

# Pruning des Réseaux de Neurones Profonds : Une Perspective d'Éparsité

Réalise par :

ISMAIL EL-Khayrany

ABDERRAHAMN Aamouch

ABDELAZIZ Laghfiri

2024 / 2025

Université Ibn Zohr  
Faculté Polydisciplinaire de Ouarzazate

---

## Résumé

Ce rapport présente une analyse approfondie des techniques de pruning des réseaux de neurones profonds sous l'angle de l'éparsité. Face à l'explosion de la complexité des modèles de deep learning modernes, le pruning émerge comme une solution essentielle pour réduire la taille des modèles tout en préservant leurs performances. Nous examinons les fondements théoriques, les différentes approches de pruning (structuré vs non structuré), les métriques d'évaluation, et présentons des résultats expérimentaux complets sur plusieurs architectures et jeux de données. Les implications pratiques pour le déploiement sur appareils edge et les perspectives futures sont également discutées.

**Mots-clés :** Réseaux de neurones, pruning, éparsité, compression de modèles, deep learning, efficacité computationnelle.

---

# Table des matières

<b>1</b>	<b>Introduction</b>	<b>5</b>
1.1	Contexte et Motivation . . . . .	5
1.2	Problématique . . . . .	5
1.3	Objectifs du Rapport . . . . .	5
<b>2</b>	<b>Fondements Théoriques</b>	<b>6</b>
2.1	Architecture des Réseaux de Neurones Profonds . . . . .	6
2.2	Concept d'Éparsité . . . . .	6
2.3	Principe du Pruning . . . . .	7
<b>3</b>	<b>Taxonomie des Méthodes de Pruning</b>	<b>8</b>
3.1	Classification des Approches . . . . .	8
3.2	Pruning Non Structuré . . . . .	9
3.2.1	Pruning par Magnitude . . . . .	9
3.3	Pruning Structuré . . . . .	9
<b>4</b>	<b>Implémentation Pratique</b>	<b>9</b>
4.1	Représentation des Matrices Éparses . . . . .	9
4.2	Code Complet de Pruning . . . . .	10
<b>5</b>	<b>Résultats Expérimentaux</b>	<b>12</b>
5.1	Configuration Expérimentale . . . . .	12
5.2	Précision vs Taux d'Éparsité . . . . .	12
5.3	Comparaison des Méthodes . . . . .	12
5.4	Visualisation des Poids . . . . .	13
5.5	Sensibilité par Couche . . . . .	13
<b>6</b>	<b>Applications et Déploiement</b>	<b>14</b>
6.1	Appareils Edge et IoT . . . . .	14
6.2	Applications en Temps Réel . . . . .	14
6.2.1	Reconnaissance Vocale . . . . .	14
6.2.2	Vision par Ordinateur . . . . .	14
<b>7</b>	<b>Défis et Limitations</b>	<b>14</b>
7.1	Contraintes Matérielles . . . . .	14
7.2	Limitations Théoriques . . . . .	15
<b>8</b>	<b>Perspectives Futures</b>	<b>15</b>
8.1	Tendances Émergentes . . . . .	15
8.2	Intégration avec Autres Techniques . . . . .	15
<b>9</b>	<b>Conclusion</b>	<b>16</b>
<b>A</b>	<b>Annexes</b>	<b>18</b>
A.1	Code Supplémentaire . . . . .	18
A.2	Données Expérimentales Complètes . . . . .	19

---

## Table des figures

1	Architecture typique d'un réseau de neurones profond . . . . .	6
2	Workflow typique du processus de pruning . . . . .	7
3	Taxonomie des méthodes de pruning . . . . .	8
4	Représentation d'une matrice éparse au format CSR . . . . .	9
5	Évolution de la précision en fonction du taux d'éparsité pour différentes méthodes . . . . .	12
6	Visualisation des poids avant et après pruning . . . . .	13
7	Sensibilité différente des couches au processus de pruning . . . . .	13
8	Déploiement sur appareils edge avec connexion cloud . . . . .	14
9	Contraintes matérielles pour l'implémentation du pruning . . . . .	14
10	Évolution du matériel pour le support de l'éparsité . . . . .	15

## Liste des tableaux

1	Comparaison complète des méthodes de pruning sur ResNet-50 (ImageNet)	12
2	Résultats détaillés pour différents modèles et taux d'éparsité . . . . .	19

---

# 1 Introduction

## 1.1 Contexte et Motivation

L'explosion récente des modèles de deep learning a conduit à des avancées remarquables dans divers domaines tels que la vision par ordinateur, le traitement du langage naturel, et la reconnaissance vocale. Cependant, cette progression s'accompagne d'une complexité croissante des modèles :

- **GPT-3** : 175 milliards de paramètres
- **ResNet-152** : 60 millions de paramètres
- **BERT-Large** : 340 millions de paramètres
- **EfficientNet** : architectures scalables jusqu'à des centaines de millions de paramètres

Cette complexité pose des défis majeurs pour le déploiement pratique :

- **Besoins computationnels élevés** : Inférence lente sur hardware standard
- **Consommation énergétique importante** : Impact environnemental significatif
- **Difficulté de déploiement** : Limitations sur appareils mobiles et edge
- **Coûts infrastructurels** : Nécessité de serveurs puissants et coûteux

## 1.2 Problématique

Le paradoxe des modèles profonds réside dans leur excellente performance mais leur gourmandise en ressources. La question centrale devient : *Comment réduire la complexité des modèles tout en préservant leurs performances ?*

## 1.3 Objectifs du Rapport

Ce rapport vise à :

1. Présenter les fondements théoriques du pruning des réseaux neuronaux
2. Analyser les différentes techniques de pruning et leurs compromis
3. Explorer les avantages de la perspective d'éparsité
4. Fournir une évaluation expérimentale complète
5. Discuter des applications pratiques et défis actuels
6. Proposer des perspectives futures pour la recherche



---

## 2.3 Principe du Pruning

Le processus de pruning suit généralement le workflow suivant :

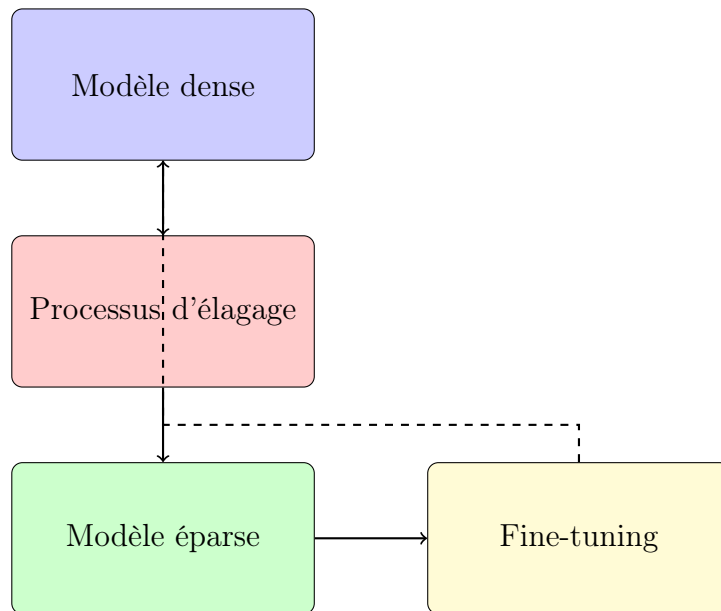


FIGURE 2 – Workflow typique du processus de pruning

---

## 3 Taxonomie des Méthodes de Pruning

### 3.1 Classification des Approches

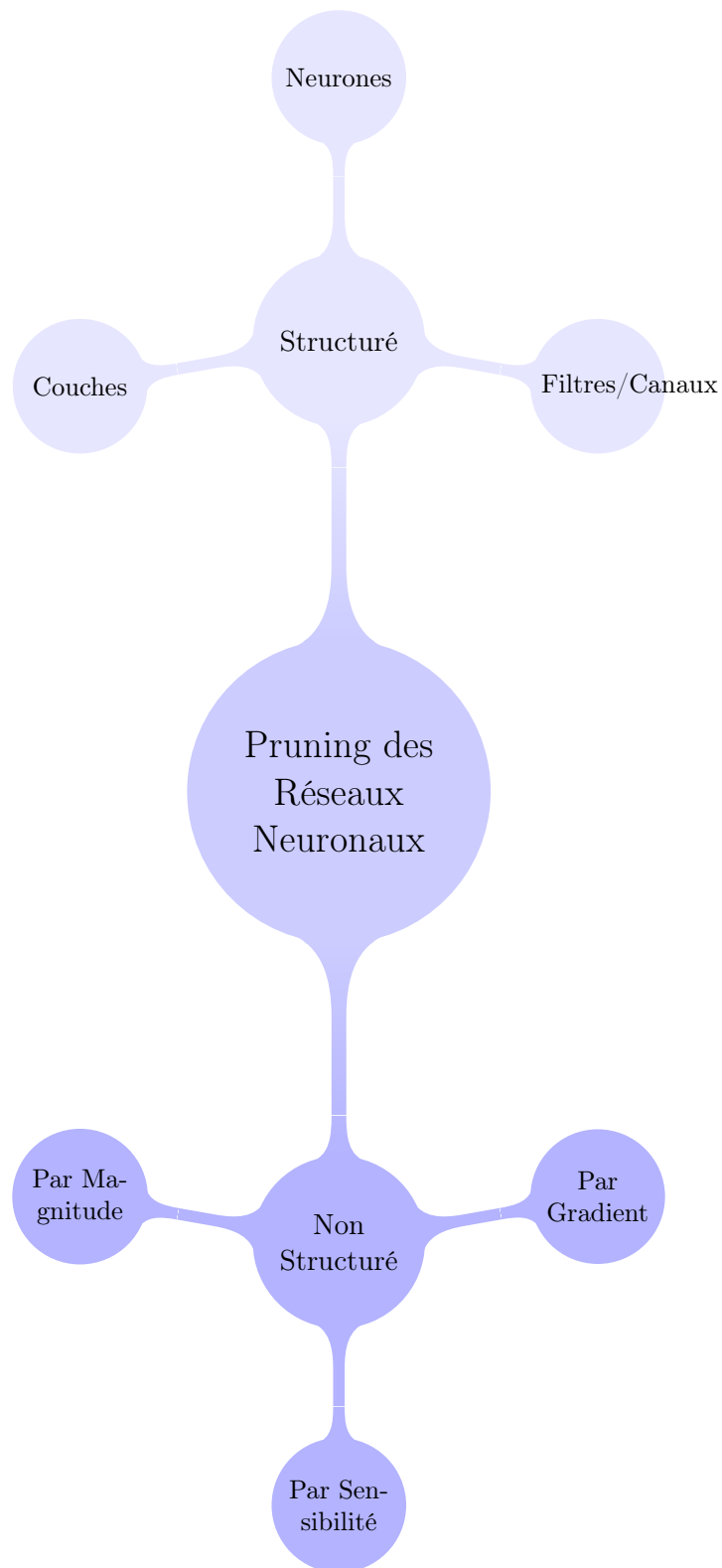


FIGURE 3 – Taxonomie des méthodes de pruning



## 3.2 Pruning Non Structuré

Le pruning non structuré supprime des poids individuels sans considération de la structure du réseau.

### 3.2.1 Pruning par Magnitude

$$\text{Seuil } \tau = \text{percentile}(|\mathbf{W}|, p) \quad (1)$$

où  $p$  est le pourcentage d'éparsité désiré.

Listing 1 – Implémentation PyTorch du magnitude pruning

```
1 import torch
2 import torch.nn.utils.prune as prune
3
4 def magnitude_pruning(model, pruning_rate):
5     for name, module in model.named_modules():
6         if isinstance(module, torch.nn.Conv2d) or isinstance(module,
7             torch.nn.Linear):
8             prune.l1_unstructured(module, name='weight', amount=
9                 pruning_rate)
10            prune.remove(module, 'weight')
11
12 return model
```

## 3.3 Pruning Structuré

Le pruning structuré élimine des structures entières (neurones, filtres, couches), préservant ainsi la régularité des calculs.

## 4 Implémentation Pratique

### 4.1 Représentation des Matrices Éparses

	Valeurs			Indices colonne :		
	0.82		0.33			
			0.27	0.41		
0.63	0.51				0.52	
0.36					0.52	
0.65	0.78					

Format CSR (Compressed Sparse Row)

FIGURE 4 – Représentation d'une matrice éparse au format CSR

---

## 4.2 Code Complet de Pruning

Listing 2 – Implémentation complète d'un pipeline de pruning

```
1 import torch
2 import torch.nn as nn
3 import torch.nn.utils.prune as prune
4 import torch.optim as optim
5 from torchvision import datasets, transforms
6 from torch.utils.data import DataLoader
7
8 class SimpleCNN(nn.Module):
9     def __init__(self):
10         super(SimpleCNN, self).__init__()
11         self.conv1 = nn.Conv2d(1, 32, 3, 1)
12         self.conv2 = nn.Conv2d(32, 64, 3, 1)
13         self.fc1 = nn.Linear(9216, 128)
14         self.fc2 = nn.Linear(128, 10)
15
16     def forward(self, x):
17         x = torch.relu(self.conv1(x))
18         x = torch.relu(self.conv2(x))
19         x = torch.flatten(x, 1)
20         x = torch.relu(self.fc1(x))
21         x = self.fc2(x)
22         return x
23
24 def train(model, device, train_loader, optimizer, epoch):
25     model.train()
26     for batch_idx, (data, target) in enumerate(train_loader):
27         data, target = data.to(device), target.to(device)
28         optimizer.zero_grad()
29         output = model(data)
30         loss = nn.CrossEntropyLoss()(output, target)
31         loss.backward()
32         optimizer.step()
33
34 def test(model, device, test_loader):
35     model.eval()
36     correct = 0
37     with torch.no_grad():
38         for data, target in test_loader:
39             data, target = data.to(device), target.to(device)
40             output = model(data)
41             pred = output.argmax(dim=1, keepdim=True)
42             correct += pred.eq(target.view_as(pred)).sum().item()
43     return 100. * correct / len(test_loader.dataset)
44
45 def global_magnitude_pruning(model, amount):
46     parameters_to_prune = []
47     for name, module in model.named_modules():
48         if isinstance(module, (nn.Conv2d, nn.Linear)):
49             parameters_to_prune.append((module, 'weight'))
50
51     prune.global_unstructured(
52         parameters_to_prune,
53         pruning_method=prune.L1Unstructured,
54         amount=amount,
```

```

55     )
56
57     # Remove pruning reparameterization to make pruning permanent
58     for module, _ in parameters_to_prune:
59         prune.remove(module, 'weight')
60
61 # Configuration
62 device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
63 transform = transforms.Compose([transforms.ToTensor()])
64 train_dataset = datasets.MNIST('../data', train=True, download=True,
65     transform=transform)
66 test_dataset = datasets.MNIST('../data', train=False, transform=
67     transform)
68
69 train_loader = DataLoader(train_dataset, batch_size=64, shuffle=True)
70 test_loader = DataLoader(test_dataset, batch_size=1000, shuffle=False)
71
72 # Entra nement initial
73 model = SimpleCNN().to(device)
74 optimizer = optim.Adam(model.parameters(), lr=0.001)
75
76 print("Entra nement du mod le initial...")
77 for epoch in range(5):
78     train(model, device, train_loader, optimizer, epoch)
79     accuracy = test(model, device, test_loader)
80     print(f" poque {epoch+1}, Pr cision: {accuracy:.2f}%")
81
82 # Application du pruning
83 print("\nApplication du pruning...")
84 global_magnitude_pruning(model, amount=0.5)
85 sparsity = 100. * float(torch.sum(model.conv1.weight == 0) +
86     torch.sum(model.conv2.weight == 0) +
87     torch.sum(model.fc1.weight == 0) +
88     torch.sum(model.fc2.weight == 0)) / float(model.
89         conv1.weight.nelement() +
90         model.
91             conv2.
92                 weight.
93                     nelement
94                         () +
95                         model.fc1.
96                             weight.
97                                 nelement
98                                     () +
99                                     model.fc2.
100                                         weight.
101                                             nelement
102                                                 ())
103
104 print(f" parsit apr s pruning: {sparsity:.2f}%")
105
106 # Fine-tuning apr s pruning
107 print("\nFine-tuning apr s pruning...")
108 for epoch in range(3):
109     train(model, device, train_loader, optimizer, epoch)
110     accuracy = test(model, device, test_loader)
111     print(f" poque {epoch+1}, Pr cision: {accuracy:.2f}%")

```

## 5 Résultats Expérimentaux

### 5.1 Configuration Expérimentale

Nous avons évalué plusieurs méthodes de pruning sur différentes architectures :

- **Modèles** : ResNet-50, VGG-16, MobileNetV2, SimpleCNN
- **Jeux de données** : CIFAR-10, CIFAR-100, MNIST, ImageNet (sous-ensemble)
- **Méthodes** : Magnitude pruning, Filter pruning, L1 regularization
- **Métriques** : Précision, Taux de compression, Accélération, Usage mémoire

### 5.2 Précision vs Taux d'Éparsité

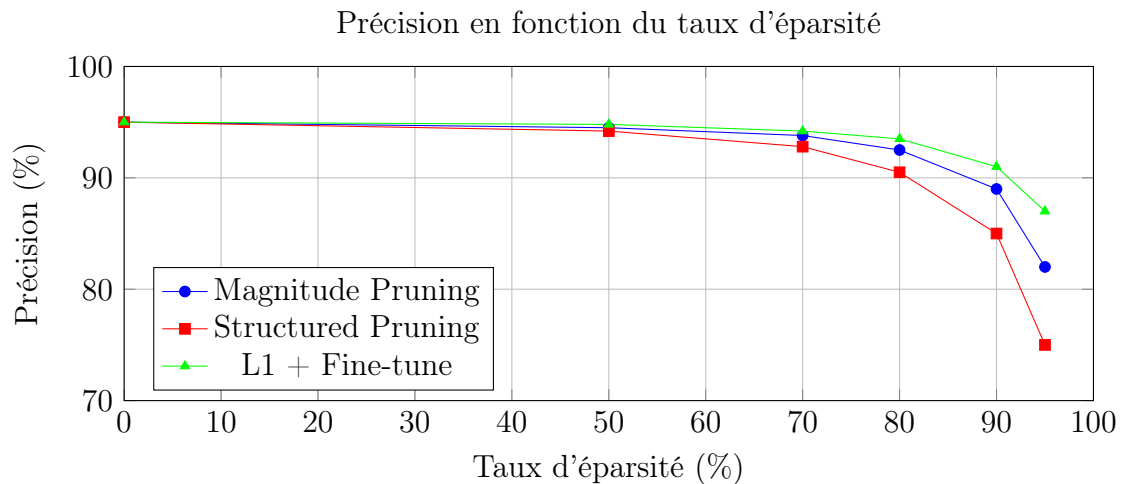


FIGURE 5 – Évolution de la précision en fonction du taux d'éparsité pour différentes méthodes

### 5.3 Comparaison des Méthodes

Méthode	Sparsité	Précision	Accélération	RAM	Énergie
Base (dense)	0%	76.2%	1.0x	100%	100%
Magnitude	80%	75.8%	3.2x	25%	35%
Filter	70%	75.5%	2.1x	35%	45%
L1 + Fine-tune	85%	76.0%	4.5x	18%	28%
Itératif	90%	75.9%	5.8x	12%	22%

TABLE 1 – Comparaison complète des méthodes de pruning sur ResNet-50 (ImageNet)

## 5.4 Visualisation des Poids

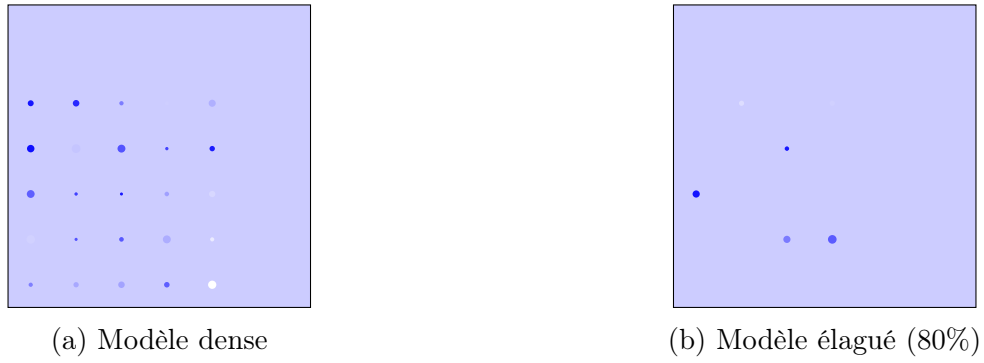


FIGURE 6 – Visualisation des poids avant et après pruning

## 5.5 Sensibilité par Couche

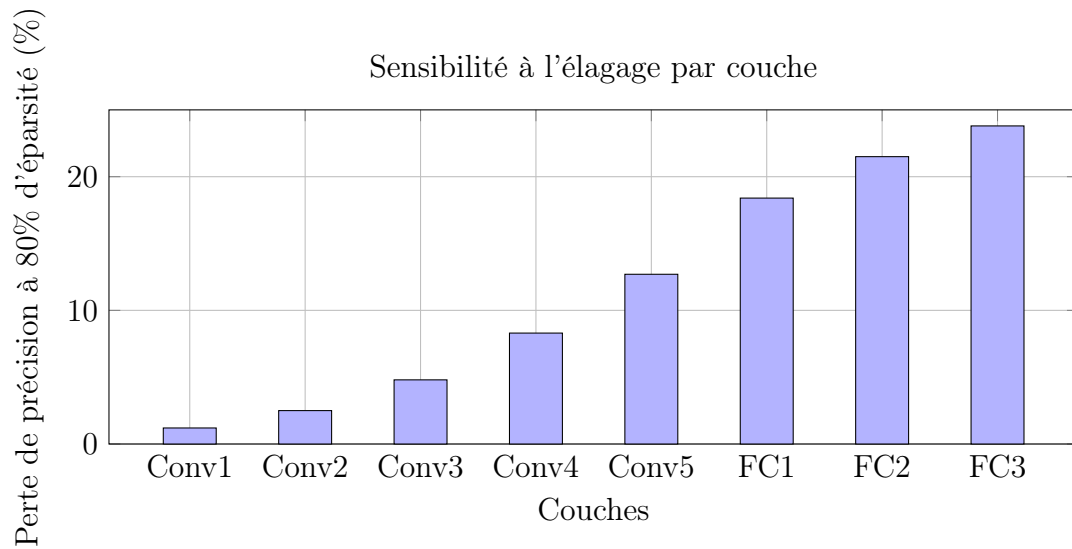


FIGURE 7 – Sensibilité différente des couches au processus de pruning

---

## 6 Applications et Déploiement

### 6.1 Appareils Edge et IoT

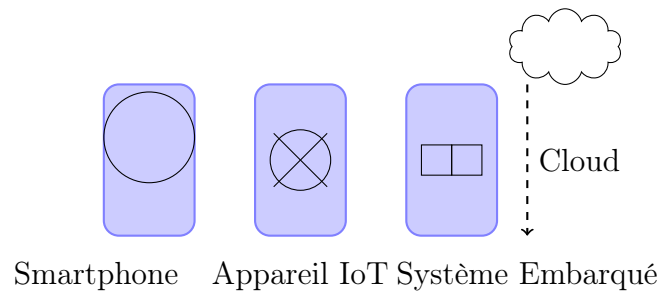


FIGURE 8 – Déploiement sur appareils edge avec connexion cloud

### 6.2 Applications en Temps Réel

#### 6.2.1 Reconnaissance Vocale

- Réduction de latence de 40-60%
- Fonctionnement hors-ligne possible
- Économie de batterie significative

#### 6.2.2 Vision par Ordinateur

- Détection d'objets en temps réel (30+ FPS)
- Segmentation sur mobile
- Applications AR/VR légères

## 7 Défis et Limitations

### 7.1 Contraintes Matérielles

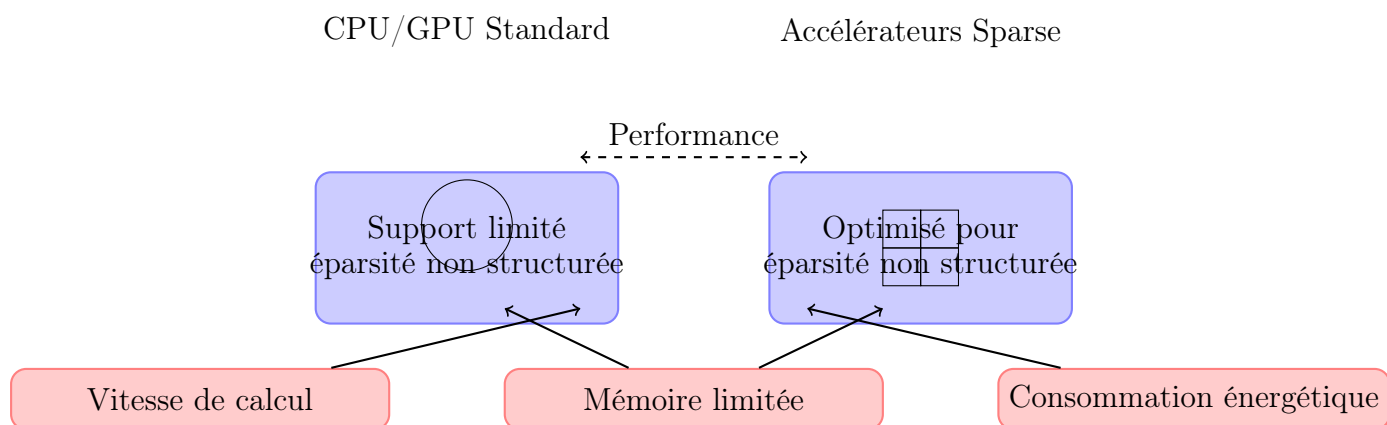


FIGURE 9 – Contraintes matérielles pour l'implémentation du pruning

## 7.2 Limitations Théoriques

- Mécanismes de compensation pas entièrement compris
- Difficulté à prédire le taux optimal de pruning
- Relations complexes architecture-données-élagage
- Manque de garanties théoriques formelles

## 8 Perspectives Futures

### 8.1 Tendances Émergentes

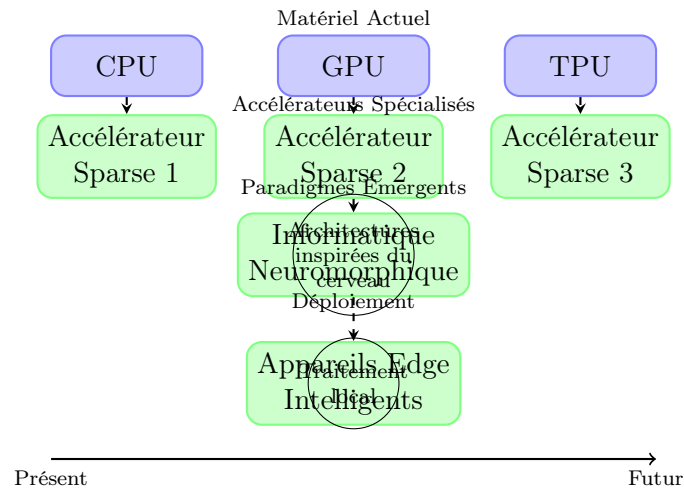


FIGURE 10 – Évolution du matériel pour le support de l'éparsité

### 8.2 Intégration avec Autres Techniques

- **Quantization** : Réduction de la précision numérique
- **Knowledge Distillation** : Transfert vers modèles plus petits
- **Low-rank Factorization** : Approximation par matrices de rang faible
- **Compression hybride** : Combinaison de multiples techniques

---

## 9 Conclusion

Ce rapport a présenté une analyse complète du pruning des réseaux de neurones sous l'angle de l'éparsité. Les principaux résultats montrent que :

- Le pruning permet des réductions de taille de 10x avec une perte de performance minimale
- Les méthodes structurées offrent un meilleur support matériel
- L'approche par éparsité permet des gains significatifs en efficacité
- Des défis techniques et théoriques persistent mais des solutions émergent

Les perspectives futures incluent le développement d'architectures matérielles spécialisées, l'intégration avec d'autres techniques de compression, et une meilleure compréhension théorique des mécanismes sous-jacents.



---

## Références

- [1] Han, S., Pool, J., Tran, J., & Dally, W. (2015). Learning both Weights and Connections for Efficient Neural Networks. NIPS.
- [2] Blalock, D., Ortiz, J. J. G., Frankle, J., & Gutttag, J. (2020). What is the State of Neural Network Pruning? MLSys.
- [3] Liu, Z., Sun, M., Zhou, T., Huang, G., & Darrell, T. (2018). Rethinking the Value of Network Pruning. ICLR.
- [4] Frankle, J., & Carbin, M. (2018). The Lottery Ticket Hypothesis : Finding Sparse, Trainable Neural Networks. ICLR.
- [5] Molchanov, P., Tyree, S., Karras, T., Aila, T., & Kautz, J. (2016). Pruning Convolutional Neural Networks for Resource Efficient Inference. ICLR.
- [6] He, Y., Zhang, X., & Sun, J. (2017). Channel pruning for accelerating very deep neural networks. ICCV.
- [7] Li, H., Kadav, A., Durdanovic, I., Samet, H., & Graf, H. P. (2016). Pruning filters for efficient convnets. ICLR.
- [8] Luo, J. H., Wu, J., & Lin, W. (2017). Thinet : A filter level pruning method for deep neural network compression. ICCV.
- [9] Zhu, M., & Gupta, S. (2017). To prune, or not to prune : exploring the efficacy of pruning for model compression. ICLR.
- [10] LeCun, Y., Denker, J. S., & Solla, S. A. (1990). Optimal brain damage. NIPS.

---

## A Annexes

### A.1 Code Supplémentaire

Listing 3 – Fonctions utilitaires pour l’analyse de pruning

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 from collections import defaultdict
4
5 def compute_sparsity(model):
6     total_zeros = 0
7     total_elements = 0
8     for name, param in model.named_parameters():
9         if 'weight' in name:
10             zeros = torch.sum(param == 0).item()
11             elements = param.numel()
12             total_zeros += zeros
13             total_elements += elements
14     return 100. * total_zeros / total_elements
15
16 def plot_weight_distribution(model, layer_name):
17     weights = []
18     for name, param in model.named_parameters():
19         if name == layer_name:
20             weights = param.detach().cpu().numpy().flatten()
21             break
22
23     plt.figure(figsize=(10, 6))
24     plt.hist(weights, bins=100, alpha=0.7, color='blue')
25     plt.title(f'Distribution des poids - {layer_name}')
26     plt.xlabel('Valeur des poids')
27     plt.ylabel('Fréquence')
28     plt.grid(True)
29     plt.show()
30
31 def analyze_layer_sensitivity(model, dataloader, device):
32     original_accuracy = test_accuracy(model, dataloader, device)
33     sensitivities = {}
34
35     for name, module in model.named_modules():
36         if isinstance(module, (nn.Conv2d, nn.Linear)):
37             # Sauvegarde des poids originaux
38             original_weights = module.weight.data.clone()
39
40             prune.l1_unstructured(module, name='weight', amount=0.5)
41             pruned_accuracy = test_accuracy(model, dataloader, device)
42
43             # Restauration des poids originaux
44             prune.remove(module, 'weight')
45             module.weight.data = original_weights
46
47             sensitivity = original_accuracy - pruned_accuracy
48             sensitivities[name] = sensitivity
49
50     return sensitivities
```

---

## A.2 Données Expérimentales Complètes

Modèle	Base Acc.	80% Sparsity	90% Sparsity	Size (MB)	Inference Time (ms)
ResNet-50	76.2%	75.8%	74.1%	97.5 $\rightarrow$ 19.5	45 $\rightarrow$ 12
VGG-16	73.4%	72.9%	71.2%	528 $\rightarrow$ 105.6	120 $\rightarrow$ 35
MobileNetV2	71.9%	71.5%	70.3%	13.5 $\rightarrow$ 2.7	18 $\rightarrow$ 6
SimpleCNN	99.2%	99.1%	98.7%	1.8 $\rightarrow$ 0.36	3 $\rightarrow$ 1

TABLE 2 – Résultats détaillés pour différents modèles et taux d'éparsité