**Project 4**

# SäxyLLM: Saxon LLM construction

Sächsisch is sexy, but current open models don't speak it. Build an LLM that speaks Saxonian by default.

## Project Description

Research real-world applications of LMs trained to speak specific dialects. Understand the methodologies used to train such models and the effects of fine-tuning on small datasets.

Explore various strategies for adapting LMs to dialects, such as: **rule-based transformation** (applying pre-defined linguistic rules), **text-based fine-tuning** (using dialect-specific corpora), etc. Based on your chosen approach, collect and preprocess Saxon dialect data. Choose a suitable pre-trained language model (e.g. **LLaMA**) and fine-tune it on your dataset.

Evaluate your model on the following tasks:

1. Translate the "About ScaDS.AI Dresden/Leipzig" Text in Sächsisch:
   *"Our team consists of more than 180 people, including renowned international researchers as well as highly skilled professionals in administrative and communicative roles. With more than 60 principal investigators, two Humboldt Professorships and up to twelve planned AI Professorships, we support excellence in research and teaching in Leipzig and Dresden. Promoting young talent is also an important part of our work, therefore we have established four Junior Research Groups that meaningfully complement our current research topics. Furthermore, we are welcoming Associated Members who contribute their expertise to our center."*
2. Ask your model to generate three texts containing at least 1000 characters based on the following prompts:
   a. What is the history of TU Dresden?
   b. Describe how LLMs work.
   c. Write a short story about two Saxons in the Deutsche Bahn.

Explain your evaluation process and chosen metrics (e.g., BLEU, perplexity, accuracy). Analyze your results and assess the quality of the dialect-specific outputs.

Discuss the benefits and risks of localized LMs. Does training on data in a specific dialect introduce biases?

## Bonus Tasks

1. Explore possible ethical issues related to dialect-specific LMs, such as misuse, stereotypes, or marginalization. Propose and implement a solution for those problems.
2. Deploy your fine-tuned model online, making it accessible for testing and demonstration.