

Rapport TP5 IDM

Génération de flux parallèles de nombres pseudo-aléatoires
Bibliothèque CLHEP

Année : 2024-2025

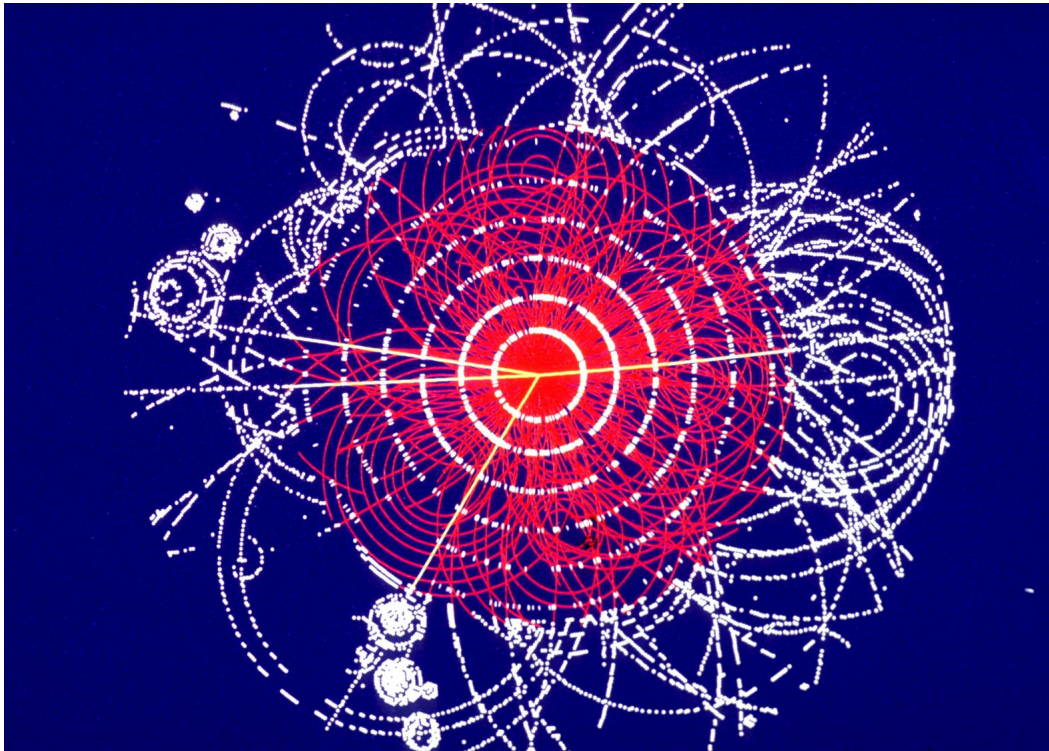


FIGURE 1 – Simulation de la désintégration d'un boson de Higgs

Table des matières

1	Introduction	3
2	Installation de la bibliothèque	4
2.1	Compilation séquentielle	4
2.2	Compilation en parallèle	4
2.3	Observations	4
3	Tests des générateurs de la librairie	5
4	Utilisation de la librairie	7
4.1	Analyse de <code>saveStatus()</code> et <code>restoreStatus()</code>	7
5	Simulation par Monte Carlo	10
5.1	Approche séquentielle	10
5.2	Approche en parallélisation	11
5.2.1	Sequence splitting	11
6	Indicateurs statistiques	14
7	Différentes techniques et bibliothèques de parallélisation	16
7.1	OpenMP (Open Multi-Processing)	16
7.1.1	Simulation d'un neutron	16
7.2	pthread	18
8	Application de la génération aléatoire de séquences en bioinformatique avec Mersenne Twister (MT)	19

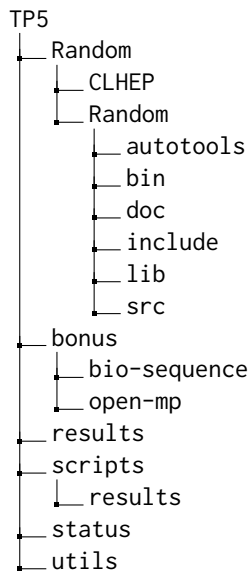
Table des figures

1	Simulation de la désintégration d'un boson de Higgs	1
2	Exemple de test (<code>testRandom</code>) sur le générateur Mersenne Twister	5
3	Restauration de l'état du générateur Mersenne Twister entre deux générations successives d'une séquence de 10 nombres aléatoires	8
4	Visualisation de la progression et sauvegarde des états du générateur Mersenne Twister	9
5	Visualisation de la progression et résultats de l'estimation du volume d'une sphère par la méthode de Monte Carlo	10
7	Indicateurs statistiques pour la simulation en parallèle	15
6	Indicateurs statistiques pour la simulation en séquentielle	15
8	Recherche d'une séquence cible (" GATTACA ") avec Mersenne Twister (approche parallèle)	19

1 Introduction

Le TP est organisé de la manière suivante. Les résultats sont disponibles dans les fichiers intitulés "**results**" (TP5-A et TP5-B). La compilation s'effectue à l'aide du Makefile : pour la première partie du TP, utilisez la commande `make`, et pour la deuxième partie, utilisez la commande : `make tp5b`.

Pour exécuter la partie B du tp5, exécuter le script `runParallelSim.sh`.



- **Random** : bibliothèque CLHEP
- **bonus** : exercices supplémentaires du TP
- **results** : résultats des simulations (TP5a et TP5b)
- **scripts** : scripts pour l'exécution des simulations
- **utils** : utilitaires pour les simulations et le traitement des données

Bibliothèque CLHEP (Class Library for High Energy Physics)

La bibliothèque CLHEP (Class Library for High Energy Physics) est une collection de bibliothèques écrites en C++ principalement utilisées dans le domaine de la physique des hautes énergies.

CLHEP fournit des outils et des classes utilitaires largement utilisés pour la simulation et l'analyse des données. Elle est intégrée dans de nombreux projets scientifiques, y compris des frameworks comme Geant4.

2 Installation de la bibliothèque

Q1.

2.1 Compilation séquentielle

Nous testons dans un premier temps la compilation séquentielle (sans parallélisation avec les jobs).

```
1 time make
```

Voici les résultats que nous obtenons respectivement pour les temps réels, utilisateurs et systèmes :

```
1 real    1m19,202s
2 user    1m11,357s
3 sys     0m7,600s
```

2.2 Compilation en parallèle

La compilation est réalisée sur ma machine personnelle, qui possèdent 12 cœurs (commande `nproc`) sous Linux (Ubuntu).

```
1 nproc
2 >> 12
```

Nous spécifions donc une compilation parallèle avec 12 cœurs logiques.

```
1 time make -j12
```

Voici les résultats que nous obtenons respectivement pour les temps réels, utilisateurs et systèmes :

```
1 real    0m29,562s
2 user    4m31,260s
3 sys     0m22,291s
```

2.3 Observations

Les résultats en mode séquentiel et en parallèle, montrent clairement l'avantage de la compilation parallélisée. Nous constatons que la parallélisation optimise le temps réel de compilation, mais maximise le temps utilisateur (`user`).

En passant de **1m19s (séquentiel)** à **29s (parallèle)**, la compilation parallèle divise presque le temps réel par un facteur de **2.68**.

Le **temps utilisateur** passe de **1m11s (séquentiel)** à **4m31s (parallèle)**. Cela reflète une augmentation de la charge de travail répartie sur plusieurs cœurs.

Bien que le temps total utilisateur + système augmente (liée au traitement simultané sur plusieurs cœurs), le gain en temps réel justifie l'approche.

En calculant le ratio du temps réel en séquentiel sur le temps réel en parallèle, nous pouvons en déduire l'efficacité de la simulation lancée en parallèle :

$$\frac{time_real_seq}{time_real_parallel} = \frac{71,4}{29} = 2,46 > 1$$

L'exécution de la simulation en parallèle est donc 2 fois plus rapide que le lancement de la simulation

3 Tests des générateurs de la librairie

Une fois la librairie installée, nous procédons à l'exécution des tests (`testRandom.cc` et `testrand.cc`).

Le fichier `testRandom.cc` initialise plusieurs moteurs aléatoires, comme le générateur `HepJamesRandom`, `RandEngine` et nous pouvons choisir l'un de ces moteurs pour effectuer des tirages aléatoires sur diverses distributions (uniforme, exponentielle, gaussienne, etc.). Les résultats des tirages sont affichés à chaque étape.

```
g++ testRandom.cc -I../include ../lib/libCLHEP-Random-2.1.0.0.a -static
```

```
----- Test on MTwistEngine -----
Flat ]0,1[      : 0.176117
Flat ]0,5[      : 0.101418
Flat ]-5,3[     : 1.90606
Exp (m=1)       : 2.12904
Exp (m=3)       : 2.30604
Gauss (m=1)     : -0.85285
Gauss (m=3,v=1) : 3.14988
Wigner(1,0.2)   : 1.11898
Wigner(1,0.2,1) : 1.01213
Wigner2(1,0.2)  : 0.988002
Wigner2(1,0.2,1): 0.546862
IntFlat [0,99[  : 28
IntFlat [-99,37[: -67
Poisson (m=3.0) : 2
Binomial(n=1,p=0.5) : 0
Binomial(n=-5,p=0.3): -1
ChiSqr (a=1)    : 0.0159581
ChiSqr (a=-5)   : -1
Gamma (k=1,l=1) : 0.973209
Gamma (k=3,l=0.5) : 4.07356
StudT (a=1)     : -2.16803
StudT (a=2.5)   : -1.19112

Shooting an array of 5 flat numbers ...

0.80256 0.582024 0.362318 0.249103 0.573408
```

FIGURE 2 – Exemple de test (`testRandom`) sur le générateur Mersenne Twister

La majorité des valeurs sont dans les plages attendues pour chaque distribution statistique.

Nous testons également les générateurs avec le fichier donné en sujet :

```
1  #include <sys/types.h>
2  #include <sys/stat.h>
3  #include <fcntl.h>
4  #include <limits.h>
5  #include <unistd.h>
6
7  #include "CLHEP/Random/MTwistEngine.h"
8
9  int main ()
10 {
11     CLHEP::MTwistEngine * s = new CLHEP::MTwistEngine();
12
13     int fs;
14     double f;
15     unsigned int nbr;
16
17     fs = open("./rngbForMarsagliaTests", O_CREAT | O_TRUNC | O_WRONLY, S_IRUSR | S_IWUSR);
18
19     for(int i = 1; i < 3000000; i++)
20     {
21         f = s->flat();
22         nbr = (unsigned int) (f * UINT_MAX);
23
24         std::cout << f << std::endl;
25
26         write(fs, &nbr, sizeof(unsigned int));
27     }
28
29     close(fs);
30
31     delete s;
32
33     return 0;
34 }
```

Le programme de test génère 3 000 000 nombres pseudo-aléatoires flottants (en utilisant le moteur Mersenne Twister), les convertit en entiers non signés, les écrit dans un fichier binaire, et affiche les valeurs flottantes à la console.

Le fichier généré peut être utilisé pour des tests (par exemple, des tests de Marsaglia pour évaluer la qualité des générateurs aléatoires).

4 Utilisation de la librairie

Q2.

Le but ici est de créer une instance du générateur de nombres aléatoires **MTwistEngine** (CLHEP).

```
1 #include "Random/Random/MTwistEngine.h"
2 #include <iostream>
3
4 int statusMT() {
5     CLHEP::MTwistEngine mtEngine;
6
7     mtEngine.saveStatus();
8
9     std::cout << "First 10 generated numbers:\n";
10    for (int i = 0; i < 10; ++i) {
11        std::cout << mtEngine.flat() << std::endl;
12    }
13
14    std::cout << "Generator state saved to 'MTwist.conf'.\n";
15    mtEngine.restoreStatus();
16
17    std::cout << "Generator restored.\n";
18    std::cout << "\nRestored generator, first 10 generated numbers:\n";
19
20    for (int i = 0; i < 10; ++i) {
21        std::cout << mtEngine.flat() << std::endl;
22    }
23
24    return 0;
25 }
26
27 int main()
28 {
29     statusMT();
30 }
```

Nous sauvegardons d'abord l'état du générateur avec `saveStatus()`. Ensuite, il génère et affiche les 10 premiers nombres aléatoires.

L'état du générateur est restauré avec `restoreStatus()`, et les 10 premiers nombres aléatoires sont de nouveau générés et affichés, reproduisant ainsi exactement la même séquence qu'au début.

4.1 Analyse de `saveStatus()` et `restoreStatus()`

La méthode `saveStatus()` enregistre l'état interne du générateur dans un fichier texte, afin de pouvoir le recharger plus tard avec `restoreStatus()`.

La méthode `restoreStatus()` recharge l'état du générateur depuis un fichier. Cela permet de reprendre la génération de nombres aléatoires à partir du même point que celui où elle avait été interrompue.

```

First 10 generated numbers:
0.176117
0.0202835
0.863258
0.118952
0.463625
0.0618998
0.823987
0.571751
0.0917234
0.777519
Generator state saved to 'MTwist.conf'.
Generator restored.

Restored generator, first 10 generated numbers:
0.176117
0.0202835
0.863258
0.118952
0.463625
0.0618998
0.823987
0.571751
0.0917234
0.777519

```

FIGURE 3 – Restauration de l'état du générateur Mersenne Twister entre deux générations successives d'une séquence de 10 nombres aléatoires

Q3.

```

1  /**
2   * @brief generates and saves multiple random generator statuses.
3   *
4   * generates multiple random generator states, saving each state as a separate file
5   * in the "status" directory. Each draw is tracked using a progress bar to visualize
6   * progress.
7   *
8   * @param numStatus number of generator states to create.
9   * @param numDraws number of random draws to generate for each status.
10  */
11  void generateAndSaveStatus(int numStatus, long numDraws) {
12
13      printf("\n ==== Saving generator status ==== \n");
14      printf("Num. status : %d, Num. draws : %ld \n", numStatus, numDraws);
15
16      std::string directory = "status";
17
18      std::string name;
19      std::string filename;
20
21      if (!std::filesystem::exists(directory)) {
22          std::filesystem::create_directory(directory);
23      }
24
25      CLHEP::MTwistEngine mtEngine;
26
27      for (int statusIndex = 0; statusIndex < numStatus; ++statusIndex) {
28          for (long i = 0; i < numDraws; ++i) {
29              mtEngine.flat();
30              if ((i + 1) % (numDraws / 100) == 0) progress(i + 1, numDraws);
31          }
32
33          name = "status-" + std::to_string(statusIndex + 1) + ".conf";
34          filename = directory + "/" + name;

```



```
35     mtEngine.saveStatus(filename.c_str());
36
37     std::cout << "\n" << name << " saved." << std::endl;
38 }
39 }
```

```
==== Saving generator status ====
Num. status : 10, Num. draws : 3000000000
[=====] 100 %
status-1.conf saved.
[=====] 100 %
status-2.conf saved.
[=====] 100 %
status-3.conf saved.
[=====] 100 %
status-4.conf saved.
[=====] 100 %
status-5.conf saved.
[=====] 100 %
status-6.conf saved.
[=====] 100 %
status-7.conf saved.
[=====] 100 %
status-8.conf saved.
[=====] 100 %
status-9.conf saved.
[=====] 100 %
status-10.conf saved.
```

FIGURE 4 – Visualisation de la progression et sauvegarde des états du générateur Mersenne Twister

5 Simulation par Monte Carlo

Q4. & Q5.

5.1 Approche séquentielle

Nous allons dans un premier temps lancer la simulation en séquentielle puis comparer le temps d'exécution et les résultats obtenus par `sequence splitting` (approche parallèle).

Voici les résultats que nous obtenons pour la simulation d'un calcul de volume d'une sphère par Monte Carlo pour 10×3 milliards de tirages.

```

=== Performing Monte Carlo Simulation ===
** Start task at: 2024-12-19 21:30:22

Num. replications : 10 / Num. draws per replication : 3000000000

Estimating volume by Monte Carlo method...
[=====] 100 %
Replication 1: Estimated volume = 4.18873

Estimating volume by Monte Carlo method...
[=====] 100 %
Replication 2: Estimated volume = 4.18878

Estimating volume by Monte Carlo method...
[=====] 100 %
Replication 3: Estimated volume = 4.18882

Estimating volume by Monte Carlo method...
[=====] 100 %
Replication 4: Estimated volume = 4.18884

Estimating volume by Monte Carlo method...
[=====] 100 %
Replication 5: Estimated volume = 4.1889

Estimating volume by Monte Carlo method...
[=====] 100 %
Replication 6: Estimated volume = 4.18885

Estimating volume by Monte Carlo method...
[=====] 100 %
Replication 7: Estimated volume = 4.18882

Estimating volume by Monte Carlo method...
[=====] 100 %
Replication 8: Estimated volume = 4.18874

Estimating volume by Monte Carlo method...
[=====] 100 %
Replication 9: Estimated volume = 4.18879

Estimating volume by Monte Carlo method...
[=====] 100 %
Replication 10: Estimated volume = 4.18878

>> Average estimated volume after 10 replications: 4.1888
Time used : 534.667 seconds

```

FIGURE 5 – Visualisation de la progression et résultats de l'estimation du volume d'une sphère par la méthode de Monte Carlo

Le temps de la simulation par Monte Carlo en **séquentiel** est donc de **534.667 secondes** soit d'environ **8 minutes**.

5.2 Approche en parallélisation

5.2.1 Sequence splitting

L'approche par "sequence splitting" consiste à diviser la séquence aléatoire en plusieurs sous-séquences indépendantes. Pour cela, on utilise des états différents du générateur aléatoire pour chaque réplcation. Chaque tâche exécute le même programme, mais utilise un flux aléatoire différent initialisé par un fichier de statut unique (approche SPMD).

Voici la démarche que nous allons adopter. Tout d'abord nous allons générer les fichiers (status) d'état aléatoires (soit les fichiers `status-x.conf`).

Notre programme doit accepter le **fichier de statut** et le **nombre de points comme arguments**. Chaque réplcation (chaque exécution) doit restaurer son propre état et effectuer le calcul de manière indépendante.

Le script suivant va nous permettra de lancer **10 réplcations en parallèle selon un status pour le MT choisi** et passé en paramètre :

```

1  #!/bin/bash
2
3  num_points=3000000000
4  status_file_prefix="status-"
5
6  mkdir -p results
7
8  overall_start_time=$(date +%s%3N)
9
10 # Run the 10 replications in parallel
11 for i in {1..10}
12 do
13     start_time=$(date "+%Y-%m-%d %H:%M:%S.%3N")
14
15     echo "[Script PID: $$] [Replication $i] Starting replication at $start_time" >>
16         replication_log.txt
17
18     status_file="../status/${status_file_prefix}${i}.conf"
19
20     if [ ! -f "$status_file" ]; then
21         echo "[Script PID: $$] [Replication $i] Error: Status file $status_file does
22             not exist!" >> replication_log.txt
23         continue
24     fi
25
26     ../../tp5b $status_file $num_points > "results/output_replication_${i}.txt" &
27
28     background_pid=$!
29
30     echo "[Script PID: $$] [Replication $i] Background job started with PID:
31         $background_pid at $start_time" >> replication_log.txt
32 done

```

Nous pouvons améliorer le script pour obtenir plus d'informations sur l'exécution de la simulation (temps total mis, volume moyen estimé) :

```

1  wait
2
3  completion_time=$(date "+%Y-%m-%d %H:%M:%S.%3N")
4  overall_end_time=$(date +%s%3N)
5
6  total_time_ms=$((overall_end_time - overall_start_time))
7  total_time_seconds=$(echo "scale=3; $total_time_ms / 1000" | bc)
8
9  echo "[Script PID: $$] *** All simulations completed at $completion_time ***" >>
    replication_log.txt
10 echo "[Script PID: $$] *** Total time taken: $total_time_seconds seconds ***" >>
    replication_log.txt
11
12 # Average estimated volume
13 total_volume=0
14 count=0
15
16 for i in {1..10}
17 do
18     output_file="results/output_replication_$i.txt"
19
20     if [ -f "$output_file" ]; then
21         estimated_volume=$(grep "Volume:" "$output_file" | awk -F ':' '{print $2}' |
            xargs)
22         if [ -n "$estimated_volume" ]; then
23             total_volume=$(echo "$total_volume + $estimated_volume" | bc)
24             count=$((count + 1))
25         fi
26     fi
27 done
28
29 if [ $count -gt 0 ]; then
30     average_volume=$(echo "scale=5; $total_volume / $count" | bc)
31     echo "Average estimated volume after $count replications: $average_volume" >>
        replication_log.txt
32     echo "Average estimated volume after $count replications: $average_volume"
33 else
34     echo "No valid estimated volumes found in the results."
35     echo "No valid estimated volumes found in the results." >> replication_log.txt
36 fi
37
38 echo "All simulations completed. Total time taken: $total_time_seconds seconds."

```

La progression des différentes simulations peut-être consultée dans le dossier `output-results-x.txt`. Les logs sont disponibles dans le fichier `result`.

Voici les résultats que nous obtenons (logs) :

```

1 [Script PID: 66620] [Replication 1] Starting replication at 2024-12-19 21:56:46.325
2 [Script PID: 66620] [Replication 1] Background job started with PID: 66624 at
   2024-12-19 21:56:46.325
3
4 [Script PID: 66620] [Replication 2] Starting replication at 2024-12-19 21:56:46.330
5 [Script PID: 66620] [Replication 2] Background job started with PID: 66626 at
   2024-12-19 21:56:46.330
6
7 [Script PID: 66620] [Replication 3] Starting replication at 2024-12-19 21:56:46.335
8 [Script PID: 66620] [Replication 3] Background job started with PID: 66628 at
   2024-12-19 21:56:46.335
9
10 [Script PID: 66620] [Replication 4] Starting replication at 2024-12-19 21:56:46.341
11 [Script PID: 66620] [Replication 4] Background job started with PID: 66630 at
   2024-12-19 21:56:46.341
12
13 [Script PID: 66620] [Replication 5] Starting replication at 2024-12-19 21:56:46.346
14 [Script PID: 66620] [Replication 5] Background job started with PID: 66632 at
   2024-12-19 21:56:46.346
15
16 [Script PID: 66620] [Replication 6] Starting replication at 2024-12-19 21:56:46.352
17 [Script PID: 66620] [Replication 6] Background job started with PID: 66634 at
   2024-12-19 21:56:46.352
18
19 [Script PID: 66620] [Replication 7] Starting replication at 2024-12-19 21:56:46.358
20 [Script PID: 66620] [Replication 7] Background job started with PID: 66636 at
   2024-12-19 21:56:46.358
21
22 [Script PID: 66620] [Replication 8] Starting replication at 2024-12-19 21:56:46.363
23 [Script PID: 66620] [Replication 8] Background job started with PID: 66638 at
   2024-12-19 21:56:46.363
24
25 [Script PID: 66620] [Replication 9] Starting replication at 2024-12-19 21:56:46.368
26 [Script PID: 66620] [Replication 9] Background job started with PID: 66640 at
   2024-12-19 21:56:46.368
27
28 [Script PID: 66620] [Replication 10] Starting replication at 2024-12-19 21:56:46.374
29 [Script PID: 66620] [Replication 10] Background job started with PID: 66642 at
   2024-12-19 21:56:46.374
30
31
32 [Script PID: 66620] *** All simulations completed at 2024-12-19 21:58:30.013 ***
33
34 [Script PID: 66620] *** Total time taken: 103.694 seconds ***
35 Average estimated volume after 10 replications: 4.18880

```

Le temps total de la simulation de Monte Carlo en **parallèle** (après terminaison des 10 threads) est de **103.694 secondes** soit environ **1 minutes et 43 secondes**.

	Nombre de points	Temps total (s)	Volume estimé (moyenne)
Séquentielle	10×3 milliards	534.667 s (8 min 55 s)	4.18888
Parallèle (sequence splitting)	10×3 milliards	103.694 s (1 min 43 s)	4.18880

TABLE 1 – Comparaison des temps d'exécution et des volumes estimés par Monte Carlo

Nous pouvons calculer l'efficacité de l'approche parallèle en calculant le ratio suivant :

$$\frac{time_seq}{time_parallel} = \frac{534.667}{103.694} = 5 > 1$$

L'exécution de la simulation en parallèle est donc 5 fois plus rapide que la simulation séquentielle de l'estimation du volume par Monte Carlo et cela pour un résultat similaire (volume moyen d'environ $\frac{4}{3} \times \pi$).

6 Indicateurs statistiques

Nous pouvons également évaluer les résultats obtenus par les approches séquentielle et parallèle à l'aide d'indicateurs statistiques tels que la moyenne, la variance, l'écart-type et l'intervalle de confiance.

— **Variance** (σ^2) :

$$\sigma^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \mu)^2$$

où μ est la moyenne et x_i est le volume estimé de la i -ième réplication.

Une variance faible indiquerait que les différentes réplifications donnent des résultats proches les uns des autres, ce qui est un signe de stabilité et de précision des estimations. Au contraire, une variance élevée pourrait indiquer la nécessité d'augmenter le nombre de points tirés ou le nombre de réplifications.

— **Écart-type** (σ) :

$$\sigma = \sqrt{\sigma^2}$$

Comme l'écart-type a la même unité que la grandeur mesurée (le volume), il est utilisé pour interpréter la précision des résultats.

— **Intervalle de confiance à 95%** :

$$IC = \mu \pm t_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}}$$

où $t_{\alpha/2}$ est la valeur critique de la distribution de Student pour $n-1$ degrés de liberté, μ est la moyenne, σ est l'écart-type et n est le nombre de réplifications.

L'intervalle de confiance à 95% permet de quantifier l'incertitude sur l'estimation du volume. Il garantit, avec une probabilité de 95%, que la "vraie" valeur du volume se trouve à l'intérieur de cet intervalle. Plus l'intervalle est étroit, plus la précision de l'estimation est élevée. La largeur de l'intervalle dépend de la variance, du nombre de réplifications n et de la valeur critique $t_{\alpha/2}$ de la loi de Student.

```

1  double averageVolume = totalVolume / numReplications;
2
3  double sumSquaredDiffs = 0.0;
4  for (double volume : volumes) {
5      sumSquaredDiffs += std::pow(volume - averageVolume, 2);
6  }
7  double variance = sumSquaredDiffs / (numReplications - 1);
8  double stdDeviation = std::sqrt(variance);
9
10 double t_value = 2.262; // Student's t-value for a 95% confidence interval (n = 10
    - 1 = 9 degrees of freedom)
11 double marginOfError = t_value * (stdDeviation / std::sqrt(numReplications));
12 double lowerBound = averageVolume - marginOfError;
13 double upperBound = averageVolume + marginOfError;
```

```
● rsegerie@rsegerie-Inspiron-15-3520:~/Bureau/ZZ3/IDM/TP5-IDM/scripts$ ./runParallelSim.sh
Average estimated volume after 10 replications: 4.1888050000
Variance of estimated volumes: .0000000023
Standard deviation of estimated volumes: .0000479583
95% confidence interval for estimated volume: [4.1887752752, 4.1888347248]
All simulations completed. Total time taken: 214.567 seconds.
```

FIGURE 7 – Indicateurs statistiques pour la simulation en parallèle

```
>> Average estimated volume after 10 replications: 4.1888
>> Standard Deviation: 5.28373e-05
>> Variance: 2.79178e-09
>> 95% Confidence Interval: [4.18877, 4.18884]
Time used : 1336.78 seconds
```

FIGURE 6 – Indicateurs statistiques pour la simulation en séquentielle

7 Différentes techniques et bibliothèques de parallélisation

7.1 OpenMP (Open Multi-Processing)

OpenMP est une interface de programmation pour le calcul parallèle sur architecture à mémoire partagée. L'API prends en charge de nombreuses plateformes (incluant GNU/Linux, OS X et Windows) pour les langages de programmation tels que le **C**, **C++** et **Fortran**.

7.1.1 Simulation d'un neutron

Nous allons implémenter la simulation d'un neutron afin de manipuler la bibliothèque **OpenMP**. Pour modéliser le comportement d'un neutron dans un milieu d'eau, plusieurs paramètres clés sont définis :

1. le **libre parcours moyen** : déterminant la distance moyenne parcourue par un neutron avant qu'un événement (collision, absorption) ne se produise.
2. le **rayon de la goutte d'eau** : définissant la frontière au-delà de laquelle le neutron est considéré comme échappé.
3. le **nombre maximum de pas** : représentant la limite du nombre de déplacements d'un neutron avant d'arrêter la simulation.


```

1  #define LAMBDA 1.0
2  #define RADIUS 10.0
3  #define MAX_STEPS 10000
4
5  double random_step() {
6      return -LAMBDA * log((double)rand() / RAND_MAX);
7  }
8
9  void simulate_neutron(int dimension, int *escaped, int *absorbed) {
10     double x = 0.0, y = 0.0, z = 0.0;
11     int steps = 0;
12
13     while (steps < MAX_STEPS) {
14         steps++;
15
16         double step_length = random_step();
17
18         if (dimension == 1) {
19             x += (rand() % 2 ? 1 : -1) * step_length;
20         } else if (dimension == 2) {
21             double angle = ((double)rand() / RAND_MAX) * 2 * M_PI;
22             x += step_length * cos(angle);
23             y += step_length * sin(angle);
24         } else if (dimension == 3) {
25             double theta = ((double)rand() / RAND_MAX) * 2 * M_PI;
26             double phi = acos(1 - 2 * ((double)rand() / RAND_MAX));
27             x += step_length * sin(phi) * cos(theta);
28             y += step_length * sin(phi) * sin(theta);
29             z += step_length * cos(phi);
30         }
31
32         double distance = sqrt(x * x + y * y + z * z);
33         if (distance > RADIUS) {
34             #pragma omp atomic
35             (*escaped)++;
36             return;
37         }
38
39         if (((double)rand() / RAND_MAX) < 0.1) {
40             #pragma omp atomic
41             (*absorbed)++;
42             return;
43         }
44     }
45 }
46
47 int main(int argc, char *argv[]) {
48     if (argc != 3) {
49         printf( "Usage: ./[executable_file] {total_neutrons} {dimension}...\n" );
50         exit( 0 );
51     }
52
53     int total_neutrons = atoi(argv[1]);
54     int dimension      = atoi(argv[2]);
55
56     int escaped  = 0;
57     int absorbed = 0;
58
59     srand(time(NULL));
60
61     #pragma omp parallel for

```

```
62     for (int i = 0; i < total_neutrons; i++) {
63         int local_escaped = 0;
64         int local_absorbed = 0;
65         simulate_neutron(dimension, &local_escaped, &local_absorbed);
66
67         #pragma omp atomic
68         escaped += local_escaped;
69
70         #pragma omp atomic
71         absorbed += local_absorbed;
72     }
73
74     printf("Total neutrons: %td\n", total_neutrons);
75     printf("* escaped : %td\n", escaped);
76     printf("* absorbed: %td\n", absorbed);
77
78     return 0;
79 }
```

7.2 pthread

La bibliothèque pthread.h fait partie de la norme POSIX et **permet de créer, gérer et synchroniser des threads dans des programmes C/C++**.

Cette bibliothèque est largement utilisée pour la programmation multithread sur les systèmes Unix/Linux et macOS.

8 Application de la génération aléatoire de séquences en bioinformatique avec Mersenne Twister (MT)

L'objectif ici est de **simuler des séquences de nucléotides (A, C, G, T)** et de **rechercher une séquence cible prédéfinie**.

Nous utiliserons le générateur Mersenne Twister pour réaliser cette simulation.

Dans un premier temps, nous exécuterons plusieurs simulations séquentielles afin de déterminer le nombre de tentatives nécessaires pour obtenir la séquence cible, représentée ici par l'oligopeptide "gattaca". Par la suite, nous paralléliserons ces essais à l'aide de la technique de **sequence splitting** abordée précédemment.

Enfin, nous effectuerons une analyse statistique comprenant le calcul de l'**écart-type** et de l'**intervalle de confiance**, afin d'estimer la probabilité d'obtenir, par hasard, des séquences plus longues, comme un génome humain).

```
● rsegerie@rsegerie-Inspiron-15-3520:~/Bureau/ZZ3/IDM/TP5-IDM/bonus/BioSequence$ ./main
Starting the simulation to find the target sequence : GATTACA
Replication 20 : Attempts = 1273
Replication 23 : Attempts = 10111
Replication 32 : Attempts = 12392
Replication 24 : Attempts = 3343
Replication 13 : Attempts = 15431
Replication 33 : Attempts = 3862
Replication 34 : Attempts = 3507
Replication 14 : Attempts = 5446
Replication 38 : Attempts = 28718
Replication 26 : Attempts = 25983
Replication 9 : Attempts = 31368
Replication 15 : Attempts = 17095
Replication 25 : Attempts = 23985
Replication 39 : Attempts = 12041
Replication 40 : Attempts = 1017
Replication 17 : Attempts = 643
Replication 35 : Attempts = 41643
Replication 27 : Attempts = 21378
Replication 28 : Attempts = 1773
Replication 36 : Attempts = 17845
Replication 1 : Attempts = 19219
Replication 18 : Attempts = 30795
Replication 2 : Attempts = 7521
Replication 21 : Attempts = 59067
Replication 37 : Attempts = 13560
Replication 22 : Attempts = 4679
Replication 16 : Attempts = 40817
Replication 29 : Attempts = 69102
Replication 10 : Attempts = 51488
Replication 11 : Attempts = 3599
Replication 3 : Attempts = 17949
Replication 19 : Attempts = 27405
Replication 30 : Attempts = 15376
Replication 12 : Attempts = 21830
Replication 31 : Attempts = 14614
Replication 4 : Attempts = 10144
Replication 5 : Attempts = 71575
Replication 6 : Attempts = 10679
Replication 7 : Attempts = 13095
Replication 8 : Attempts = 3084

=== Results ===
Number of replications : 40
Target sequence GATTACA
Sequence length : 7
Mean number of attempts : 19611.30
Standard deviation of attempts 17947.58
95% Confidence Interval for attempts: [14049.29, 25173.31]
```

FIGURE 8 – Recherche d'une séquence cible ("GATTACA") avec Mersenne Twister (approche parallèle)

La probabilité de tirer **au hasard** la bonne lettre à une position donnée est :

$$p_1 = \frac{1}{4}$$

Pour générer la séquence correcte de **3 milliards de bases**, chaque position doit être correcte et indépendante des autres. La probabilité totale P de générer **par hasard** cette chaîne exacte est donc le produit des probabilités de chaque position :

$$P = \left(\frac{1}{4}\right)^N$$

où $N = 3 \times 10^9$, correspondant à 3 milliards de nucléotides.

En réécrivant, nous obtenons :

$$P = 4^{-N} = 4^{-3 \times 10^9}$$

On peut également écrire 4 sous la forme 2^2 , ce qui donne :

$$P = (2^2)^{-3 \times 10^9} = 2^{-2 \cdot 3 \times 10^9} = 2^{-6 \times 10^9}$$

```

1  ** Target sequence: AAATTTGCGTTTCGATTAG
2  sub-sequence: A
3  * attempts: 1
4  * probability (theoretical): 0.2500000000
5  =====
6  sub-sequence: AA
7  * attempts: 21
8  * probability (theoretical): 0.0625000000
9  =====
10 sub-sequence: AAA
11 * attempts: 11
12 * probability (theoretical): 0.0156250000
13 =====
14 sub-sequence: AAAT
15 * attempts: 971
16 * probability (theoretical): 0.0039062500
17 =====
18 sub-sequence: AAATT
19 * attempts: 2507
20 * probability (theoretical): 0.0009765625
21 =====
22 sub-sequence: AAATTT
23 * attempts: 2525
24 * probability (theoretical): 0.0002441406
25 =====
26 sub-sequence: AAATTTG
27 * attempts: 1451
28 * probability (theoretical): 0.0000610352
29 =====
30 sub-sequence: AAATTTGC
31 * attempts: 28257
32 * probability (theoretical): 0.0000152588
33 =====
34 sub-sequence: AAATTTGCG
35 * attempts: 797676
36 * probability (theoretical): 0.0000038147
37 =====
38 sub-sequence: AAATTTGCGT
39 * attempts: 361018
40 * probability (theoretical): 0.0000009537
41 =====
42 sub-sequence: AAATTTGCGTT
43 * attempts: 12687274
44 * probability (theoretical): 0.0000002384
45 =====
46 sub-sequence: AAATTTGCGTTC
47 * attempts: 10513771
48 * probability (theoretical): 0.0000000596
49 =====
50 sub-sequence: AAATTTGCGTTTCG
51 * attempts: 95543582
52 * probability (theoretical): 0.0000000149
53 =====
54 sub-sequence: AAATTTGCGTTTCGA
55 * attempts: 55044844
56 * probability (theoretical): 0.0000000037
57 =====
58 sub-sequence: AAATTTGCGTTTCGAT
59 * attempts: 423629810
60 * probability (theoretical): 0.0000000009
61 =====

```

Nous avons généré successivement des sous-séquences de la chaîne cible, et avons mesuré le nombre d'essais nécessaires pour obtenir chaque sous-séquence.

Les résultats montrent que, à mesure que la sous-séquence devient plus longue, le nombre d'essais nécessaires pour la générer augmente de manière exponentielle, ce qui est en accord avec la probabilité théorique. Par exemple, pour la sous-séquence "A", qui représente une seule position, il suffit d'un seul essai pour l'obtenir. Cependant, pour des sous-séquences plus longues, comme "AAATTGCGTTCGAT", le nombre d'essais nécessaires est extrêmement élevé (423,629,810 essais).