

2020-02332 - Post-Doctoral Research Visit F/M Inspecting and Debugging a VQA System through Explanations

Contract type : Fixed-term contract
Level of qualifications required : PhD or equivalent
Fonction : Post-Doctoral Research Visit

About the research centre or Inria department

The Inria Lille - Nord Europe Research Centre was founded in 2008 and employs a staff of 360, including 300 scientists working in sixteen research teams. Recognised for its outstanding contribution to the socio-economic development of the Hauts-De-France région, the Inria Lille - Nord Europe Research Centre undertakes research in the field of computer science in collaboration with a range of academic, institutional and industrial partners.

The strategy of the Centre is to develop an internationally renowned centre of excellence with a significant impact on the City of Lille and its surrounding area. It works to achieve this by pursuing a range of ambitious research projects in such fields of computer science as the intelligence of data and adaptive software systems. Building on the synergies between research and industry, Inria is a major contributor to skills and technology transfer in the field of computer science.

Context

Applications are welcome for a post-doctoral position at Inria as part of the **HyAIAI Inria challenge**. More specifically, this position is part of a collaboration between the **Lacodam** and **SequeL** Inria teams. The post-doc will share her time between the two groups (**Lacodam** in Rennes, **SequeL** in Lille). She will actively participate in the activities of HyAIAI, in particular reporting on her work in the HyAIAI meetings.

Background

Visual Question Answering (VQA) is a research area concerned with the construction of computer systems that can answer questions in natural language from the contents of an image. VQA carries potential applications in multimodal information retrieval.

Current VQA solutions rely on deep learning techniques. Being a problem on multimodal data, this implies to merge both images and questions into a common representation space. This is challenging because images and texts are very different data types, treated by means of different neural network architectures: CNNs (Convolutional Neural Networks) are the state of the art for image classification and representation, whereas text processing often resorts to RNNs (Recurrent Neural Networks). VQA solutions must orchestrate both technologies leading to systems that are extremely complex.

The complexity of existing VQA solutions makes the task of inspection and debugging very hard. In particular, neural networks are black-box models: one requires a significant amount of work and solid expertise to understand the inner-workings of the network. It becomes therefore very difficult to understand why a VQA system makes a mistake. Such a task, however, is vital for the progress of research in VQA.

Assignment

The main purpose of this post-doctoral fellowship is to port the principles of interpretable AI and ML to the domain of the visual question answering. Attaining such a goal requires us to overcome other challenges such as understanding what makes a good explanation in a multi-modal setting.

Main activities

The post-doctoral researcher in charge of this project will work on methods to explain the output of a VQA system in order to understand why it erred (or not).

In this regard, we aim at deploying post-hoc interpretability modules that can provide hints on the logic behind a VQA module for a given case. Techniques such as **LIME** [Ribeiro et. al., 2016], **SHAP** [Lundberg and Lee, 2017], or **Anchors** [Ribeiro et. al., 2018] provide the foundations to generate such explanations for any type of black-box model. That said, those techniques deliver explanations in terms of an interpretable space that depends on the data type (e.g., tabular data, texts, images) and is usually very different from the representation space of the black-box. For instance, pixel channels and feature maps in images are replaced by super-pixels (image segments) in explanations, whereas word embeddings are substituted by word occurrences. No research work until now has tried to reconcile those representations in a multimodal setting, thus the main challenge is to find a common interpretable representation for explanations that encompasses both images and texts in a VQA setting. We are particularly interested in representations that resort to semantically meaningful units such as known objects or patterns (e.g., pointy borders, a nose, limbs, vehicles) as studied in **network dissection**. This would allow us to yield explanations in terms of the presence or absence of those objects. Such explanations could be enhanced with semantic knowledge, such as a taxonomy, in order to signal interesting associations automatically, e.g., pointy borders in a bird may refer to its beak. Other recent approaches [Shi et al., 2019], [Yi et al., 2018] have focused on the extraction of knowledge from the input before defining a reasoning program to execute. This knowledge may be a start in the definition of an interpretable space for explanations of VQA systems.

Traditional post-hoc interpretability modules do not make any assumptions about the architecture of the model they try to explain, i.e., they are model-agnostic. In the context of VQA systems, a possible solution is to drop this assumption and mine neuron activation patterns that identify the instances for which the system fails. We could use contrast pattern mining techniques for this purpose.

Our research will make use of publicly available VQA datasets such as GQA, VQA-E, C-VQA, and CLEVR.

General Information

- **Theme/Domain :** Optimization, machine learning and statistical methods Information system (BAP E)
- **Town/city :** Lille
- **Inria Center :** CRI Lille - Nord Europe
- **Starting date :** 2020-06-01
- **Duration of contract :** 2 years
- **Deadline to apply :** 2020-06-30

Contacts

- **Inria Team :** SEQUEL
- **Recruiter :**
Galarraga Del Prado Luis / luis.galarraga-del-prado@inria.fr

About Inria

Inria is the French national research institute dedicated to digital science and technology. It employs 2,600 people. Its 200 agile project teams, generally run jointly with academic partners, include more than 3,500 scientists and engineers working to meet the challenges of digital technology, often at the interface with other disciplines. The Institute also employs numerous talents in over forty different professions. 900 research support staff contribute to the preparation and development of scientific and entrepreneurial projects that have a worldwide impact.

Instruction to apply

To apply for the position, the candidate must send an email to the list of contacts below. The email should include:

- A CV
- A statement letter explaining the candidate's motivations to apply for the position
- At least two recommendation letters

Contacts

Philippe Preux (philippe.preux@inria.fr)

Luis Galarraga (luis.galarraga@inria.fr)

Defence Security :

This position is likely to be situated in a restricted area (ZRR), as defined in Decree No. 2011-1425 relating to the protection of national scientific and technical potential (PPST). Authorisation to enter an area is granted by the director of the unit, following a favourable Ministerial decision, as defined in the decree of 3 July 2012 relating to the PPST. An unfavourable Ministerial decision in respect of a position situated in a ZRR would result in the cancellation of the appointment.

Recruitment Policy :

As part of its diversity policy, all Inria positions are accessible to people with disabilities.

Warning : you must enter your e-mail address in order to save your application to Inria. Applications must be submitted online on the Inria website. Processing of applications sent from other channels is not guaranteed.

Skills

We are searching for motivated candidates with a PhD degree in Computer Science and with competences in machine learning (preferably with focus on deep learning). Knowledge in data mining, e.g., sequence and itemset mining, will be also appreciated.

The candidate should be proficient in written and spoken English (at least B2 level according to the CEFR system).

Benefits package

- Subsidized meals
- Partial reimbursement of public transport costs
- Leave: 7 weeks of annual leave + 10 extra days off due to RTT (statutory reduction in working hours) + possibility of exceptional leave (sick children, moving home, etc.)
- Possibility of teleworking (after 6 months of employment) and flexible organization of working hours
- Professional equipment available (videoconferencing, loan of computer equipment, etc.)
- Social, cultural and sports events and activities
- Access to vocational training
- Social security coverage

Remuneration

Gross monthly salary (before taxes) : 2653 €