# Splitting approach and noise reduction for PyRAT's neural network propagation

**Keywords**: PyRAT, Neural Network, Abstract Interpretation, Zonotopes

## Institution

The French Alternative Energies and Atomic Energy Commission (CEA) is a key player in research, development, and innovation. Drawing on the widely acknowledged expertise gained by its 16,000 staff spanned over 9 research centers with a budget of 4.1 billion Euros, CEA actively participates in more than 400 European collaborative projects with a large number of academic (notably as a member of Paris-Saclay University) and industrial partners. Within the CEA Technological Research Division, the CEA List institute addresses the challenges coming from smart digital systems.

Among other activities, CEA List's Software Safety and Security Laboratory (LSL) research teams design and implement automated analysis in order to make software systems more trustworthy, to exhaustively detect their vulnerabilities, to guarantee conformity to their specifications, and to accelerate their certification. Recently the field of activity of the laboratory has been extended to artificial intelligence safety and security verification. In particular, PyRAT and CAISAR are two tools developed in that context to verify safety properties on neural networks using abstract interpretation techniques.

## Objectives

As mentioned, PyRAT leverages the principles of abstract interpretation to propagate abstract domains (input) through abstract operations representing the layers of the neural network and in order to assess the reachable states (output). In comparison to classical software verification, PyRAT works directly on the weights, biases, and parameters of a neural network model thus making PyRAT lighter and faster to use for neural network analyses. PyRAT is developed in Python as it is a widely used language for neural network frameworks such as Keras, Pytorch or Tensorflow. As of now, the primary use of PyRAT is to assess robustness w.r.t. some perturbation around inputs on small neural networks. However, on larger neural networks or on larger inputs a simple pass with PyRAT will lack precision.

To compensate the loss of precision in such cases, we introduced a domain splitting approach inside PyRAT. This iterative approach splits the input domain into smaller ones until PyRAT can prove the property on them or find a counterexample. This also allows to use PyRAT to also prove some safety properties, i.e. bigger intervals of input on certain neural networks.

For this internship, we focus on the noise introduced during a PyRAT analysis. Noise is introduced by PyRAT during its propagation through the network when it encounters non-linear functions. In neural networks such functions can be ReLUs, sigmoid... In such cases, PyRAT will overapproximate the function with a linear function and add some noise which in turn decreases the precision. To counter this, a first path to explore is to develop a splitting mechanism on certain function which are linear by pieces to avoid adding noise and handle them linearly in parallel. However, such splitting is exponential so a careful study of when to use it should be made. Additionally, this approach on domains such as the zonotopes that we are using in PyRAT will entail to add constraints on the zonotopes complexifying their concretization. A second idea to develop in this internship in order to prevent the number of noises to slow down processing, is to reduce periodically the number of noises. Here heuristics need to be planned to understand when the best timing to reduce the noise symbols is during propagation.

## Qualifications

- **Minimal**
- Master student or 2nd or 3rd year of engineering school
- knowledge of Python
- notions of AI and neural networks, and of AI frameworks (TF, Keras, Pytorch, ...)
- ability to work in a team

- **Preferred**

- some knowledge of abstract interpretation or formal method

## Characteristics

- **Duration:** 5 to 6 months from early 2022
- **Location:** CEA Nano-INNOV, Paris-Saclay Campus, France
- **Compensation:**
  - €700 to €1300 monthly stipend (determined by CEA compensation grids)
  - maximum €229 housing and travel expense monthly allowance (in case a relocation is needed)
  - CEA buses in Paris region and 75% refund of transit pass
  - subsidized lunches

## Application

If you are interested in this internship, please send to the contact persons an application containing:

- your resume;
- a cover letter indicating how your curriculum and experience match the qualifications expected and how you would plan to contribute to the project;
- your bachelor and master 1 transcripts;
- the contact details of two persons (at least one academic) who can be contacted to provide references.

Applications are welcomed until the position is filled. Please note that the administrative processing may take up to 3 months.

## Contact persons

For further information or details about the internship before applying, please contact:

- Augustin Lemesle (augustin.lemesle@cea.fr)
- Maxime Jacquemin Maxime.Jacquemin@cea.fr
- Zakaria Chihani (zakaria.chihani@cea.fr)