



SORBONNE UNIVERSITÉ

RAPPORT

---

# PLDAC : Enhancing SPLADE with Large Vocabularies in Neural Information Retrieval

---

*Étudiants :*

Yuan LIU

Herve NGUYEN

*Encadré par :*

Benjamin PIWOWARSKI

14 mai 2024

## Résumé

SPLADE est un modèle neuronal pour la recherche d'information, qui se sert d'une représentation sparse pour avoir plusieurs avantages par rapport aux modèles denses tel que :

- La possibilité d'utiliser un dictionnaire inversé de manière efficace.
- Avoir une correspondance lexicale explicite entre la requête et les résultats.
- Une meilleure interprétabilité.

Notre objectif ici est d'optimiser la partie backpropagation en tirant partie de la représentation sparse des données d'entraînement.

Le projet consiste d'abord à déterminer des expressions du backward de l'apprentissage de SPLADE et utiliser ces expressions pour optimiser la backpropagation avec des exécutions en parallèle sur GPU.

**Repository GitHub :** [https://github.com/Mosakana/PLDAC\\_Splade\\_GPU\\_Calculate](https://github.com/Mosakana/PLDAC_Splade_GPU_Calculate)

# Table des matières

<b>1 Définitions</b>	<b>3</b>
1.1 Tenseurs et dimensions . . . . .	3
1.2 Forward . . . . .	3
<b>2 Expressions des gradients</b>	<b>3</b>
2.1 Gradient de $W$ . . . . .	3
2.2 Gradient de $X$ . . . . .	4
2.3 Gradient de $b$ . . . . .	5
<b>3 Concept du calcul GPU</b>	<b>6</b>
3.1 L'optimisation sur le calcul du gradient de $W$ . . . . .	6
3.2 L'optimisation sur le calcul du gradient de $X$ . . . . .	6
3.3 L'optimisation sur le calcul du gradient de $b$ . . . . .	6
<b>4 Implémentation au niveau du code</b>	<b>6</b>
4.1 Vérification des résultats . . . . .	6
4.2 Récupération des indices à élément non nuls . . . . .	7
4.3 Utilisation de la librairie triton . . . . .	8
4.4 Évaluation des performances . . . . .	10
<b>5 Perspectives</b>	<b>11</b>

# 1 Définitions

## 1.1 Tenseurs et dimensions

Pour aborder le sujet, nous devons d'abord définir les tenseurs de base que nous allons manipuler.

- Le tenseur  $X$  qui contient les données d'entraînement,  $X \in \mathbb{R}^{B \times L \times D}$ . Avec  $B$  la taille du batch,  $L$  la longueur maximum d'une séquence et  $D$  la dimension d'un emdedding d'un mot.
- Le tenseur  $W$  de poids,  $W \in \mathbb{R}^{D \times V}$ . Avec  $D$  la dimension d'un embedding et  $V$  la taille du vocabulaire.
- Le vecteur  $b$  du biais,  $b \in \mathbb{R}^V$ . Avec  $V$  la taille du vocabulaire.
- Le masque  $M$ ,  $M \in \mathbb{R}^{B \times L}$  dont les éléments prennent soit la valeur  $-\infty$  ou 0. Avec  $B$  la taille du batch,  $V$  la taille du vocabulaire.
- $XW + b + M$  est donc de dimension  $B \times L \times D$ .

## 1.2 Forward

Le forward de SPLADE prends l'expression suivante :

$$f(X, W, b) = \text{ReLU}(\max(XW + b + M)) \quad (1)$$

Avec  $\max$  ici le max sur la dimension  $L$ , sachant que  $XW + b + M$  est de dimension  $B \times L \times V$ . Alors  $\max(XW + b + M)$  et sa  $\text{ReLU}$  est donc de dimension  $B \times V$ .

# 2 Expressions des gradients

Afin de pouvoir essayer d'améliorer le backward de SPLADE avec du code parallélisé, nous avons besoin de l'expression des gradients de  $W$ ,  $X$ ,  $b$  par rapport à la loss.

Soit  $\nabla \in B \times V$  le gradient de la loss par rapport à  $f(X, W, b)$ ,  $\nabla_{st}$  est un élément de  $\nabla$  à la position respective.

## 2.1 Gradient de $W$

On a d'abord, avec la règle de la chaîne :

$$\frac{\partial L}{\partial W_{dv}} = \sum_{b \in \{1, \dots, B\}, v' \in \{1, \dots, V\}} \frac{\partial L}{\partial f_{b'v'}} \frac{\partial f_{b'v'}}{\partial W_{dv}} \quad (2)$$

$$= \sum_{b' \in \{1, \dots, B\}, v' \in \{1, \dots, V\}} \nabla_{b'v'} \frac{\partial f_{b'v'}}{\partial W_{dv}} \quad (3)$$

Le gradient de  $W_{dv}$  par rapport à  $f_{st}$  prend la forme :

$$\frac{\partial f_{b'v'}}{\partial W_{dv}} = \frac{\partial \text{ReLU}(\max(XW + b + M))_{b'v'}}{\partial W_{dv}} \quad (4)$$

$$= \frac{\partial \text{ReLU}(\max(X_{b'\bullet} W_{\bullet v'} + b_{v'} + M_{b'}))}{\partial W_{dv}} \quad (5)$$

Soit  $l^* \in \mathbb{N}^{B \times V}$  la matrice des indices sur la dimension  $L$  de  $XW + b + M$ , où  $l_{b,v}^*$  est donc le  $l$  sur  $L$  où on obtient la valeur du max dans  $\max(XW + b + M)$  pour un couple  $(b, v)$  donné.

On remarque donc que  $\frac{\partial f_{b'v'}}{\partial W_{dv}}$  peut prendre différentes valeurs possible selon le cas.

$$\frac{\partial f_{b'v'}}{\partial W_{dv}} = \begin{cases} 0 & \text{si } v' \neq v \\ 0 & \text{si } l \neq l_{b',v'}^*, \forall l \in [0, L[ \\ \text{sinon} & (6) \end{cases} \quad (6)$$

Si nous ne sommes pas les deux premiers cas, nous avons donc (le masque ici étant à 0) :

$$\frac{\partial f_{b'v'}}{\partial W_{dv}} = \frac{\partial (X_{b'l_{b',v'}^*} \bullet W_{\bullet v} + b_w)}{\partial W_{dv}} \quad (7)$$

$$= X_{b'l_{b',v'}^*} d \quad (8)$$

Par conséquent selon (3), (6), (8) :

$$\frac{\partial L}{\partial W_{dv}} = \sum_{b' \in \{1, \dots, B\}} \nabla_{b'v} X_{b'l_{b',v'}^*} d \quad (9)$$

Et pour une colonne de  $W$ , l'expression est :

$$\frac{\partial L}{\partial W_{\bullet v}} = \sum_{b' \in \{1, \dots, B\}} \nabla_{b'v} X_{b'l_{b',v'}^*} \bullet \quad (10)$$

## 2.2 Gradient de X

Comme avec  $W$ , avec la règle de la chaîne :

$$\frac{\partial L}{\partial X_{bld}} = \sum_{b' \in \{1, \dots, B\}, v' \in \{1, \dots, V\}} \frac{\partial L}{\partial f_{b'v'}} \frac{\partial f_{b'v'}}{\partial X_{bld}} \quad (11)$$

$$= \sum_{b' \in \{1, \dots, B\}, v' \in \{1, \dots, V\}} \nabla_{b'v'} \frac{\partial f_{b'v'}}{\partial X_{bld}} \quad (12)$$

Le gradient de  $X_{bld}$  par rapport à  $f_{st}$  prend la forme :

$$\frac{\partial f_{b'v'}}{\partial X_{bld}} = \frac{\partial \text{ReLU}(\max(XW + b + M))_{b'v'}}{\partial X_{bld}} \quad (13)$$

$$= \frac{\partial \text{ReLU}(\max(X_{b'\bullet\bullet} W_{\bullet v'} + b_{v'} + M_{b'}))}{\partial X_{bld}} \quad (14)$$

Donc :

$$\frac{\partial f_{b'v'}}{\partial X_{bld}} = \begin{cases} 0 & \text{si } b' \neq b \\ 0 & \text{si } l \neq l_{b',v'}^* \\ \text{sinon} & (15) \end{cases} \quad (15)$$

Si  $b = b$  et  $l = l_{bv'}^*$ , nous avons donc :

$$\frac{\partial f_{b'v'}}{\partial X_{bld}} = \frac{X_{bl^*} \cdot W_{\bullet v'} + b_{v'}}{X_{bld}} \quad (16)$$

$$= W_{dv'} \quad (17)$$

Par conséquent, selon (12), (16), (17) :

$$\frac{\partial L}{\partial X_{bld}} = \sum_{v' \in \{1, \dots, V\} \setminus l=l_{bv'}^*} \nabla_{bv'} W_{dv'} \quad (18)$$

Pour notre projet on utilisera donc :

$$\frac{\partial L}{\partial X_{bl^*}} = \sum_{v' \in \{1, \dots, V\} \setminus l=l_{bv'}^*} \nabla_{bv'} W_{\bullet v'} \quad (19)$$

### 2.3 Gradient de $b$

Avec la règle de la chaîne :

$$\frac{\partial L}{\partial b_v} = \sum_{b' \in \{1, \dots, B\}, v' \in \{1, \dots, V\}} \frac{\partial L}{\partial f_{b'v'}} \frac{\partial f_{b'v'}}{\partial b_v} \quad (20)$$

$$= \sum_{b' \in \{1, \dots, B\}, v' \in \{1, \dots, V\}} \nabla_{b'v'} \frac{\partial f_{b'v'}}{\partial b_v} \quad (21)$$

Le gradient de  $b_v$  par rapport à  $f_s t$  s'exprime de la manière suivante :

$$\frac{\partial f_{b'v'}}{\partial b_v} = \frac{\partial \text{ReLU}(\max(XW + b + M))_{b'v'}}{\partial b_v} \quad (22)$$

$$= \frac{\partial \text{ReLU}(\max(X_{b'..} W_{\bullet v'} + b_{v'} + M_{b'}))}{\partial b_v} \quad (23)$$

$$= \begin{cases} 0 & \text{si } v \neq v' \\ 0 & \text{si } M_{b'} = -\infty \\ 1 & \text{sinon} \end{cases} \quad (24)$$

D'après (21) et (24) :

$$\frac{\partial L}{\partial b_v} = \sum_{b' \in \{1, \dots, B\} \setminus M_{b'}=0} \nabla_{b'v} \quad (25)$$

$$\frac{\partial L}{\partial b_{\bullet}} = \sum_{b' \in \{1, \dots, B\} \setminus M_{b'}=0} \nabla_{b'\bullet} \quad (26)$$

### 3 Concept du calcul GPU

#### 3.1 L'optimisation sur le calcul du gradient de $W$

Selon la formule (10) du gradient de  $W$  :

$$\frac{\partial L}{\partial W_{\bullet v}} = \sum_{b'} \nabla_{b'v} X_{b'l_{bv}^*} \bullet$$

Il est possible de calculer cette formule en parallèle par rapport à la dimension  $D$ . Chaque processus s'occupe un  $d$  et parcourt la liste qui contient des triplets sous la forme  $(b, v, l_{bv}^*)$  en accumulant la valeur de  $\nabla_{bv} X_{bl_{bv}^* d}$  à la position  $(d, v)$  de la matrice du gradient de  $W$ .

Il existe une autre approche. La tâche de calcul peut être divisé par la dimension  $V$ . Cependant, pour cette approche, la liste de triplet doit être transformée à une liste 2D qui regroupe les triplets de la même valeur de  $v$  dans une sous-liste. Chaque processus itère une sous-liste de triplets de  $v$  correspondant. À chaque itération, les valeurs de  $\nabla_{bv} X_{bl_{bv}^*} \bullet$  par rapport au triplet pris  $(b, v, l_{bv}^*)$  sera sauvegardé à la position de  $(\bullet, v)$  de la matrice du gradient de  $W$ .

#### 3.2 L'optimisation sur le calcul du gradient de $X$

La formule du gradient de  $X$  (cf. (19)) :

$$\frac{\partial L}{\partial X_{bl\bullet}} = \sum_{v' \setminus l=l_{bv}^*} \nabla_{bv'} W_{\bullet v'}$$

Ce calcul peut aussi être divisé par la dimension  $D$ . Les processus parallèles chacun travaille sur un  $d$  différent, accumule le produit de  $\nabla_{bv} W_{dv}$  à l'adresse  $(b, l_{bv}^*, d)$  de la matrice du gradient de  $X$  par rapport au triplet  $(b, v, l_{bv}^*)$  récupéré.

Cependant, cette tâche ne peut pas être coupée par la dimension  $V$ . Si deux processus avec  $v$  différent, mais ils ont les même valeur de  $b$  et  $l^*$  par une coïncidence. Lorsqu'ils exécutent les calculs, il est possible de se produit un conflit de I/O, celui-ci provoque possiblement des erreurs en resultat.

Pour le gradient de  $X$ , il est possible de découper le calcul par la dimension de  $B$ . Pourtant, généralement, les données d'entraînement ont une taille de  $B$  qui est plus petite que celle de  $V$  ou  $D$ . Donc il n'est pas très efficace de diviser la tâche par rapport à la dimension  $B$ .

#### 3.3 L'optimisation sur le calcul du gradient de $b$

Ici, il y a finalement peu d'intérêt d'implémenter le calcul sous Triton car `scatter_add_` de Pytorch déjà optimisé existe pour cette opération.

## 4 Implémentation au niveau du code

### 4.1 Vérification des résultats

Premièrement, pour vérifier nos résultats. Nous avons écrit un module test nommé `check_grad_splade` qui est exécutable pour tester nos expressions par rapport à ce que

torch obtient. (avec torch.gradcheck)

## 4.2 Récupération des indices à élément non nuls

Dû à la nature sparse de  $X$ , de la ReLU et du masque, il est pertinent de vouloir seulement manipuler les éléments non-nuls parmi de nombreux éléments nuls.

Ainsi, au niveau du forward on récupère les indices du résultat à la sortie du forward des éléments non nuls.

```
1     relu = ReLU()
2     maximum, max_indice = torch.max(output, 1)
3     result = relu(maximum)
4
5     indice_not_zero = torch.nonzero(result, as_tuple=True)
6
7     effective_indice = []
8     for b, v in zip(*indice_not_zero):
9         effective_indice.append((b, v, max_indice[b, v]))
```

Listing 1 – Récupération des indices effectives



### 4.3 Utilisation de la librairie triton

Pour implémenter les équations de calcul directement avec le GPU, nous avons utilisé la bibliothèque Triton de OpenAI qui permet d'écrire et compiler des kernel GPU en python sans passer par le code CUDA.

Le flag "@triton.jit" permet de déclarer à l'interpréteur python de considérer la fonction en dessous en tant que fonction kernel à exécuter dans le GPU pour un processus donné.

À l'intérieur de la fonction on se retrouve donc à manipuler des pointeurs dans une plage de données, dans un point de vue très bas niveau.

Ci dessous le processus correspondant va itérer sur une rangée pour calculer une partie locale du gradient.

```

1  @triton.jit
2  def gradient_w_kernel(index_ptr, grad_ptr, delta_ptr, x_ptr,
3                        D, N_INDEX,
4                        stride_index_dim0, stride_index_dim1,
5                        stride_grad_d, stride_grad_v,
6                        stride_delta_b, stride_delta_v,
7                        stride_x_b, stride_x_l, stride_x_d,
8                        BLOCK_SIZE_D: tl.constexpr):
9
10     pid = tl.program_id(axis=0)
11     start_point = pid * BLOCK_SIZE_D
12     offsets = start_point + tl.arange(0, BLOCK_SIZE_D)
13     mask = offsets < D
14
15     for i in range(N_INDEX):
16         b = tl.load(index_ptr + (i * stride_index_dim0) + 0 *
17                     stride_index_dim1)
18         v = tl.load(index_ptr + (i * stride_index_dim0) + 1 *
19                     stride_index_dim1)
20         lmax = tl.load(index_ptr + (i * stride_index_dim0) + 2 *
21                       stride_index_dim1)
22
23         delta = tl.load(delta_ptr + (b * stride_delta_b + v *
24                                     stride_delta_v))
25         x = tl.load(x_ptr + (b * stride_x_b + lmax * stride_x_l +
26                             offsets * stride_x_d), mask=mask)
27         grad = tl.load(grad_ptr + (offsets * stride_grad_d + v *
28                                    stride_grad_v), mask=mask)
29
30         grad += delta * x
31
32         tl.store(grad_ptr + (offsets * stride_grad_d + v *
33                             stride_grad_v), grad, mask=mask)

```

Listing 2 – Fonction à compiler en tant que kernel pour calculer le gradient de W

Pour assigner utiliser des processus du GPU et calculer une partie du gradient on fait donc appel à une fonction wrapper qui appeler la fonction kernel correspondante. Dans l'exemple ci-dessous on lance donc les kernels GPU sur l'ensemble des processus dans la grille de thread (grid).

```

1  def compute_gradient_w(index, delta, x, grad_weight):
2      N_INDEX = index.shape[0]
3      D = x.shape[2]
4
5      BLOCK_SIZE_D = triton.next_power_of_2(D)
6      num_warps = 4
7      if BLOCK_SIZE_D >= 2048:
8          num_warps = 8
9      if BLOCK_SIZE_D >= 4096:
10         num_warps = 16
11
12     grid = lambda meta: (triton.cdiv(D, meta['BLOCK_SIZE_D']), )
13     gradient_w_kernel[grid](index, grad_weight, delta, x,
14                             D, N_INDEX,
15                             index.stride(0), index.stride(1),
16                             grad_weight.stride(0), grad_weight.
17                                 stride(1),
18                             delta.stride(0), delta.stride(1),
19                             x.stride(0), x.stride(1), x.stride(2),
20                             num_warps=num_warps,
21                             BLOCK_SIZE_D=BLOCK_SIZE_D)
22
23     return grad_weight.clone().detach().requires_grad_(True)

```

Listing 3 – Fonction wrapper qui permet d'appeler les fonction kernels

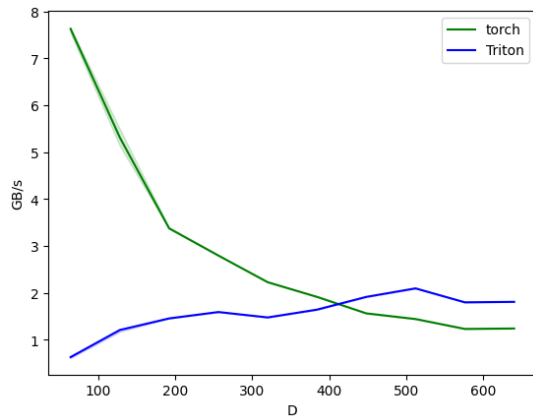
## 4.4 Évaluation des performances

Pour les performances nous avons tenté de comparer la version avec la bibliothèque torch et notre améliorée version sur Triton.

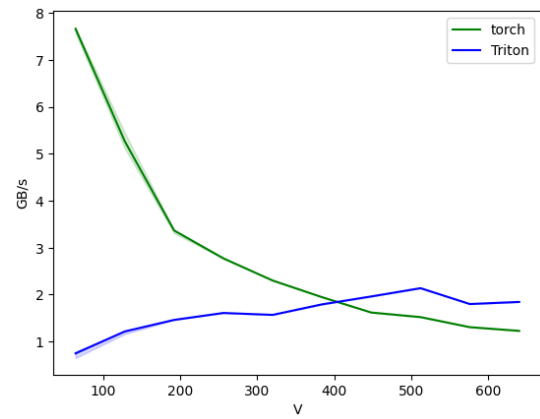
Pour l'optimisation de GPU calcul, 2 versions est fait pour une comparaison d'efficace parmi les choix de la dimension coupée. La première version est de la coupe en dimension  $D$  sur  $W$  et la coupe en dimension  $D$  sur  $X$ . La deuxième version est de la coupe en dimension  $V$  sur  $W$  et la coupe en dimension  $D$  sur  $X$ .

Les tenseurs sont initialisés de manière aléatoire avec `torch.randn`.

La métrique mesurée ici est le débit de calcul en Gigaoctet par seconde.



(a) coupe en  $D$  sur  $W$



(b) coupe en  $V$  sur  $W$

FIGURE 1 – La comparaison de performance entre différente coupe en dimension

Nous avons alors mesuré de débit de traitement / de calcul entre les deux versions. Les performances de Triton sont inférieures jusqu'à 400 éléments dans un batch, mais surpassent celles de Torch au-delà.

Parmi les deux versions, les tendances des variations de performance sont similaires. Mais, la performance de la coupe en  $V$  sur  $W$  est légèrement plus haute que celle de la coupe en  $D$  sur  $W$ .

En raison de la limite de l'environnement d'exécution (RTX 4070 12Go), nous ne pouvons pas aller plus loin que la taille de 640 de chaque dimension. Il existe aussi une possibilité que la différence minuscule de performance entre ces deux versions élargisse en un environnement d'une capacité de calcul plus élevée. (eg. Nvidia A100)

## 5 Perspectives

Jusqu'à maintenant, nous n'avons que tenté d'optimiser le backward de SPLADE. Cependant, il y a aussi de l'intérêt à prendre avantage de la nature parsimonieuses des données pour également optimiser le forward.

D'ailleurs, les tests effectués ici sont des tests isolés à petite échelle. Il aurait été aussi intéressant de voir si nos gains de performance se répliquent lorsque l'on entraîne SPLADE sur le vrai dataset MSMARCO, ce qui n'est pas possible avec les équipements que nous avons à disposition.