Section for Refs</b>	<b>12</b>
<b>8</b>	<b>Appendix: EFT Algorithms</b>	<b>13</b>

## 1 Introduction

In computer aided geometric design, polynomials are usually expressed in Bernstein form. Polynomials in this form are usually evaluated by the de Casteljau algorithm. This algorithm has a round-off error bound which grows only linearly with degree, even though the number of arithmetic operations grows quadratically. The Bernstein basis is optimally suited ([FR87]) for polynomial evaluation; it is typically more accurate than the monomial basis, for example in Figure 1 evaluation via Horner’s method produces a jagged curve for points near a tripl root, but the de Casteljau algorithm produces a smooth curve. Nevertheless the de Casteljau algorithm returns results arbitrarily less accurate than the working precision  $\mathbf{u}$  when evaluating  $p(s)$  is ill-conditioned. The relative accuracy of the computed evaluation with the de Casteljau algorithm (DeCasteljau) satisfies ([MP99]) the following a priori bound:

$$\frac{|p(s) - \text{DeCasteljau}(p, s)|}{|p(s)|} \leq \text{cond}(p, s) \times \mathcal{O}(\mathbf{u}). \quad (1)$$

In the right-hand side of this inequality,  $\mathbf{u}$  is the computing precision and the condition number  $\text{cond}(p, s) \geq 1$  only depends on  $s$  and the Bernstein coefficients of  $p$  — its expression will be given further.

For ill-conditioned problems, such as evaluating  $p(s)$  near a multiple root, the condition number may be arbitrarily large, i.e.  $\text{cond}(p, s) > 1/\mathbf{u}$ , in which case most or all of the computed digits will be incorrect. In some cases, even the order of magnitude of the computed value of  $p(s)$  can be incorrect.

To address ill-conditioned problems, error-free transformations (EFT) can be applied in *compensated algorithms* to account for roundoff. Error-free transformations were studied in great detail in [ORO05] and open a large number of applications. In [LGL06], a compensated Horner’s algorithm was described to

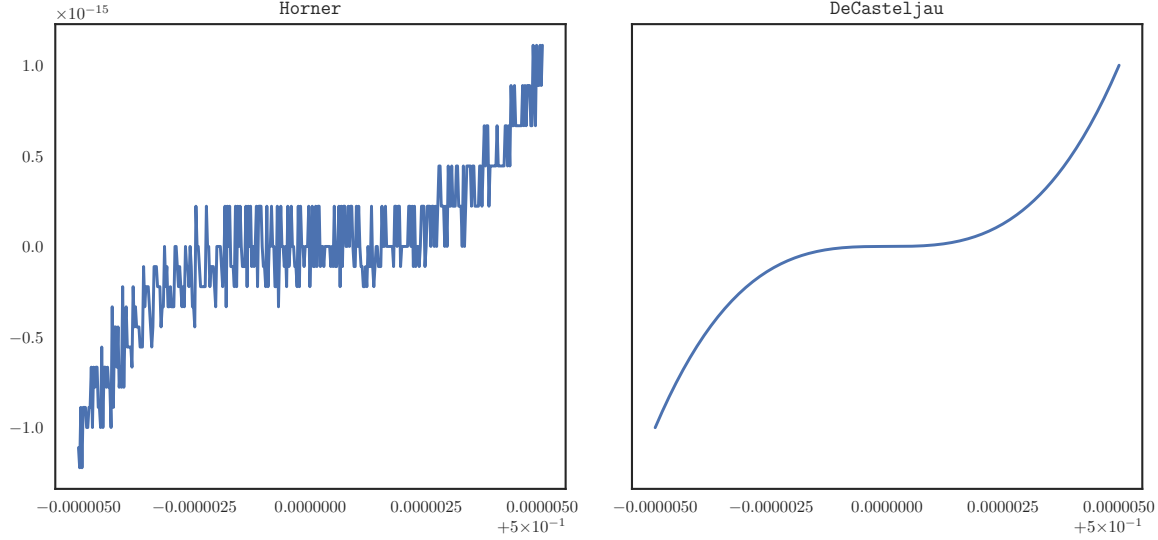


Figure 1: Comparing Horner's method to the de Casteljau method for evaluating  $p(s) = (2s - 1)^3$  in the neighborhood of its multiple root  $1/2$ .

evaluate a polynomial in the monomial basis. In [JLCS10], a similar method was described to perform a compensated version of the de Casteljau algorithm. In both cases, the  $\text{cond}(p, s)$  factor is moved from  $\mathbf{u}$  to  $\mathbf{u}^2$  and the computed value is as accurate as if the computations were done in twice the working precision. For example, the compensated de Casteljau algorithm (**CompDeCasteljau**) satisfies

$$\frac{|p(s) - \text{CompDeCasteljau}(p, s)|}{|p(s)|} \leq \mathbf{u} + \text{cond}(p, s) \times \mathcal{O}(\mathbf{u}^2). \quad (2)$$

For problems with  $\text{cond}(p, s) < 1/\mathbf{u}^2$ , the relative error is  $\mathbf{u}$ , i.e. accurate to full precision, aside from rounding to the nearest floating point number. Figure 2 shows this shift in relative error from **DeCasteljau** to **CompDeCasteljau**.

In [GLL09], the authors generalized the compensated Horner's algorithm to produce a method for evaluating a polynomial as if the computations were done in  $K$  times the working precision for any  $K \geq 2$ . This result motivates this paper, though the approach there is somewhat different than ours. They perform each computation with error-free transformations and interpret the errors as coefficients of new polynomials. They then evaluate the error polynomials, which (recursively) generate second order error polynomials and so on. This recursive property causes the number of operations to grow exponentially in  $K$ . Here, we instead have a fixed number of error groups, each corresponding to roundoff from the group above it. For example, when  $(1 - s)b_j^{(n)} + sb_{j+1}^{(n)}$  is computed in floating point, any error is filtered down to the error group below it.

As in (1), the accuracy of the compensated result (2) may be arbitrarily bad for ill-conditioned polynomial evaluations. For example, as the condition number grows in Figure 2, some points have relative error exactly equal to 1; this indicates that  $\text{CompDeCasteljau}(p, s) = 0$ , which is a complete failure to evaluate the order of magnitude of  $p(s)$ . For root-finding problems  $\text{CompDeCasteljau}(p, s) = 0$  when  $p(s) \neq 0$  can cause premature convergence and incorrect results. We describe how to defer rounding into progressively smaller error groups and improve the accuracy of the computed result by a factor of  $\mathbf{u}$  for every error group added. So we derive **CompDeCasteljauK**, a  $K$ -fold compensated de Casteljau algorithm that satisfies the following a priori bound for any arbitrary integer  $K$ :

$$\frac{|p(s) - \text{CompDeCasteljauK}(p, s)|}{|p(s)|} \leq \mathbf{u} + \text{cond}(p, s) \times \mathcal{O}(\mathbf{u}^K). \quad (3)$$

This means that the computed value with **CompDeCasteljauK** is now as accurate as the result of the de Casteljau algorithm performed in  $K$  times the working precision with a final rounding back to the working precision.

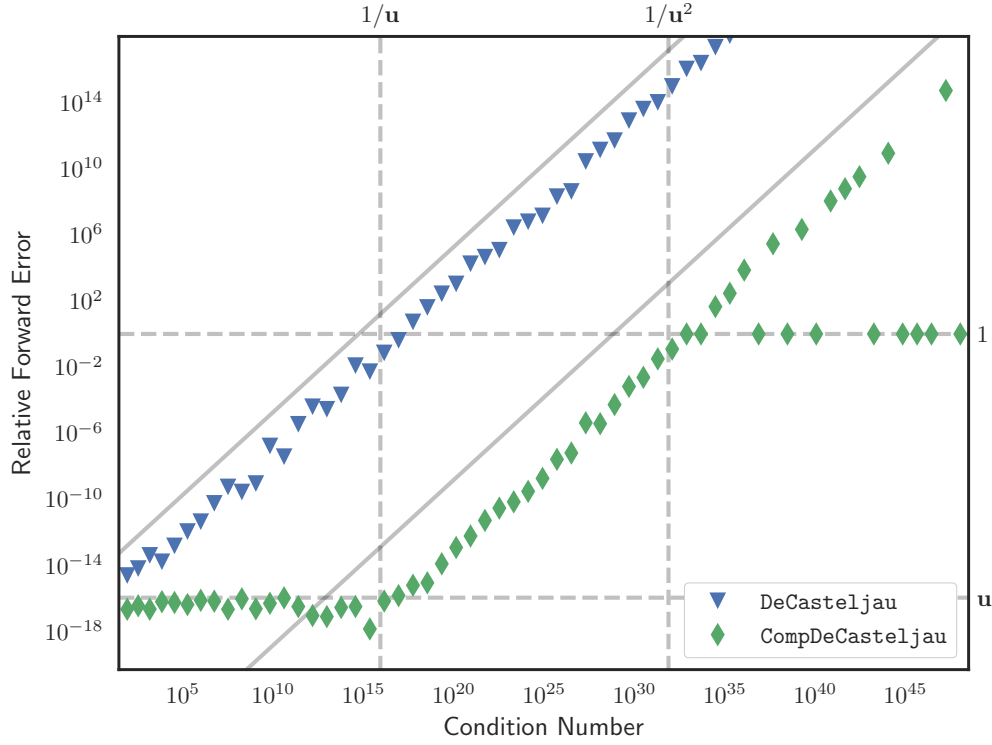


Figure 2: Evaluation of  $p(s) = (s - 1)(s - 3/4)^7$  represented in Bernstein form.

The paper is organized as follows. Section 2 introduces some basic notation and results about floating point arithmetic, error-free transformations and the de Casteljau algorithm. In Section 3, the compensated algorithm for polynomial evaluation from [JLCS10] is reviewed and notation is established for the expansion. In Section 4, the  $K$  compensated algorithm is provided and a forward error analysis is performed. Finally, in Section 5 we give numerical tests to illustrate the practical efficiency of our algorithms.

## 2 Basic notation and results

Throughout this paper we assume that the computation in floating point arithmetic obeys the model

$$a \star b = \text{fl}(a \circ b) = (a \circ b)(1 + \varepsilon_1) = (a \circ b)/(1 + \varepsilon_2) \quad (4)$$

where  $\star \in \{\oplus, \ominus, \otimes, \oslash\}$ ,  $\circ \in \{+, -, \times, \div\}$  and  $|\varepsilon_1|, |\varepsilon_2| \leq \mathbf{u}$ . The symbol  $\mathbf{u}$  is the unit round-off and  $\star$  represents the floating point computation, e.g.  $a \oplus b = \text{fl}(a + b)$ . (For IEEE-754 floating point double precision,  $\mathbf{u} = 2^{-53}$ .) We also assume that the computed result of  $\alpha \in \mathbf{R}$  in floating point arithmetic is denoted by  $\hat{\alpha}$  or  $\text{fl}(\alpha)$  and  $\mathbf{F}$  denotes the set of all floating point numbers (see [Fig02] for more details). Following [Fig02], we will use the following classic properties in error analysis.

1. If  $\delta_i \leq \mathbf{u}$ ,  $\rho_i = \pm 1$ , then  $\prod_{i=1}^n (1 + \delta_i)^{\rho_i} = 1 + \theta_n$ ,
2.  $|\theta_n| \leq \gamma_n := n\mathbf{u}/(1 - n\mathbf{u})$ ,
3.  $(1 + \theta_k)(1 + \theta_j) \leq (1 + \theta_{k+j})$ ,
4.  $\gamma_k + \gamma_j + \gamma_k\gamma_j \leq \gamma_{k+j}$ ,
5.  $(1 + \mathbf{u})^j \leq 1/(1 - j\mathbf{u})$ .

Now let us introduce some results concerning error-free transformation (EFT). For a pair floating point numbers  $a, b \in \mathbf{F}$ , there exists a floating point number  $y$  satisfying  $a \circ b = x + y$  where  $x = \text{fl}(a \circ b)$  and  $\circ \in \{+, -, \times\}$ . The transformation  $(a, b) \rightarrow (x, y)$  is regarded as an error-free transformation. The error-free transformation algorithms of the sum and product of two floating point numbers used later in this paper are the **TwoSum** algorithm by Knuth [Knu97] and **TwoProd** algorithm by Dekker [Dek71], respectively. The following theorem exhibits the important properties of **TwoSum** and **TwoProd** algorithms (see [ORO05]).

**Theorem 1** ([ORO05]). For  $a, b \in \mathbf{F}$  and  $P, \pi, S, \sigma \in \mathbf{F}$ , **TwoSum** and **TwoProd** satisfy

- $[S, \sigma] = \text{TwoSum}(a, b)$ ,  $x = \text{fl}(a + b)$ ,  $S + \sigma = a + b$ ,  $\sigma \leq \mathbf{u}|S|$ ,  $\sigma \leq \mathbf{u}|a + b|$
- $[P, \pi] = \text{TwoProd}(a, b)$ ,  $P = \text{fl}(a \times b)$ ,  $P + \pi = a \times b$ ,  $\pi \leq \mathbf{u}|P|$ ,  $\pi \leq \mathbf{u}|a \times b|$ .

The letters  $\sigma$  and  $\pi$  are used to indicate that the errors came from sum and product, respectively. See the appendix in Section 8 for implementation details.

Next, we recall the de Casteljau algorithm:

---

**Algorithm 1** *de Casteljau algorithm for polynomial evaluation.*

---

```

function res = DeCasteljau( $b, s$ )
     $n = \text{length}(b) - 1$ 
     $\hat{r} = 1 \ominus s$ 

    for  $j = 0, \dots, n$  do
         $\hat{b}_j^{(n)} = b_j$ 
    end for

    for  $k = n - 1, \dots, 0$  do
        for  $j = 0, \dots, k$  do
             $\hat{b}_j^{(k)} = (\hat{r} \otimes \hat{b}_j^{(k+1)}) \oplus (s \otimes \hat{b}_{j+1}^{(k+1)})$ 
        end for
    end for

    res =  $\hat{b}_0^{(0)}$ 
end function

```

---

**Theorem 2** ([MP99]). If  $p(s) = \sum_{j=0}^n b_j B_{j,n}(s)$  and  $\text{DeCasteljau}(p, s)$  is the value computed by the de Casteljau algorithm then

$$|p(s) - \text{DeCasteljau}(p, s)| \leq \gamma_{2n} \sum_{j=0}^n |b_j| B_{j,n}(s). \quad (5)$$

The relative condition number of the evaluation of  $p(s) = \sum_{j=0}^n b_j B_{j,n}(s)$  in Bernstein form used in this paper is (see [MP99], [FR87]):

$$\text{cond}(p, s) = \frac{\tilde{p}(s)}{|p(s)|} = \frac{\sum_j |b_j| B_{j,n}(s)}{|p(s)|}, \quad (6)$$

where  $B_{j,n} \geq 0$ .

### 3 Compensated de Casteljau

In this section we review the compensated de Casteljau algorithm from [JLCS10]. In order to track the local errors at each update step, we use four EFTs:

$$[\widehat{r}, \rho] = \text{TwoSum}(1, -s) \quad (7)$$

$$[P_1, \pi_1] = \text{TwoProd}(\widehat{r}, \widehat{b}_j^{(k+1)}) \quad (8)$$

$$[P_2, \pi_2] = \text{TwoProd}(s, \widehat{b}_{j+1}^{(k+1)}) \quad (9)$$

$$[\widehat{b}_j^{(k)}, \sigma_3] = \text{TwoSum}(P_1, P_2) \quad (10)$$

With these, we can exactly describe the local error between the exact update and computed update:

$$\ell_{1,j}^{(k)} = \pi_1 + \pi_2 + \sigma_3 + \rho \widehat{b}_j^{(k+1)} \quad (11)$$

$$(1-s) \cdot \widehat{b}_j^{(k+1)} + s \cdot \widehat{b}_{j+1}^{(k+1)} = \widehat{b}_j^{(k)} + \ell_{1,j}^{(k)}. \quad (12)$$

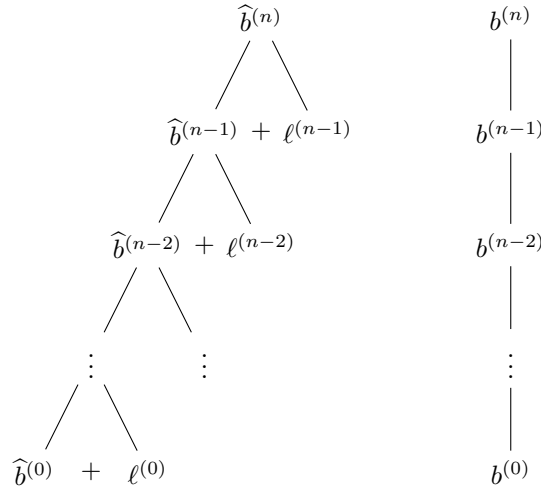


Figure 3: Local round-off errors

By defining the global errors at each step

$$\partial b_j^{(k)} = b_j^{(k)} - \widehat{b}_j^{(k)} \quad (13)$$

we can see that the local errors accumulate in  $\partial b^{(k)}$ :

$$\partial b_j^{(k)} = (1-s) \cdot \partial b_j^{(k+1)} + s \cdot \partial b_{j+1}^{(k+1)} + \ell_{1,j}^{(k)}. \quad (14)$$

Hence, when computed in exact arithmetic, we have

$$p(s) = \widehat{b}_0^{(0)} + \partial b_0^{(0)}. \quad (15)$$

The idea behind the compensated de Casteljau algorithm is to compute both the local error and the updates of the global error with floating point operations:

---

**Algorithm 2** *Compensated de Casteljau algorithm for polynomial evaluation.*

---

```

function res = CompDeCasteljau( $b, s$ )
     $n = \text{length}(b) - 1$ 
     $[\widehat{r}, \rho] = \text{TwoSum}(1, -s)$ 

    for  $j = 0, \dots, n$  do
         $\widehat{b}_j^{(n)} = b_j$ 
         $\widehat{\partial b}_j^{(n)} = 0$ 
    end for

    for  $k = n - 1, \dots, 0$  do
        for  $j = 0, \dots, k$  do
             $[P_1, \pi_1] = \text{TwoProd}(\widehat{r}, \widehat{b}_j^{(k+1)})$ 
             $[P_2, \pi_2] = \text{TwoProd}(s, \widehat{b}_{j+1}^{(k+1)})$ 
             $[\widehat{b}_j^{(k)}, \sigma_3] = \text{TwoSum}(P_1, P_2)$ 
             $\widehat{\ell}_{1,j}^{(k)} = \pi_1 \oplus \pi_2 \oplus \sigma_3 \oplus (\rho \otimes \widehat{b}_j^{(k+1)})$ 
             $\widehat{\partial b}_j^{(k)} = \left( \widehat{r} \otimes \widehat{\partial b}_j^{(k+1)} \right) \oplus \left( s \otimes \widehat{\partial b}_{j+1}^{(k+1)} \right) \oplus \widehat{\ell}_{1,j}^{(k)}$ 
        end for
    end for

    res =  $\widehat{b}_0^{(0)} \oplus \widehat{\partial b}_0^{(0)}$ 
end function

```

---

When comparing this computed error to the exact error, the difference depends only on  $s$  and the Bernstein coefficients of  $p$ :

**Theorem 3** ([JLCS10]). If no underflow occurs and  $n \geq 2$ , then

$$\left| \partial b_0^{(0)} - \widehat{\partial b}_0^{(0)} \right| \leq 2\gamma_{3n+1}\gamma_{3(n-1)} \sum_{j=0}^n |b_j| B_{j,n}(s). \quad (16)$$

Using the bound on the error in  $\partial b$ , the algorithm can be shown to be as accurate as if the computations were done in twice the working precision:

**Theorem 4** ([JLCS10]). If no underflow occurs,  $n \geq 2$  and  $s \in [0, 1]$

$$\frac{|p(s) - \text{CompDeCasteljau}(p, s)|}{|p(s)|} \leq \mathbf{u} + 2\gamma_{3n}^2 \text{cond}(p, s). \quad (17)$$

## 4 $K$ Compensated de Casteljau

In order to raise from twice the working precision to  $K$  times the working precision, we continue using EFTs when computing  $\widehat{\partial b}^{(k)}$ . By tracking the round-off from each floating point evaluation via an EFT, we can form a cascade of global errors:

$$b_j^{(k)} = \widehat{b}_j^{(k)} + \partial b_j^{(k)} \quad (18)$$

$$\partial b_j^{(k)} = \widehat{\partial b}_j^{(k)} + \partial^2 b_j^{(k)} \quad (19)$$

$$\partial^2 b_j^{(k)} = \widehat{\partial^2 b}_j^{(k)} + \partial^3 b_j^{(k)} \quad (20)$$

$\vdots$

In the same way local error can be tracked when updating  $\widehat{b}_j^{(k)}$ , it can be tracked for updates that happen down the cascade:

$$(1-s) \cdot \widehat{b}_j^{(k+1)} + s \cdot \widehat{b}_{j+1}^{(k+1)} = \widehat{b}_j^{(k)} + \ell_{1,j}^{(k)} \quad (21)$$

$$(1-s) \cdot \widehat{\partial b}_j^{(k+1)} + s \cdot \widehat{\partial b}_{j+1}^{(k+1)} + \ell_{1,j}^{(k)} = \widehat{\partial b}_j^{(k)} + \ell_{2,j}^{(k)} \quad (22)$$

$$(1-s) \cdot \widehat{\partial^2 b}_j^{(k+1)} + s \cdot \widehat{\partial^2 b}_{j+1}^{(k+1)} + \ell_{2,j}^{(k)} = \widehat{\partial^2 b}_j^{(k)} + \ell_{3,j}^{(k)} \quad (23)$$

$\vdots$

In **CompDeCasteljau**, after a single stage of error filtering we “give up” and use  $\widehat{\partial b}^{(k)}$  instead of  $\partial b^{(k)}$  (without keeping around any information about the round-off error). In order to obtain results that are as accurate as if computed in  $K$  times the working precision, we must continue filtering errors down  $(K-1)$  times, so that at the final level we compute  $\widehat{\partial^{K-1} b}^{(k)}$  and don’t worry about round-off error.

When computing  $\widehat{\partial^F b}$  (i.e. the error after  $F$  stages of filtering) there will be several sources of round-off. In particular, there will be

- errors when computing  $\widehat{\ell}_{F,j}^{(k)}$  from the terms in  $\ell_{F,j}^{(k)}$
- an error for the “missing”  $\rho \cdot \widehat{\partial^F b}_j^{(k+1)}$  in  $(1-s) \cdot \widehat{\partial^F b}_j^{(k+1)}$
- an error from the product  $\widehat{r} \otimes \widehat{\partial^F b}_j^{(k+1)}$
- an error from the product  $s \otimes \widehat{\partial^F b}_{j+1}^{(k+1)}$
- two errors from the two  $\oplus$  when combining the three terms in  $\left[ \widehat{r} \otimes \widehat{\partial^F b}_j^{(k+1)} \right] \oplus \left[ s \otimes \widehat{\partial^F b}_{j+1}^{(k+1)} \right] \oplus \widehat{\ell}_{F,j}^{(k)}$

For example, in (11):

$$\ell_{1,j}^{(k)} = \underbrace{\pi_1}_{P_1 = \widehat{r} \otimes \widehat{b}_j^{(k+1)}} + \underbrace{\pi_2}_{P_2 = s \otimes \widehat{b}_{j+1}^{(k+1)}} + \underbrace{\sigma_3}_{P_1 \oplus P_2} + \underbrace{\rho \cdot \widehat{b}_j^{(k+1)}}_{(1-s) \widehat{b}_j^{(k+1)}} \quad (24)$$

After each stage, we’ll always have

$$\ell_{F,j}^{(k)} = e_1 + e_2 + \dots + e_M + \rho \cdot \widehat{\partial^{F-1} b}_j^{(k+1)}$$

where the terms  $e_1, \dots, e_M$  come from using **TwoSum** and **TwoProd** when computing  $\widehat{\partial^{F-1} b}_j^{(k)}$  and the  $\rho$  term comes from the roundoff in  $1 \ominus s$  when multiplying  $(1-s)$  by  $\widehat{\partial^{F-1} b}_j^{(k+1)}$ . With this in mind, we can define an EFT (**LocalErrorEFT**) that computes  $\widehat{\ell}$  and tracks all round-off errors generated in the process:

---

**Algorithm 3** *EFT for computing the local error.*

---

**function**  $[\eta, \widehat{\ell}] = \text{LocalErrorEFT}(e, \rho, \delta b)$

$M = \text{length}(e)$

$[\widehat{\ell}, \eta_1] = \text{TwoSum}(e_1, e_2)$

**for**  $j = 3, \dots, M$  **do**

$[\widehat{\ell}, \eta_{j-1}] = \text{TwoSum}(\widehat{\ell}, e_j)$

**end for**

$[P, \eta_M] = \text{TwoProd}(\rho, \delta b)$

$[\widehat{\ell}, \eta_{M+1}] = \text{TwoSum}(\widehat{\ell}, P)$

**end function**

---

With this EFT helper in place, we can perform  $(K - 1)$  error filtrations:

---

**Algorithm 4**  $K$ -compensated de Casteljau algorithm.

---

```

function res = CompDeCasteljau( $b, s$ )
     $n = \text{length}(b) - 1$ 
     $[\widehat{r}, \rho] = \text{TwoSum}(1, -s)$ 

    for  $j = 0, \dots, n$  do
         $\widehat{b}_j^{(n)} = b_j$ 
        for  $F = 1, \dots, K - 1$  do
             $\widehat{\partial^F b_j}^{(n)} = 0$ 
        end for
    end for

    for  $k = n - 1, \dots, 0$  do
        for  $j = 0, \dots, k$  do
             $[P_1, \pi_1] = \text{TwoProd}(\widehat{r}, \widehat{b}_j^{(k+1)})$ 
             $[P_2, \pi_2] = \text{TwoProd}(s, \widehat{b}_{j+1}^{(k+1)})$ 
             $[\widehat{b}_j^{(k)}, \sigma_3] = \text{TwoSum}(P_1, P_2)$ 

             $e = [\pi_1, \pi_2, \sigma_3]$ 
             $\delta b = \widehat{b}_j^{(k+1)}$ 

            for  $F = 1, \dots, K - 2$  do
                 $[\eta, \ell] = \text{LocalErrorEFT}(e, \rho, \delta b)$ 
                 $M = \text{length}(\eta)$ 

                 $[P_1, \eta_{M+1}] = \text{TwoProd}(\widehat{r}, \widehat{\partial^F b_j}^{(k+1)})$ 
                 $[P_2, \eta_{M+2}] = \text{TwoProd}(s, \widehat{\partial^F b_{j+1}}^{(k+1)})$ 
                 $[S_3, \eta_{M+3}] = \text{TwoSum}(P_1, P_2)$ 
                 $[\widehat{\partial^F b_j}^{(k)}, \eta_{M+4}] = \text{TwoSum}(S_3, \ell)$ 

                 $e = \eta$ 
                 $\delta b = \widehat{\partial^F b_j}^{(k+1)}$ 
            end for

             $\widehat{\ell} = \text{LocalError}(e, \rho, \delta b)$ 
             $\widehat{\partial^{K-1} b_j}^{(k)} = \left( \widehat{r} \otimes \widehat{\partial^{K-1} b_j}^{(k+1)} \right) \oplus \left( s \otimes \widehat{\partial^{K-1} b_{j+1}}^{(k+1)} \right) \oplus \widehat{\ell}$ 
        end for
    end for

    res =  $\widehat{b}_0^{(0)}$ 
    for  $F = 1, \dots, K - 1$  do
        res = res  $\oplus \widehat{\partial^F b_0}^{(0)}$ 
    end for
end function

```

---



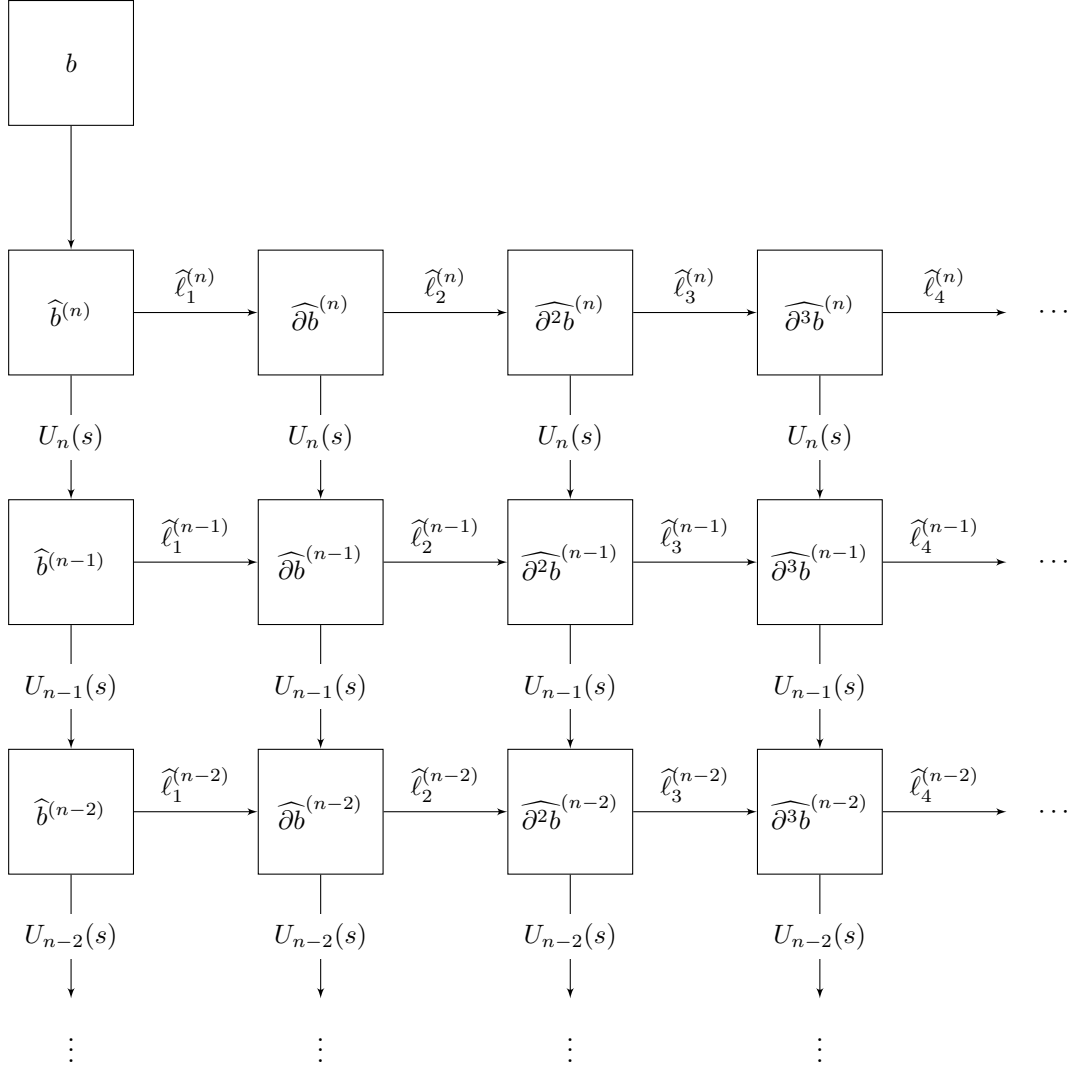


Figure 4: Filtering errors

## 5 Numerical experiments

Figure 5 and Figure 6 show how things improve.

## 6 de Casteljau's Method

Consider de Casteljau's method to evaluate a degree  $n$  polynomial in Bernstein-Bézier form with control points  $b_j$ :

$$b_j^{(n)} = b_j \tag{25}$$

$$b_j^{(k)} = (1-s)b_j^{(k+1)} + sb_{j+1}^{(k+1)} \tag{26}$$

$$b(s) = b_0^{(0)}. \tag{27}$$

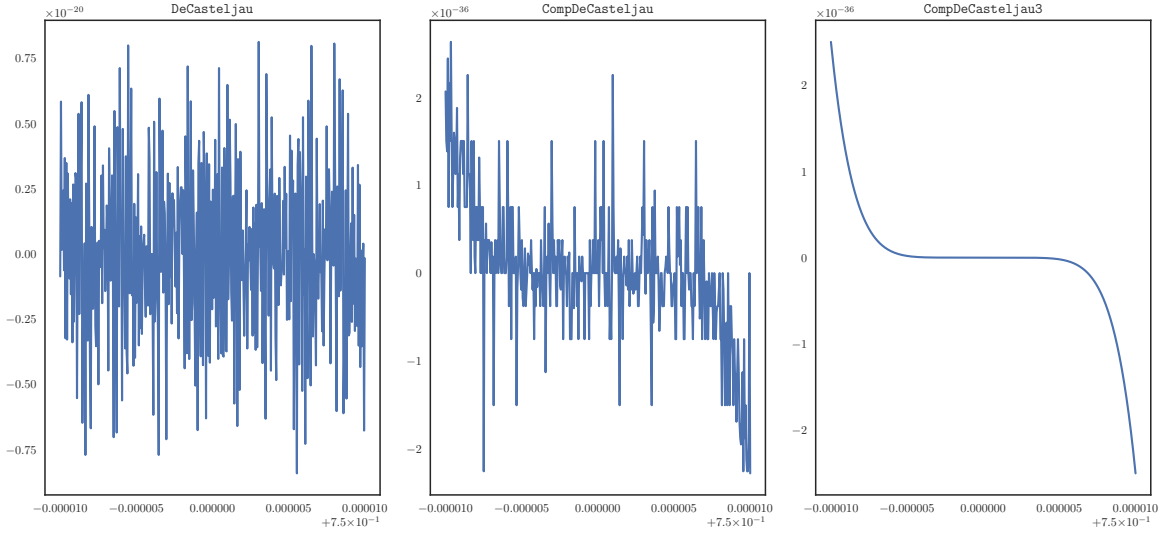


Figure 5: Evaluation of  $p(s) = (s - 1)(s - 3/4)^7$  in the neighborhood of its multiple root  $3/4$ .

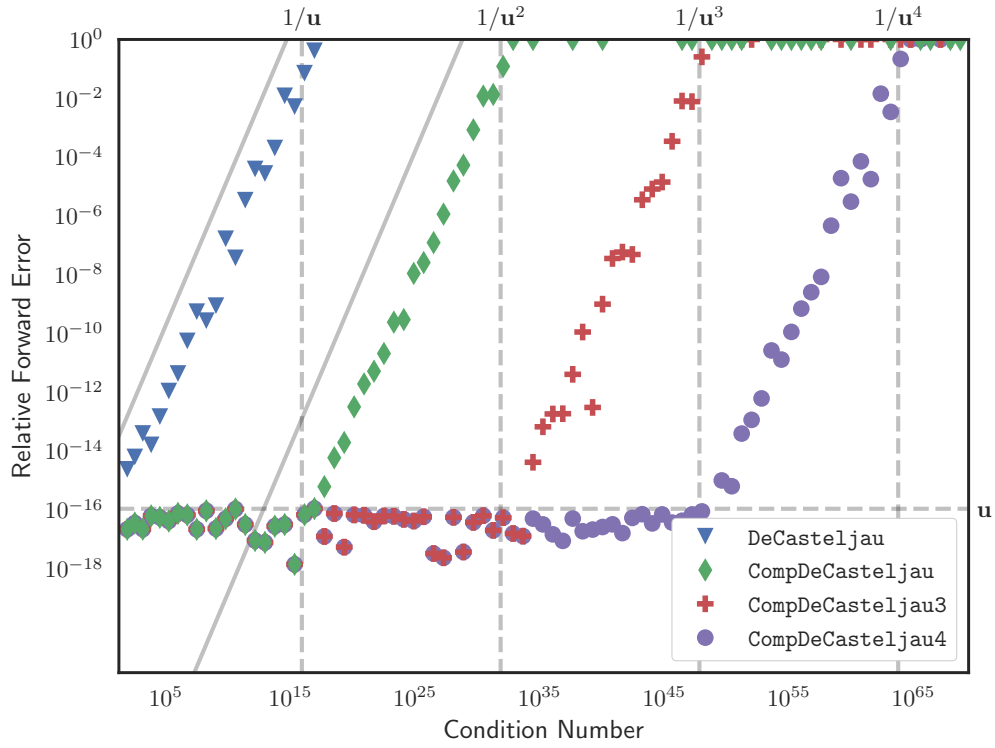


Figure 6: Accuracy of evaluation of  $p(s) = (s - 1)(s - 3/4)^7$  represented in Bernstein form.

## 6.1 Condition Number

For a polynomial  $p(x)$  in the power basis, we have ([LGL06]):

$$\text{cond}(p(x)) = \frac{\tilde{p}(|x|)}{|p(x)|} = \frac{\sum_j |a_j| |x|^j}{|p(x)|}. \quad (28)$$

In particular, this means that if  $x \geq 0$  and each  $a_j \geq 0$  we must necessarily have  $\text{cond}(p(x)) = 1$ . To see an example in use, consider  $p(x) = (x-1)^n$  and input values of the form  $x = 1 + \delta$  (with  $|\delta| \ll 1$ ). Since  $a_j = \binom{n}{j}(-1)^{n-j}$  we have  $\tilde{p}(x) = (x+1)^n$  hence

$$\text{cond}(p(1+\delta)) = \frac{(2+\delta)^n}{|\delta|^n} = \left|1 + \frac{2}{\delta}\right|^n. \quad (29)$$

As  $\delta \rightarrow 0$ , this value approaches  $\infty$  (as expected).

For a polynomial  $p(s)$  in Bernstein form, we have ([JLCS10]):

$$\text{cond}(p(s)) = \frac{\tilde{p}(s)}{|p(s)|} = \frac{\sum_j |b_j| |B_{j,n}(s)|}{|p(s)|}. \quad (30)$$

The Bernstein form is suited for  $s \in [0, 1]$ , which means  $B_{j,n}(s) \geq 0$  typically. If  $s \in [0, 1]$  and each  $b_j \geq 0$  we must necessarily have  $\text{cond}(p(s)) = 1$ . To see an example in use, consider

$$p(s) = (1-2s)^n = [(1-s)-s]^n = \sum_j \binom{n}{j} (1-s)^{n-j} (-s)^j = \sum_j (-1)^j B_{j,n}(s) \quad (31)$$

and input values of the form  $x = \frac{1}{2} + \delta$  (with  $|\delta| \ll \frac{1}{2}$ ). Since  $b_j = (-1)^j$  we have  $\tilde{p}(s) = [(1-s)+s]^n = 1$

$$\text{cond}\left(p\left(\frac{1}{2} + \delta\right)\right) = \frac{1}{|2\delta|^n}. \quad (32)$$

As  $\delta \rightarrow 0$ , this value approaches  $\infty$  (as expected).

## 6.2 Example of Compensated de Casteljau Failing

Consider

$$p(s) = (4s-3)^3(8s+7) = (-189)B_{0,4}(s) + (-54)B_{1,4}(s) + 57B_{2,4}(s) + (-32)B_{3,4}(s) + 15B_{4,4}(s) \quad (33)$$

at the point  $s = \frac{3}{4} + 800\mathbf{u}$ . In exact arithmetic, we have  $p(s) = 13(3200\mathbf{u})^3 + \mathcal{O}(\mathbf{u}^4)$ . However, using the compensated de Casteljau algorithm results in  $\widehat{b}_0^{(0)} \oplus \widehat{\partial} b_0^{(0)} = 0$ :

$k$	$j$	$\widehat{b}_j^{(k)}$	$\widehat{\partial} b_j^{(k)}$	$\widehat{\partial^2} b_j^{(k)}$	$\widehat{\partial^3} b_j^{(k)}$
3	0	$-87\frac{3}{4} + 108032\mathbf{u}$	$-32\mathbf{u}$	0	0
3	1	$29\frac{1}{4} + 88768\mathbf{u}$	$32\mathbf{u}$	0	0
3	2	$-9\frac{3}{4} - 71200\mathbf{u}$	0	0	0
3	3	$3\frac{1}{4} + 37600\mathbf{u}$	0	0	0
2	0	$187200\mathbf{u}$	$-15360000\mathbf{u}^2$	0	0
2	1	$-62408\mathbf{u}$	$8\mathbf{u} - 128000000\mathbf{u}^2$	0	0
2	2	$20800\mathbf{u}$	$87040000\mathbf{u}^2$	0	0
1	0	$-6\mathbf{u} - 199688192\mathbf{u}^2$	$6\mathbf{u} - 99831808\mathbf{u}^2$	$-90112000000\mathbf{u}^3$	0
1	1	$-2\mathbf{u} + 66568192\mathbf{u}^2$	$2\mathbf{u} + 33271808\mathbf{u}^2$	$172032000000\mathbf{u}^3$	0
0	0	$-3\mathbf{u} + 7296\mathbf{u}^2$	$3\mathbf{u} - 7296\mathbf{u}^2$	$13(3200\mathbf{u})^3 + 3048(512\mathbf{u})^4$	$962(128\mathbf{u})^4$

### 6.3 Selection of Test Cases

From [DP15] (end of Section 3):

We can observe that, in this case, the algorithm with a good behavior everywhere is the de Casteljau algorithm

In the same paper (when referring to [Bez13] at the beginning of Section 2):

assuming that all control points are positive. This assumption avoided ill-conditioned polynomials. In this section, we shall show that this is a natural assumption in Computer Aided Geometric Design (from now on, C.A.G.D.) and that it permits to assure high relative precision for the evaluation through a large family of representations in C.A.G.D.

From the same author, in [MP05] (towards the end of Section 5, at the bottom of page 109):

Let us observe that in this case, the de Casteljau algorithm presents better stability properties for the evaluation near the roots. In fact, the de Casteljau algorithm has good behaviour even when using simple precision, although the running error bound is not so accurate in points close to the roots.

### 6.4 Operation Counts

After implementing for  $K = 2, 3, \dots, 12$  and instrumenting all relevant floating point operations, the  $K$ -fold Horner requires

$$(5 \cdot 2^K - 8)n + ((K + 8)2^K - 12K - 6) = \mathcal{O}((n + K)2^K) \quad (34)$$

flops to evaluate a degree  $n$  polynomial (this only applies when  $n \geq K - 1$ ). As a comparison, the non-compensated form of Horner requires  $2n$  flops. Of these,  $(2^{K-1} - 1)n - 2^{K-1}(K - 3) - 2$  are FMA (fused-multiply-add) instructions.

After implementing for  $K = 2, 3, 4, 5$  and instrumenting all relevant floating point operations, the  $K$ -fold de Casteljau requires

$$(15K^2 - 34K + 26)T_n + K + 5 = \mathcal{O}(n^2K^2) \quad (35)$$

flops to evaluate a degree  $n$  polynomial. (Here  $T_n$  is the  $n$ th triangular number.) As a comparison, the non-compensated form of de Casteljau requires  $3T_n + 1$  flops. Of these,  $(3K - 4)T_n$  are FMA instructions. On hardware that doesn't support FMA, every FMA will be exchanged for 10  $\ominus$ 's and 6  $\otimes$ 's so the count will increase by  $(10 + 6 - 1)(3K - 4)T_n$ .

## 7 Bogus Section for Refs

Here they are, for now

- Newton with compensated Horner ([Gra08])

## References

- [Bez13] Licio Hernanes Bezerra. Efficient computation of Bézier curves from their Bernstein–Fourier representation. *Applied Mathematics and Computation*, 220:235–238, Sep 2013.
- [Dek71] T. J. Dekker. A floating-point technique for extending the available precision. *Numerische Mathematik*, 18(3):224–242, jun 1971.
- [DP15] Jorge Delgado and J.M. Peña. Accurate evaluation of Bézier curves and surfaces and the Bernstein–Fourier algorithm. *Applied Mathematics and Computation*, 271:113–122, Nov 2015.
- [FR87] R.T. Farouki and V.T. Rajan. On the numerical condition of polynomials in Bernstein form. *Computer Aided Geometric Design*, 4(3):191–216, Nov 1987.

- [GLL09] Stef Graillat, Philippe Langlois, and Nicolas Louvet. Algorithms for accurate, validated and fast polynomial evaluation. *Japan Journal of Industrial and Applied Mathematics*, 26(2-3):191–214, Oct 2009.
- [Gra08] Stef Graillat. Accurate simple zeros of polynomials in floating point arithmetic. *Computers & Mathematics with Applications*, 56(4):1114–1120, Aug 2008.
- [Hig02] Nicholas J. Higham. *Accuracy and Stability of Numerical Algorithms*. Society for Industrial and Applied Mathematics, jan 2002.
- [JLCS10] Hao Jiang, Shengguo Li, Lizhi Cheng, and Fang Su. Accurate evaluation of a polynomial and its derivative in Bernstein form. *Computers & Mathematics with Applications*, 60(3):744–755, Aug 2010.
- [Knu97] Donald E. Knuth. *Art of Computer Programming, Volume 2: Seminumerical Algorithms (3rd Edition)*. Addison-Wesley Professional, 1997.
- [LGL06] Philippe Langlois, Stef Graillat, and Nicolas Louvet. Compensated Horner Scheme. In Bruno Buchberger, Shin'ichi Oishi, Michael Plum, and Sigfried M. Rump, editors, *Algebraic and Numerical Algorithms and Computer-assisted Proofs*, number 05391 in Dagstuhl Seminar Proceedings, Dagstuhl, Germany, 2006. Internationales Begegnungs- und Forschungszentrum für Informatik (IBFI), Schloss Dagstuhl, Germany.
- [MP99] E. Mainar and J.M. Peña. Error analysis of corner cutting algorithms. *Numerical Algorithms*, 22(1):41–52, 1999.
- [MP05] E. Mainar and J. M. Peña. Running Error Analysis of Evaluation Algorithms for Bivariate Polynomials in Barycentric Bernstein Form. *Computing*, 77(1):97–111, Dec 2005.
- [ORO05] Takeshi Ogita, Siegfried M. Rump, and Shin'ichi Oishi. Accurate Sum and Dot Product. *SIAM Journal on Scientific Computing*, 26(6):1955–1988, Jan 2005.

## 8 Appendix: EFT Algorithms

Find here concrete implementation details on the EFTs described in Theorem 1. They do not use branches, nor access to the mantissa that can be time-consuming

---

**Algorithm 5** *EFT of the sum of two floating point numbers.*

---

```

function  $[S, \sigma] = \text{TwoSum}(a, b)$ 
     $S = a \oplus b$ 
     $z = S \ominus a$ 
     $\sigma = (a \ominus (S \ominus z)) \oplus (b \ominus z)$ 
end function
```

---

In order to avoid branching to check which among  $|a|, |b|$  is largest, `TwoSum` uses 6 flops rather than 3.

---

**Algorithm 6** *Splitting of a floating point number into two parts.*

---

```

function  $[h, \ell] = \text{Split}(a)$ 
     $z = a \otimes (2^r + 1)$ 
     $h = z \ominus (z \ominus a)$ 
     $\ell = a \ominus h$ 
end function
```

---

For IEEE-754 double precision floating point number,  $r = 26$  so  $2^r + 1$  will be known before `Split` is called. In all, `Split` uses 4 flops.

---

**Algorithm 7** *EFT of the product of two floating point numbers.*


---

```

function  $[P, \pi] = \text{TwoProd}(a, b)$ 
   $P = a \otimes b$ 
   $[a_h, a_\ell] = \text{Split}(a)$ 
   $[b_h, b_\ell] = \text{Split}(b)$ 
   $\pi = a_\ell \otimes b_\ell \ominus (((P \ominus a_h \otimes b_h) \ominus a_\ell \otimes b_h) \ominus a_h \otimes b_\ell)$ 
end function

```

---

This implementation of **TwoProd** requires 17 flops. For processors that provide a fused-multiply-add operator (FMA), **TwoProd** can be rewritten to use only 2 flops:

---

**Algorithm 8** *EFT of the sum of two floating point numbers with a FMA.*


---

```

function  $[P, \pi] = \text{TwoProdFMA}(a, b)$ 
   $P = a \otimes b$ 
   $\pi = \text{FMA}(a, b, -P)$ 
end function

```

---