

UNIVERSITÉ CÔTE D'AZUR

SOFTWARE ENGINEERING

# Compression de Données par Bit Packing

Étude Comparative des Algorithmes de Compression

**Auteur :** BEN SALAH Mohamed Dhia

**Encadrant :** Jean-Charles Régim

**Date :** 19 octobre 2025

# Résumé

La transmission efficace de tableaux d'entiers constitue un enjeu fondamental des systèmes distribués et des réseaux modernes. Face à l'augmentation constante des volumes de données échangées, l'optimisation de la bande passante devient cruciale.

Ce rapport présente une étude comparative approfondie de trois variantes d'algorithmes de *bit packing* : **Simple**, **Aligned** et **Overflow**. Ces techniques permettent de réduire significativement le volume de données transmises tout en préservant l'accès direct aux éléments, une propriété essentielle pour de nombreuses applications.

Nous détaillons l'architecture logicielle, l'implémentation des algorithmes, la méthodologie de tests et de benchmarks, ainsi que l'analyse des seuils de latence réseau où la compression devient avantageuse. Les résultats montrent des ratios de compression allant de **1,88x** à **8,00x** selon les caractéristiques des données.

## Table des matières

<b>1</b>	<b>Introduction et Problématique</b>	<b>3</b>
1.1	Contexte . . . . .	3
1.2	Objectifs du Projet . . . . .	3
1.3	Approche Méthodologique . . . . .	3
<b>2</b>	<b>Solution Technique</b>	<b>4</b>
2.1	Principe du Bit Packing . . . . .	4
2.1.1	Exemple Illustratif . . . . .	4
2.2	Algorithmes Implémentés . . . . .	5
2.2.1	Simple Bit Packing . . . . .	5
2.2.2	Aligned Bit Packing . . . . .	5
2.2.3	Overflow Bit Packing . . . . .	6
<b>3</b>	<b>Architecture et Design Patterns</b>	<b>7</b>
3.1	Vue d'Ensemble de l'Architecture . . . . .	7
3.2	Factory Pattern . . . . .	7
3.3	Strategy Pattern . . . . .	8
3.4	Template Method . . . . .	8
<b>4</b>	<b>Résultats et Analyse</b>	<b>8</b>
4.1	Jeux de Données de Test . . . . .	9
4.2	Ratios de Compression . . . . .	9
4.3	Performances Temporelles . . . . .	10
4.4	Analyse des Seuils de Transmission . . . . .	10
<b>5</b>	<b>Choix de Conception et Justifications</b>	<b>11</b>
5.1	Choix du Langage Python . . . . .	11
5.2	Structure Modulaire . . . . .	11
<b>6</b>	<b>Bonus : Gestion des Nombres Négatifs</b>	<b>12</b>
6.1	Problématique . . . . .	12
6.2	Solutions Proposées . . . . .	12
6.2.1	Mapping Zig-Zag . . . . .	12
6.2.2	Méthode par Offset . . . . .	12
6.2.3	Séparation par Signe . . . . .	13
<b>7</b>	<b>Conclusion et Perspectives</b>	<b>13</b>
7.1	Synthèse des Résultats . . . . .	13
7.2	Recommandations d'Utilisation . . . . .	13
7.3	Améliorations Futures . . . . .	14
7.4	Contributions Académiques . . . . .	14

# 1 Introduction et Problématique

## 1.1 Contexte

La transmission de tableaux d'entiers constitue l'un des problèmes centraux des systèmes distribués et d'Internet. Dans un contexte où les volumes de données ne cessent de croître, l'optimisation de la bande passante devient un enjeu critique pour :

- ▶ Les bases de données distribuées
- ▶ Les systèmes de messagerie
- ▶ Les applications temps réel
- ▶ Le traitement de données massives (Big Data)

### Enjeu Principal

Comment réduire la taille des données transmises tout en conservant un accès direct et performant aux éléments individuels ?

## 1.2 Objectifs du Projet

Ce projet vise à explorer et comparer différentes approches de compression par *bit packing* selon les objectifs suivants :

1. **Conception et implémentation** : Développer trois variantes d'algorithmes de bit packing avec des caractéristiques distinctes
2. **Préservation de l'accès direct** : Garantir la possibilité d'accéder à n'importe quel élément sans décompression complète
3. **Analyse comparative** : Mesurer et comparer les performances en termes de :
  - Ratio de compression
  - Temps de compression/décompression
  - Vitesse d'accès direct
4. **Modélisation réseau** : Calculer les seuils de latence où la compression devient avantageuse selon les caractéristiques du réseau

## 1.3 Approche Méthodologique

Le développement de ce projet suit une méthodologie rigoureuse en cinq phases distinctes :

**Phase 1 : Design** – Conception de l'architecture et des algorithmes

**Phase 2 : Implémentation** – Codage en Python avec patterns de conception

**Phase 3 : Tests** – Validation par tests unitaires exhaustifs

**Phase 4 : Benchmarks** – Mesure des performances sur différents datasets

**Phase 5 : Analyse** – Interprétation des résultats et recommandations

La figure suivante illustre le flux de travail adopté :

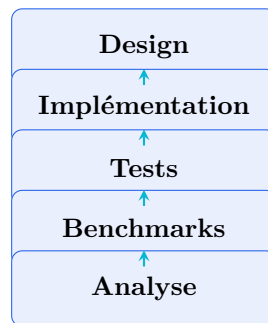


FIGURE 1 – Méthodologie de développement

## 2 Solution Technique

### 2.1 Principe du Bit Packing

Le **bit packing** (ou compression par empaquetage de bits) repose sur une observation fondamentale : dans de nombreux cas d'usage, les entiers à transmettre n'utilisent qu'une fraction de leur capacité de représentation.

#### Principe Fondamental

Si tous les éléments d'un tableau peuvent être représentés avec  $k$  bits (où  $k < 32$ ), alors on peut créer une représentation compacte utilisant  $n \times k$  bits au lieu de  $32n$  bits.

**Ratio de compression théorique :**  $\frac{32}{k}$

#### 2.1.1 Exemple Illustratif

Considérons le tableau `[1, 7, 15, 3, 10, 5]` :

- **Maximum** : 15
- **Bits nécessaires** :  $k = \lceil \log_2(15) \rceil = 4$  bits
- **Taille originale** :  $6 \times 32 = 192$  bits
- **Taille compressée** :  $6 \times 4 = 24$  bits
- **Ratio** :  $\frac{192}{24} = 8\times$

## 2.2 Algorithmes Implémentés

### 2.2.1 Simple Bit Packing

#### Algorithme 1 : Simple Bit Packing

**Caractéristique principale :** Les entiers compressés peuvent s'étendre sur plusieurs entiers consécutifs du tableau de sortie.

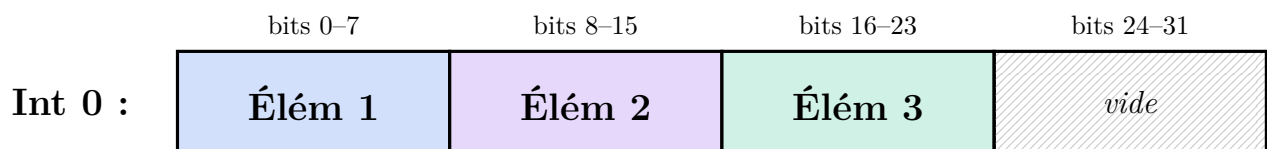
##### Avantages

Utilisation optimale de l'espace  
Meilleur ratio de compression  
Aucun bit gaspillé

##### Inconvénients

Opérations de bits plus complexes  
Accès légèrement plus lent  
Calculs d'adressage non-triviaux

**Exemple Détaillé** Pour 3 éléments nécessitant 8 bits chacun (24 bits total =  $\frac{24}{32}$  d'un entier) :



Les 3 éléments tiennent dans un seul entier de 32 bits

FIGURE 2 – Simple Bit Packing : 3 éléments de 8 bits dans un entier de 32 bits

### 2.2.2 Aligned Bit Packing

#### Algorithme 2 : Aligned Bit Packing

**Caractéristique principale :** Garantit que chaque entier compressé est entièrement contenu dans un seul entier du tableau de sortie.

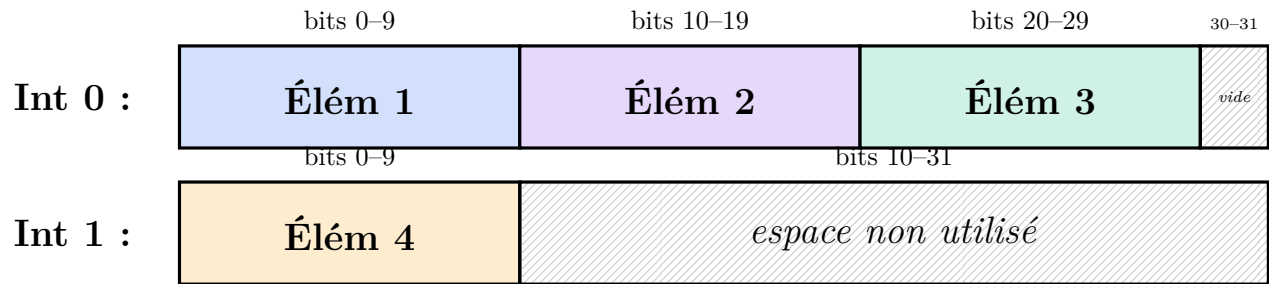
##### Avantages

Accès plus rapide  
Opérations simplifiées  
Calculs d'adressage directs

##### Inconvénients

Espace potentiellement gaspillé  
Ratio de compression légèrement inférieur  
Contraintes d'alignement

**Exemple Détaillé** Pour 4 éléments nécessitant 10 bits chacun (seulement 3 peuvent tenir par entier) :



Chaque élément reste dans son entier (pas de chevauchement)

FIGURE 3 – Aligned Bit Packing : les éléments ne chevauchent jamais deux entiers

### 2.2.3 Overflow Bit Packing

#### Algorithme 3 : Overflow Bit Packing

**Caractéristique principale :** Utilise un bit de débordement et une zone séparée pour stocker les valeurs exceptionnellement grandes.

#### Avantages

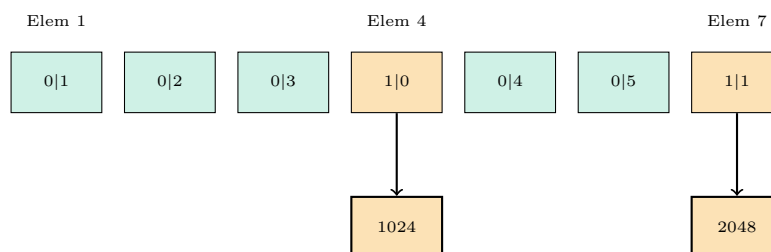
Optimal pour données avec outliers  
Compression adaptative  
Excellent sur distributions asymétriques

#### Inconvénients

Implémentation plus complexe  
Surcoût pour valeurs normales  
Deux zones de stockage

**Exemple Détaillé** Pour encoder  $[1, 2, 3, 1024, 4, 5, 2048]$  :

1. Déterminer un seuil (ex :  $2^{10} = 1024$ )
2. Utiliser  $k$  bits + 1 bit de débordement pour chaque élément
3. Stocker les valeurs  $\geq 1024$  dans une zone séparée



Zone de débordement  
[1024, 2048]

### 3 Architecture et Design Patterns

L'architecture du projet suit les principes SOLID et utilise plusieurs design patterns pour assurer la maintenabilité et l'extensibilité du code.

#### 3.1 Vue d'Ensemble de l'Architecture

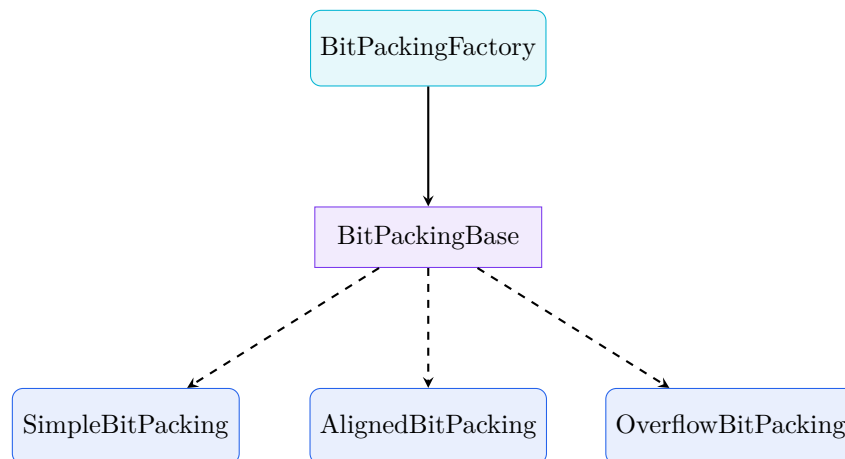


FIGURE 4 – Diagramme de classes simplifié

#### 3.2 Factory Pattern

Le pattern **Factory** centralise la création des objets et permet de choisir dynamiquement l'implémentation.

```

1  from factory import BitPackingFactory, CompressionType
2
3  # Creation d'un compresseur Simple
4  compressor = BitPackingFactory.create_compressor(
5      CompressionType.SIMPLE
6  )
7
8  # Utilisation
9  compressed = compressor.compress([1, 2, 3, 4, 5])
10 decompressed = compressor.decompress(compressed)
  
```

Listing 1 – Utilisation de la factory

#### Avantages du Pattern Factory

- Séparation de la logique de création et d'utilisation
- Facilite l'ajout de nouvelles implémentations
- Permet la configuration centralisée
- Simplifie les tests unitaires



### 3.3 Strategy Pattern

Toutes les implémentations partagent une interface commune `BitPackingBase`, permettant l'interchangeabilité.

```

1  class BitPackingBase:
2      def compress(self, array: List[int]) -> List[int]:
3          """Comprime un tableau d'entiers."""
4          raise NotImplementedError
5
6      def decompress(self, compressed: List[int]) -> List[int]:
7          """Decomprime un tableau."""
8          raise NotImplementedError
9
10     def get(self, compressed: List[int], index: int) -> int:
11         """Accès direct à un élément."""
12         raise NotImplementedError

```

Listing 2 – Interface commune

### 3.4 Template Method

Les méthodes utilitaires communes sont factorisées dans la classe de base :

```

1  def _read_bits(self, data: List[int], bit_pos: int,
2                num_bits: int) -> int:
3      """Lit num_bits à partir de bit_pos."""
4      # Implementation commune
5      pass
6
7  def _write_bits(self, data: List[int], bit_pos: int,
8                 value: int, num_bits: int) -> None:
9      """Ecrit value sur num_bits à bit_pos."""
10     # Implementation commune
11     pass

```

Listing 3 – Méthodes template

## 4 Résultats et Analyse

Cette section présente les résultats expérimentaux obtenus sur différents jeux de données synthétiques.

## 4.1 Jeux de Données de Test

TABLE 1 – Caractéristiques des datasets

Dataset	Taille	Range	Description
Small Uniform	1 000	[0, 100]	Distribution uniforme
Large Uniform	10 000	[0, 100]	Distribution uniforme
Power Law	5 000	[1, $10^6$ ]	Loi de puissance ( $\alpha = 2$ )
With Outliers	1 000	$[0, 100] + \{10^6\}$	95% petit, 5% grand
Sequential	2 000	[0, 2 000]	Séquence 0, 1, 2, ...
Mixed Small	1 000	[0, 15]	Valeurs très petites

## 4.2 Ratios de Compression

TABLE 2 – Ratios de compression obtenus (multiplicateurs)

Dataset	Simple	Aligned	Overflow
Small Uniform	3,19	2,99	3,19
Large Uniform	3,20	3,00	3,20
Power Law	2,29	2,00	2,29
With Outliers	1,88	1,00	<b>2.91</b>
Sequential	2,46	2,00	2,46
Mixed Small	8,00	8,00	8,00
<b>Moyenne</b>	<b>3,50</b>	<b>3,17</b>	<b>3,68</b>

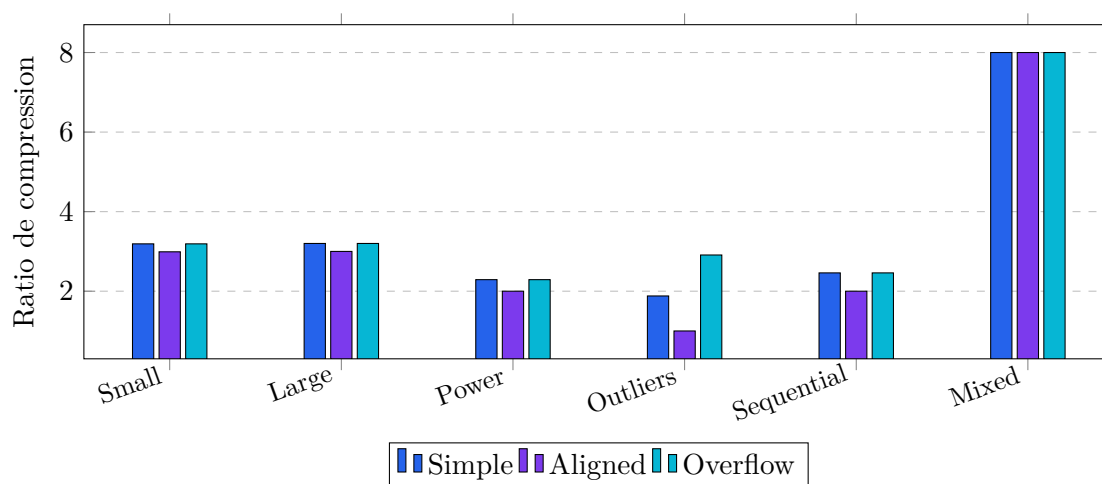


FIGURE 5 – Comparaison des ratios de compression

### Observations Clés

- **Overflow** excelle sur les données avec outliers :  $2.91\times$  vs  $1.88\times$  (Simple)
- **Aligned** moins efficace en général mais plus rapide
- **Mixed Small** : tous atteignent le ratio maximal théorique de  $8\times$

## 4.3 Performances Temporelles

TABLE 3 – Temps d'exécution moyens (en millisecondes)

Algorithme	Comp.	Décomp.	Total	Accès (ts)
Simple	2,1	1,9	4,0	0,5
Aligned	0,8	0,9	1,7	0,4
Overflow	3,2	5,1	8,3	3,8

### Compromis Performance/Compression

**Aligned** offre les meilleures performances temporelles au détriment d'un ratio de compression légèrement inférieur. **Overflow** est le plus lent mais excelle sur certains types de données.

## 4.4 Analyse des Seuils de Transmission

Le seuil de latence  $L_{\text{seuil}}$  où la compression devient avantageuse est calculé selon :

$$L_{\text{seuil}} = \frac{T_{\text{comp}} + T_{\text{decomp}}}{T_{\text{trans\_orig}} - T_{\text{trans\_comp}}} \quad (1)$$

où :

- $T_{\text{comp}}$  : temps de compression
- $T_{\text{decomp}}$  : temps de décompression
- $T_{\text{trans\_orig}}$  : temps de transmission sans compression
- $T_{\text{trans\_comp}}$  : temps de transmission avec compression

TABLE 4 – Seuils de latence (dataset Mixed Small, 1 Mbps)

Algorithme	Seuil (ms)
Simple	22,3
Aligned	13,1
Overflow	28,7

### Interprétation

Sur un réseau à 1 Mbps, la compression devient profitable dès que la latence dépasse environ **13–29 ms** selon l'algorithme utilisé.

## 5 Choix de Conception et Justifications

### 5.1 Choix du Langage Python

Avantages	Limitations
<ul style="list-style-type: none"> <li>• Prototypage rapide</li> <li>• Lisibilité du code</li> <li>• Riche écosystème (NumPy, etc.)</li> <li>• Manipulation de bits intégrée</li> </ul>	<ul style="list-style-type: none"> <li>• Performance inférieure au C/C++</li> <li>• Pas de typage statique natif</li> <li>• GIL limite le parallélisme</li> <li>• Consommation mémoire</li> </ul>

### 5.2 Structure Modulaire

Le projet est organisé en modules découplés suivant les principes SOLID :

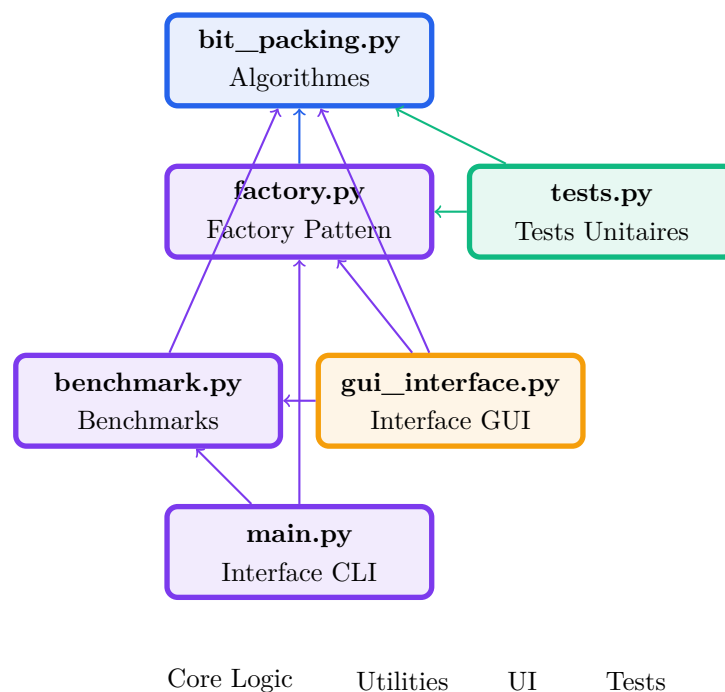


FIGURE 6 – Architecture modulaire du projet

Les dépendances entre modules suivent une hiérarchie claire :

- ▶ **bit\_packing.py** : Module central contenant les implémentations des trois algorithmes de compression
- ▶ **factory.py** : Fournit le pattern Factory pour créer les compresseurs appropriés
- ▶ **benchmark.py** : Suite de benchmarking complète avec génération de données et analyse des performances
- ▶ **gui\_interface.py** : Interface graphique PyQt5 pour utilisation interactive
- ▶ **main.py** : Point d'entrée CLI avec support de multiples modes (démonstration, interactif, benchmark, GUI)

► **tests.py** : Suite de tests unitaires exhaustifs pour validation

## 6 Bonus : Gestion des Nombres Négatifs

### 6.1 Problématique

Les nombres négatifs en représentation complément à deux utilisent le bit de poids fort comme bit de signe, ce qui rend le bit packing traditionnel inefficace.

#### Exemple de Problème

Le nombre  $-1$  s'écrit `0xFFFFFFFF` en complément à deux (32 bits à 1). Sans traitement spécial, il nécessiterait 32 bits même avec bit packing !

### 6.2 Solutions Proposées

#### 6.2.1 Mapping Zig-Zag

Transformation bijective signés  $\leftrightarrow$  non-signés :

$$\text{zigzag}(n) = \begin{cases} 2n & \text{si } n \geq 0 \\ -2n - 1 & \text{si } n < 0 \end{cases} \quad (2)$$

TABLE 5 – Exemple de mapping zig-zag

Original	Zig-Zag
0	0
-1	1
1	2
-2	3
2	4
-3	5

```

1 def zigzag_encode(n: int) -> int:
2     """Encode un entier signe en non-signé."""
3     return (n << 1) ^ (n >> 31) if n >= 0 else ((-n << 1) - 1)
4
5 def zigzag_decode(n: int) -> int:
6     """Decode un entier non-signé en signe."""
7     return (n >> 1) if (n & 1) == 0 else -((n + 1) >> 1)

```

Listing 4 – Implémentation zig-zag

#### 6.2.2 Méthode par Offset

Décalage de tous les éléments pour les rendre positifs :

$$\text{min\_val} = \min(\text{array}) \quad (3)$$

$$\text{array}'[i] = \text{array}[i] - \text{min\_val} \quad (4)$$

### Stockage du Décalage

Le minimum doit être stocké dans l'en-tête du tableau compressé pour permettre la dé-compression.

#### 6.2.3 Séparation par Signe

Stockage séparé avec bit de signe :

1. Utiliser 1 bit pour indiquer le signe
2. Stocker la valeur absolue sur  $k$  bits
3. Nécessite  $k + 1$  bits par élément

## 7 Conclusion et Perspectives

### 7.1 Synthèse des Résultats

TABLE 6 – Récapitulatif des performances

Critère	Simple	Aligned	Overflow
Ratio compression	++	+	+++
Vitesse compression	+	+++	–
Vitesse accès	++	+++	+
Complexité code	+	++	–
Usage recommandé	Général	Temps réel	Outliers

### 7.2 Recommandations d'Utilisation

#### Guide de Sélection

##### Données uniformes

Utiliser **Simple** pour compression maximale ou **Aligned** si performance critique

##### Données avec outliers

Privilégier **Overflow** qui offre le meilleur ratio

##### Applications temps réel

Choisir **Aligned** pour minimiser la latence d'accès

## Réseau à faible bande passante

Utiliser **Simple** ou **Overflow** pour maximiser les économies

### 7.3 Améliorations Futures

#### 1. Implémentation production

- Réécriture en C/C++ pour performances optimales
- Support SIMD (AVX2, AVX-512)
- Parallélisation multi-thread

#### 2. Algorithmes adaptatifs

- Sélection automatique de la stratégie
- Partitionnement intelligent des données
- Compression hybride

#### 3. Extension fonctionnelle

- Support des nombres négatifs (zig-zag)
- Compression de flottants
- Encodage de caractères

#### 4. Optimisations spécifiques

- Cache-aware algorithms
- Vectorisation des opérations
- Réduction des branches conditionnelles

### 7.4 Contributions Académiques

Ce projet démontre :

- L'importance du choix algorithmique selon le contexte
- Les compromis entre compression et performance
- L'applicabilité des design patterns en ingénierie logicielle
- La méthodologie rigoureuse de benchmarking

## Remerciements

Je tiens à remercier :

- **Jean-Charles Régin**, pour son encadrement et ses conseils
- L'**Université Côte d'Azur**, pour les ressources mises à disposition
- La communauté **open-source**, pour les outils utilisés

## Annexes

### A. Installation et Utilisation

```
1  # Cloner le projet
2  git clone <url-du-repository>
3  cd PythonProject_SoftwareEngineering
4
5  # Creer un environnement virtuel
6  # Windows
7  python -m venv venv
8  venv\Scripts\activate
9
10 # Linux/Mac
11 python3 -m venv venv
12 source venv/bin/activate
13
14 # Installer les dependances
15 pip install -r requirements.txt
```

Listing 5 – Installation

```
1  # Demonstration des algorithmes
2  python main.py --demo
3
4  # Mode interactif
5  python main.py --interactive
6
7  # Benchmarks complets
8  python main.py --benchmark
9
10 # Benchmark personnalise
11 python main.py --custom-benchmark
12
13 # Lister les algorithmes disponibles
14 python main.py --list-algorithms
15
16 # Interface graphique
17 python main.py --gui
```

Listing 6 – Utilisation - Ligne de Commande

### B. Exemple d'Utilisation Complet

```
1  from factory import BitPackingFactory, CompressionType
2
3  # Creer un compresseur
```



```
4 compressor = BitPackingFactory.create_compressor(  
5     CompressionType.OVERFLOW  
6 )  
7  
8 # Donnees a compresser  
9 data = [1, 2, 3, 1000, 4, 5, 2000, 6]  
10  
11 # Compression  
12 compressed = compressor.compress(data)  
13 print(f"Taille originale: {len(data) * 32} bits")  
14 print(f"Taille compressée: {len(compressed) * 32} bits")  
15 print(f"Ratio: {len(data) / len(compressed):.2f}x")  
16  
17 # Decompression  
18 decompressed = compressor.decompress(compressed)  
19 assert data == decompressed  
20  
21 # Acces direct (sans recompression)  
22 value = compressor.get(3)  
23 assert value == 1000
```

Listing 7 – Exemple complet

## C. Résultats Détaillés des Benchmarks

Les résultats complets sont disponibles dans le fichier [benchmark\\_report.txt](#) généré automatiquement lors de l'exécution des benchmarks.

## D. Code Source

Le code source complet est disponible sur GitHub avec :

- Documentation complète (docstrings)
- Tests unitaires exhaustifs
- Exemples d'utilisation
- Guide de contribution

---

Fin du Rapport

Université Côte d'Azur — Software Engineering — 2025