



## Extension GL : **Trigonométrie**

Groupe GL 35

Janvier 2023

## Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
1.1	Les flottants Deca: flottant simple précision . . . . .	4
1.2	Stratégie adoptée . . . . .	4
1.3	Méthode de validation . . . . .	5
<b>2</b>	<b>ULP</b>	<b>6</b>
2.1	Principe . . . . .	6
2.2	Résultats . . . . .	6
2.3	Graphes . . . . .	6
<b>3</b>	<b>Cos et Sin</b>	<b>8</b>
3.1	Choix d'implémentation . . . . .	8
3.2	Difficultés et amélioration . . . . .	8
3.3	Résultats sin . . . . .	9
3.4	Résultats cos . . . . .	9
3.5	Graphes . . . . .	9
<b>4</b>	<b>Asin</b>	<b>12</b>
4.1	Choix d'implémentation . . . . .	12
4.2	Difficultés et amélioration . . . . .	13
4.3	Résultats . . . . .	13
4.4	Graphes . . . . .	13
<b>5</b>	<b>Atan</b>	<b>15</b>
5.1	Choix d'implémentation . . . . .	15
5.2	Difficultés et amélioration . . . . .	15
5.3	Résultats . . . . .	15
5.4	Graphes . . . . .	16

## 1 Introduction

L'objectif de cette partie est de créer une classe Maths dans laquelle on implémente des méthodes trigonométriques, tout en obtenant une bonne précision pour chaque fonction, en différenciant selon les sous-domaines du domaine de définition suivantes :

- `float ulp(float f)`
- `float sin(float f)`
- `float cos(float f)`
- `float asin(float f)`
- `float atan(float f)`

Nous avons du implémenté d'autres fonctions qu'on a pensées nécessaire pour l'implémentation des méthodes demandées, à savoir :

- `float puiss(float a, int b)`
- `float abs(float f)`
- `float fact(float f)`
- `float sqrt(float f)`

Cette partie d'extension sera abordé comme suit, d'abord une implémentation de la fonction ULP, qui sera nécessaire pour le calcul de précision tout au long du travail sur l'extension. Ensuite un choix justifié d'algorithmes pour implémenter les différentes méthodes trigonométriques, en cherchant le bon compromis entre précision et quantité de calculs de la machine.

### 1.1 Les flottants Deca: flottant simple précision

D'après la norme IEEE 754, un nombre flottant simple précision est stocké dans un mot de 32 bits : 1 bit de signe, 8 bits pour l'exposant et 23 pour la mantisse.

Cette représentation impose que l'exposant d'un nombre normalisé va donc de  $-126$  à  $+127$ .



Figure 1: Format simple précision 32bits

### 1.2 Stratégie adoptée

Notre stratégie consiste à implémenter dans un premier temps nos algorithmes dans java, afin de les tester et d'avoir des idées concrètes sur les précisions attendues, tout en respectant les limitations de l'Ima. Et aussi pour ne pas avoir des retards concernant les délais de rendus, car attendre que le compilateur soit prêt ,ne va pas nous laisser beaucoup de temps pour tester.

Nous avons pensé aussi à faire une deuxième couche de tests après la fin de la partie essentielle du compilateur avec objet. En effet, dès que le compilateur était prêt, nous avons implémenté les algorithmes des fonctions trigonométriques sur Deca, tout en les adaptant pour qu'ils restent fonctionnels.

Tout au long des tests, nous avons priorisé le traçage des courbes, car l'analyse des graphes permet de mieux détecter les problèmes. Les courbes qu'on a tracées ont été créées par Latex, après création de fichiers (.csv) sous java, dans lesquelles on a stockés nos résultats.

#### Remarque:

L'écriture flottante de java est limitée à 8 chiffres après la virgule, contrairement à l'écriture hexadécimale du Deca qui n'est pas limitée. Ce qui donne une meilleure évaluation de la précision sur Deca par les mêmes algorithmes.

### 1.3 Méthode de validation

Nous tenons à préciser que notre objectif est d'avoir une fonction ULP qui marche parfaitement, pour ne pas influencer le calcul des précisions qui vient après.

Nous visons aussi des précisions de quelques ULP pour les fonctions trigonométriques.

On a aussi tenu compte dans nos implémentations de réduire les calculs au maximum, en faisant par exemple beaucoup de tests qui nous ont permis de choisir le bon intervalle de réduction, ou encore en diminuant au maximum les degrés de polynômes utilisés en développement en série entière.

Pour le calcul d'erreur des fonctions implémentés, nous avons calculé la différence en valeur absolue entre la valeur calculé par notre méthode et celle de la classe Math du Java, le tout divisé par Ulp de la valeur réelle (celle calculée par la classe du Java).

## 2 ULP

### 2.1 Principe

ULP: Unit in the Last Place, est l'écart entre deux nombres consécutifs à virgule flottante, c'est-à-dire la valeur que représente le chiffre le moins significatif (chiffre le plus à droite) s'il est égal à 1. Elle est utilisée comme mesure de précision dans les calculs numériques.

La méthode ULP implémenté dans ce projet, est inspirée de l'algorithme donné pour un double dans la thèse de Monsieur Jean-Michel Muller:[these] qu'on a adapté aux flottants, et auquel on a ajouté des modifications, après de nombreux tests, pour avoir les résultats attendus.

Cet algorithme vérifie d'abord que le flottant est bien dans l'intervalle  $[2^{-126}; 2^{127}]$ , sinon, s'il est plus grand il retourne la valeur max  $2^{103}$ , et s'il est plus petit, il retourne la valeur du nombre dé-normalisé non nul le plus proche de zéro qui est :  $2^{-149}$ .

Ensuite, il utilise une recherche binaire afin d'encadrer le flottant comme suit :  $2^i \leq x \leq 2^{i+1}$  pour trouver l'exposant le plus proche du nombre donné, et finalement il renvoie 2 à la puissance l'exposant trouvé.

Notre compilateur ne supporte pas la déclaration de variables flottantes plus grandes que du  $2^{31} - 1$  en écriture de  $2^x$ , donc on a du écrire en dur toutes les valeurs qui dépassent cette puissance à l'aide de leurs écriture scientifiques.

### 2.2 Résultats

Les résultats obtenus avec la fonction ULP sont parfaitement précis, vu que l'erreur avec l' ULP de la classe Math de Java est nulle partout.

Le tableau suivant résume ces résultats.

Intervalle	Pas	Nombre d'échantillons	Erreur entre notre ulp et Math.ulp
$[0; 10]$	$3.8146973 * 10^{-6}$	2621440	0
$[2^{-40}; 2^{-7}]$	$2^{-28}$	2097152	0
$[2^{-7}; 2^{12}]$	$9.765625 * 10^{-4}$	4194296	0
$[2^{12}; 2^{40}]$	$9.765625 * 10^{-4}$	4194296	0

### 2.3 Graphes

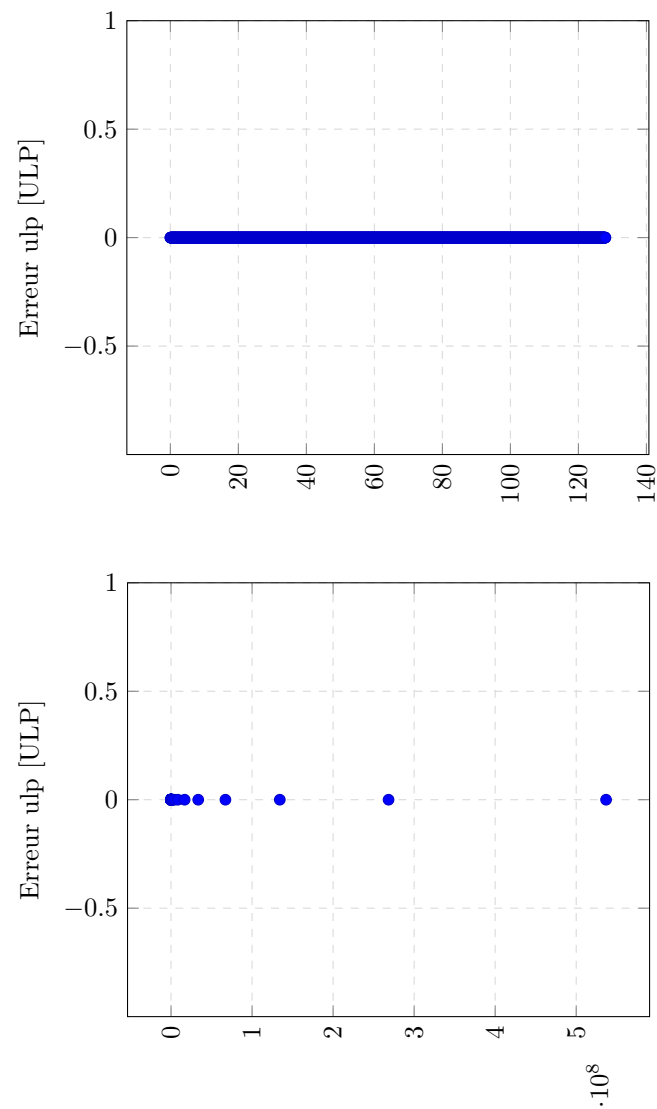


Figure 2: ULP plots

### 3 Cos et Sin

#### 3.1 Choix d'implémentation

Pour ces deux fonctions trigonométriques, on a choisi de les implémenter à l'aide de séries de Taylor. La définition de la série de Taylor est comme suit : Soit  $f$  une fonction d'une variable réelle ou complexe, indéfiniment dérivable en un point  $a$ . La série de Taylor de  $f$  en  $a$  est la série de fonctions :

$$f(a) + \frac{f'(a)}{1!}(x-a) + \frac{f''(a)}{2!}(x-a)^2 + \frac{f^{(3)}(a)}{3!}(x-a)^3 + \dots$$

**Choix du nombre de coefficients :** La courbe de la fonction cos implémentée par 8 coefficients de la série de Taylor coïncide parfaitement avec la courbe de cos sur l'intervalle de  $[-\pi; \pi]$  (et ça va de même pour la courbe de sin). Ce qui nous permet de travailler dans un premier temps sur cet intervalle, et de l'élargir à  $[-2\pi; 2\pi]$  à l'aide des propriétés de cos et de sin suivantes :

- $\cos(x+\pi) = -\cos(x)$
- $\sin(x+\pi) = -\sin(x)$

On élargit l'intervalle après sur  $\mathbb{R}$ , en retranchant à chaque fois la valeur  $2\pi$  des nombres hors notre intervalle.

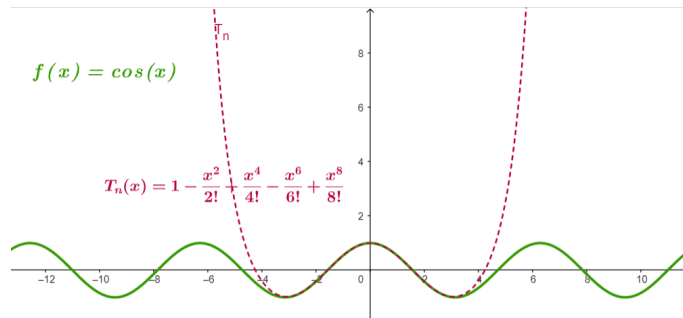


Figure 3: Graphe représentant la fonction : cos et sa série de Taylor à l'ordre 8

#### 3.2 Difficultés et amélioration

Après cette première implémentation, on s'est avéré que prendre un intervalle de départ de  $[-\pi; \pi]$ , qui est assez grand entraîne des ULP d'une valeur de 2 de temps en temps. Alors on a réduit l'intervalle encore plus pour ne travailler que sur  $[0; \pi/4]$ , et déduire les valeurs du reste de l'intervalle  $[-2\pi; 2\pi]$  à l'aide des relations trigonométriques :



- $\cos(\pi/2-x) = \sin(x)$
- $\sin(\pi/2-x) = \cos(x)$
- $\cos(\pi-x) = -\cos(x)$
- $\sin(\pi-x) = \sin(x)$

Alors que pour le reste de l'intervalle, on fait la même chose qu'à la première démarche.

### 3.3 Résultats sin

D'après les tests effectués pour la fonction sin sur Java, on trouve:

Intervalle	step	erreur en ULP	max erreur	nombres d'éléments
$[pi/4; \pi/2]$	$2^{-23}$	0.43120888	3.0	6588398
$[pi/4; \pi/2]$	$2^{-18}$	0.42268613	2.0	205888
$[pi/2; \pi]$	$2^{-10}$	0.22364366	0.75	1609
$[pi; 2\pi]$	$2^{-10}$	0.22142512	0.75	3217
$[2\pi; 4\pi]$	$2^{-10}$	0.20667103	0.5	6434
$[50; 50 + 2\pi]$	$2^{-10}$	0.14266211	0.328125	6434

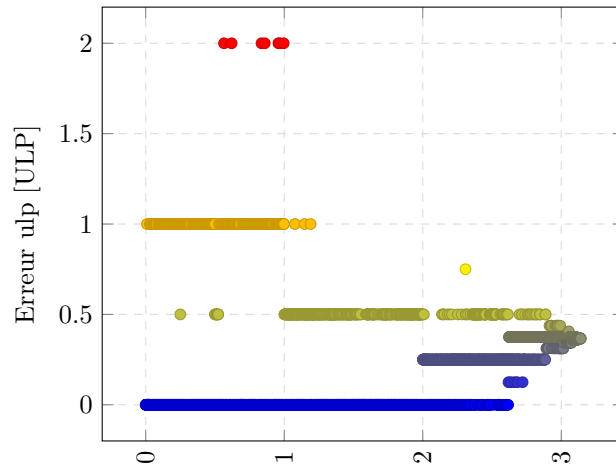
### 3.4 Résultats cos

D'après les tests effectués pour la fonction cos sur Java, on trouve:

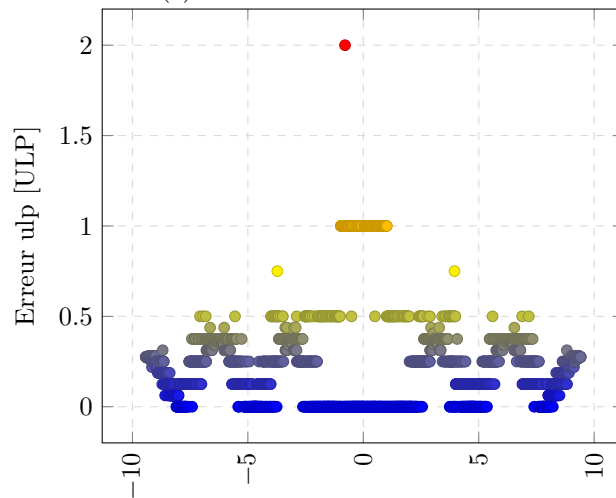
Intervalle	step	erreur en ULP	max erreur	nombres d'éléments
$[0; \pi/4]$	$2^{-18}$	0.33533767	2.0	205888
$[pi/4; \pi/2]$	$2^{-23}$	0.43120888	3.0	6588398
$[pi/2; \pi]$	$2^{-10}$	0.20397955	0.75	1609
$[pi; 2\pi]$	$2^{-10}$	0.20091482	0.75	3217
$[2\pi; 4\pi]$	$2^{-10}$	0.21499515	0.625	6434
$[50; 50 + 2\pi]$	$2^{-10}$	0.16506171	0.37304688	6434

### 3.5 Graphes

Ci-dessous, les graphes représentant les résultats des deux fonctions Cos et Sin, avant et après optimisations.

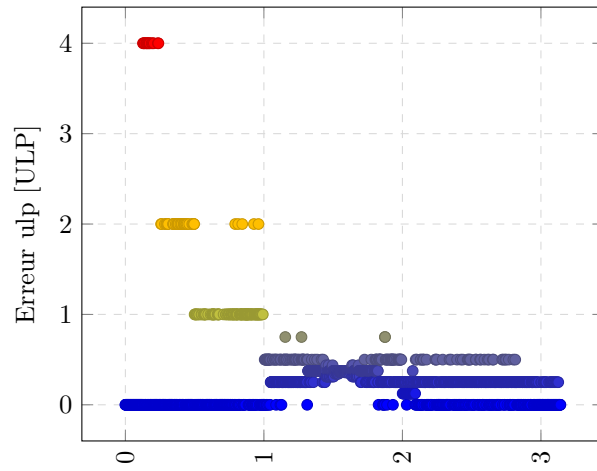


(a) Sin Plot avant correction

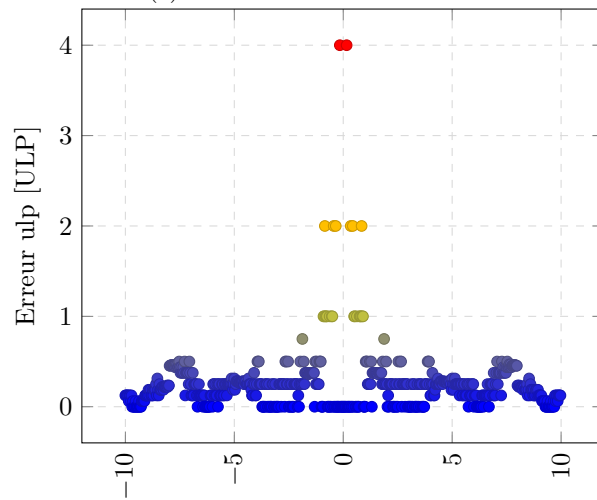


(b) Sin Plot Final

Figure 4: Sin Plots



(a) Cos Plot avant correction



(b) Cos Plot Final

Figure 5: Cos Plots

## 4 Asin

### 4.1 Choix d'implémentation

Le choix d'algorithme pour cette fonction, a demandé beaucoup de tests de différents algorithmes. On en cite un qui était le plus performant entre eux et qu'on a changé à la fin.

Pour cette première méthode, on a utilisé 2 implémentations différentes selon l'intervalle de départ. Sur l'intervalle  $[0; 0.6]$  on a utilisé un développement en série entière à l'ordre 21 de la fonction *arcsin*, alors que sur l'intervalle restant de  $[0.6; 1]$ , on a opté pour une approximation polynomiale, en s'aidant aussi du polynôme de Horner pour réduire le cout des opérations. Et on profite de l'impairité de l'*arcsin*, pour déduire les valeurs de l'*arcsin* sur l'intervalle  $[-1; 0]$ :

$$\text{arcsin}(x) = -\text{arcsin}(-x);$$

Mais on s'est aperçu que cette méthode donnait des erreurs qui s'élevaient à du  $10^{-5}$  autour de  $1/\sqrt{2}$ , ce qui représente un très grand nombre d'ULP(6a).

On a donc décidé de procéder autrement, en n'utilisant qu'un développement en série de Taylor sur l'intervalle  $[0; 1/\sqrt{2}]$  et de profiter de la relation trigonométrique de l'*arcsin*, pour déduire les valeurs sur l'intervalle  $[1/\sqrt{2}; 1]$ , et ainsi ne calculer que des valeurs proches de 0 (6b):

$$\text{arcsin}(x) + \text{arcsin}(\sqrt{1-x^2}) = \pi/2;$$

Ce qui donne des meilleurs résultats, mais encore une fois, le problème persiste un peu au alentour de  $1/\sqrt{2}$ .

## 4.2 Difficultés et amélioration

Ce problème a été dû au manque de termes dans notre développement en série entière, qui s'élevait à l'ordre **21**, mais qu'on a trouvé insuffisant lors du traçage du graphe, vu qu'on remarquait une divergence de l'erreur autour de  $1/\sqrt{2}$ .

Afin de le corriger, on a poussé un peu plus le développement en série entière, vers un polynôme de degré 31, bien que cela ajoute plus de calcul (7). Et on a laissé l'intervalle  $[0; 1/\sqrt{2}]$  comme intervalle d'étude initial.

## 4.3 Résultats

Ci-dessous, les résultats obtenus par les tests effectués sur Java:

Intervalle	step	erreur en ULP	max erreur	nombres d'éléments
$[0; 1/\sqrt{2}]$	$2^{-23}$	<b>0.55730605</b>	<b>5.0</b>	<b>5931642</b>
$[1/\sqrt{2}; 1]$	$2^{-23}$	<b>0.5891706</b>	<b>7.0</b>	<b>3439330</b>

## 4.4 Graphes

Les courbes ci-après résument les résultats de la fonction Asin avant et après optimisations.

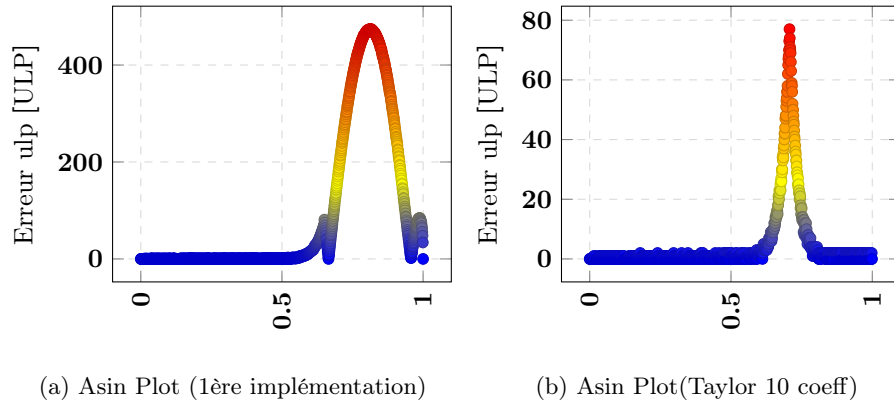


Figure 6: Asin Plots

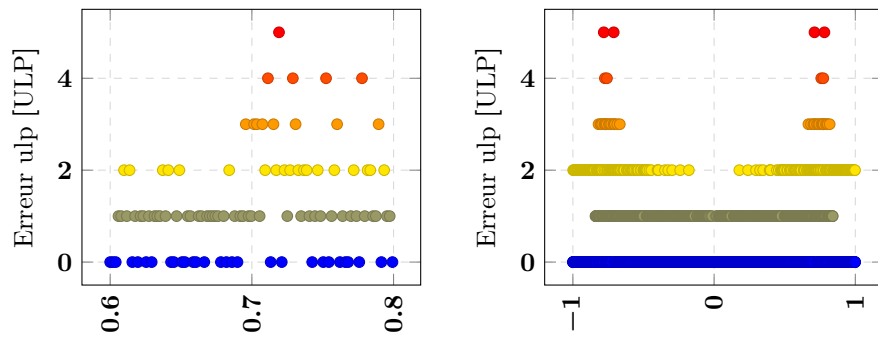


Figure 7: Asin Plots finals(15 coefficients Taylor)

## 5 Atan

### 5.1 Choix d'implémentation

Pour la méthode **atan**, on a choisi comme implémentation les polynômes d'Hermite. C'est une autre séquence à décrire des polynômes dans  $\mathbb{Q}[x]$  qui s'approche d'**arctan**(x) uniformément sur  $[0, 1]$  et qui le fait beaucoup plus rapidement que la séquence de polynômes de Taylor centrée sur 0. Nous notons qu'une telle séquence d'approximation fournit une approximation sur tout  $\mathbb{R}$ , grâce aux identités suivantes :

$$\mathbf{arctan}(x) = -\mathbf{arctan}(-x) = \pi/2 - \mathbf{arctan}(1/x)$$

On a utilisé le polynôme d'Hermite  $h_3(x)$  qui est un polynôme d'ordre 15, et qui d'après l'article [5], donne une erreur de l'ordre de  $10^{-7}$ .

Afin de justifier notre choix d'implémentation, on a effectué de nombreux tests sur l'intervalle  $[0; 1]$  pour comparer l'implémentation de l'**arctan** avec Séries de Taylor à l'ordre 9 et l'implémentation avec polynômes d'Hermite. Et on a remarqué qu'au voisinage de 0, les deux méthodes ont une précision de  $10^{-8}$ , mais plus on s'approche de 1, et précisément après 0,6, l'erreur de la méthode avec Taylor commence à augmenter pour atteindre du  $10^{-4}$  dans l'intervalle de  $[0, 65 ; 0, 7]$ .

**Exemple:**

Intervalle	Pas	Nombre d'échantillons	Résultats par Hermite	Résultats par Taylor
$[0, 65; 0, 7]$	$2^{-23}$	419431	$3, 2338175 * 10^{-8}$	$8, 9091645 * 10^{-4}$

### 5.2 Difficultés et amélioration

L'implémentation avec les polynômes d'Hermite nous a donné de bons résultats (7), mais on a quelques grands valeurs de l'ULP après 0,6.

Pour régler ce problème, on s'aide de la relation trigonométrique de **arctan**:

$$\mathbf{arctan}(x) + \mathbf{arctan}(y) = \mathbf{arctan}((x + y)/(1 - xy)) \iff xy < 1.$$

Et donc si on prend  $y = 1$ , on peut choisir l'intervalle qu'on veut pour x. Ici c'est  $[0, 6; 1[$  qui nous pose problème.

### 5.3 Résultats

D'après les tests effectués sur Java on trouve les résultats suivants:

Intervalle	step	erreur en ULP	max erreur	nombres d'éléments
[0.0; 0.59]	$2^{-23}$	<b>1.434709</b>	<b>8.0</b>	<b>4949279</b>
[0.59; 1]	$2^{-23}$	<b>0.5891706</b>	<b>7.0</b>	<b>3439330</b>
[1; 1000]	$2^{-7}$	<b>0.36701545</b>	<b>3.0</b>	<b>127872</b>

#### Quelques valeurs sur Deca:

Après implémentation de la fonction **arctan** sur Deca, on a refait quelques tests pour s'assurer des précisions trouvées.

On trouve les mêmes ordre de grandeurs. Avec des fois des précisions meilleures; comme on l'a indiqué au début de la documentation, cela est dû à la restriction de l'écriture flottante de Java qui est limité à 8 chiffres après la virgule, contrairement à l'écriture hexadécimale du Deca qui n'est pas limitée.

valeur	Résultats sur Deca(hex)	résultat atan java(hex))	erreur	erreur en ulp
<b>1</b>	<b>0x1.921fb6p - 1</b>	<b>0x1.921fb6p - 1</b>	<b>0.0</b>	<b>0.0</b>
<b>0.0001</b>	<b>0x1.a36e2ep - 14</b>	<b>0x1.a36e26p - 14</b>	<b><math>6 * 10^{-14}</math></b>	<b>8.24634e - 03</b>
<b>5656</b>	<b>0x1.92142p + 0</b>	<b>0x1.92141ep + 0</b>	<b><math>10^{-9}</math></b>	<b>8.38861e - 03</b>
<b>0,005</b>	<b>0x1.47ad62p - 8</b>	<b>0x1.47ad62p - 8</b>	<b>0.0</b>	<b>0.0</b>
<b>0.0000009</b>	<b>0x1.e32f0ep - 21</b>	<b>0x1.e32f0ep - 21</b>	<b>0.0</b>	<b>0.0</b>

## 5.4 Graphes

Les courbes ci-après résument les résultats de la fonction Atan avant et après optimisations.



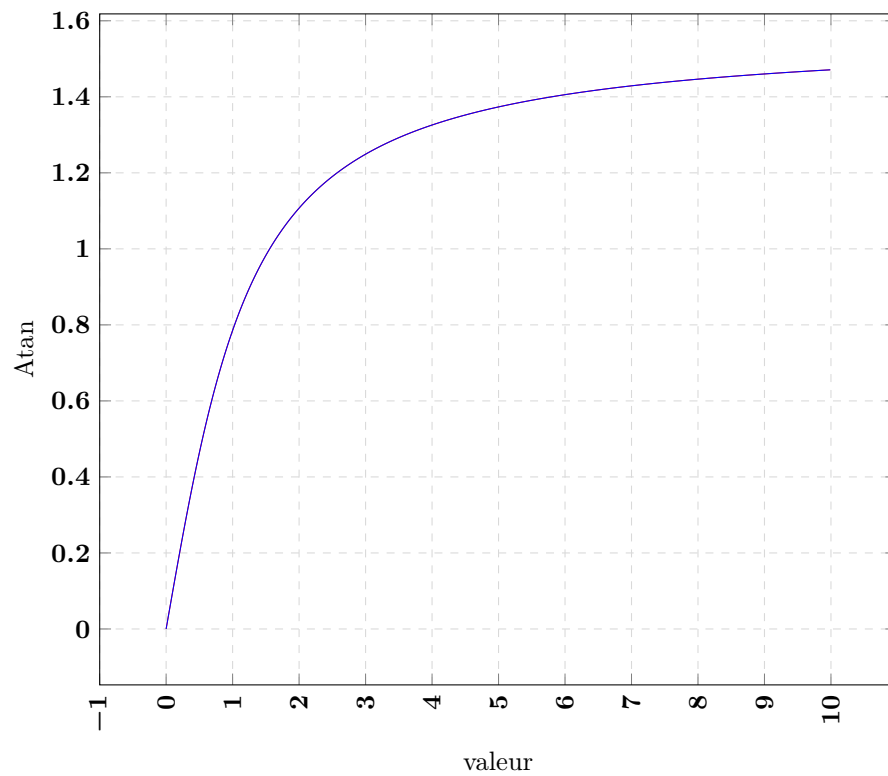


Figure 8: Atan Deca (en bleue )sur Atan Java (en rouge )

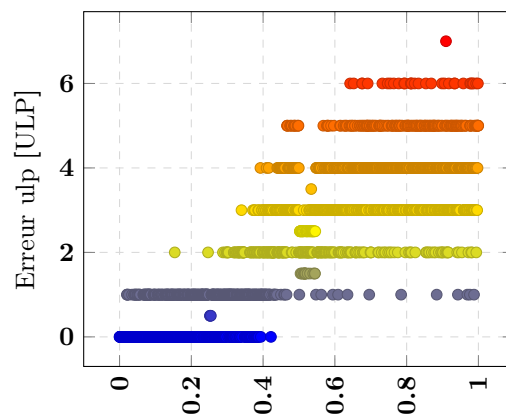
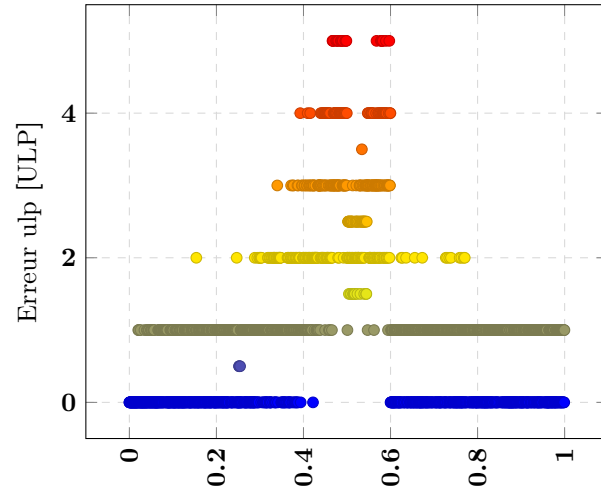
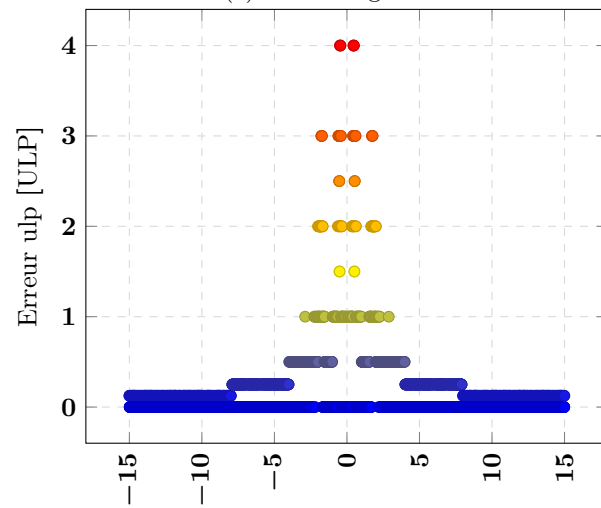


Figure 9: atan avant correction



(a) atan corrigé



(b) atan corrigé sur un grand intervalle

Figure 10: Atan corrigé

## References

- [1] [https://en.wikipedia.org/wiki/Unit\\_in\\_the\\_last\\_place](https://en.wikipedia.org/wiki/Unit_in_the_last_place)
- [2] <https://citeseerx.ist.psu.edu/document?repid=rep1&type=pdf&doi=bc07e8d52dbbac5bef9551bbc111b95d82dc4766>
- [3] [https://fr.wikipedia.org/wiki/S%C3%A9rie\\_de\\_Taylor](https://fr.wikipedia.org/wiki/S%C3%A9rie_de_Taylor)
- [4] <https://www.geogebra.org/m/s9SkCsvC>
- [5] [https://www.academia.edu/5490106/A\\_SEQUENCE\\_OF\\_HERMITE\\_INTERPOLATING\\_LIKE\\_POLYNOMIALS\\_FOR\\_APPROXIMATING\\_ARCTANGENT](https://www.academia.edu/5490106/A_SEQUENCE_OF_HERMITE_INTERPOLATING_LIKE_POLYNOMIALS_FOR_APPROXIMATING_ARCTANGENT)
- [6] <https://stackoverflow.com/questions/15393573/square-root-with-babylonian-or-herons-method>