# STATEC

# *Text classification with transformers*

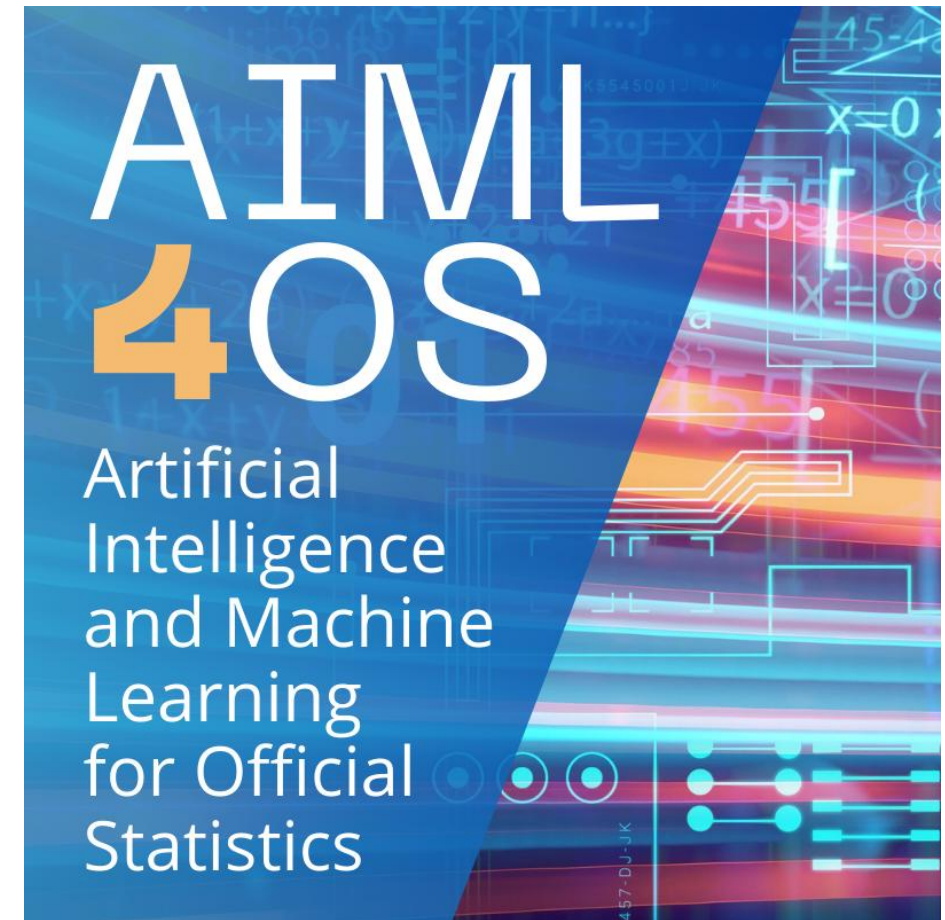**STATEC Datalab Workshop**

16/09/2025

# Introduction

- **Official statistics rely on standardized classifications (e.g., NACE, ISCO or COICOP codes)**

- **Text-to-code problems : assign classification codes to textual descriptions**

    *"I work as a software developer"* $\Rightarrow$ *ISCO 21: Science and Engineering Professionals*

- **Manual coding can be time consuming**

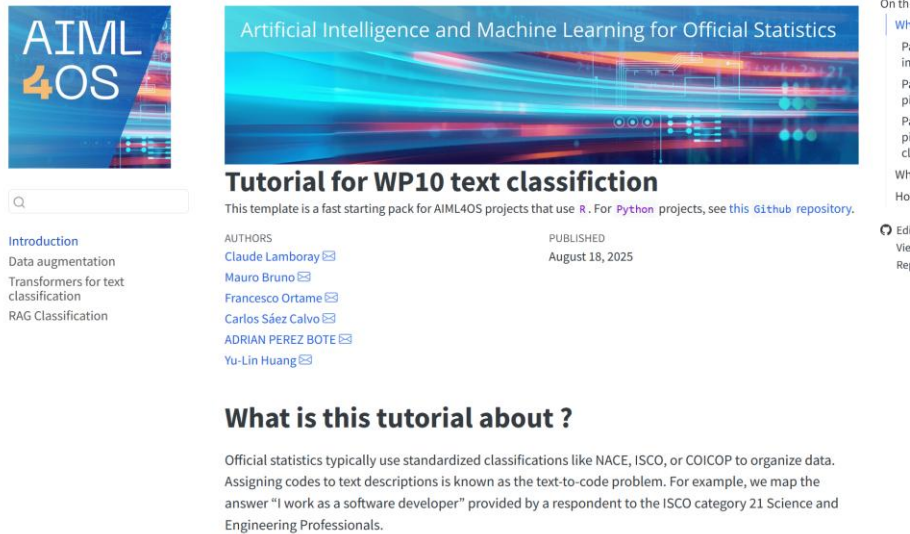- **Tutorials to demonstrate how transformer models can be applied to solve this problem more efficiently**

STATEC
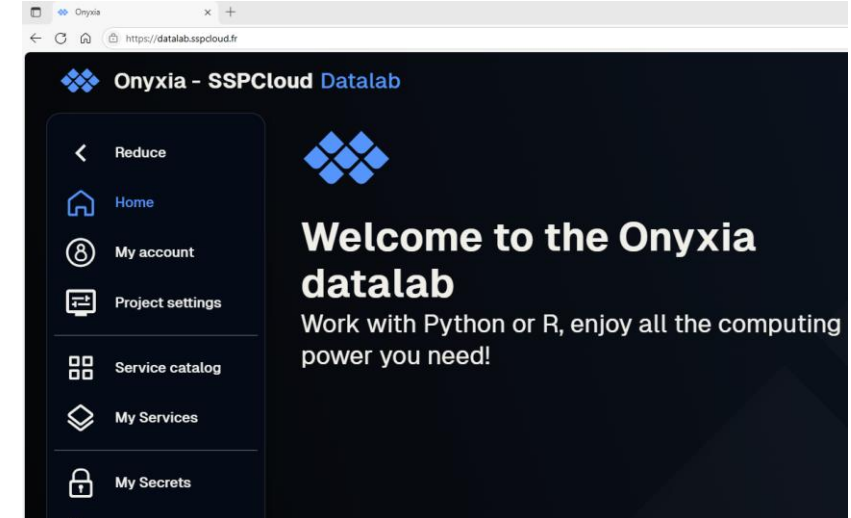
# Introduction

- **European project AIML4OS**

  https://cros.ec.europa.eu/dashboard/aiml4os

# Organization of the tutorial



**Launch Notebook
on SSP Cloud**

**Notebooks on github website
https://aiml4os.github.io/WP10_tutorial_text_classification/**

**SSP Cloud as datalab platform
https://datalab.sspcloud.fr/**

*Need to create an account on SSP cloud
Contact innovation@insee.fr if email domain not covered*

STATEC

# Agenda

1.  **Synthetic data generation to improve the training data set**

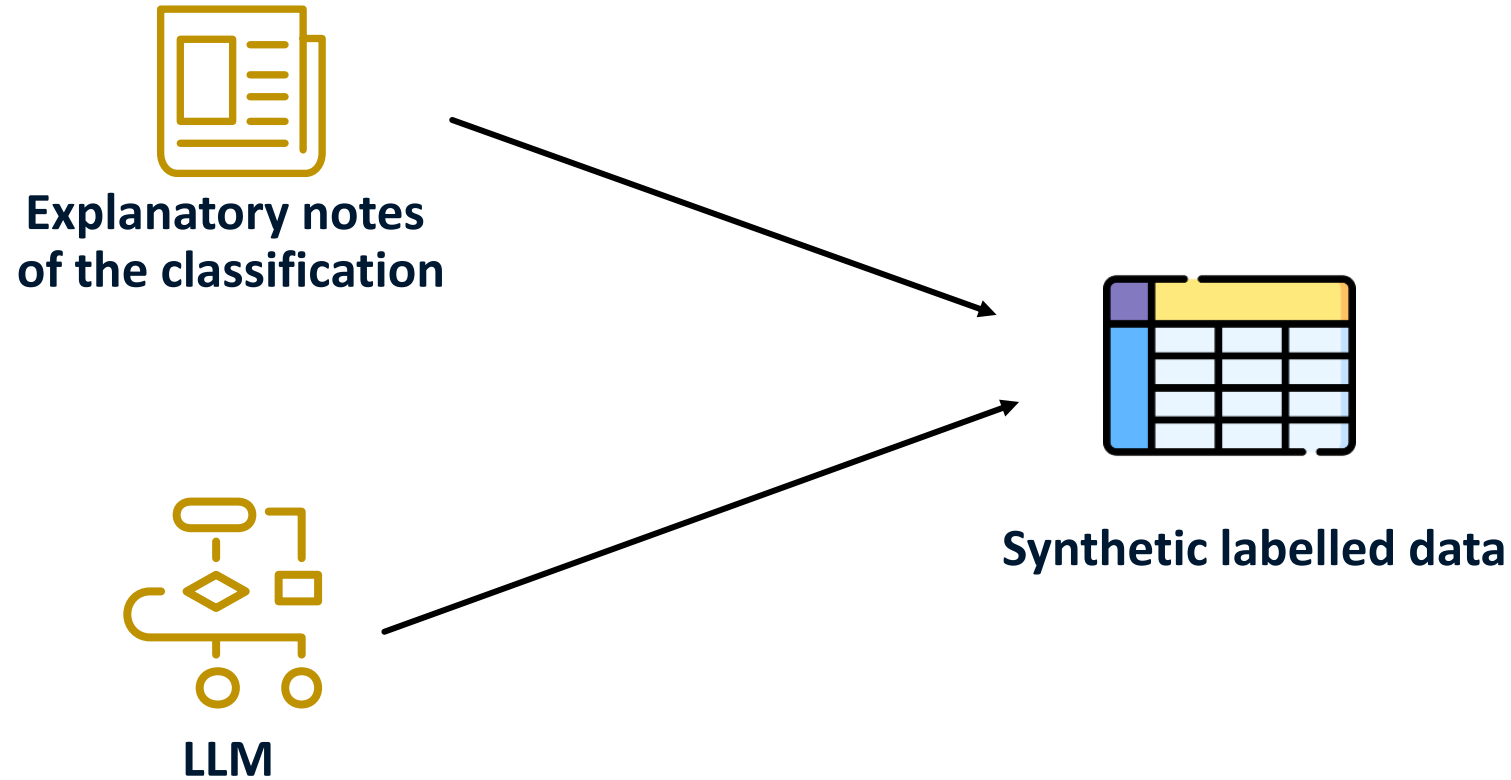    Carlos Sáez Calvo, INE Spain

2.  **Fine-tuning a transformer pipeline for text classification**

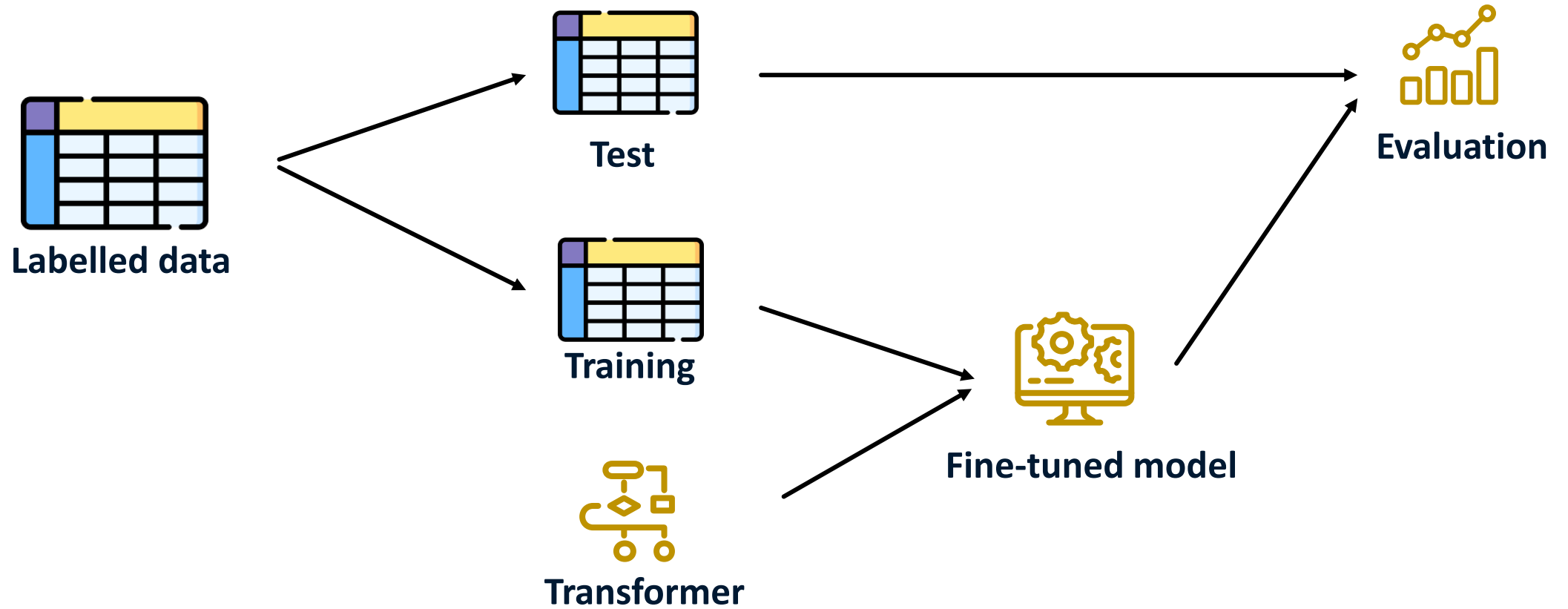    Yu-Lin Huang, STATEC and SnT (University of Luxembourg)

3.  **Setting up a simple RAG pipeline for automatic text classification**

    Francesco Ortame, ISTAT

**STATEC**

# 1. Synthetic data generation to improve the training data set

**Explanatory notes
of the classification**

**LLM**

**Synthetic labelled data**

# 2. Fine-tuning a transformer pipeline for text classification



Labelled data → Test → Evaluation

Labelled data → Training → Fine-tuned model

Transformer → Fine-tuned model

Fine-tuned model → Evaluation

# 3. Setting up a simple RAG pipeline for automatic text classification



Explanatory notes of the classification

*"I work as a software developer"*

**Text string**

**Top-k classes**

**LLM**

**Labelled data**