

Compression d'Entiers par Bit Packing

Projet de Software Engineering

LOUNIS Elies

Master 1 Informatique
Année Universitaire 2024-2025



Repository GitHub :
<https://github.com/Elieslns/bit-packing-compression>

Table des matières

1 Introduction	3
1.1 Objectifs du projet	3
2 Problématique	3
2.1 Contexte de transmission	3
2.2 Observation clé	3
2.3 Contraintes	3
2.4 Défis techniques	4
3 Solutions Techniques	4
3.1 Bit Packing - Principe fondamental	4
3.2 Variante Consecutive	4
3.3 Variante Non-Consecutive	4
3.4 Overflow Areas	5
3.5 Encodage des nombres négatifs	5
3.5.1 Problème	5
3.5.2 Solution : Offset Encoding	5
4 Choix d'Implémentation et Démarche Personnelle	6
4.1 Langage : Python	6
4.1.1 Hésitation initiale	6
4.1.2 Pourquoi Python finalement	6
4.2 Architecture : Factory Pattern	6
4.2.1 Comment j'en suis arrivé là	6
4.2.2 Ma solution avec le Factory	7
4.2.3 Ce que ça m'apporte concrètement	7
4.3 Structure des données compressées	7
4.3.1 Problème initial et réflexion	7
4.3.2 Implémentation détaillée	8
4.3.3 Choix du marqueur 0xFFFFFFFF	8
4.3.4 Coût de cette approche	8
4.4 Opérations binaires : La partie la plus difficile	8
4.4.1 Mes galères avec les bits	8
4.4.2 Extraction de bits : Comment je m'y suis pris	9
4.4.3 Insertion de bits : Encore plus compliqué	9
4.4.4 Pourquoi j'ai pas utilisé de bibliothèque	9
4.4.5 Comment j'ai testé tout ça	10
5 Difficultés Rencontrées et Solutions Personnelles	10
5.1 Le cauchemar du mode Non-Consecutive avec Overflow	10
5.1.1 Où j'ai bloqué pendant 2 jours	10
5.1.2 Ma première tentative (ratée)	10
5.1.3 La solution (après avoir beaucoup cherché)	10
5.2 Le problème des nombres négatifs (que j'avais oublié)	11
5.2.1 Quand j'ai réalisé le problème	11
5.2.2 Ce que j'ai regardé comme solutions	11
5.3 Mesure de performance : Mes galères avec les timings	11

5.3.1	Mes mesures étaient n'importe quoi	11
5.3.2	Ce que j'ai mis en place	11
6	Résultats et Performances	12
6.1	Ratios de compression	12
6.2	Temps d'exécution	12
6.3	Break-even Latency	12
6.4	Comparaison des variantes	12
7	Réflexions Personnelles et Apprentissages	12
7.1	Ce que j'ai appris techniquement	12
7.1.1	Manipulation binaire approfondie	12
7.1.2	Importance de la cohérence algorithme	13
7.2	Ce que j'ai appris en ingénierie logicielle	13
7.2.1	L'importance des tests	13
7.2.2	Architecture et extensibilité	13
7.3	Trade-offs et compromis acceptés	13
7.3.1	Performance vs Lisibilité	13
7.3.2	Métadonnées vs Taille	13
7.4	Ce que j'ai appris et ce que je ferais différemment	14
7.4.1	Les bonnes décisions	14
7.4.2	Mes erreurs	14
7.5	Applications futures	14
8	Conclusion	14
8.1	Objectifs atteints	14
8.2	Résultats clés	15
8.3	Applications pratiques et utilité réelle	15
8.3.1	Cas d'usage que j'ai identifiés	15
8.3.2	Quand NE PAS utiliser cette solution	15
8.4	Impact du projet sur ma compréhension	15
8.5	Bilan personnel	16
Annexes		17

1 Introduction

Dans le contexte de la transmission de données sur un réseau, le temps de transmission est directement proportionnel à la quantité de données à envoyer. Pour des applications nécessitant l'envoi fréquent de tableaux d'entiers (bases de données, analytics, IoT, etc.), cette contrainte peut devenir un goulot d'étranglement majeur.

Ce projet explore l'utilisation du **Bit Packing** comme méthode de compression pour réduire la taille des données transmises, tout en maintenant une contrainte forte : la possibilité d'accéder directement à un élément du tableau sans avoir à décompresser l'ensemble des données.

1.1 Objectifs du projet

- Implémenter une compression efficace d'arrays d'entiers
- Garantir un accès direct aux éléments (méthode `get(i)`)
- Gérer les cas particuliers : outliers et nombres négatifs
- Mesurer les performances et calculer le seuil de rentabilité
- Fournir une architecture extensible (Factory Pattern)

2 Problématique

2.1 Contexte de transmission

Lors de la transmission d'un tableau d'entiers standard sur un réseau, chaque entier est typiquement encodé sur 32 bits (4 octets). Pour un tableau de n éléments, la taille totale est donc de $n \times 32$ bits.

Exemple : Un tableau de 1000 entiers représente $1000 \times 4 = 4000$ octets à transmettre.

2.2 Observation clé

Dans de nombreux cas pratiques, les valeurs contenues dans le tableau n'utilisent pas la totalité des 32 bits disponibles. Par exemple :

- Pour représenter des valeurs de 0 à 100, seulement 7 bits sont nécessaires
- Pour des valeurs de 0 à 1000, seulement 10 bits suffisent

Il existe donc un potentiel significatif de compression en utilisant uniquement le nombre minimum de bits requis.

2.3 Contraintes

Trois contraintes majeures doivent être respectées :

1. **Compression sans perte** : Toutes les données doivent être récupérables exactement
2. **Accès direct** : Possibilité de lire le i -ème élément sans décompresser l'intégralité du tableau
3. **Performance** : Le temps de compression/décompression ne doit pas annuler le gain de transmission

2.4 Défis techniques

- **Outliers** : Quelques valeurs très grandes peuvent forcer l'utilisation de nombreux bits pour tous les éléments
- **Nombres négatifs** : En complément à deux, -1 s'écrit comme 111...111, impossible à compresser efficacement
- **Métadonnées** : Les informations nécessaires à la décompression doivent être stockées avec les données compressées

3 Solutions Techniques

3.1 Bit Packing - Principe fondamental

Le Bit Packing consiste à :

1. Analyser le tableau pour déterminer le nombre minimum de bits k nécessaire pour représenter la valeur maximale
2. Encoder chaque valeur sur exactement k bits (au lieu de 32)
3. Empaqueter ces valeurs bit par bit dans des entiers 32 bits

Exemple mathématique :

Pour un tableau $[1, 2, 3, 4, 5]$:

- $\max = 5 = 101_2$
- Bits nécessaires : $k = \lceil \log_2(5) \rceil + 1 = 3$ bits
- Taille non compressée : $5 \times 32 = 160$ bits
- Taille compressée : $5 \times 3 = 15$ bits
- Ratio : $160/15 \approx 10.67 \times$

3.2 Variante Consecutive

Dans cette variante, une valeur compressée peut s'étendre sur deux entiers 32 bits consécutifs.

Avantages :

- Utilisation optimale de l'espace (aucun bit gaspillé)
- Meilleur ratio de compression

Inconvénients :

- Calcul de position plus complexe pour l'accès direct
- Lecture potentiellement dans 2 entiers pour une seule valeur

Formule d'accès : Pour accéder au i -ème élément compressé sur k bits :

$$\begin{aligned} \text{position_bit} &= i \times k \\ \text{index_entier} &= \lfloor \text{position_bit}/32 \rfloor \\ \text{offset_bit} &= \text{position_bit} \bmod 32 \end{aligned}$$

3.3 Variante Non-Consecutive

Dans cette variante, une valeur ne s'étend jamais sur deux entiers. Si elle ne tient pas dans l'espace restant, elle commence au prochain entier.

Avantages :

- Calcul de position simplifié
- Lecture toujours dans un seul entier

Inconvénients :

- Bits inutilisés (padding)
- Ratio de compression légèrement moins bon

3.4 Overflow Areas

Pour gérer efficacement les outliers, une zone de débordement séparée est utilisée.

Principe :

1. Identifier les valeurs nécessitant significativement plus de bits
2. Les stocker séparément dans une "overflow area" (32 bits complets)
3. Utiliser un bit de drapeau pour chaque élément :
 - flag = 0 : valeur normale (k bits)
 - flag = 1 : index vers overflow area

Exemple :

Pour [1, 2, 3, 1000000, 4, 5] :

- Sans overflow : tous en 20 bits $\rightarrow 6 \times 20 = 120$ bits
- Avec overflow :
 - 5 valeurs normales : $5 \times (1 + 3) = 20$ bits (flag + valeur)
 - 1 overflow : $(1 + 1) + 32 = 34$ bits (flag + index + valeur complète)
 - Total : 54 bits \rightarrow Gain de 2.2×

Structure du tableau compressé avec overflow :

[données_bit_packed] [0xFFFFFFFF] [métadonnées] [overflow_area]

3.5 Encodage des nombres négatifs

3.5.1 Problème

En représentation complément à deux, les nombres négatifs ont tous leurs bits de poids fort à 1 :

- $-1 = 1111111\dots111_2$ (32 bits à 1)
- Impossible à compresser avec le bit packing standard

3.5.2 Solution : Offset Encoding

L'encodage par décalage consiste à ajouter un offset pour rendre les valeurs négatives positives :

Algorithme :

```

1 def encode(value, k_bits):
2     if value < 0:
3         return 2***(k_bits - 1) + abs(value)
4     else:
5         return value
6
7 def decode(encoded, k_bits):
8     threshold = 2***(k_bits - 1)
9     max_positive = threshold - 1

```

```

10     if encoded > max_positive:
11         return -(encoded - threshold)
12     else:
13         return encoded

```

Exemple avec 4 bits (plage : -7 à +7) :

Valeur	Encodé	Binaire
-5	$8 + 5 = 13$	1101
-3	$8 + 3 = 11$	1011
-1	$8 + 1 = 9$	1001
0	0	0000
1	1	0001
3	3	0011
5	5	0101

4 Choix d'Implémentation et Démarche Personnelle

4.1 Langage : Python

4.1.1 Hésitation initiale

Au début du projet, j'ai hésité entre Java et Python. Le sujet autorisait les deux langages, et dans mon cursus on a beaucoup utilisé Java. Mais honnêtement, je suis beaucoup plus à l'aise avec Python - c'est le langage que j'utilise le plus souvent pour mes projets personnels et mes TPs.

4.1.2 Pourquoi Python finalement

Plusieurs raisons m'ont convaincu de rester sur Python :

1. **Je connais bien le langage** : J'ai fait pas mal de projets en Python l'année dernière, donc je savais que je pourrais avancer rapidement. Avec Java, j'aurais passé plus de temps à me rappeler la syntaxe des classes abstraites et tout ça.
2. **Les opérations binaires sont lisibles** : Même si au début j'ai galéré avec les shifts et les masks, une fois qu'on comprend le principe, c'est assez clair à écrire en Python. Par exemple (`value > start_bit`) & `mask` c'est direct.
3. **Pas besoin de bibliothèques externes** : J'ai juste utilisé `math`, `time` et `statistics` qui sont dans la lib standard. Ça évite les problèmes de dépendances.
4. **Focus sur l'algo** : Le but c'était de comprendre le bit packing, pas de faire le code le plus rapide du monde. Python m'a permis de me concentrer sur la logique plutôt que sur les pointeurs ou la gestion mémoire.

4.2 Architecture : Factory Pattern

4.2.1 Comment j'en suis arrivé là

Au début, j'avais juste fait 4 classes différentes : `BitPackingConsecutive`, `BitPackingNonConsecutive` etc. Et dans mon `main.py`, j'instanciais directement les classes. Ça marchait, mais c'était assez désorganisé.

Le problème c'est que quand j'ai voulu faire des benchmarks, je devais importer toutes les classes et les instancier une par une. Mon code ressemblait à ça :

```

1 from bit_packing import BitPackingConsecutive,
2     BitPackingNonConsecutive
# ... imports qui n'en finissent pas
3
4 comp1 = BitPackingConsecutive()
5 comp2 = BitPackingNonConsecutive()
# etc.

```

C'était vraiment pas pratique. J'ai cherché sur Google "python create objects dynamically" et je suis tombé sur le Factory Pattern. J'avoue que je connaissais pas trop ce pattern avant (on l'a vu en cours mais j'avais pas tout compris).

4.2.2 Ma solution avec le Factory

J'ai créé une classe BitPackingFactory qui me permet de faire :

```

1 factory = BitPackingFactory()
2 compressor = factory.create('consecutive')

```

Ça simplifie vraiment le code. Dans mes benchmarks, je peux juste boucler sur une liste de strings et créer les compresseurs à la volée.

4.2.3 Ce que ça m'apporte concrètement

- **Plus simple à utiliser** : Je tape juste 'consecutive' au lieu de me rappeler du nom exact de la classe
- **Moins de bugs** : Quand j'ai ajouté l'overflow, j'ai juste modifié le factory. Avant j'aurais dû toucher à plein de fichiers différents et j'aurais sûrement oublié quelque chose
- **Tests plus faciles** : Mon code de test fait une boucle sur ['consecutive', 'non_consecutive', ...] et ça teste tout automatiquement

Bon, j'ai conscience que mon factory est assez basique (c'est juste un gros if/elif), mais ça fait le job pour ce projet.

4.3 Structure des données compressées

4.3.1 Problème initial et réflexion

Au début, j'avais deux options pour gérer les métadonnées nécessaires à la décompression :

Option A : Passer les métadonnées séparément

```

1 compressed_data, metadata = compress(array)
2 result = decompress(compressed_data, metadata)

```

Option B : Inclure les métadonnées dans le tableau compressé

```

1 compressed = compress(array) # Tout est dedans
2 result = decompress(compressed) # Autonome

```

J'ai choisi l'Option B pour plusieurs raisons importantes :

1. **Autonomie** : Le tableau compressé peut être transmis, sauvegardé, ou décompressé sans contexte additionnel. C'est crucial pour un système réel où les données peuvent transiter par plusieurs systèmes.
2. **Robustesse** : Impossible de "perdre" les métadonnées ou de les mélanger avec celles d'un autre tableau. Tout est auto-contenu.
3. **Simplicité d'utilisation** : L'API est plus simple. L'utilisateur manipule un seul objet (la liste compressée) au lieu de deux.

4.3.2 Implémentation détaillée

Structure : [données_bit_packed...] [0xFFFFFFFF] [métadonnées]

Métadonnées (6+ entiers) :

- [0] : 0xFFFFFFFF (marqueur de début des métadonnées)
- [1] : original_length (nombre d'éléments originaux)
- [2] : overflow_size (taille de l'overflow area)
- [3] : value_bits (bits par valeur normale)
- [4] : bits_for_overflow_index (bits pour indexer overflow)
- [5] : nb_indices (nombre d'éléments en overflow)
- [6+] : overflow_area (valeurs complètes 32 bits)

4.3.3 Choix du marqueur 0xFFFFFFFF

Le choix de 0xFFFFFFFF comme marqueur n'est pas anodin :

- C'est une valeur facilement reconnaissable (tous les bits à 1)
- Elle a peu de chances d'apparaître naturellement dans les données bit-packed (qui utilisent moins de 32 bits)
- En cas de corruption, le marqueur absent signale immédiatement un problème

4.3.4 Coût de cette approche

Oui, ajouter 6+ entiers de métadonnées augmente la taille compressée. Mais le surcoût est fixe (24 octets) indépendamment de la taille du tableau. Pour 1000 éléments, cela représente seulement 0.6% de la taille originale. Le compromis est largement acceptable pour gagner en autonomie et robustesse.

4.4 Opérations binaires : La partie la plus difficile

4.4.1 Mes galères avec les bits

Honnêtement, les opérations binaires c'est la partie où j'ai le plus galéré. J'ai passé des heures à déboguer des cas où j'avais inversé un bit ou oublié un masque. Le pire c'est qu'une petite erreur et tout le tableau est corrompu, donc c'est vraiment difficile de localiser le problème.

Au début je visualisais pas bien ce qui se passait. J'ai dû écrire plein de tests avec des petits nombres pour bien comprendre les shifts et les masks.

4.4.2 Extraction de bits : Comment je m'y suis pris

Après avoir lu pas mal de docs et regardé quelques tutos, j'ai fini par comprendre qu'il fallait faire en deux étapes :

```
1 def _extract_bits(self, value, start_bit, num_bits):
2     mask = (1 << num_bits) - 1
3     return (value >> start_bit) & mask
```

Mon raisonnement :

1. D'abord je décale avec `>>` pour mettre les bits que je veux en position basse 2. Ensuite je masque avec `&` pour virer les bits que je veux pas

Un exemple que j'ai testé :

Imaginons j'ai `value = 0b11010110` (214) et je veux les 3 bits à partir de la position 2 :

- Après `>> 2` : `0b00110101`
- Mon mask : `(1 << 3) - 1 = 0b00000111`
- Résultat : `0b00110101 & 0b00000111 = 0b00000101 = 5`

Ça marche ! Mais j'avoue que j'ai mis du temps à capter le principe du mask avec `(1 << num_bits) - 1`. Au début je faisais des trucs bizarres qui marchaient pas.

4.4.3 Insertion de bits : Encore plus compliqué

Pour insérer des bits, c'est encore plus complexe. J'ai dû m'y reprendre à plusieurs fois :

```
1 def _set_bits(self, target, value, start_bit, num_bits):
2     mask = (1 << num_bits) - 1
3     value &= mask                      # Nettoyer la valeur
4     target &= ~(mask << start_bit)    # Faire de la place
5     return target | (value << start_bit) # Insérer
```

Les 3 étapes (que j'ai compris après plusieurs bugs) :

1. D'abord je nettoie la valeur avec un mask pour être sûr qu'elle tient dans le nombre de bits voulu
2. Ensuite je mets à zéro les bits dans `target` où je vais écrire. Le `~` ça inverse le mask (j'ai mis du temps à comprendre ça)
3. Enfin je combine avec `|` pour insérer la valeur décalée

J'avoue que cette fonction m'a donné du fil à retordre. J'ai d'abord essayé sans l'étape 2 et évidemment ça écrasait n'importe quoi.

4.4.4 Pourquoi j'ai pas utilisé de bibliothèque

J'ai vu qu'il y avait `bitstring` et `bitarray` mais :

- Je voulais comprendre comment ça marche vraiment (sinon j'apprends rien)
- Je sais pas si vous aurez accepté des libs externes
- Mes fonctions sont assez simples, pas besoin d'une bibliothèque complexe pour ça

4.4.5 Comment j'ai testé tout ça

Pour être sûr que ça marche, j'ai fait plein de petits tests avec des nombres faciles à visualiser en binaire. Par exemple extraire 3 bits de 0b11010110 et vérifier que je récupère bien ce que j'attends.

J'ai aussi fait un test où j'extrais puis je réinsère la même valeur, elle doit rester identique. Ce test m'a permis de trouver plusieurs bugs au début.

5 Difficultés Rencontrées et Solutions Personnelles

5.1 Le cauchemar du mode Non-Consecutive avec Overflow

5.1.1 Où j'ai bloqué pendant 2 jours

Clairement la partie la plus difficile du projet. J'ai passé presque 2 jours là-dessus et j'ai failli abandonner l'overflow pour le mode non-consecutive.

Le problème c'est que chaque élément a une taille variable : soit 1 flag + `value_bits` pour les valeurs normales, soit 1 flag + `bits_for_overflow_index` pour les overflow. Et en mode non-consecutive, je dois m'assurer qu'un élément ne traverse jamais la frontière des 32 bits.

Mon bug principal :

Au début mon code de décompression était vraiment très lent (environ 200ms pour 1000 éléments alors que ça devrait être < 2ms). En analysant j'ai vu que je recalculais la position depuis le début pour chaque élément. Du coup c'était du $O(n^2)$ au lieu de $O(n)$.

5.1.2 Ma première tentative (ratée)

J'ai d'abord essayé d'utiliser une taille fixe maximum :

```
1 max_size = 1 + max(value_bits, bits_for_overflow_index)
2 if position + max_size > 32:
3     # Passer au prochain entier
```

Mais ça marchait pas parce qu'un petit élément overflow qui aurait pu tenir était déplacé inutilement. Mes tests échouaient.

5.1.3 La solution (après avoir beaucoup cherché)

En fait il fallait que je maintienne `bit_pos` et `int_idx` entre les itérations au lieu de tout recalculer. Ça paraît évident maintenant mais sur le moment j'avais pas vu.

```
1 bit_pos = 0 # Position courante, pas recalculée !
2 int_idx = 0
3
4 for i in range(n):
5     # Lire le flag pour connaître la taille REELLE
6     flag = read_flag(bit_pos, int_idx)
7
8     # Calculer si ça passe ou pas
9     if pos_before + taille_reelle > 32:
10        # Aller au prochain entier
```

Après avoir compris ça, je suis passé de 200ms à 1ms.

5.2 Le problème des nombres négatifs (que j'avais oublié)

5.2.1 Quand j'ai réalisé le problème

En fait au début j'avais complètement oublié de gérer les nombres négatifs. J'ai codé tout le système avec juste des positifs. C'est seulement quand j'ai fait des tests avec des valeurs négatives que j'ai vu que mon ratio de compression était quasi nul (environ 1.1x au lieu de 4x).

J'ai mis un moment à comprendre pourquoi. En fait en complément à deux, -1 c'est 0xFFFFFFFF, tous les bits à 1. Donc impossible à compresser avec ma méthode.

5.2.2 Ce que j'ai regardé comme solutions

J'ai cherché sur internet "how to compress negative integers" et j'ai trouvé plusieurs trucs :

1. **ZigZag encoding** : Apparemment c'est ce qu'utilise Protocol Buffers. Ça a l'air très efficace mais honnêtement ça semblait assez compliqué à implémenter. J'ai pas bien compris la formule avec les XOR.
2. **Bit de signe séparé** : Stocker le signe à part et compresser juste la valeur absolue. Simple mais il faut une structure en plus et je savais pas trop comment gérer ça avec mon système d'accès direct.
3. **Offset encoding** : Ajouter un offset pour tout rendre positif. Ça m'a semblé le plus simple.

J'ai choisi l'offset parce que c'était le plus simple à implémenter et ça suffit pour le projet.

5.3 Mesure de performance : Mes galères avec les timings

5.3.1 Mes mesures étaient n'importe quoi

Au début je faisais juste un `time.time()` avant et après, et j'affichais la différence. Le problème c'est que d'une exécution à l'autre j'avais des valeurs complètement différentes. Genre 0.3ms puis 1.2ms puis 0.5ms pour exactement le même code.

J'ai cherché pourquoi et j'ai compris que y'a plein de trucs qui font varier les mesures : le cache, le garbage collector de Python, les autres programmes qui tournent, etc.

5.3.2 Ce que j'ai mis en place

J'ai fait un système plus sérieux :

- **Warm-up** : Je lance 3 fois avant de mesurer, pour "chauffer" le système
- **Mesures multiples** : Je mesure 10 fois et je prends la médiane (pas la moyenne parce que la médiane est moins sensible aux valeurs bizarres)
- **Timer précis** : J'utilise `time.perf_counter()` au lieu de `time.time()` parce que c'est plus précis

Avec ça mes mesures sont beaucoup plus stables. Avant ça variait de 50%, maintenant c'est < 5%.

6 Résultats et Performances

6.1 Ratios de compression

Tests effectués sur différents types de données :

Données	Taille originale	Taille compressée	Ratio
[1-100] (1000 éléments)	4000 octets	1000 octets	4.0x
[1-1000] (1000 éléments)	4000 octets	1250 octets	3.2x
Avec outliers	4000 octets	800 octets	5.0x
Nombres négatifs [-100,100]	4000 octets	1000 octets	4.0x

6.2 Temps d'exécution

Mesures effectuées avec le protocole de timing (10 runs, médiane) :

Opération	100 éléments	1000 éléments	10000 éléments
Compression	0.05 ms	0.4 ms	4.2 ms
Décompression	0.04 ms	0.3 ms	3.8 ms
Accès direct (get)	0.002 ms	0.002 ms	0.002 ms

6.3 Break-even Latency

La compression devient rentable lorsque :

$$t_{\text{compression}} + t_{\text{décompression}} + t_{\text{transmission_compressé}} < t_{\text{transmission_non_compressé}}$$

Pour un tableau de 1000 entiers avec ratio 4.0x et bande passante de 100 Mbps :

- Overhead compression/décompression : 0.7 ms
- Gain en transmission : 3.0 ms
- Bénéfice net : +2.3 ms

Conclusion : La compression est rentable dès qu'il y a une latence réseau significative.

6.4 Comparaison des variantes

Variante	Ratio	Temps compress.	Complexité get()
Consecutive	4.0x	0.4 ms	$O(1)$ calculé
Non-Consecutive	3.8x	0.4 ms	$O(1)$ calculé
Overflow Consecutive	5.0x	0.5 ms	$O(n)$ pire cas
Overflow Non-Consecutive	4.8x	0.5 ms	$O(n)$ pire cas

7 Réflexions Personnelles et Apprentissages

7.1 Ce que j'ai appris techniquement

7.1.1 Manipulation binaire approfondie

Avant ce projet, je comprenais les opérations binaires en théorie. Maintenant, je les maîtrise en pratique. J'ai appris que :

- Le débogage d'opérations binaires nécessite de visualiser les nombres en binaire

- Un seul bit mal placé peut corrompre toute une structure de données
- Les masques binaires sont puissants mais il faut les construire avec soin
- Python gère les entiers de taille arbitraire, ce qui peut masquer des bugs (dépassement de 32 bits)

7.1.2 Importance de la cohérence algorithme

La plus grande leçon : en compression, l'encodeur et le décodeur doivent être **parfaitemnt symétriques**. Chaque décision prise lors de la compression (passer au prochain entier, ajouter du padding, etc.) doit être exactement répliquée lors de la décompression.

Cette contrainte m'a forcé à :

- Documenter précisément chaque décision
- Créer des tests exhaustifs pour valider la symétrie
- Refactorer pour extraire la logique commune

7.2 Ce que j'ai appris en ingénierie logicielle

7.2.1 L'importance des tests

Sans une suite de tests complète, ce projet aurait été impossible. J'ai créé :

- Tests unitaires pour chaque fonction binaire
- Tests d'intégration pour chaque variante de compression
- Tests de propriétés : `decompress(compress(x)) == x`
- Tests de performance pour détecter les régressions

Ces tests m'ont sauvé des dizaines de fois lors des refactorisations.

7.2.2 Architecture et extensibilité

Le Factory Pattern et l'héritage m'ont permis d'ajouter les variantes overflow sans réécrire le code existant. Cette extensibilité n'était pas gratuite : elle a nécessité de :

- Identifier les abstractions correctes (`BitPackingBase`)
- Séparer les responsabilités (`compression`, métadonnées, `overflow`)
- Penser à l'interface avant l'implémentation

Sans cette architecture, ajouter l'overflow aurait nécessité de dupliquer tout le code.

7.3 Trade-offs et compromis acceptés

7.3.1 Performance vs Lisibilité

J'aurais pu optimiser davantage (`bitarray`, Cython, C extensions). Mais pour un projet académique, j'ai privilégié :

- Code lisible et compréhensible
- Pas de dépendances externes
- Facilité de modification et d'extension

Le ratio performance/complexité était optimal pour ce projet.

7.3.2 Métadonnées vs Taille

Inclure les métadonnées dans le tableau compressé augmente la taille de 24 octets. Pour 100 éléments (400 octets), c'est 6%. Pour 10000 éléments (40 Ko), c'est 0.06%.

J'ai accepté ce surcoût pour gagner en robustesse et autonomie. C'est un compromis conscient et justifié.

7.4 Ce que j'ai appris et ce que je ferais différemment

7.4.1 Les bonnes décisions

Avec du recul, y'a des trucs que je suis content d'avoir fait :

- Le Factory Pattern : Au début je trouvais ça un peu overkill mais finalement ça m'a vraiment simplifié la vie pour les tests
- Les tests : Même si j'aime pas trop ça, les tests m'ont sauvé plusieurs fois quand j'ai modifié le code
- Faire les variantes simples d'abord : J'ai commencé par consecutive/non-consecutive avant de faire l'overflow. Sinon j'aurais été perdu

7.4.2 Mes erreurs

Si je devais refaire ce projet, je changerais plusieurs choses :

- **Commencer par les tests** : J'ai écrit beaucoup de tests après avoir codé. C'était une erreur, j'aurais dû faire du TDD dès le début. Ça m'aurait évité de déboguer pendant des heures.
- **Penser aux négatifs plus tôt** : J'ai complètement oublié les nombres négatifs au début. J'ai dû modifier plein de trucs après coup.
- **Mesurer les perfs dès le début** : J'ai découvert le bug $O(n^2)$ très tard dans le projet. Si j'avais mesuré les performances dès le début je l'aurais vu tout de suite.
- **Commenter en codant** : J'ai ajouté tous les commentaires à la fin. Du coup j'avais oublié pourquoi j'avais fait certains trucs d'une certaine manière.

Bon, c'était mon premier vrai projet de compression donc c'est normal de faire des erreurs. Au moins j'ai appris plein de choses !

7.5 Applications futures

Ce projet m'a donné des idées d'extensions :

- **Compression adaptative** : Analyser les données et choisir automatiquement la meilleure variante
- **Compression par blocs** : Diviser un grand tableau en blocs et les compresser indépendamment pour paralléliser
- **Delta encoding** : Pour des séries temporelles où les différences sont petites
- **Compression hybride** : Combiner bit packing avec d'autres techniques (RLE, dictionnaires)

8 Conclusion

8.1 Objectifs atteints

Ce projet a permis d'implémenter avec succès :

- ✓ Compression efficace par Bit Packing (ratios de 3.2x à 5.0x)
- ✓ Accès direct aux éléments sans décompression totale
- ✓ Gestion des outliers via overflow areas

- ✓ Support des nombres négatifs par offset encoding
- ✓ Architecture extensible avec Factory Pattern
- ✓ Suite complète de tests et benchmarks

8.2 Résultats clés

1. La compression Bit Packing est particulièrement efficace pour des données avec des valeurs dans une plage limitée
2. Les overflow areas améliorent significativement le ratio pour des données avec quelques outliers
3. Le temps de compression/décompression est négligeable comparé au gain de transmission
4. L'accès direct est maintenu avec une complexité $O(1)$ pour les variantes standard

8.3 Applications pratiques et utilité réelle

8.3.1 Cas d'usage que j'ai identifiés

Cette solution est particulièrement adaptée pour :

1. **IoT et capteurs** : Un capteur de température envoie des valeurs entre 0-50°C (6 bits). Sur 1000 mesures, gain de 5.3x.
2. **Analytics** : Les compteurs web (vues, clics) sont souvent < 10000 (14 bits). Pour des dashboards temps réel avec milliers de métriques, la compression réduit significativement la latence.
3. **Bases de données colonnes** : Dans une colonne "age" (0-120), chaque valeur utilise 7 bits au lieu de 32. Pour 1 million de lignes, économie de 3 Mo.
4. **Gaming** : Synchronisation d'états (positions de joueurs) sur réseau. Les coordonnées bornées se compressent très bien.

8.3.2 Quand NE PAS utiliser cette solution

Il est important de reconnaître les limitations :

- **Données très variables** : Si les valeurs utilisent déjà tous les bits (nombres aléatoires, hashs), le gain est nul.
- **Tableaux très petits** : Pour < 50 éléments, le surcoût des métadonnées (24 octets) peut annuler le gain.
- **Mises à jour fréquentes** : Modifier un élément nécessite de décompresser, modifier, recompresser. Pas adapté aux structures mutables.

8.4 Impact du projet sur ma compréhension

Ce projet m'a permis de comprendre concrètement :

- **Pourquoi** les bases de données colonnes (Parquet, ORC) utilisent le bit packing
- **Comment** Protocol Buffers et autres formats de sérialisation optimisent l'espace
- **Où** se situe le compromis entre compression et temps CPU
- **L'importance** des détails d'implémentation pour la performance

8.5 Bilan personnel

Au-delà des aspects techniques, ce projet m'a enseigné :

- La rigueur nécessaire pour implémenter correctement un algorithme bas niveau
- L'importance de tester exhaustivement, surtout pour du code binaire
- La valeur d'une architecture propre pour faciliter l'évolution
- Le besoin de compromis réfléchis entre différents objectifs (performance, lisibilité, maintenabilité)

C'est un projet que je pourrais réellement utiliser dans un contexte professionnel, moyennant quelques optimisations. Cette applicabilité concrète rend le travail fourni d'autant plus satisfaisant.

Annexes

A. Commandes d'exécution

Exécuter les exemples :

```
python main.py
```

Lancer les benchmarks :

```
python benchmark.py
```

B. Repository GitHub

Tout le code source est disponible sur :

<https://github.com/Elieslns/bit-packing-compression>

C. Structure des fichiers

```
.  
bit_packing.py          # Implémentation de base  
bit_packing_overflow.py # Compression avec overflow  
factory.py              # Factory pattern  
timing.py               # Mesures de performance  
benchmark.py            # Suite de benchmarks  
main.py                 # Point d'entrée  
README.md               # Documentation
```