

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

Large Language Models

Point clés

- ▶ Etat de l'art pour le traitement de langage naturel.
- ▶ Réseaux de Neurones avec une architecture basé sur le transformer¹ (annexe 1)
- ▶ Taille : entre 1 et 405 Milliards de neurones

¹Vaswani et al, Attention is all you need,2017

Auto-attention

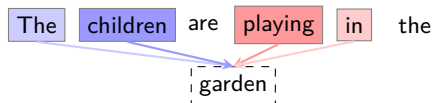


Figure: Illustration du mécanisme d'auto-attention

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

Travaux connexes

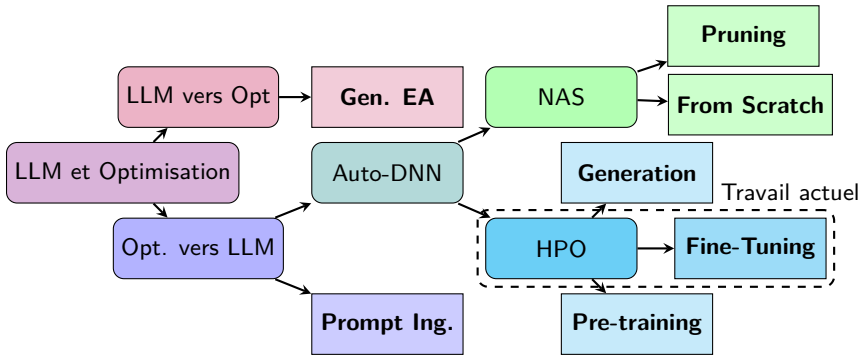


Figure: Classification des travaux similaires

Sommaire

2. Design et Implémentation

Strategie de Recherche : Optimisation Bayésienne par Process Gaussien

Principe

Utiliser un substitut moins cher à optimiser pour explorer l'espace de recherche

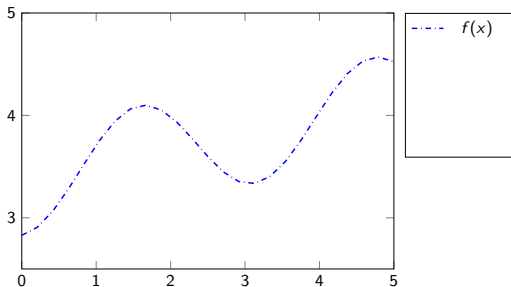


Figure: Exemple d'un surrogate sur une fonction en 1D

Strategie de Recherche : Optimisation Bayésienne par Process Gaussien

Principe

Utiliser un substitut moins cher à optimiser pour explorer l'espace de recherche

Algorithme

- Echantillon de n Points (LHS)

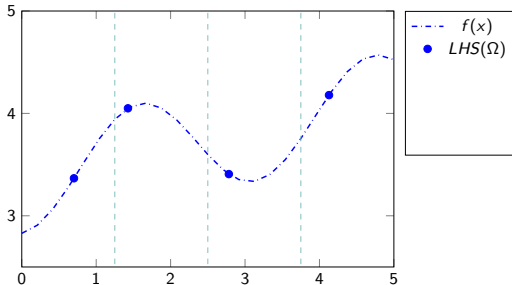
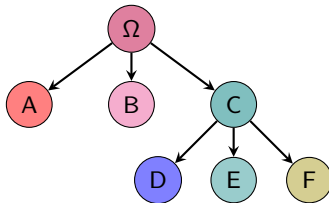


Figure: Exemple d'un surrogate sur une fonction en 1D



Search Strategy : BaMSOO (TO DO)

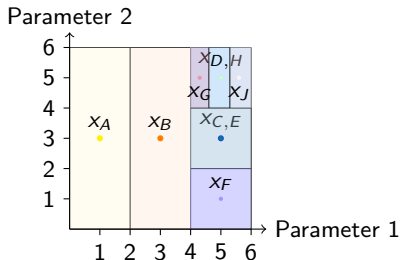


Figure: Partition de l'espace de recherche par SOO

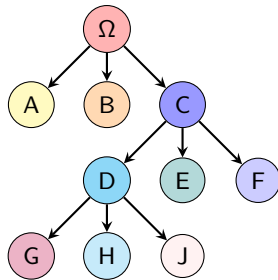


Figure: Arbre correspondant à S00

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
- ▶ Intégralement open-source²

²https://github.com/Kiwy3/B0_PBO_HPO_LLM

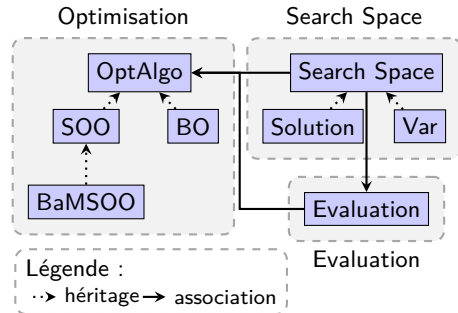


Figure: Diagramme de l'implémentation

Sommaire

1. Introduction

2. Design et Implémentation

3. Résultats et Analysis

4. Conclusion

Résultats des 3 algorithms

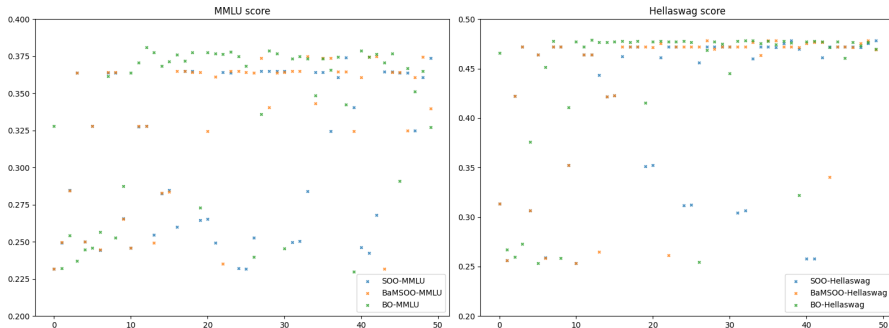


Figure: Résumé des résultats par sampling. Détails par algorithme dans les annexes 4, 5 et 6

Perspectives

Poursuite du travail

- ▶ Retour sur l'article et présentation en conférence (si validation)
- ▶ Elargissement de l'espace de recherche
- ▶ Diversification sur les modèles/données

Généralisation hors LLM

- ▶ Parallélisation d'algorithme d'optimisation de fonction coûteuse pour Exascale (project Exa-MA)
- ▶ Intégration de substitut dans les algorithme parallèle à Partition ³

³Partition-based Parallel Bayesian Optimisation

Sommaire

4. Conclusion

Résultats du stage

Une conclusion

Apprentissage

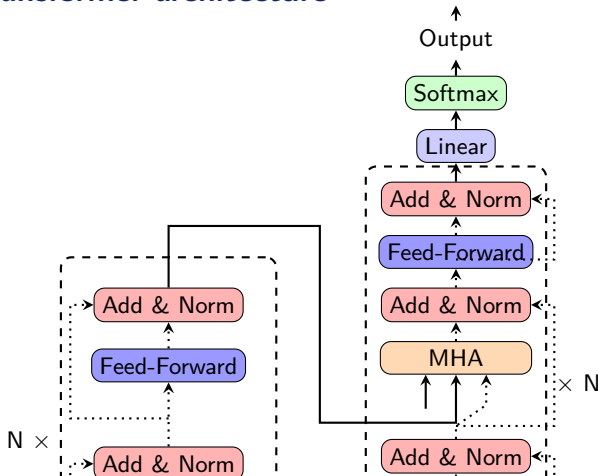
Une conclusion

Poursuite du projet professionnel

Une conclusion

Merci.

Annexe 1 : Transformer architecture



Annexe 2 : Multi-Head Attention

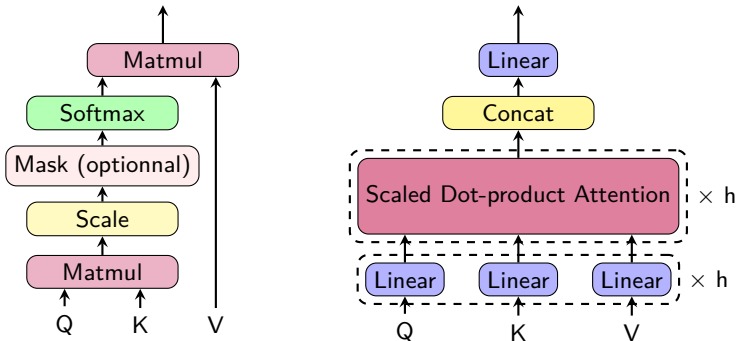


Figure: Illustration du mécanisme d'auto-attention : A droite le mécanisme complet, a gauche le *Scaled Dot-product Attention*

Annexe 3 : Low Rank Adaptation (LoRA)

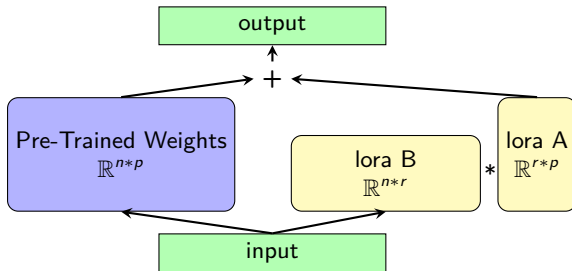


Figure: Illustration de l'application du Low Rank Adaptation (LoRA)

Annexe 5 : Résultats pour S00

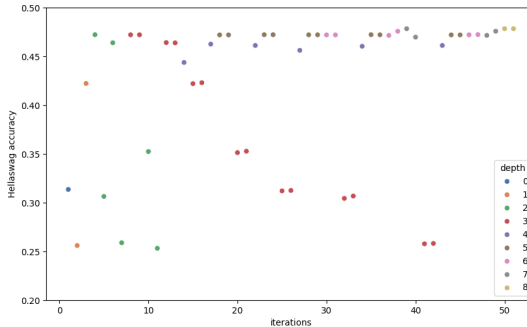


Figure: Evolution des score lors de l'expérience SOO

