![Aalto University logo]

**Aalto-yliopisto**
**Aalto-universitetet**
**Aalto University**

ELEC-E7852
**Computational Interaction and Design**

# Assignment A4b

*Email Assistant*

**Aitor Urruticoechea Puig**
aitor.urruticoecheapuig@aalto.fi
Student N°101444219
November 2024

## Acknowledgement

## Contents

## List of Figures

## 1   Introduction and Approach

In this assignment, the object of study is alignment. As AI systems become more complex, it is im‑perative that, as engineers, we learn how to adapt their behaviour to human values and personalized goals. This will be explored this time around in a simplified manner by developing a generative email assistant. This assistant should be able to adapt to the writing style of the user, and identify the context in which it is writing (topic, recipient, etc), in order to generate a coherent output which is aligned with what would be expected of the user. To do so, a training dataset of 15 emails is provided from a single user, which include different exchanges with different people.

The proposed approach uses this limited dataset to construct a "knowledge database" which, after being curated during "training", can be used as a standalone source for a Large Language Model (LLM) to generate emails. This "knowledge database" is to be concise, so as to be useful for a LLM to generate emails without being too wordy (to avoid going over the token limit these models tend to have) nor convoluted (to avoid echo‑chambers or prompt‑breaking "knowledge"). To curate this knowledge, the idea is to have the LLM iterate over its creation as more of the training emails are fed to it. See the proposed training loop:

```
1  conversation_history = [
2      {
3          "role": "system",
4          "content": """You are a master mail writer. You get to know your users as well as
           ↪  possible, to the point where you understand in-depth their way of expressing
           ↪  themselves, and how this are adapted to each situation within their social
           ↪  circles.
5      This allows you to write mail in their place, perfectly adapted in character and
           ↪  style to the way they write to that concrete recipient."""
6      },
7  ]
```

```
1  def generate_context(dataset:list, datapath:str):
2      knowledge = ""
3      for email in dataset:
4          print("Email number: " + str(dataset.index(email)+1) + " out of " +
           ↪  str(len(dataset)), end="\r")
5          training_conversation = [
6              conversation_history[0]
7          ,{
8              "role": "user",
9              "content": """You are getting to know a person better in order to write their
               ↪  emails. So far you have gathered this data : \n"""
10             + knowledge + """\n Use the following example email, written by your user, to
               ↪  further curate your knowledge database:\n"""
11             + email + """\n\n Return the updated knowledge database; which should always
               ↪  have the following format:
12             - Name : [name]
13             - Relationships: [list of relationships, their nature, and other observations]
14             - Writing style: [list of writing style characteristics, such as formality,
               ↪  humor, etc, and how they change depending on context]
15             - Other: [any other relevant information like idiosyncrasies, interests,
               ↪  typical expressions that they use, etc]"""
16         }]
17         knowledge = generate_response(training_conversation)
18     print("Training completed! Read " + str(len(dataset)) + " emails.")
19     with open(datapath, "w") as f:
20         f.write(knowledge)
21     return knowledge
```

This is heavily reliant in prompt engineering, yet as the knowledge generation for each email happens independently (as the conversation is a "new" one every time) we allow the model to correct itself: previously generated "knowledge" is only given as a guideline, so corrections are expected with the new information coming from this new email. In a way, one can think of this approach of the same step performing the user model and task model described in the assignment instructions at the same time (as it itself corrects its assumptions on the information given). This is added to the relative freedom for format in the curation of this "knowledge". Using the proposed approach, we are only making sure key areas are analysed and stored, but both in the Other container and in the contents themselves, the LLM can choose how to fill them up according to the nature of the provided emails. The chosen key areas for analysis are Relationships (their nature, and other observations), Writing style (characteristics such as formality, humour, etc, and how they change depending on context), and Other (including any other relevant information like idiosyncrasies, interests, typical expressions that they use, etc].

All in all, this "knowledge database" is used for complementing the prompting for the actual email generation, which by itself is uncomplicated. By including this compendium of traits from the user in the prompt, it is expected that the LLM will be able to generate emails more in line with what would be expected from the user.

```
1  def email_generator(recipient, subject, knowledge):
2      email = generate_response([
3      conversation_history[0],{
4          "role": "user",
5          "content": """You are writing an email to """ + recipient + """ as your user. The
           ↪  subject of the email is: """ + subject + """.
6          The email should be written in the style of your user, using the knowledge you have
           ↪  gathered about them.
7          Here is the knowledge you have gathered so far: \n""" + knowledge + """\n\n Write
           ↪  the email:"""
8      }])
9      return email
```

## 2 Results and Conclusions

This approach, while simple, has yielded very compelling results. The first interesting thing to notice is how the "knowledge database" has taken form. Models like GPT‑4 (the one used) are remarkably good at the task of extracting insights, and this is proven once more here. Not only things like typical spelling errors or frequently used expressions are recorded, but also tones and intentions (like Tom trying to impress Bob, or how close he is to his mother). This proves that the approach is at least useful in generating a complex enough profile from the user, despite it consisting of 15 emails. More interestingly, this approach should be easy to generalize beyond this case, as a "knowledge database" like this will still work for other relation types or contexts.

The generated emails seem to also prove the hypothesis that, given reliable‑enough background information, the alignment problem facing email generation is, to an extent, solvable. For an example, see Figure 1 with three emails with the same subject ("Coffee on Friday") yet addressed to the three different potential recipients included in the database. The LLM is able to discern from the "knowledge database" the type of relationship the user has with each of the recipients, and creates a compelling story in line with it; to the point that I had to double‑check the original "ground truth" for some references the LLM was making to ensure they were not hallucinated (Tom does indeed mention his Dad in passing, though only as a secondary thing).

In short, the approach where the two functions (user and task) are combined in the training by having the model generate a reusable database with profile insights has proven quite successful. More testing would be needed with a more diverse dataset to solidify these findings, yet it is clear that as an email generator its alignment is light years ahead of what would be achieved with traditional/simple prompt engineering or even by having the LLM scheme through a couple of example mails; without having to sacrifice extreme amounts of tokens in feeding extra context.

(a) Example mail to "Mom"



(b) Example mail to "Emilia".



(c) Example mail to "Bob".

Figure 1: Example of generated emails, with the same subject, to the three possible recipients.