

# Offre n°2024-07921

# Doctorant F/H Des Grands Modèles de Language 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 relsoudre les plus grands problelmes scientifiques en un temps raisonnable, les applications sont parallelliseles et lanceles sur des supercalculateurs. Cependant, ces supercalculateurs sont de plus en plus complexes et puissants, ce qui entraine une elvolution des applications (ex., nouveaux algorithmes pour le passage all l'elchelle, combinaison de modelles de programmation parallelle). Cette elvolution lelve de nombreux delfis de programmation et un relel besoin d'outils et techniques pour aider les delveloppeurs all utiliser au mieux les diffellentes machines et architectures all leur disposition. En effet, all grande elchelle, les delveloppeurs d'applications font face all de nouvelles erreurs, lielles au parallellisme, souvent difficiles all analyser et corriger. Aujourd'hui, s'assurer que les applications parallelles s'exellcutent correctement devient aussi important que d'obtenir de bonnes performances.

Les grands modelles de langage (LLMs) sont un sujet de recherche en pleine elvolution. En particulier, leurs relicents succells pour gellnellrer du texte pertinent et relipondre all des questions en font des candidats attrayants dans le domaine de la vellrification.

## Objectif:

L'objectif de cette thèse est d'exploiter et d'adapter les Grands Mode 🛮 les de Langage pour identifier et corriger les erreurs dans les programmes paralle 🖺 les. Pour cela, nous proposons d'entrai 🗈 ner des mode 🖺 les sur des ensembles de donne 🖺 es soigneusement ge 🖺 ne 🖺 re 🖺 s et e 🖺 tiquete 🖺 s gra 🖺 ce a 🖺 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: CreDation d'un jeu de donneDes

Un jeu de donne des de haute qualite dest une condition ne dessaire pour cre des mode des predicis. Pour cre der notre jeu de donne des, nous nous appuierons sur deux sources compledimentaires contenant des codes corrects et incorrects. Dans un premier temps, nous exploiterons la base de donne des git d'EasyPAP [1], une plateforme qui enseigne la programmation paralle de. Bien que limite den taille, le code soumis par les editudiants est repredisentatif des erreurs que font les dedibutants. Nous explorerons ensuite Github via son API intedigred pour collecter des codes redels et plus consediquents en taille. Les projets seront sedlectionne des selon les issues, pull requests et descriptions des commits. Nous red cupe d'ereons le code avant et apre des commits pertinents.

#### Axe 2: Labellisation

Une fois le jeu de donne les cre le l, l'e la tape cruciale est d'e liqueter les programmes, c'est-a l-dire d'associer chaque programme avec un label (erreur pre l'sente dans le code ou corrige le). Pour cela, on utilisera des techniques de NLP. Les descriptions des commits et toute me la-information associe le (e.g., CI) seront analyse les avec TF-IDF (ou optionnellement des textes d'embeddings a la word 2 vec). Les vecteurs obtenus seront traite la avec NMF [2] pour en extraire les diffe l'entes classes d'erreurs que nous e l'tudierons. En paralle le, nous pourrons e l'galement directement utiliser des LLMs (e.g., Chat GPT) sur les commit pour les grouper. De plus, nous analyserons l'embedding des codes avant et apre la chaque commit [3] : les vecteurs obtenus seront clusterise les pour grouper des changements similaires. A terme, nous unifirons les deux classifications pour cre le run processus de labellisation plus ge l'ne l'al.

#### Axe 3: Entrai nement des mode les

Nous visons deux types de mode les. Nous commencerons par cre le r des mode les supervise ls (Code2 Error) qui prennent le code source (ou une repre le sentation du compilateur, e.g., LLVM IR) d'un programme et pre l'disent la cate le gorie d'erreur associe le au programme (base le sur la labellisation). Ces mode les permettront de classer les codes incorrects et d'enrichir les descriptions des proble l'es met les descriptions des proble l'es mode l'es permettront de classer les codes incorrects et d'enrichir les descriptions des proble l'es. En de l'atil, Code les permettront de classer les codes incorrects et d'enrichir les descriptions des proble l'es. En de l'atil, Code les descriptions des proble l'es mode l'es upervise le description des vecteurs, a l'es partir des codes, auxquels nous appliquerons un mode le supervise le (e.g., arbre de de l'esion) pour decider du label. Les codes avant et apre le commit serviront a l'antore la version incorrecte et sa correction. Nous avons de l'alle l'une version pre l'iminaire (ir l'evec & arbre de decision) sur 2000 codes tests de l'die le pour la ve l'ification MPI et souhaitons passer ce mode l'e l'e l'elle sur de vrais codes.

Ensuite, nous utiliserons des LLMs (Code2Fix). Pour chaque erreur (et donc groupe de commits associe\(\Pi\)s, nous entrai\(\Pi\)nerons un LLM specialise\(\Pi\). Ce LLM recevra les codes corrects et incorrects associe\(\Pi\)s a\(\Pi\) une certaine erreur. Nous utiliserons ici les codes sources (car plus utile pour l'utilisateur) et entrai\(\Pi\)nerons (fine tuning) le LLM pour passer de la version errone\(\Pi\)e a\(\Pi\) la version correcte. Nous pourrons appliquer Code2Error sur un programme inconnu pour identifier le type d'erreur qu'il contient, et appeller le LLM Code2Fix associe\(\Pi\) a\(\Pi\) l'erreur pour essayer de la re\(\Pi\)soudre. Notre intuition est qu'un LLM spe\(\Pi\)cialise\(\Pi\) par erreur sera plus efficace. Enfin, on pourra explorer la granularite\(\Pi\) du code pour Code2Error \(\&\Circ\) Code2Fix. De petites granularite\(\Pi\)s seront faciles a\(\Pi\) ge\(\Pi\)rer pour le mode\(\Pi\)le et donc pour trouver la localisation de l'erreur au moment de la correction mais pourront manquer de contexte pour traiter certaines erreurs complique\(\Pi\)es. Ce sera un compromis a\(\Pi\) explorer.

#### Axe 4: Disse mination

Les diffe I rents mode I les (Code2Fix, Code2Error) seront applique Is a I des projets existants pour chercher et corriger des erreurs existantes. Nous validerons e II galement nos mode I les sur des erreurs que nous aurons exclues du jeu de donne I les pour l'apprentissage afin de mettre en avant la ge I ne II ralisation de notre me II thode et estimer a II quel point deux erreurs sont similaires (si nous pouvons pre II dire une erreur avec des informations provenant d'une autre erreur, il est probable qu'elles soient lie I les). Les experts en outils de ve II rification pourront utiliser cette information pour de II finir de nouvelles topologies d'erreurs. Enfin, nous envisageons d'e II tendre notre ensemble de donne I les avec de nouveaux codes ge I ne II re II sautomatiquement par le biais des LLMs (Dataset 2 Code).

#### Références:

[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.

# Compétences

- Motivation
- Curiosite II et capacite II a II apprendre de nouveaux concepts

- Expellrience avec l'ellcriture de scripts (ex., Python)
- Mai I trise 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

# 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 canditate 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 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.