# Leveraging Tool-Calling to enhance LLM Capabilities

## Student

**Moritz Schiesser**

**Introduction:** Recent advancements in generative natural language processing (NLP) models have lead companies and academic research institutions alike to develop and deploy increasingly sophisticated models for a variety of applications.
These models have shown to be capable of generating grammatically correct and coherent human-like text.
Since these models not only learn to produce grammatically correct text but also learn to recall and reproduce information from their training data, they are often unable to talk about topics that are not present in their training data.
While they are able to generalize well to unseen data, when it comes to factual real-world and real-time information, they often fall short.
This can be bypassed by providing the model information about a topic by extending the prompt with this information.
Further, LLMs cannot directly interact with other systems.
Tool-calling can be used to take action as a consequence of the user's input, be it providing information to the model, or taking action on the user's behalf.

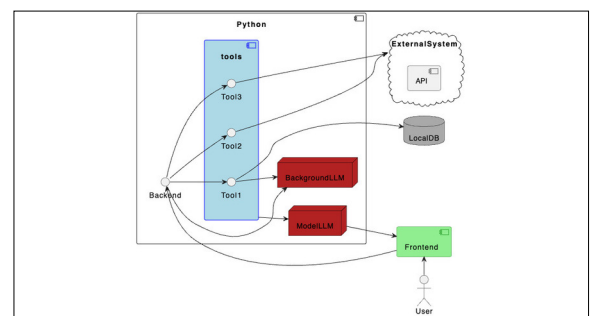This project explores how tool-calling can be leveraged to enhance the capabilities of large language models.

**Approach / Technology:** LLMs do not have access to real-time information, which makes it difficult for them to generate up-to-date information, neither can they interact with other systems directly. Tool-calling can be used to enhance the capabilities of LLMs.
A tool, as referred to in this context, is a fragment of code that interacts with other systems and provides relevant information to the LLM. For example, a tool can be used to fetch real-time information from a database or an API and provide this information to the model. To decide, which tool should be invoked, an LLM is presented with the available tools and their descriptions, alongside the user prompt. Based on the selection, required information is extracted by querying the LLM again, which is used to invoke the tool. The result of the tool is fed back to the LLM, alongside the original user prompt, to generate the final output. This project explores and illustrates this by building a chatbot system that covers parts of the `themovieDB` API.
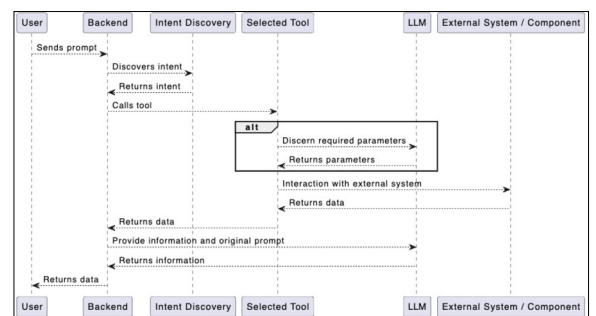
**Result:** The sample application successfully demonstrates how tool-calling can be leveraged to enhance the capabilities of various models.
The approach of using generative AI to select which programmatic tool to call, and then using the output of that tool to help the LLM generate the final output, is shown to be viable and effective.
While the selection of the generative text model has a great impact on the reliability and performance of the

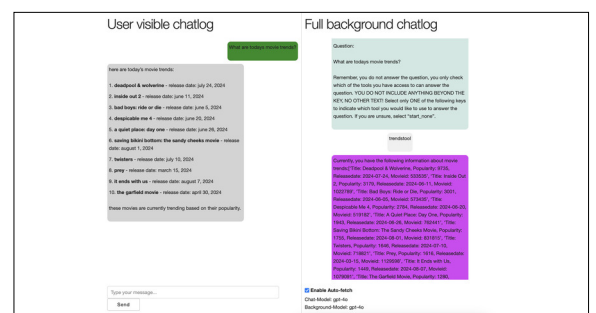application, even comparatively small models can be used to great effect.

**An overview over the components of this approach, and how they interact with each other**
Own presentment



**Sequence diagram of tool-calling**
Own presentment



**The UI of the sample application with the typical start of a conversation between user and system**
Own presentment

## Advisor
**Prof. Dr. Mitra Purandare**

## Subject Area
**Data Science**

OST

Eastern Switzerland University of Applied Sciences | Project Theses 2024 | Master of Science in Engineering | Technik und IT