

# Optimization des Hyperparamètres appliquée au Fine Tuning de LLM

Basé sur l'article : *Bayesian and Partition-Based Optimization for Hyperparameter Optimization of LLM Fine-Tuning*

Nathan Davouse

# Sommaire

## 1. Introduction

## 2. Design et Implémentation

## 3. Résultats et Analysis

## 4. Conclusion

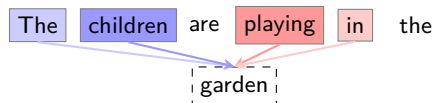
## Large Language Models

## Point clés

- ▶ Etat de l'art pour le traitement de langage naturel.
- ▶ Réseaux de Neurones avec une architecture basé sur le transformer<sup>a</sup> (annexe 1)
- ▶ Taille : entre 1 et 405 Milliards de neurones

<sup>a</sup>Vaswani et al, Attention is all you need, 2017

## Auto-attention



**Figure:** Illustration du mécanisme d'auto-attention

L'auto-attention est la clé du LLM, en permettant de comprendre le contexte



## Optimisation des Hyperparamètres (OHP)

## Hyperparamètres

Paramètres qui ne sont pas entraînés par le modèle  
(learning rate, dropout ...)

## Objectifs

- ▶ Meilleure performance qu'en manuel
- ▶ Retirer le besoin d'expertise

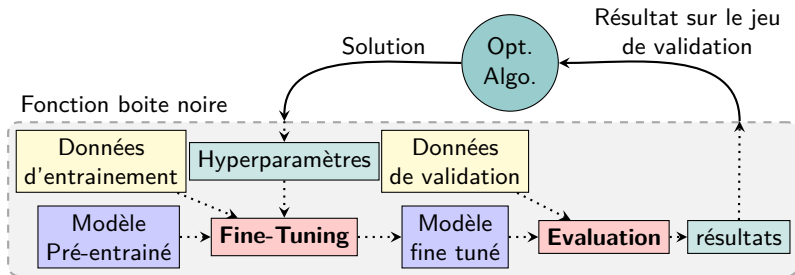


Figure: Fonctionnement général de l'optimisation des hyperparamètres



## Travaux connexes

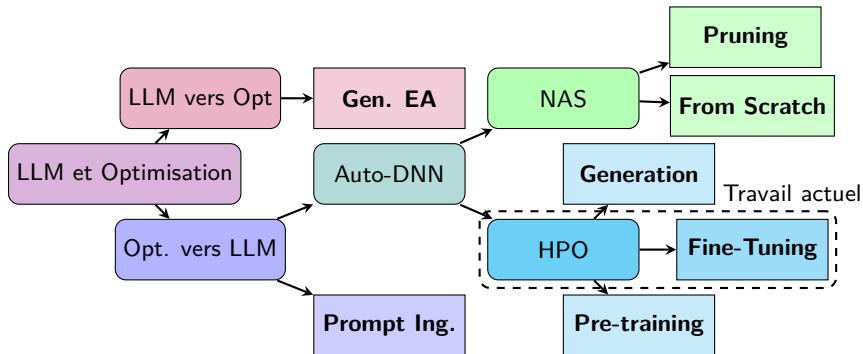


Figure: Classification des travaux similaires

# Sommaire

1. Introduction

**2. Design et Implémentation**

3. Résultats et Analysis

4. Conclusion



## Espace de Recherche

Hyperparamètres	Plage d'Optimisation		Type	Conversion
	Borne Inf.	Borne Sup.		
Learning Rate	−10	−1	log.	$f(x) = 10^x$
LoRA rank	1	64	ent.	$f(x) = \text{round}(x)$
LoRA scale	1	64	ent.	$f(x) = \text{round}(x)$
Dropout	0	0.5	cont.	$f(x) = x$
Weight Decay	−3	−1	log.	$f(x) = 10^x$

### Table: Résumé de l'espace de recherche

- Variables mixtes : étape de conversion nécessaire

## Strategie de Recherche : Optimisation Bayésienne par Process Gaussien

## Algorithm

- ▶ Echantillon de  $n$  Points (LHS)
- ▶ Evaluer ces  $n$  points
- ▶ Jusqu'à fin du budget
  - Entraîner le Process Gaussien (GP)
  - Optimiser ce GP pour obtenir un nouveau Point
  - Evaluer ce nouveaux point

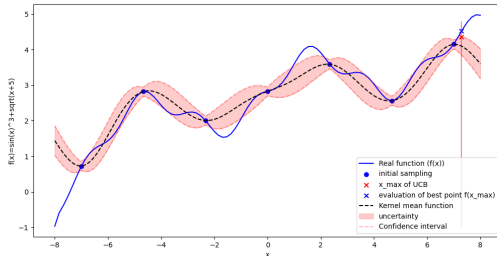


Figure: Exemple d'un surrogate sur une fonction en 1D





# Stratégie d'Evaluation de Solutions

## Implémentation

### ► Fine Tuning

- modèle : LLaMa-3.2-1B
- 

### ► Evaluation

- librairie lm\_eval
- Evaluation par la précision sur des jeu de données Benchmark : Hellaswag et MMLU

## Implémentation

- Programmation Orienté Object en Python
- Travail de documentation : *readme*, indication de type...
- Objectif : permettre le réusage
- Utilisable en ligne de commande pour Grid5000

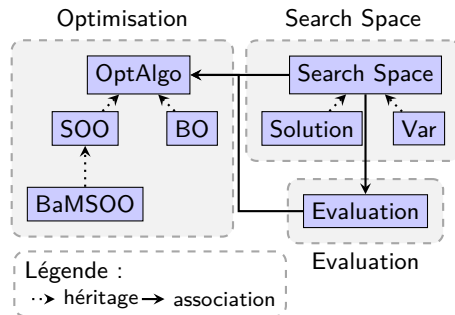


Figure: Diagramme de l'implémentation

# Sommaire

1. Introduction

2. Design et Implémentation

3. Résultats et Analysis

4. Conclusion

## Expérimentation



## LHS : Résultats

## Résultats des 3 algorithms

## Analyse

## Prospectives

# Sommaire

1. Introduction

2. Design et Implémentation

3. Résultats et Analysis

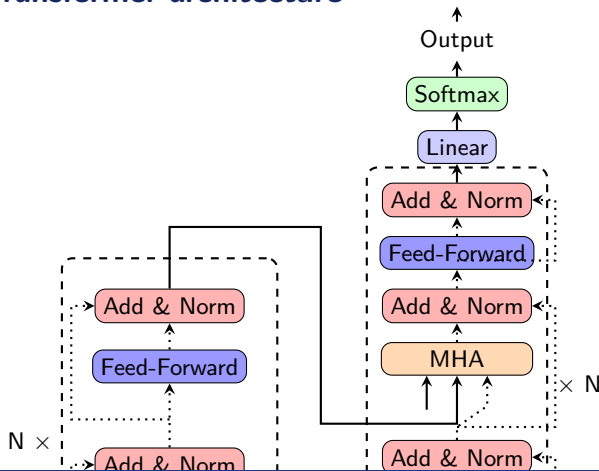
4. Conclusion

## Conclusion

Une conclusion

*Merci.*

## Annexes 1 : Transformer architecture





## Annexes 2 : Multi-Head Attention

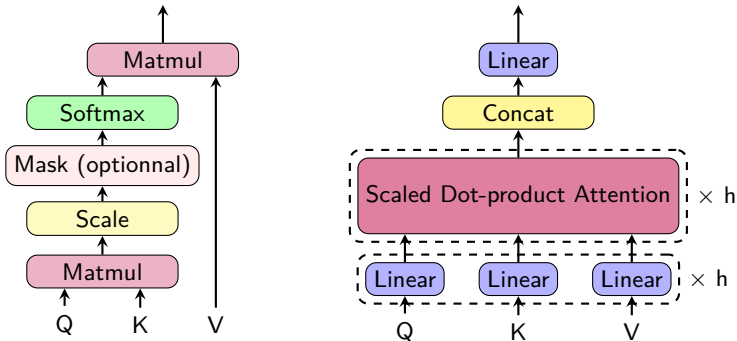


Figure: Illustration du mécanisme d'auto-attention : A droite le mécanisme complet, a gauche le *Scaled Dot-product Attention*

## Annexes 3 : Low Rank Adaptation (LoRA)

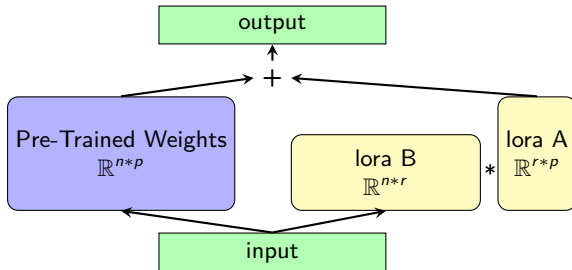


Figure: Illustration de l'application du Low Rank Adaptation (LoRA)