

INRAE - Unité MIA Paris Saclay & IATE

Offre de Stage - Master 2 / Ingénieur - 6 mois

Sujet : Complétion automatique de définitions d'ontologies à l'aide de modèles de deep learning (LLM)

Contexte

Les ontologies sont essentielles pour structurer et partager les connaissances dans les sciences de l'agriculture, l'alimentation et l'environnement. Elles permettent une description standardisée et cohérente des concepts, facilitant l'interopérabilité des données et leur réutilisation à travers divers projets de recherche.

Dans ce contexte, l'ontologie **TransformON** (Weber et al. 2023) a été développée pour couvrir l'ensemble des domaines de connaissance produits par le département Transform d'INRAE. Basée sur la core ontologie **PO²** (Munch et al. 2022), TransformON structure le vocabulaire autour des itinéraires de production des aliments et des bioproduits, en lien avec les bénéfices et les risques pour la santé humaine et l'environnement. Cependant, certaines définitions de concepts manquent ou sont incomplètes, ce qui peut nuire à l'utilisation optimale de l'ontologie.

Ce stage a pour objectif de combler ces lacunes en utilisant des techniques avancées de deep learning, notamment des grands modèles de langage (LLM), pour générer automatiquement des définitions pertinentes et complètes.

Objectifs

L'objectif du stage est de développer une méthode permettant de compléter automatiquement les définitions manquantes de concepts dans des ontologies en s'appuyant sur l'ontologie TransformON. Le travail sera structuré en deux étapes principales :

1. **Fine-tuning** : Adapter un LLM existant pour la tâche spécifique de génération de définitions d'ontologies en se basant sur la structure et les définitions de concepts existante de l'ontologie (Erbacher et al. 2024).
2. **Renforcement avec feedback humain (RLHF)** : Mettre en place un processus d'apprentissage par renforcement avec feedback humain (Reinforcement Learning with Human Feedback) pour améliorer la pertinence et la qualité des définitions générées, en intégrant des retours d'experts du domaine (Ouyang et al. 2022).

Tâches

Les principales missions du stage sont les suivantes :

1. **État de l'art** : Réaliser une revue de littérature sur les méthodes actuelles de génération de texte avec des LLM pour les tâches de complétion d'ontologies et de génération de définitions.
2. **Préparation des données** : Sélectionner et préparer un corpus de définitions d'ontologies pour le fine-tuning du modèle.

3. **Fine-tuning du modèle** : Choisir et adapter un modèle de langage pré-entraîné pour la tâche de complétion de définitions, en ajustant les hyperparamètres et en optimisant les performances sur des exemples annotés.
4. **RLHF** : Intégrer le feedback humain dans le modèle en utilisant des méthodes de RLHF pour ajuster la génération de définitions en fonction de critères qualitatifs définis avec les experts du domaine.
5. **Évaluation** : Évaluer la qualité des définitions générées, comparer avec les modèles de base et réaliser des tests d'acceptabilité par les utilisateurs finaux.
6. **Documentation et présentation des résultats** : Produire un rapport de stage complet et préparer une présentation finale pour l'équipe de recherche.

Compétences requises

- Formation en informatique, intelligence artificielle, ou domaine connexe (Master 2, école d'ingénieur).
- Bonne compréhension des concepts de machine learning et deep learning.
- Connaissances des modèles de langage (LLM) et techniques de fine-tuning.
- Familiarité avec le Python et les bibliothèques de machine learning (TensorFlow, PyTorch, Hugging Face).
- Intérêt pour le traitement du langage naturel (NLP) et les ontologies.
- Connaissance des bases en renforcement par feedback humain (RLHF) serait un plus.

Références

Magalie Weber, Patrice Buche, Liliana Ibanescu, Stéphane Dervaux, Hervé Guillemin, et al.. PO2/TransformON, an ontology for data integration on food, feed, bioproducts and biowaste engineering. *npj Science of Food*, 2023, 7, pp.47. [10.1038/s41538-023-00221-2](https://doi.org/10.1038/s41538-023-00221-2). [hal-04197618](https://hal.archives-ouvertes.fr/hal-04197618)

Mélanie Munch, Patrice Buche, Stéphane Dervaux, Juliette Dibie, Liliana Ibanescu, et al.. Combining ontology and probabilistic models for the design of bio-based product transformation processes. *Expert Systems with Applications*, 2022, 203, pp.117406. [10.1016/j.eswa.2022.117406](https://doi.org/10.1016/j.eswa.2022.117406). [hal-03662183](https://hal.archives-ouvertes.fr/hal-03662183)

Ouyang, L., Wu, J., Jiang, X., Almeida, D., Wainwright, C., Mishkin, P., ... & Lowe, R. (2022). Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35, 27730-27744.

Erbacher, P., Falissard, L., Guigue, V., & Soulier, L (2024) Navigating Uncertainty: Optimizing API Dependency for Hallucination Reduction in Closed-Book Question Answering. In ECIR 2024

Encadrement

Le stagiaire sera encadré par l'équipe de l'unité MIA Paris-Saclay, en collaboration avec des experts en ontologies et traitement du langage naturel.

Conditions

- Lieu de travail : Campus Agro Paris-Saclay, 22 place de l'Agronomie, 91120 Palaiseau
 - Laboratoire d'accueil : MIA Paris Saclay, équipe EKINOCS
 - Durée : 6 mois.
 - Gratification : environ 570 euros/mois
-

Pour postuler :

Merci d'envoyer votre CV et une lettre de motivation à stephane.dervaux@inrae.fr, liliana.ibanescu@agroparistech.fr, patrice.buche@inrae.fr, vincent.quigue@agroparistech.fr