

Éléments de correction

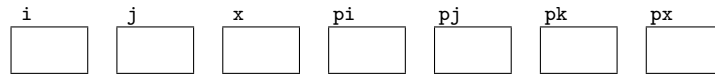
A Notions fondamentales

Exercice 1 : Introduction aux pointeurs **A**

Notions de base sur les pointeurs

exo1-pointeurs.c	exo1-pointeurs.f90
1 #include <stdio.h>	1 program pointeur
2 #include <stdlib.h>	2
3	3 implicit none
4 #define ULI (unsigned long int) // conversion	4
5	5 integer, target :: i=1, j=2
6 int main(void) {	6 integer, pointer :: pi=>null(), pj=>null(), pk
7 int i=1, j=2;	7 real, target :: x
8 int *pi=NULL, *pj=NULL, *pk;	8 real, pointer :: px=>null()
9 float x, *px=NULL;	9
10	10 pi => i
11 pi = &i;	11 pj => j
12 pj = &j;	12 write (*,*) "i =", i, ", j =", j
13 printf("i=%d, j=%d, *pi=%d, *pj=%d\n",	13 write (*,*) "pi =", pi, ", pj =", pj
14 i, j, *pi, *pj);	14 write (*, *) "pi => i :", associated(pi, i)
15 printf("pi=%lu, pj=%lu\n", ULI pi, ULI pj);	15 stop ! fin etape 1
16 exit(EXIT_SUCCESS); // fin etape 1	16
17	17 i = i*10
18 i *= 10;	18 write (*,*) "i =", i, ", pi =", pi
19 printf("i=%d, *pi=%d, pi=%lu\n", i, *pi, ULI pi);	19 pi = pi*10
20 *pi *= 10;	20 write (*,*) "i =", i, ", pi =", pi
21 printf("i=%d, *pi=%d, pi=%lu\n", i, *pi, ULI pi);	21 stop ! fin etape 2
22 exit(EXIT_SUCCESS); // fin etape 2	22
23	23
24 printf("pk=%lu\n", ULI pk);	24 pk => pj
25 pk = pj; // memes types	25 if (associated(pk, j)) then
26 printf("pk=%d, pk=%lu\n", *pk, ULI pk);	26 write (*,*) "pk pointe vers j"
27 if (pk == &j) {	27 write (*,*) "pk =", pk, "j =", j
28 printf("pk pointe vers j\n");	28 endif
29 }	29 stop ! fin etape 3
30 exit(EXIT_SUCCESS); // fin etape 3	30
31	31 pk = pi ! memes types
32 *pk = *pi;	32 write (*,*) "i =", i, " pk =", pk
33 printf("i=%d, *pk=%d\n", i, *pk);	33 stop ! fin etape 4
34 exit(EXIT_SUCCESS); // fin etape 4	34
35	35 x = 3.5
36 x = 3.5f;	36 px => x
37 px = &x;	37 write (*,*) "x =", x, "px =", px
38 printf("x=%g, *px=%g, px=%lu\n", x, *px, ULI px);	38 px = pj
39 *px = *pj;	39 write (*,*) "j =", j, " x =", x, "px =", px
40 printf("j=%d, *px=%g, px=%lu\n", j, *px, ULI px);	40 stop ! fin etape 5
41 exit(EXIT_SUCCESS); // fin etape 5	41
42 }	42 end program pointeur

Le code **fourni** dans les fichiers `exo1-pointeurs.c` et `exo1-pointeurs.f90` permet d'expérimenter les opérations sur les pointeurs en dehors de tout passage de paramètre dans des procédures. Exécuter ce code étape par étape en passant en commentaire au fur et à mesure les instructions d'arrêt (`exit` en C et `stop` en fortran). À chaque étape du programme, représenter schématiquement les zones mémoire des différentes variables et noter leur valeur ; préciser vers quelle zone pointent les variables pointeurs ainsi que les valeurs des variables pointées.



Précautions à prendre avec les pointeurs

Reprendre la même démarche avec les fichiers `exo2-pointeurs.c` et `exo2-pointeurs.f90`. Expliquer avec précision quelles instructions provoquent des erreurs ou des avertissements lors de la compilation. Certains comportements peuvent apparaître aléatoires à l'exécution. Modifier les programmes pour les éviter (quitte à supprimer certaines instructions).

exo2-pointeurs.c	exo2-pointeurs.f90
<pre> 1 #include <stdio.h> 2 #include <stdlib.h> 3 4 #define ULI (unsigned long int) 5 6 int main(void) { 7 int i=1; 8 int *pi=NULL, *pj; 9 float x, *px=NULL; 10 11 printf("pi=%lu\n", ULI pi); 12 *pi = 4; // essayer d'exécuter 13 exit(EXIT_SUCCESS); // fin etape 1 14 15 pi = &i; 16 printf("pi=%lu\n", ULI pi); 17 *pi *= 10; 18 printf("i=%d, *pi=%d, pi=%lu\n", i, *pi, ULI pi); 19 exit(EXIT_SUCCESS); // fin etape 2 20 21 printf("pj=%lu\n", ULI pj); 22 // resultat aleatoire 23 *pj = 5; // essayer d'exécuter 24 pj = NULL; 25 printf("pj=%lu\n", ULI pj); 26 // resultat certain 27 exit(EXIT_SUCCESS); // fin etape 3 28 29 x = 3.5f; 30 px = &x; 31 printf("x=%g , *px=%g, px=%lu\n", x, *px, ULI px); 32 *px = *pi; 33 printf("i=%d , *px=%g, px=%lu\n", i, *px, ULI px); 34 exit(EXIT_SUCCESS); // fin etape 4 35 36 px = pi; // essayer pour etape 5 37 printf("i=%d , *px=%g, px=%lu\n", i, *px, ULI px); 38 exit(EXIT_SUCCESS); // fin etape 5 39 } </pre>	<pre> 1 program pointeur 2 3 implicit none 4 5 integer, target :: i=1 6 integer, pointer :: pi=>null(), pj 7 real, target :: x 8 real, pointer :: px=>null() 9 10 write(*,*) " pi assoc :", associated(pi) 11 pi = 4 ! essayer d'exécuter 12 stop ! fin etape 1 13 14 pi => i 15 write(*,*) " pi assoc :", associated(pi) 16 pi = 4 17 write (*,*) "i =",i,", pi =",pi 18 stop ! fin etape 2 19 20 write(*,*) " pj assoc :", associated(pj) 21 ! resultat aleatoire 22 pj = 5 ! essayer d'exécuter 23 nullify(pj) 24 write(*,*) " pj assoc :", associated(pj) 25 ! resultat certain 26 stop ! fin etape 3 27 28 x = 3.5 29 px => x 30 write (*,*) "x =", x, "px =", px 31 px = pi 32 write (*,*) "i =", i, " x =", x, "px =", px 33 stop ! fin etape 4 34 35 px => pi ! essayer de compiler 36 write (*,*) "i =", i, "px =", px 37 stop ! fin etape 5 38 39 end program pointeur </pre>

```

1  #include <stdio.h>
2  #include <stdlib.h>
3
4  #define ULI (unsigned long int)
5
6  int main(void) {
7      int i=1;
8      int *pj, *pi=NULL;
9      float x, *px=NULL;
10
11     printf("pi=%lu\n", ULI pi);
12     /*pi = 4; // essayer d'executer
13     //exit(EXIT_SUCCESS); // fin etape 1
14
15     pi = &i;
16     printf("pi=%lu\n", ULI pi);
17     *pi *= 10;
18     printf("i=%d, *pi=%d, pi=%lu\n", i, *pi, ULI pi);
19     //exit(EXIT_SUCCESS); // fin etape 2
20
21     printf("pj=%lu\n", ULI pj);
22     // resultat aleatoire
23     /*pj = 5; // essayer d'executer, peut echouer
24     pj = NULL;
25     printf("pj=%lu\n", ULI pj);
26     // resultat certain
27     //exit(EXIT_SUCCESS); // fin etape 3
28
29     x = 3.5f;
30     px = &x;
31     printf("x=%g , *px=%g, px=%lu\n", x, *px, ULI px);
32     *px = *pi; // garder pour voir l'effet ! (warning)
33     printf("i=%d , *px=%g, px=%lu\n", i, *px, ULI px);
34     //exit(EXIT_SUCCESS); // fin etape 4
35
36     // des entiers en octal correspondant
37     // en float a des "valeurs speciales"
38     i = 017740000000; // +Inf en float
39     i = 037740000000; // -Inf en float
40     i = 017740000001; // +NaN (quiet)
41     i = 037740000001; // -NaN (quiet)
42     i = 037750000001; // -NaN (signaling)
43     px = pi; // essayer pour etape 5
44     printf("i=%d , *px=%g, px=%lu\n", i, *px, ULI px);
45     exit(EXIT_SUCCESS); // fin etape 5
46 }

```

En fortran, l'association entre pointeurs de types différents provoque une erreur dès la compilation.

```

px => pi ! essayer de compiler
1
Error: Different types in pointer assignment at (1)

```

Cela prévient des erreurs graves qui pourraient en résulter. Il faut donc supprimer cette instruction pour réussir à compiler y compris pour mettre en œuvre les étapes précédentes.

En langage C, au contraire, elle ne produit qu'un warning à la compilation :

```

warning: assignment from incompatible pointer type

```

En C, on peut donc par exemple pointer une variable de type entier via un pointeur de `float` : le motif de 32 bits représentant un entier signé ne peut pas toujours être interprété en tant que `float` valide. Il correspond soit à des valeurs très différentes de la valeur de l'entier, soit à des codes de nombres réels dénormalisés (de valeur absolue inférieure à `FLT_MIN`, cf. annexe du polycopié fortran sur la norme IEEE), voire des codes spéciaux pour les NaN, +Inf ou -Inf. Le programme `exo2-pointeurs+.c` permet de stocker les valeurs octales d'un entier signé qui codent ces valeurs non numériques en `float`.

En fortran, on peut affecter des motifs de bits à des variables réelles via des lectures en format B (norme fortran 2008, avec `gfortran` ou `ifort`). Mais il n'est pas possible de faire pointer un pointeur de réel vers un entier : une solution pour lire le motif binaire comme un entier est d'utiliser la fonction intrinsèque `transfer`.

Exercice 2 : Fonctions avec retour sans effet de bord A

Éditer les fichiers suivants, `produit.c` en C, puis `produit.f90` en fortran 90. À chaque question, prendre soin de faire une copie en incrémentant le nom de façon à conserver toutes les versions.

<code>produit.c</code>	<code>produit.f90</code>
<pre> 1 #include <stdio.h> 2 #include <stdlib.h> 3 4 /* declaration : prototype */ 5 int prod(int x, int y); 6 7 /* definition : les instructions */ 8 int prod(int x, int y) { 9 int z; /* locale */ 10 z = x * y; 11 return z; /* valeur de retour */ 12 } 13 14 int main(void) { 15 int a1, b1, a2, b2; 16 int p; 17 printf("Saisir 4 entiers a1,b1,a2,b2\n"); 18 scanf("%d %d %d %d", 19 &a1, &b1, &a2, &b2); 20 /* deux appels dans une expression */ 21 p = prod(a1, b1) + prod(a2, b2); 22 printf("a1*b1 + a2*b2 = %d\n", p); 23 exit(EXIT_SUCCESS); 24 }</pre>	<pre> 1 module m_prod ! => nom du fichier .mod 2 implicit none 3 contains ! procédures de module 4 function prod(x, y) 5 integer :: prod 6 integer, intent(in) :: x, y 7 integer :: z ! locale 8 z = x * y 9 prod = z ! valeur de retour 10 return ! sans valeur 11 end function prod 12 end module m_prod 13 14 program produit 15 use m_prod ! relit m_prod.mod 16 ! ce qui rend visible l'interface 17 integer :: a1, b1, a2, b2 18 integer :: p 19 write(*,*) "Saisir 4 entiers a1,b1,a2,b2" 20 read(*,*) a1, b1, a2, b2 21 ! deux appels dans une expression 22 p = prod(a1, b1) + prod(a2, b2) 23 write(*,*) "a1*b1 + a2*b2 = ", p 24 end program produit</pre>

- 1) Ces programmes font appel à une fonction `prod` qui ne modifie pas ses arguments et rend une valeur à l'appelant. Compiler dans les deux langages, observer les fichiers créés. Tester ces programmes. Modifier la fonction `prod` pour qu'elle n'utilise plus de variable locale `z` et tester.

B Essayer d'introduire une modification de la valeur de `x` dans `prod` en fortran. Expliquer.

Solution

En fortran, création du fichier `m_prod.mod` qui contient l'interface des procédures définies dans le module `m_prod`. L'attribut `intent(in)` de `x` interdit toute modification, par exemple :

In file `produit.f90`:7

```

x = x + 1
1
```

Error: Can't assign to INTENT(IN) variable 'x' at (1)

- 2) Déclarer la variable `p` en tant que réel `float` en C (ajuster aussi le format d'affichage) et `real` en fortran dans le programme principal. Expliquer les comportements à la compilation et à l'exécution. Quelle conversion est effectuée ?

Solution

Le produit, calculé en entier, est converti en réel, ce qui modifie l'affichage en fortran. Noter que le problème de l'éventuel dépassement de capacité se pose dans le type entier dans la fonction : le produit est limité en valeur absolue à `INT_MAX` en C et `HUGE(1)` en fortran, c'est à dire $2^{31} - 1$ sur une machine 32 bits.

- 3) En conservant `p` réel dans chaque langage, déclarer les variables `a1`, `b1`, `a2` et `b2` en tant que réels dans le programme principal (et adapter le format de saisie en C). Expliquer les comportements à la compilation et à l'exécution. En langage C, reprendre la compilation en ajoutant l'option `-Wconversion` et conclure.

Solution

La fonction reste définie avec des paramètres formels entiers, mais elle est appelée avec des paramètres effectifs réels. Le fortran refuse alors de convertir et si l'interface de `prod` est visible de l'appelant (ce

qui est le cas avec le USE `m_prod`), l'erreur est détectée dès la compilation :

```
p = prod(a1, b1) + prod(a2, b2)
      1
```

Error: Type mismatch in parameter 'x' at (1). Passing REAL(4) to INTEGER(4)
et de même pour `a1`.

Le C fait naturellement les conversions (à cause du passage par copie de valeur). Le compilateur n'avertit pas sauf avec l'option `-Wconversion`.

```
gcc-mni produit-fl+int.c
```

```
produit-fl+int.c: In function 'main':
```

```
produit-fl+int.c:19: attention : conversion to 'int' from 'float' may alter its value
produit-fl+int.c:19: attention : conversion to 'int' from 'float' may alter its value
produit-fl+int.c:19: attention : conversion to 'int' from 'float' may alter its value
produit-fl+int.c:19: attention : conversion to 'int' from 'float' may alter its value
produit-fl+int.c:19: attention : conversion to 'float' from 'int' may alter its value
```

Un détail en C : Jusqu'à gcc 4.2, l'option `-Wconversion` n'était pas recommandée en usage normal où le prototype est déclaré, car elle émettait un avertissement même dans le cas où les arguments formels de `prod` sont déclarés `float` ainsi que les arguments lors de l'appel. Il n'y a alors plus de conversion, mais `-Wconversion` signalait un comportement différent de celui où on ne déclare pas de prototype pour `prod` (où, pour des raisons historiques les `float` sont systématiquement promus en `double`).

```
produit-float.c:19: warning: passing argument 1 of 'prod' as integer rather
than floating due to prototype
```

Depuis gcc 4.3, ce comportement est celui de `-Wtraditional-conversion` alors que `-Wconversion` avertit des vrais problèmes potentiels de conversion sans référence à un type implicite... Donc, l'option `-Wtraditional-conversion` est de peu d'intérêt.

Attention, en nov 2010, `sappli` est encore sous gcc-3.4.6 (et les mandriva 2007.1 sous gcc-4.1.2!) Donc on évitera `-Wconversion` sur le serveur. Heureusement, les mandriva 2009.1 utilisent gcc-4.3.2, donc `-Wconversion` est conseillé. C'est un cas où le serveur et les machines virtuelles n'ont pas le même comportement.

Extrait du man gcc 4.3 ou supérieur :

```
-Wtraditional-conversion
```

```
Warn if a prototype causes a type conversion that is different from what would happen to the
same argument in the absence of a prototype. This includes conversions of fixed point to
floating and vice versa, and conversions changing the width or signedness of a fixed point
argument except when the same as the default promotion.
```

```
-Wconversion
```

```
between real and integer, like abs (x) when x is double; conversions between signed and
unsigned, like unsigned ui = -1; and conversions to smaller types, like sqrtf (M_PI). Do not
warn for explicit casts like abs ((int) x) and ui = (unsigned) -1, or if the value is not
changed by the conversion like in abs (2.0). Warnings about conversions between signed and
unsigned integers can be disabled by using -Wno-sign-conversion.
```

- 4) Modifier la fonction `prod` de manière à ce que le programme se compile sans avertissement lorsque `a1`, `b1`, `a2` et `b2` sont déclarés réels dans le programme principal. Tester de nouveau.

Solution

Il faut déclarer réels les paramètres d'entrée et aussi de préférence le résultat (sinon double conversion).

Exercice 3 : Procédures sans valeur de retour A

Ce sont les sous-programmes (*subroutine*) en fortran et les fonctions sans retour (résultat de type *void*) en C. Elles sont censées produire des effets de bord : entrées/sorties ou modification des paramètres passés.

- 1) Copier, puis éditer, tester et commenter les programmes `produit_et_somme` suivants, en C puis en fortran 90 :

```

1  _____ produit_et_somme.c _____
2  #include <stdio.h>
3  #include <stdlib.h>
4  /* declaration : prototype */
5  void prod_som(int x, int y);
6  /* definition */
7  void prod_som(int x, int y){
8      /* x et y : parametres d'entree */
9      /* variables locales */
10     int prod, som;
11     prod = x * y ;
12     som = x + y ;
13     printf("produit = %d\n", prod);
14     printf("somme = %d\n", som);
15     return;
16 }
17 int main(void){
18     int a, b;
19     printf("Saisir a, b:\n");
20     scanf("%d %d", &a, &b);
21     prod_som(a, b);
22     exit(EXIT_SUCCESS);
23 }

```

```

1  _____ produit_et_somme.f90 _____
2  module m_prod_som
3  implicit none
4  contains
5      subroutine prod_som(x, y)
6          ! parametres d'entree
7          integer, intent(in) :: x, y
8          ! variables locales
9          integer                :: prod, som
10         prod = x*y
11         som = x+y
12         write(*,*) 'produit =', prod
13         write(*,*) 'somme = ', som
14         return
15     end subroutine prod_som
16 end module m_prod_som
17 program produit_et_somme
18     use m_prod_som
19     integer :: a, b
20     write(*,*) 'Saisir a, b:'
21     read(*,*) a, b
22     call prod_som(a, b)
23 end program produit_et_somme

```

- 2) On souhaite modifier le programme `produit_et_somme` de manière à ce que les valeurs du produit et de la somme soient renvoyées au programme principal qui les imprimera : on déclarera donc les variables `prod` et `som` dans le programme principal et la procédure `prod_som` comportera alors quatre paramètres.

En C, les arguments sont **passés par copie de valeur** et les modifications effectuées sur les paramètres dans la fonction appelée n'affectent que sa copie locale des paramètres. Il est donc nécessaire de communiquer à la fonction appelée l'adresse du paramètre pour le partager avec l'appelant : le type du paramètre formel recevant copie de cette adresse est donc un pointeur sur le type de la variable à modifier, elle-même déclarée côté appelant.

À titre d'exemple, expliquer la différence de syntaxe entre les fonctions `printf` et `scanf`. Modifier le code `produit_et_somme.c` de manière à passer en argument l'adresse des variables à modifier.

En fortran, les arguments sont par défaut **passés par référence** ; les zones mémoires des paramètres sont alors communes à l'appelé et à l'appelant. Les modifications effectuées sur les paramètres dans la procédure sont donc visibles par l'appelant.

Effectuer la modification en précisant l'attribut `intent` pour tous les paramètres du sous-programme.

```

1 #include <stdio.h>
2 #include <stdlib.h>
3
4 void prod_som(int x, int y, int *p_prod, int *p_som);
5
6 void prod_som(int x, int y, int *p_prod, int *p_som){
7     /* pointeurs sur les entiers */
8     *p_prod = x * y;
9     *p_som = x + y;
10
11     printf("fonction: produit = %d\n", *p_prod);
12     printf("fonction: somme = %d\n", *p_som);
13     return;
14 }
15
16 int main(void){
17     int a, b;
18     int prod, som;
19
20     printf("Saisir 2 entiers a, b:\n");
21     scanf("%d %d", &a, &b);
22     printf("a = %d b = %d\n", a, b);
23     prod_som(a, b, &prod, &som);
24
25     printf("main: produit = %d\n", prod);
26     printf("main: somme = %d\n", som);
27
28     exit(EXIT_SUCCESS);
29 }

```

```

1 module m_prod_som
2 implicit none
3 contains
4     subroutine prod_som(x, y, prod, som)
5         ! parametres d'entree
6         integer, intent(in) :: x, y
7         integer, intent(out) :: prod, som
8         prod = x * y
9         som = x + y
10        write(*,*) 'ss-pgm : produit =', prod
11        write(*,*) 'ss-pgm : somme = ', som
12        return
13    end subroutine prod_som
14 end module m_prod_som
15
16 program produit_et_somme
17 use m_prod_som
18 integer :: a, b
19 integer :: s, p
20 write(*,*) 'Saisir 2 entiers a, b:'
21 read(*,*) a, b
22 write(*,*) "a= ", a, " b= ", b
23 call prod_som(a, b, p, s)
24 write(*,*) 'pgm ppal: produit =', p
25 write(*,*) 'pgm ppal: somme = ', s
26 end program produit_et_somme

```

- 3) En Fortran 90 et en C, écrire une procédure `echange` qui admet deux arguments entiers `x` et `y`, et qui les intervertit. Les variables `x` et `y` sont à la fois des arguments d'entrée et de sortie. Appeler cette procédure depuis le programme principal et imprimer les valeurs avant et après l'échange (on pourra vérifier ici la nécessité d'utiliser des pointeurs en C).

```

1 #include <stdio.h>
2 #include <stdlib.h>
3
4 void vrai_echange(int *p_x, int *p_y); /* declaration */
5
6 void vrai_echange(int *p_x, int *p_y){ /* definition */
7     int c;
8     c = *p_y;
9     *p_y = *p_x;
10    *p_x = c;
11    printf("vrai_echange: p_x et p_y = %p %p \n", p_x, p_y);
12    return;
13 }
14
15 void faux_echange(int x, int y); /* declaration */
16
17 void faux_echange(int x, int y){ /* definition */
18     int c;
19     c = y;
20     y = x;
21     x = c;
22     printf("faux_echange: x et y = %d %d \n", x, y);
23     return;
24 }
25
26 int main(void){
27     int a, b;
28     printf("Saisir 2 entiers a et b:\n");
29     scanf("%d %d", &a, &b);
30     printf("main: a et b avant vrai_echange = %d %d\n", a, b);
31     vrai_echange(&a, &b);
32     printf("main: a et b apres vrai_echange = %d %d\n", a, b);
33     faux_echange(a, b); /* echange des copies */
34     printf("main: a et b apres faux_echange = %d %d\n", a, b);
35     exit(EXIT_SUCCESS);
36 }

```

```

1 module m_echange
2 implicit none
3 contains
4     subroutine echange(x, y)
5         integer, intent(inout) :: x, y
6         integer :: z ! locale
7         z = x
8         x = y
9         y = z
10        write(*,*) "dans echange : a = ", x, " b = ", y
11    end subroutine echange
12 end module m_echange
13
14 program ppal
15 use m_echange
16 integer :: a, b
17 write(*,*) "Saisir 2 entiers a et b"
18 read(*,*) a, b
19 write(*,*) "ppal avant echange : a = ", a, " b = ", b
20 call echange(a, b)
21 write(*,*) "ppal apres echange : a = ", a, " b = ", b
22 end program ppal

```

B Applications

Exercice 4 : Fonction de Planck A

Soit la fonction de Planck

$$\text{Planck}(x) = \frac{x^3}{e^x - 1} \quad \text{où } x = \frac{h\nu}{kT} \quad \text{est un argument sans dimension}$$

- 1) Écrire un programme `planck` qui, en vue d'un tracé, échantillonne sur un intervalle fini la fonction de Planck et affiche sur la sortie standard les valeurs de x et de $\text{Planck}(x)$ (à raison d'un point par ligne). Dans un premier temps, on décrira l'intervalle $[0, 5; 10]$ par pas de 0,5
- 2) Déplacer le calcul de la fonction de Planck dans une fonction `f_planck` à un argument x renvoyant la valeur de $\text{Planck}(x)$. La fonction `f_planck` sera appelée à chaque passage dans la boucle décrivant l'intervalle.
- 3) Pour $|x|$ tendant vers zéro, l'expression initiale conduit à une indétermination. Lever cette indétermination par un développement de $\exp(x)$ au premier ordre. Modifier la fonction `f_planck` pour prendre en compte cette expression approchée quand $|x|$ est petit. On comparera $|x|$ avec :
 - `EPSILON(x)` en fortran (fonction intrinsèque)
 - `FLT_EPSILON` ou `DBL_EPSILON` en C (macro-constantes définies dans `float.h`)
 Reprendre l'échantillonnage à partir de 0 pour un pas quelconque.
- 4) Tracer la fonction de Planck en fonction de x . Pour cela, rediriger les données de la sortie standard par défaut vers un fichier `planck1.dat` et utiliser `gnuplot` pour tracer les courbes (commande `gnuplot` puis `plot "planck1.dat" using ($1):($2) with points` où $\$1$ et $\$2$ sont les champs correspondant aux données que vous voulez tracer).
- 5) B Comparer les deux types d'erreurs dominantes commises au voisinage de $x = 0$ selon la méthode de calcul :
 - **erreur systématique de troncature** (liée au nombre limité de termes) en utilisant le développement limité de l'expression autour de 0 d'une part : elle diminue quand on s'approche de 0.
 - **erreur aléatoire d'arrondi** (en considérant le facteur le plus imprécisément évalué) avec l'expression initiale de $\text{Planck}(x)$ d'autre part : elle augmente quand on s'approche de 0.
 Déterminer pour quelle valeur de $|x|$ ces deux erreurs sont égales et on doit changer de méthode de calcul afin de minimiser l'erreur totale. Modifier la fonction `f_planck` en conséquence.

Solution

```

1 #include <stdio.h>
2 #include <stdlib.h>
3 #include <tgmath.h> // fonctions generiques
4 #include <float.h>
5 /* standard c99 pour sqrtf expf et fabsf */
6
7 float planck(float x); /* déclaration */
8
9 float planck(float x){
10     float y;
11     if(fabs(x) > sqrt(2.*FLT_EPSILON)) {
12         // appelle fabsf et sqrtf grace a tgmath.h
13         y = x*x*x / (exp(x) - 1.f);
14     }
15     else{
16         y = x * x ;
17     }
18     return y;
19 }
20
21 int main(void){
22     int n, i;
23     float x, xmin, xmax, dx ;
24     float y ;
25     xmin = 0.f;
26     xmax = 10.f;
27     dx = .5f;
28     n = 1 + (int) ((xmax - xmin) / dx) ;
29     for (i =1; i<= n; i++){
30         x = xmin + (i -1) * dx ;
31         y = planck(x) ;
32         printf("%g %g\n", x, y);
33     }
34     exit(EXIT_SUCCESS) ;
35 }

```

```

1 module m_planck
2     implicit none
3     contains
4     function planck(x)
5         real :: planck
6         real, intent(in) :: x
7         ! en comparant :
8         ! l'erreur d'arrondi sur l'expression complète
9         ! et l'erreur de troncature sur le développement
10        ! au 1er ordre, on doit changer de mode de calcul
11        ! pour x/2 = \epsilon / x
12        if(abs(x) > sqrt(2.*epsilon(x))) then
13            planck = x**3 / (exp(x) - 1.)
14        else
15            planck = x**2
16        endif
17    end function planck
18 end module m_planck
19
20 program tabule_planck
21     use m_planck
22     real :: x, y, xmin, xmax, dx
23     integer :: n, i
24
25     xmin = 0.
26     xmax = 10.
27     dx = .5
28     n = 1 + int((xmax - xmin) / dx)
29     do i = 1, n
30         x = xmin + (i -1) * dx
31         y = planck(x)
32         write(*,*) x, y
33     end do
34 end program tabule_planck

```


Quand on doit évaluer $f(x) = \frac{x^3}{\exp(x) - 1}$ au voisinage de zéro, il y a une indétermination que l'on peut lever par un développement limité du dénominateur $d(x) \approx 1 + x + x^2/2 + \dots - 1 \approx x(1 + x/2 + \dots)$. Ainsi $f(x) \approx \frac{x^2}{1 + x/2 + \dots} \approx x^2(1 - x/2 + \dots)$. On peut tronquer le développement limité de $f(x)$ pour ne garder qu'un terme avec $f_1(x) = x^2$: l'erreur relative de troncature associée avec ce calcul approché est de $|x/2|$.

Ainsi, pour évaluer la fonction de Planck, on doit distinguer :

- (a) les très faibles valeurs de x pour lesquelles on calcule $f_1(x) = x^2$: on commet alors essentiellement une erreur de troncature.
- (b) les autres valeurs de x pour lesquelles on utilise la formule initiale $f_2(x) = \frac{x^3}{\exp(x) - 1}$ qui présente une erreur d'arrondi qui peut être importante quand $|x|$ est faible.

Si on note ε le plus grand flottant positif tel que $1 + \varepsilon$ ne soit pas différent de 1, l'erreur relative maximale de représentation d'un flottant est ε . Dans une addition ou une soustraction, chacun des termes est représenté approximativement : le résultat est donc entaché d'une erreur dite d'arrondi dont la valeur absolue est majorée par la somme des erreurs de représentation des deux termes. Dans le cas d'une soustraction de termes proches, l'erreur relative d'arrondi peut donc être importante. Dans le calcul du dénominateur pour x faible, l'erreur absolue sur $\exp(x)$ (mais aussi sur $\exp(x) - 1$ car 1 est représenté exactement) est majorée par ε . Mais comme le dénominateur est de l'ordre de x , l'erreur relative sur ce dénominateur est de l'ordre de ε/x . Si on néglige l'erreur d'arrondi dans le calcul de x^3 et dans le rapport, l'erreur relative sur $f_2(x)$ est aussi de l'ordre de ε/x .

Le basculement entre les deux calculs doit se faire lorsque les erreurs relatives sont égales, soit $\varepsilon/x = x/2$. Ainsi,

- (a) pour $|x| \leq \sqrt{2\varepsilon}$, on calcule $f(x) \approx x^2$ et l'erreur dominante est déterministe et due à la troncature
- (b) pour $|x| \geq \sqrt{2\varepsilon}$, on calcule $f(x) = \frac{x^3}{\exp(x) - 1}$ et l'erreur dominante est aléatoire et due à l'arrondi.

N.-B. : Dans la description d'un intervalle par découpage à pas constant, il est conseillé de calculer le nombre de pas à effectuer et de faire une boucle avec compteur. De plus, on évitera d'accumuler les erreurs d'arrondi par additions successives ($x = x + dx$ à chaque pas) en effectuant une multiplication par le numéro du pas ($x = x_{\min} + (i-1)*dx$ où i est l'entier servant de compteur).

Exercice 5 : Calcul du nombre d'or par dichotomie A

- 1) Montrer que la racine positive de l'équation $x^2 = x + 1$ est le nombre d'or ϕ .
- 2) Résolution de l'équation $x^2 = x + 1$ par dichotomie.

Bref rappel de la méthode de dichotomie : on se donne un intervalle initial de recherche dont les bornes encadrent la racine ; puis on divise cet intervalle par le milieu et on le remplace par le sous-intervalle deux fois plus petit dans lequel se trouve la racine, et on itère le processus. La longueur de l'intervalle dont les bornes encadrent la racine diminue donc en 2^{-n} en fonction de l'ordre n d'itération.

- (a) Vérifier que cette équation admet une solution dans l'intervalle $[1; 2]$.
- (b) Écrire une fonction `dichotomie` qui accepte comme arguments les deux bornes de l'intervalle d'entrée (soit $[a; b]$). Cette fonction remplace l'un des arguments par le milieu de l'intervalle de départ (par exemple $[\frac{a+b}{2}; b]$), de telle sorte que la solution de l'équation soit toujours contenue dans l'intervalle de sortie.
- (c) Demander la précision souhaitée à l'utilisateur du programme et appeler la fonction `dichotomie` autant de fois que nécessaire pour trouver la solution de l'équation à la précision adéquate.
- (d) B Indiquer ce qui limite la précision et comment modifier le programme pour calculer par exemple avec 10 chiffres significatifs.

Solution

La largeur de l'intervalle est divisée par 2 à chaque étape, donc il est possible de connaître a priori le nombre d'étapes nécessaires pour une précision donnée.

$$b_n - a_n = \frac{b_{n-1} - a_{n-1}}{2} = \frac{b_0 - a_0}{2^n}$$

On peut donc utiliser une boucle avec compteur au lieu de la boucle `while` présentée dans les codes ci-dessous.

Fortran (sous-type réel à choisir avec `prec`)

```

----- dichotomie-or.f90 précision ajustable -----
1 module rect ! Le module
2   implicit none
3   ! variable déclarée au niveau du module
4   !integer, parameter :: prec=selected_real_kind(p=6) ! simple préc.
5   !integer, parameter :: prec=selected_real_kind(p=15) !double préc.
6   private f ! => pas de visibilité hors du module
7   contains
8     function f(x)
9       real(kind=prec), intent(in) :: x
10      real(kind=prec) :: f
11      f = x**2-x-1._prec
12    end function f
13
14    subroutine dichotomie(x1,x2)
15      real(kind=prec),intent(inout) :: x1, x2
16      real(kind=prec)                :: x3
17      x3 = (x1+x2)/2._prec
18      if (f(x1)*f(x3) < 0) then
19        x2 = x3
20      else
21        x1 = x3
22      endif
23    end subroutine dichotomie
24 end module rect
25
26 program nombredor
27   use rect
28   implicit none
29   real(kind=prec) :: p, g, eps
30   integer :: compteur
31   ! Determination du nombre d'or par dichotomie dans [1,2]
32   write(*,*) "Entrer la précision choisie eps:"
33   write(*,*) "pas mieux que ", EPSILON(1._prec)
34   read(*,*) eps
35   p = 1.
36   g = 2.
37   compteur = 0
38   do while (abs(g-p) > eps)
39     call dichotomie(p,g)
40     compteur = compteur + 1
41   enddo
42   write(*,*) "Le nombre d'iterations nécessaire pour &
43     & une precision de ", eps, "est", compteur
44   write(*,*) "Nombre d'or =", (g+p)/2._prec
45   write(*,*) "Nombre d'or =", (1._prec+sqrt(5._prec))/2._prec
46 end program nombredor

```

Langage C

dichotomie-or.c version float		dichotomie-or-dbl.c version double	
1	#include <stdio.h>	1	#include <stdio.h>
2	#include <stdlib.h>	2	#include <stdlib.h>
3	#include <math.h>	3	#include <tgmath.h>
4	#include <limits.h>	4	#include <limits.h>
5	#include <float.h>	5	#include <float.h>
6		6	
7	float f(float x);	7	double f(double x);
8		8	
9	float f(float x){	9	double f(double x){
10	return (x*x -x -1.f);	10	return (x*x -x -1.);
11	}	11	}
12		12	
13	void dichotomie(float *p_x1, float *p_x2);	13	void dichotomie(double *p_x1, double *p_x2);
14		14	
15	void dichotomie(float *p_x1, float *p_x2){	15	void dichotomie(double *p_x1, double *p_x2){
16	float x3;	16	double x3;
17	x3 = (*p_x1 + *p_x2)/2.f;	17	x3 = (*p_x1 + *p_x2)/2.;
18	if (f(*p_x1)*f(x3) < 0){	18	if (f(*p_x1)*f(x3) < 0){
19	*p_x2 = x3;	19	*p_x2 = x3;
20	}	20	}
21	else{	21	else{
22	*p_x1 = x3;	22	*p_x1 = x3;
23	}	23	}
24	return;	24	return;
25	}	25	}
26		26	
27	int main(void){	27	int main(void){
28	float p, g, eps;	28	double p, g, eps;
29	int compteur;	29	int compteur;
30	/* Determination du nombre d'or par dichotomie */	30	/* Determination du nombre d'or par dichotomie */
31	printf("Entrer la precision choisie eps: \n");	31	printf("Entrer la precision choisie eps: \n");
32	printf("pas mieux que %g\n", FLT_EPSILON);	32	printf("pas mieux que %lg\n", DBL_EPSILON);
33	scanf("%g", &eps);	33	scanf("%lg", &eps);
34	p = 1.;	34	p = 1.;
35	g = 2.;	35	g = 2.;
36	compteur = 0;	36	compteur = 0;
37	while(fabs(g-p) > eps){	37	while(fabs(g-p) > eps){
38	dichotomie(&p, &g);	38	dichotomie(&p, &g);
39	compteur++;	39	compteur++;
40	}	40	}
41	printf("Le nombre d'iterations necessaire "	41	printf("Le nombre d'iterations necessaire "
42	"pour une precision de %g est %d \n", eps, compteur);	42	"pour une precision de %g est %d \n", eps, compteur);
43	printf("Nombre d'or = %f \n", (g+p)/2.);	43	printf("Nombre d'or = %.16lf \n", (g+p)/2.);
44	printf("Nombre d'or = %f \n", (1.+sqrt(5.f))/2.f);	44	printf("Nombre d'or = %.16lf \n", (1.+sqrt(5.))/2.);
45	exit(EXIT_SUCCESS);	45	exit(EXIT_SUCCESS);
46	}	46	}

N.-B. : En C, ne pas confondre la fonction `abs` d'argument entier, avec `fabs` d'argument réel. Utiliser un test sur `abs` avec un argument de type `float` ou `double` arrêtera la boucle dès qu'il est inférieur à 1.

Exercice 6 : Approximation de π

On cherche à obtenir une estimation de π par différentes méthodes.

Première méthode : série alternée A

On utilise le fait que $\arctan(1) = \frac{\pi}{4}$ et le développement en série entière de $\arctan(x)$:

$$\arctan(x) = \sum_{k=0}^{\infty} (-1)^k \frac{x^{2k+1}}{2k+1} \quad \text{pour } |x| \leq 1$$

- 1) Comment calculer de proche en proche $(-1)^k x^{2k+1}$ sans faire appel à des calculs de puissance ?
- 2) Écrire une fonction `mon_arctan` qui calcule la somme des n premiers termes de la série entière précédente. Les arguments de cette procédure seront x et n . Utiliser cette procédure pour obtenir une estimation de π . Contrôler la précision relative de l'approximation en comparant avec la valeur de π calculée avec par exemple la fonction `atan`.
- 3) Sachant que, dans une série alternée, le reste peut être majoré par la valeur absolue du premier terme négligé dans la somme finie, intégrer dans le programme le calcul du nombre de termes suffisant pour limiter l'erreur de troncature à une valeur choisie par l'utilisateur (en ignorant les erreurs d'arrondi).
- 4) B Dans quel ordre doit on effectuer la somme pour minimiser l'erreur d'arrondi ?

Solution

- 1) Le calcul de la série entière nécessite une alternance de signe qui peut être programmée avec l'opérateur d'élevation à la puissance en fortran mais en C, on doit soit appeler la fonction `pow` soit utiliser l'opérateur `%` de modulo. On peut aussi explicitement programmer l'alternance de façon itérative en évitant l'élevation à la puissance grâce à une variable intermédiaire `u` :

```

u = x
s = u
do k=1, n
  u = -u * x * x
  s = s + u / (2*k + 1)
end do

```

```

u = x ;
s = u ;
for (k=1; k<=n; k++){
  u = -u * x * x ;
  s = s + u / (2*k + 1) ;
}

```

- 2) Les codes proposés pages 13 et 14 n'utilisent pas ces méthodes car ils commencent par sommer les plus petits termes (plus grands indices).
- 3) Si $|x| \leq 1$, la série est alternée et le reste est majoré par la valeur absolue du premier terme négligé. L'erreur absolue est donc majorée par $\frac{x^{2n+3}}{2n+3}$ c'est à dire $\frac{1}{2n+3}$ pour $x = 1$. Pour que l'erreur relative soit inférieure à ε , il faut sommer jusqu'à $k = n = \frac{2}{\pi\varepsilon} - 1, 5$. L'erreur de troncature décroît linéairement avec le nombre de termes. Mais la précision relative finale est de toute façon limitée par l'erreur d'arrondi, soit `EPSILON(x)` en `REAL` et `FLT_EPSILON` en `float`, tous deux de l'ordre de 10^{-7} .
- 4) Encore faut-il limiter l'**accumulation** des erreurs d'arrondi. Pour cela, il est préférable de sommer en commençant par les plus petits termes. Dans le cas général $|x| \leq 1$, la mise en œuvre de la sommation dans cet ordre est facilitée par une factorisation qui s'inspire du schéma de Horner.

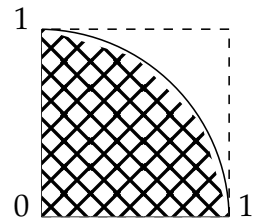
$$\sum_{k=0}^n (-1)^k \frac{x^{2k+1}}{2k+1} = x \left(1 - \frac{1}{3}x^2 \left(1 - \frac{3}{5}x^2 \left(1 - \frac{5}{7}x^2 \left(\dots \left(1 - \frac{2n-1}{2n+1}x^2 \right) \right) \right) \right) \right) \quad (7.1)$$

Deuxième méthode : Monte-Carlo **AB**

L'estimation de π peut aussi être obtenue par un calcul d'aire. Dans la figure ci-contre, le rapport de l'aire du secteur angulaire de 90 degrés à l'aire du carré vaut :

$$\frac{\text{Aire en grisé}}{\text{Aire du carré}} = \frac{\pi}{4}$$

L'aire peut être estimée par une méthode de Monte-Carlo dont l'algorithme est le suivant. On tire aléatoirement n points **répartis uniformément** dans le carré $[0; 1] \times [0; 1]$. Le rapport des aires est estimé par le rapport du nombre de points situés à l'intérieur du secteur angulaire au nombre total de points tirés.



On utilisera les générateurs de nombres pseudo-aléatoires uniformes fournis pour chaque langage :

- en fortran, le sous-programme intrinsèque standard `random_number` permet des tirer des valeurs équidistribuées sur l'intervalle $[0; 1]$ via `CALL RANDOM_NUMBER(x)` où `x` est un **réel**.
- en C, la fonction `rand()` qui rend un **entier** entre 0 et `RAND_MAX` (voir man 3 `rand`).

- 1) Rajouter au module précédent une procédure effectuant l'intégration de Monte-Carlo.
- 2) **B** Comment évolue la précision avec le nombre de tirages ? Conclure sur la méthode qui vous semble la plus efficace.

Solution (en imposant le nombre de termes de la série et en commençant par les plus grands indices)

Pour le programme principal :

```

1 program pi
2 use m_approx
3 implicit none
4 integer :: n
5 real :: x, mon_pi
6 write(*,*) 'Entrer un nombre reel x tel que |x| <=1'
7 read (*,*) x

```

```

8 write(*,*) 'Entrer le nombre de termes du developpement limite'
9 read (*,*) n
10 write(*,*) "Le developpement limite a l'ordre ", 2*n+1, " &
11     &de 4 * atan(", x, ") est ", 4 * mon_arctan(x,n)
12 ! mon_arctan est une fonction
13 write(*,*) "Avec la fonction intrinsèque "
14 write(*,*) " 4 * atan(", x,") = ", 4 * atan(x)
15 write(*,*) ! saut de ligne
16 write(*,*) 'Entrer le nombre de points dans la methode de Monte-Carlo'
17 read (*,*) n
18 call montecarlo(n, mon_pi)
19 ! montecarlo est un sous-programme
20 write(*,*) 'Valeur approchee de pi = ', mon_pi
21 end program pi

```

Et pour le module :

```

1 module m_approx
2 implicit none
3 contains
4   function mon_arctan(x, n) ! calcul de atan(x) à l'ordre n
5     real, intent(in) :: x
6     integer, intent(in) :: n
7     real :: mon_arctan
8     integer :: k
9     mon_arctan = 0.
10    do k=n, 0, -1 ! faire la somme en commençant par les plus petits
11      mon_arctan = mon_arctan + (-1)**k * x**(2*k + 1) / (2*k + 1)
12    enddo
13  end function mon_arctan
14
15  subroutine montecarlo(n, resul)
16    integer, intent(in) :: n
17    real, intent(out) :: resul
18    ! version "vectorielle" sans boucle
19    real, dimension(n) :: tirage_x, tirage_y, modul
20    call random_number(tirage_x) ! on tire n nb aleatoires compris dans [0;1[
21    call random_number(tirage_y) ! idem
22    modul = tirage_x**2 + tirage_y**2 ! module carré de (x,y)
23    ! compter le nombre de points dans le quart de disque
24    resul = 4. * count(modul <= 1) / real(n)
25  end subroutine montecarlo
26 end module m_approx

```

En C (version float)

```

1 #include <stdio.h>
2 #include <stdlib.h>
3 #include <math.h>
4
5 float mon_arctan(float x, int n);
6 float mon_arctan(float x, int n){ /* calcul de atan(x) à l'ordre n */
7   float somme;
8   int k ;
9   somme = 0.f;
10  for(k=n; k>=0; k--) { /* faire la somme en commençant par les plus petits */
11    somme += (1 - 2*(k%2)) * pow(x, (float) (2 * k + 1)) /
12              (float) (2 * k + 1) ;
13  }
14  return somme ;

```

```

15 }
16
17 void montecarlo(int n, float *result);
18 void montecarlo(int n, float *result){
19     int compt = 0;
20     int k;
21     float x, y ;
22     for (k = 1; k<=n ; k++){
23         x = (float) rand()/ ((float) RAND_MAX + 1.);
24         y = (float) rand()/ ((float) RAND_MAX + 1.);
25         if( (x*x + y*y) <= 1 ) {
26             compt++;
27         }
28     }
29     *result = 4.f * (float) compt / (float) n ;
30 }
31
32 int main(void){
33     int n;
34     float x, mon_pi ;
35     x = 1.f;
36     printf("Entrer le nombre de termes du développement limité\n");
37     scanf("%d", &n);
38     printf("Le développement limité à l'ordre %d de 4 * atan(%f)"
39           " est %f\n", 2*n +1, x, 4.f * mon_arctan(x, n) );
40     printf("Et 4 * atan(%f) = %f\n", x, 4.f * atan(x));
41     printf("Erreur relative par la série %e\n",
42           (mon_arctan(x,n) - atan(x))/ atan(x));
43     printf("Entrer le nombre de points utilisés par la méthode de Monte-Carlo\n");
44     scanf("%d", &n);
45     montecarlo(n, &mon_pi) ;
46     printf("Valeur approchée de pi = %f\n", mon_pi);
47     printf("Erreur relative par MonteCarlo %e\n", (mon_pi/4.f - atan(x))/ atan(x));
48     exit(EXIT_SUCCESS);
49 }

```

La vitesse de convergence de la méthode de Monte Carlo (comptage de points) est faible : l'écart relatif est proportionnel à $1/\sqrt{N}$, alors qu'il évolue en $1/N$ dans le calcul par série alternée.

N.-B. : En langage C, on prendra soin de ne pas calculer `RAND_MAX + 1` en type `int`, sous peine de dépassement de capacité qui donne une valeur négative.