

# Large Language Models (LLM)

PA154 Language Modeling (12.2)

Pavel Rychlý

pary@fi.muni.cz

May 14, 2024

## prompt engineering

- simple prompt:  
Q&A  
Q: {question}  
A:
- general knowledge:  
Generate some knowledge about the concepts in the input.  
Input: {question} Knowledge:
- task specific:  
If {premise} is true, is it also true that {hypothesis}? ||| {entailed}.

## Generated trained data

- there are never enough text with the right Q&As
- generate data from a pattern especially for chain-of-thought
- variation of variables (numbers in math, ...)
- generated using LLM

## From LM to Chat

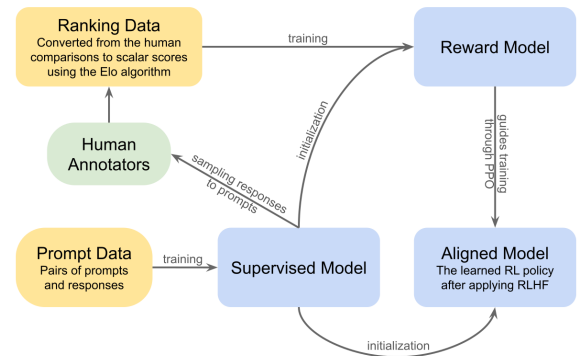
- LM generates prompt continuation
- prompt engineering
- fine tune on chat texts
- some Q&A/chat texts are in training data

## Chain-of-thought

- direct answer is not correct for complex questions
- solve a problem in steps
- Q: {question}  
A: Let's think step by step.

## RLHF

Reinforcement Learning from Human Feedback



## Alignment

- annotated data by humans
- training using RLHF
- eliminate toxicity, bias

## Foundation models

- LLM without fine tuning
- can be used to adapt on a new domain/language
- fine-tuned on a specific task
  - chat
  - question answering
  - summarization

## LoRA

- Low-Rank Adaptation of Large Language Models
- fine-tune only a small fraction of parameters
- usually only attention matrices

