



Offre n°2024-07921

## Doctorant F/H Des Grands Modèles de Langage pour la détection et la correction des erreurs dans les applications HPC

Type de contrat : CDD

Niveau de diplôme exigé : Bac + 5 ou équivalent

Fonction : Doctorant

### A propos du centre ou de la direction fonctionnelle

Le centre Inria de l'université de Bordeaux est un des neuf centres d'Inria en France et compte une vingtaine d'équipes de recherche. Le centre Inria est un acteur majeur et reconnu dans le domaine des sciences numériques. Il est au cœur d'un riche écosystème de R&D et d'innovation : PME fortement innovantes, grands groupes industriels, pôles de compétitivité, acteurs de la recherche et de l'enseignement supérieur, laboratoires d'excellence, institut de recherche technologique...

### Contexte et atouts du poste

Nous proposons un contrat de thèse sur une durée de 3 ans dans l'équipe Storm (<https://team.inria.fr/storm/>) du centre Inria de l'Université de Bordeaux.

### Mission confiée

Afin de résoudre les plus grands problèmes scientifiques en un temps raisonnable, les applications sont parallélisées et lancées sur des supercalculateurs. Cependant, ces supercalculateurs sont de plus en plus complexes et puissants, ce qui entraîne une évolution des applications (ex., nouveaux algorithmes pour le passage à l'échelle, combinaison de modèles de programmation parallèle). Cette évolution implique de nombreux défis de programmation et un réel besoin d'outils et techniques pour aider les développeurs à utiliser au mieux les différentes machines et architectures à leur disposition. En effet, à grande échelle, les développeurs d'applications font face à de nouvelles erreurs, liées au parallélisme, souvent difficiles à analyser et corriger. Aujourd'hui, s'assurer que les applications parallèles s'exécutent correctement devient aussi important que d'obtenir de bonnes performances.

Les grands modèles de langage (LLMs) sont un sujet de recherche en pleine évolution. En particulier, leurs récents succès pour générer du texte pertinent et répondre à des questions en font des candidats attrayants dans le domaine de la vérification.

#### Objectif:

L'objectif de cette thèse est d'exploiter et d'adapter les Grands Modèles de Langage pour identifier et corriger les erreurs dans les programmes parallèles. Pour cela, nous proposons d'entraîner des modèles sur des ensembles de données soigneusement générés et étiquetés grâce à une combinaison de techniques d'apprentissage et de traitement du langage naturel.

#### Collaboration :

La personne recrutée sera sous la direction d'Emmanuelle Saillard et Mihail Popov. Elle sera également en lien avec Pablo Oliveira (Université de Versailles) et Eric Petit (Intel).

### Principales activités

Le programme de recherche est découpé en 4 axes d'exploration.

## Axe 1 : Cr ation d'un jeu de donnees

Un jeu de donnees de haute qualite est une condition necessaire pour creer des modeles precis. Pour creer notre jeu de donnees, nous nous appuierons sur deux sources complementaires contenant des codes corrects et incorrects. Dans un premier temps, nous exploiterons la base de donnees git d'EasyPAP [1], une plateforme qui enseigne la programmation parallele. Bien que limitee en taille, le code soumis par les  tudiants est repr sentatif des erreurs que font les debutants. Nous explorerons ensuite Github via son API integree pour collecter des codes reels et plus conséquents en taille. Les projets seront selectionnes selon les issues, pull requests et descriptions des commits. Nous recuperons le code avant et apres les commits pertinents.

## Axe 2 : Labellisation

Une fois le jeu de donnees cree, l' tape cruciale est d'etiqueter les programmes, c'est- -dire d'associer chaque programme avec un label (erreur presente dans le code ou corrige). Pour cela, on utilisera des techniques de NLP. Les descriptions des commits et toute meta-information associee (e.g., CI) seront analysees avec TF-IDF (ou optionnellement des textes d'embeddings a la word2vec). Les vecteurs obtenus seront traites avec NMF [2] pour en extraire les differentes classes d'erreurs que nous  tudierons. En parallele, nous pourrions egalement directement utiliser des LLMs (e.g., ChatGPT) sur les commit pour les grouper. De plus, nous analyserons l'embedding des codes avant et apres chaque commit [3] : les vecteurs obtenus seront clusterises pour grouper des changements similaires. A terme, nous unifierons les deux classifications pour creer un processus de labellisation plus general.

## Axe 3 : Entraınement des modeles

Nous visons deux types de modeles. Nous commencerons par creer des modeles supervisees (Code2Error) qui prennent le code source (ou une representation du compilateur, e.g., LLVM IR) d'un programme et preddisent la categorie d'erreur associee au programme (basee sur la labellisation). Ces modeles permettront de classer les codes incorrects et d'enrichir les descriptions des problemes. En detail, Code2Error utilisera un embedding (e.g., ir2vec, code2vec) pour generer des vecteurs, a partir des codes, auxquels nous appliquerons un modele supervise (e.g., arbre de decision) pour decider du label. Les codes avant et apres le commit serviront a donner a l'arbre la version incorrecte et sa correction. Nous avons deja valide une version preliminaire (ir2vec & arbre de decision) sur 2000 codes tests dedies pour la verification MPI et souhaitons passer ce modele a l'echelle sur de vrais codes.

Ensuite, nous utiliserons des LLMs (Code2Fix). Pour chaque erreur (et donc groupe de commits associes), nous entraınerons un LLM specialise. Ce LLM recevra les codes corrects et incorrects associes a une certaine erreur. Nous utiliserons ici les codes sources (car plus utile pour l'utilisateur) et entraınerons (fine tuning) le LLM pour passer de la version erronee a la version correcte. Nous pourrions appliquer Code2Error sur un programme inconnu pour identifier le type d'erreur qu'il contient, et appeler le LLM Code2Fix associe a l'erreur pour essayer de la resoudre. Notre intuition est qu'un LLM specialise par erreur sera plus efficace. Enfin, on pourra explorer la granularite du code pour Code2Error & Code2Fix. De petites granularites seront faciles a gerer pour le modele et donc pour trouver la localisation de l'erreur au moment de la correction mais pourront manquer de contexte pour traiter certaines erreurs complexes. Ce sera un compromis a explorer.

## Axe 4 : Dissemination

Les differents modeles (Code2Fix, Code2Error) seront appliques a des projets existants pour chercher et corriger des erreurs existantes. Nous validerons egalement nos modeles sur des erreurs que nous aurons exclues du jeu de donnees pour l'apprentissage afin de mettre en avant la generalisation de notre methode et estimer a quel point deux erreurs sont similaires (si nous pouvons preddire une erreur avec des informations provenant d'une autre erreur, il est probable qu'elles soient liees). Les experts en outils de verification pourront utiliser cette information pour definir de nouvelles topologies d'erreurs. Enfin, nous envisageons d'entendre notre ensemble de donnees avec de nouveaux codes generees automatiquement par le biais des LLMs (Dataset2Code).

## R f rences :

[1] A. Lasserre, R. Namyst, and P.-A. Wacrenier. EasyPAP : a framework for learning parallel programming. In 2020 IEEE International Parallel and Distributed Processing Symposium Workshops (IPDPSW), pages 276–283, 2020.

[2] S. Heldens, P. Hijma, B. Werkhoven, J. Maassen, A. Belloum, and R. Van Nieuwpoort. The landscape of exascale research : A data-driven literature analysis. ACM Computing Surveys, 53(2) :1–43, Mar. 2020.

[3] H. Wang, G. Ye, Z. Tang, S. H. Tan, S. Huang, D. Fang, Y. Feng, L. Bian, and Z. Wang. Combining graph-based learning with automated data collection for code vulnerability detection. Trans. Info. For. Sec., 16 :1943–1958, jan 2021.

## Comp tences

- Motivation
- Curiosite et capacite a apprendre de nouveaux concepts

- Expérience avec l'écriture de scripts (ex., Python)
- Maîtrise des bases Linux
- Des connaissances en ML est un plus

## Avantages

- Restauration subventionnée
- Transports publics remboursés partiellement
- Congés: 7 semaines de congés annuels + 10 jours de RTT (base temps plein) + possibilité d'autorisations d'absence exceptionnelle (ex : enfants malades, déménagement)
- Possibilité de télétravail et aménagement du temps de travail
- Équipements professionnels à disposition (visioconférence, prêts de matériels informatiques, etc.)
- Prestations sociales, culturelles et sportives (Association de gestion des œuvres sociales d'Inria)
- Accès à la formation professionnelle
- Sécurité sociale

## Rémunération

Montant salaire brut 1e année : 2100€

Montant salaire brut 2e et 3e année : 2190€

## Informations générales

- **Thème/Domaine** : Calcul distribué et à haute performance  
Calcul Scientifique (BAP E)
- **Ville** : Talence
- **Centre Inria** : [Centre Inria de l'université de Bordeaux](#)
- **Date de prise de fonction souhaitée** : 2024-10-01
- **Durée de contrat** : 3 ans
- **Date limite pour postuler** : 2024-07-31

## Contacts

- **Équipe Inria** : [STORM](#)
- **Directeur de thèse** :  
Saillard Emmanuelle / [emmanuelle.saillard@inria.fr](mailto:emmanuelle.saillard@inria.fr)

## A propos d'Inria

Inria est l'institut national de recherche dédié aux sciences et technologies du numérique. Il emploie 2600 personnes. Ses 215 équipes-projets agiles, en général communes avec des partenaires académiques, impliquent plus de 3900 scientifiques pour relever les défis du numérique, souvent à l'interface d'autres disciplines. L'institut fait appel à de nombreux talents dans plus d'une quarantaine de métiers différents. 900 personnels d'appui à la recherche et à l'innovation contribuent à faire émerger et grandir des projets scientifiques ou entrepreneuriaux qui impactent le monde. Inria travaille avec de nombreuses entreprises et a accompagné la création de plus de 200 start-up. L'institut s'efforce ainsi de répondre aux enjeux de la transformation numérique de la science, de la société et de l'économie.

## L'essentiel pour réussir

Le ou la candidate doit avoir un bon niveau de programmation.

De plus, la personne recrutée devra relire de la bibliographie scientifique, écrire des rapports/articles et présenter ses travaux devant la communauté. De ce fait, un bon niveau de communication en anglais sera fortement apprécié.

**Attention:** Les candidatures doivent être déposées en ligne sur le site Inria. Le traitement des candidatures adressées par d'autres canaux n'est pas garanti.

## Consignes pour postuler

Si vous êtes intéressés, merci de bien vouloir candidater via le site [jobs.inria](https://jobs.inria.fr) avec les documents suivants :

- CV
- lettre de motivation
- notes master
- lettre de recommandation le cas échéant

**Sécurité défense :**

Ce poste est susceptible d'être affecté dans une zone à régime restrictif (ZRR), telle que définie dans le décret n°2011-1425 relatif à la protection du potentiel scientifique et technique de la nation (PPST). L'autorisation d'accès à une zone est délivrée par le chef d'établissement, après avis ministériel favorable, tel que défini dans l'arrêté du 03 juillet 2012, relatif à la PPST. Un avis ministériel défavorable pour un poste affecté dans une ZRR aurait pour conséquence l'annulation du recrutement.

**Politique de recrutement :**

Dans le cadre de sa politique diversité, tous les postes Inria sont accessibles aux personnes en situation de handicap.