# Reinforcement Learning from Human Feedback in Story Generation

- **Supervisor**: Abdullatif Köksal
- **Examiner**: Prof. Schütze
- **BSc, MSc, Open**: MSc
- **Title**: Reinforcement Learning from Human Feedback in Story Generation
- **Summary**: The recent advancements in large language models have led to an increased interest in the use of reinforcement learning (RL) in text generation. OpenAI has suggested that reinforcement learning can significantly improve the quality of text generation, even a 1.3B parameters model with RL generates more coherent and plausible text than a 175B parameter model without RL [1].
The project will be focused on story generation with a dataset of short stories collected from Reddit, which includes information about the number of upvotes for each story for a given prompt. This dataset was proposed in a recent master's thesis [2], which demonstrated that the upvote information can be used to improve the quality of generated stories without RL. In this project, we will evaluate a set of reinforcement learning algorithms from human feedback (i.e., the number of upvotes) using the trlx library (https://github.com/CarperAI/trlx) to compare the quality of generated stories.
- **Prerequisites**: Reinforcement Learning, Pytorch/HuggingFace
⚠️ Prior knowledge in reinforcement learning is required for this project, as implementing RL in NLP can be challenging.

**References:**
**[1]** Training language models to follow instructions with human feedback, https://arxiv.org/abs/2203.02155
**[2]** Contact me for a copy of the thesis
**Recommended Readings:**
**1.** https://huggingface.co/blog/rlhf