

BEACHELOR THESIS
KOGNITIVE INFORMATIK

SYSTEMATIC
GENERATION OF
UNIT-TESTS WITH
DIFFERENT LLMS
AND PROMPT-DESIGN

RICHARD PAMIES

University of Bielefeld

SUPERVISED BY
PROF. DR.-ING. BENJAMIN PAASSEN,
M.SC. JESPER DANNAT

JULY 15, 2024

CONTENTS

Contents	ii
1 Introduction	1
2 Background and related work	2
2.1 Unit-tests	2
2.2 LLMs	2
2.3 related work	2
3 Methods	3
3.1 high level description of my approach	3
3.2 provide overview diagram illustrating my steps	3
4 Implementation	4
4.1 for each of the steps how it was implemented in the code	4
5 Experiments	5
5.1 what i want to show	5
5.2 how i want to show it	5
5.3 baselines and assumptions	5
5.4 what was my initial assumption	5
5.5 Table for experiments	5
5.6 Discussion of the results	5
6 Conclusion	6
6.1 re-telling of the paper	6
6.2 high-level overview of experimental results	6
6.3 Limitations	6
6.4 what future work could accomplish	6
6.5 final words	6
7 latex features	7

INTRODUCTION

Unit-test are really important and useful. But most often those tests don't get written because of laziness, so the automatic generation of those tests arise as a possible solution.

Large Language Models (LLMs) offer for the the possibility of creating useful unit-tests. But those generated test are flawed because those LLMs dont work with reason, but instead the statistically most probable answer.

Prior work has tried under different circumstances to solve this problem.[paper 1,2,3].

My goal in this paper is to work out if prompt design by the LLM makes sense, and which LLMs, paired with which prompt design, provide the best answers.

BACKGROUND AND RELATED WORK

2.1 UNIT-TESTS

Unit-tests are tests written for already existing code to make sure the implemented functionalities still work after further work on the code is done.

The usefulness of unit-test is not controversial. Rather the time it takes to write them and the metrics used to test their usefulness.

Code coverage and statement coverage is often mentioned when talking about validating unit-tests. –Talking about the mechanisms of both and why they are not enough–

Here come other metrics i use in my work and their pros and cons:
mutation coverage
assert count and variety

2.2 LLMS

An llm is only as good as its architecture and its training data. Generating unit-tests for mainstream code might be sufficient. But it is really easy to write code which covers niche subjects. And in those cases the LLM would need critical thinking, understanding and reasoning to fully understand the code and write useful and encompassing unit-tests.

2.2.1 *GPT*

2.3 RELATED WORK

METHODS

3.1 HIGH LEVEL DESCRIPTION OF MY APPROACH

I used the API-endpoints of different LLMs to send a structured prompt to those and receive code which should be the unit-test. Then i execute the unit

The sequence of my approach to generate these unit-tests is:

1. refining a prompt by the LLM
2. sending a prompt to the LLM
3. receiving solution and extracting the unit-test
4. executing the unit-test
5. re-sending the prompt with appended error msg if the unit-test fails
6. re-sending the unit-test with the instruction to remove the defective line

The structured prompt could look like this:

Lemma Let the coefficients of the polynomial

$$a_0 + a_1x + a_2x^2 + \cdots + a_{m-1}x^{m-1} + x^m$$

be integers. Then any real root of the polynomial is either integral or irrational.

This high level description what it should accomplish detailed description for reproduction for experienced user use illustrating diagrams as much as possible (no source code, but perhaps pseudo code) provide a good reason for each design choice(source or math) use a running example to illustrate the steps

3.2 PROVIDE OVERVIEW DIAGRAM ILLUSTRATING MY STEPS

IMPLEMENTATION

4.1 FOR EACH OF THE STEPS HOW IT WAS IMPLEMENTED IN THE CODE

EXPERIMENTS

5.1 WHAT I WANT TO SHOW

5.2 HOW I WANT TO SHOW IT

5.3 BASELINES AND ASSUMPTIONS

5.4 WHAT WAS MY INITIAL ASSUMPTION

5.5 TABLE FOR EXPERIMENTS

5.6 DISCUSSION OF THE RESULTS

CONCLUSION

6.1 RE-TELLING OF THE PAPER

6.2 HIGH-LEVEL OVERVIEW OF EXPERIMENTAL RESULTS

6.3 LIMITATIONS

6.4 WHAT FUTURE WORK COULD ACCOMPLISH

6.5 FINAL WORDS

LATEX FEATURES

sparse¹

Did you know? The discrete fourier transformation shown in equation 7.1 is the backbone of the modern information-society.

$$f_m \tag{7.1}$$

For words of science, see [**botsch2010polygon**]. Unfortunately, that book has nothing about [wolves](#).

¹ This is a footnote