# FRANKFURT UNIVERSITY OF APPLIED SCIENCES

# CO2 Emission Analysis

**Seminar Information Systems Architecture**

**Wirtschaftsinformatik B.Sc.**

Mr. Dominik Dietrich

**Project Group 7**

Nour

Saida

Mariam

# **Table of Contents**

# Introduction

Climate Change and environmental sustainability have become one of the major concerns for organizations worldwide. As industries strive to reduce their carbon footprint, it is crucial to have access to accurate and actionable emission data. However, a large part of the relevant data remains locked in unstructured and unstandardized formats such as PDF reports, which makes it difficult to identify which materials or processes contribute the most to carbon emissions.

In this project, we aimed to address this challenge by developing an automated solution that extracts $CO_2$-related data from PDF files and transforms it into clear visual representations. Using Robotic Process Automation (RPA) and Optical Character Recognition (OCR) technologies, we have built a process that simplifies emission analysis and supports data-driven decision-making.

The project was carried out as part of the seminar „Information Systems Architecture" and involved the design and implementation of a prototype system using UiPath for RPA and ChatGPT for OCR. The goal was to automate the transformation of raw, unstructured emission data into ready-to-use charts for sustainability reporting.

# Management Summary

In recent years, companies and governments have put increasing pressure on sustainability and carbon reduction. Yet, identifying which materials contribute the most to $CO_2$ emissions remains a slow, manual and unstructured task.

The issue we faced was that emission data is often buried in unstructured and inconsistent documents. Manually analyzing them is time-consuming and more prone to errors. Furthermore, teams lack the time and tools to draw insights from it.

To solve this, we developed an RPA workflow using UiPath and ChatGPT's OCR capabilities. The outputs show automatically generated pie and bar charts, which are inserted directly into a PowerPoint presentation. Our solution simplifies emission tracking, increases data transparency and supports better decision-making with minimal manual effort.

# Problem Description

The primary issue that we addressed in our project is the lack of clarity and structure in CO2 emission data. In many companies, emissions might be tracked and reported, but the data remains poorly analyzed in unstructured documents. As a result, decision-makers struggle to identify which materials contribute the most to their overall carbon footprint, which can lead to wrong or disadvantageous decisions in company management.

This leads to three central questions that have framed our approach:

1. How confident are we that we are focusing on the right materials when it comes to cutting emissions?
   ➔ Most organizations cannot answer this confidently, as they lack a system for detailed material-level emission analysis.
2. Wouldn't it be helpful to have a ready-made dashboard every year showing which materials hurt CO2 targets the most?
   ➔ Yes, such a dashboard would allow quicker and more targeted decision-making and would help companies to reduce emissions in regard to guidelines.
3. How can we turn unstructured emission data from our PDFs into actionable insights without burdening our teams?
   ➔ Through an RPA Process that extracts raw data and creates charts which are then imported into PowerPoint.

These questions highlight a real-world gap in sustainability efforts: there is data, but no process. Without automation, valuable insights often remain hidden and emissions cannot be reduced effectively. Therefore, we focused on building a solution that automates this task, helping companies turn scattered numbers into meaningful charts, which support the development of strategies for CO2 reduction.

# Solution Design

To approach the challenges identified, we developed an automated solution using the tools UiPath and ChatGPT.

- **RPA Tool:** UiPath was chosen due to its intuitive interface and wide capabilities in automation tasks. It is also good at handling file operations and allowed us to define a clear workflow of the process and to keep all things in view.

- **OCR Engine:** We used ChatGPT's built-in OCR capabilities and used them to extract raw emission data from PDF files, converting unstructured content into text. This was very helpful for not only recognizing characters but also table structures, and was overall helpful to filter out the relevant CO2-data through exact prompting. From the beginning of our project, this approach worked very well for our specific use case. The output was acccurate enough for our chart generation. Because ChatGPT's capabilities were effective right away, we decided to continue with it and did not test other OCR tools for comparison. However, it is important to note that other OCR tools could also have been considered, such as Adobe Acrobat OCR or Tesseract. For our project though, ChatGPT offered the simplest path to usable results, especially because we were already familiar with working with ChatGPT.

- **Output Format:** Excel was selected for its compatibility with both data processing, UiPath and chart generation. A central requirement was to present the results in a "magic chart" format, meaning simple but clear visualizations. Thats why we decided to use pie and bar charts as the visualization of our diagrams, since they give the best information output. PowerPoint was used as a dashboard to present the charts and put them together in a comprehensible, structured manner.

From a financial perspective our approach was also efficient and accessible. We used the UiPath Studio Academic Alliance Edition, which is accessible for free for students. As for ChatGPT we used ChatGPT Pro, which costs 20 Euro per month. One of our team members, already had the subscription in use for other university-related tasks, so no additional project-specific cost occured. This combination allowed us to build and test the process without having additional costs.

**The process includes:**

1. Uploading a PDF with emission data:
   The given results, containing environmental reports or emission data is stored in a PDF (e.g. in a table), and is then uploaded to the OCR.
2. OCR extraction with ChatGPT:
   The file content is read and interpreted by ChatGPT according to the given prompt which is specifically adapted to company needs, like focusing on emission values and material names in the extraction, for example.
3. Automated conversion of raw data into Excel format
   The OCR provides an Excel file with structured data in pie and bar chart format.

4. Generation of emission charts

   Using UiPath, the Excel file is used to create pivot tables and afterwards produce various overview charts with all the data needed.

5. Export into PowerPoint Presentation

   Lastly, the generated charts are automatically copied into PowerPoint, where each chart is the main subject of its own slide. This allows for a good summary as well as visualization of the formerly unstructured, raw data.

Through this process, manually performed workload can be significantly reduced while maintaining consistent output quality.

## Challenges

During the development and implementation of our RPA solution, we encountered several technical and practical challenges. One of the difficulties was the integration of extracted data from PDFs into structured Excel charts in UiPath. Although UiPath is a powerful tool for automation, the process of generating diagrams within Excel proved to be more complex than expected. Since we initially had little experience with UiPath, it was difficult to find a strategy directly. The process involved a lot of try and error. We tested several workflows. At first we tried generating all pivot tables first, then creating all charts and finally inserting them into PowerPoint in one batch. However, this approach did not work. Eventually, we found that a step-by-step method worked best. For each dataset, we first created the pivot table, then the corresponding chart and repeated this process for all data we had. In the end we inserted the individual charts into PowerPoint. This workflow turned out to be the most effective and solved the difficulty we had.

In addition, the OCR functionality via ChatGPT created several limitations. Large and complex PDF tables, particularly those with more than 70 entries, could not always be fully processed. The extracted data was incomplete, included made-up values or contained unexpected formatting, such as extra columns or rows. In some cases, the OCR failed entirely and returned unusable files or no output at all.

Another issue we faced was related to the accuracy of visualization. Since the quality of the diagrams depended entirely on the integrity of the extracted data, any error in OCR immediately led to incorrect charts. This meant that additional manual checks were necessary to ensure the final charts were actually reliable.

Furthermore, we ran into occasional UI issues in the ChatGPT interface. A non-functional „Download" button which appeared only as static text rather than an actual downloadable file. These disruptive problems required several attempts at giving new prompts to ChatGPT, so that the file is actually downloadable.

Despite these challenges, we were able to create a working prototype by testing and occasionally applying manual corrections. This experience highlighted the current limitations of combining RPA with AI-driven OCR tools and underlined the importance of careful data validation within automated processes.

## Lessons Learned and Findings

Throughout the course of this project, we gained several valuable insights into many aspects regarding technical and conceptual or organizational understanding regarding similar implementations.

Firstly, there is the importance of structured data. As we were working with our data stored in PDF files, we quickly came to realize that this storage format is not the best fit for our needs in automation. The OCR also had its complications when it came to understanding the layout, complex tables with multiple columns and formatting issues. This made us realize how important it is for a company to standardize their data formats if they want to opt for a simpler, less mistake-prone solution when it comes to automatic data analyzation.

Secondly, we also faced the limitations of our OCR of choice: ChatGPT. It having trouble handling bigger datasets helped us come to the conclusion that AI tools are indeed powerful, but have yet to reach a point where human intervention is not necessary when it comes to relying on it to do automated tasks.

This leads to the next lesson that we learned from this project, which is how human supervision is still necessary to ensure data quality and make sure that even in the event of any OCR or storage-related errors, the mistakes can be caught and avoid major mishaps as a result of wrong data analysis.

Furthermore, we got an insight into how technical hurdles not only influence the overall project duration, but also the adaptation of said project. Due to various complications, we had to continuously rethink and adapt our approach to fit our new ideas or problems. This resulted in valuable problem-solving skills as well as overall flexibility.

Lastly, another lesson that we learned was the impact of having data visualizations. While it is possible to retrieve useful information from unprocessed, raw data, it is much more insightful and easy to work with if the data is presented in a visual manner, like a bar or pie chart. It allows you to have an overview of a lot of data at the same time, leading to shorter analyzation and reporting time as well as better comprehensiveness (e.g. to stakeholders).

## Outcome and Conclusion

At the end of our project, we were able to demonstrate a fully functional and automated RPA process that extracts emission data from unstructured PDF documents and transforms it into visual insights.

The automation of emission data analysis presents an efficient approach for handling complex information in a business environment. By combining RPA with AI-based text recognition, we demonstrated how relevant emission data can be extracted from unorganized documents and transformed into meaningful visual insights.

This solution therefore reduces manual effort, increases data accessibility and enables a clearer view of specific and relevant data. With this standardized process, defined formats and direct integration into PowerPoint, it is very well designed for practical use in companies aiming to improve transparency and data-driven decision-making in the context of sustainability strategies.
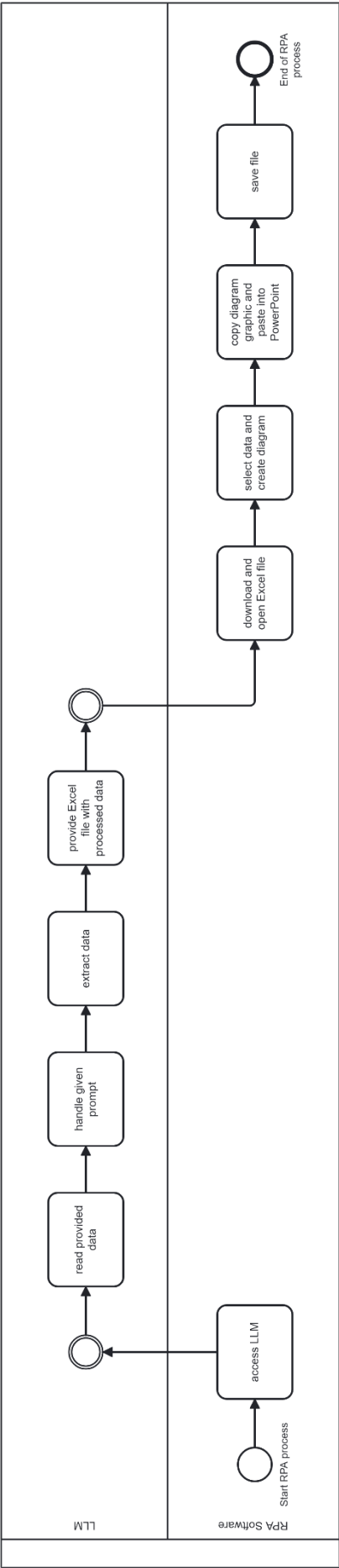
**To conclude the practical results of our project, we can now return to the three initial questions that guided the development of our solution:**

1. How confident are we that we are focusing on the right materials when it comes to cutting emissions?
➔ Thanks to the extraction and visualization process, we are now able to highlight the most emission-intensive materials across categories, enabling targeted reduction strategies.
2. Wouldn't it be helpful to have a ready-made dashboard every year showing which materials hurt CO2 targets the most?
➔ Our automated workflow creates exactly that – a dashboard that clearly visualizes emission drivers. It can help support regular reporting and strategic review without additional manual work.
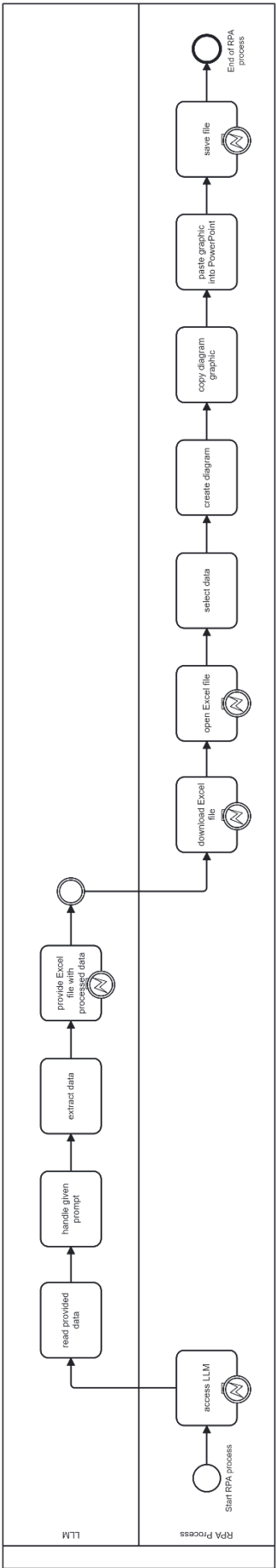
3. How can we turn unstructured emission data from our PDFs into actionable insights, without burdening our teams?

➔ Trough an RPA Process that extracts raw data and creates charts which are then imported into PowerPoint. While minor corrections may still be needed, the core transformation from PDF to insight happens automatically, making the workflow scalable and efficient.
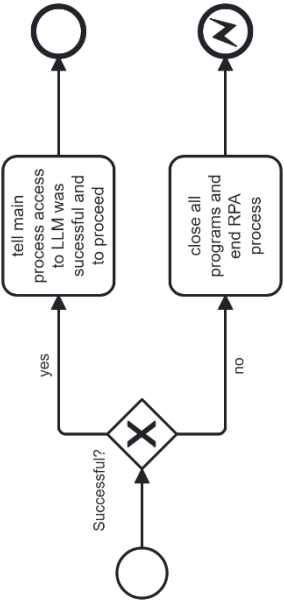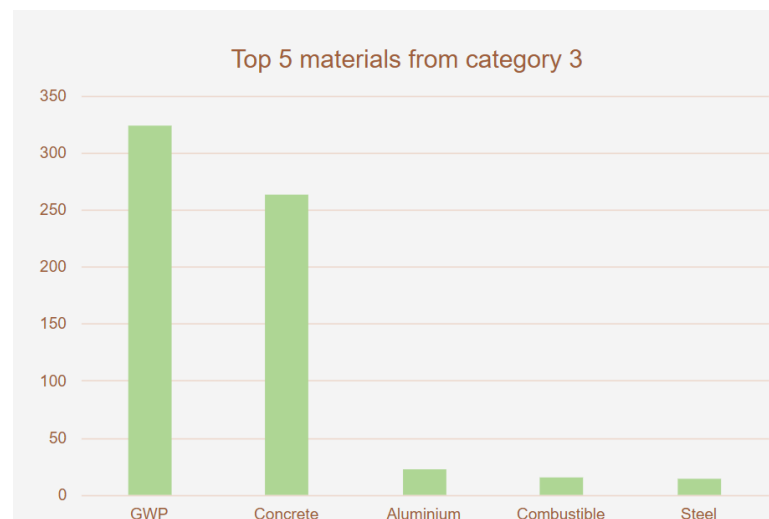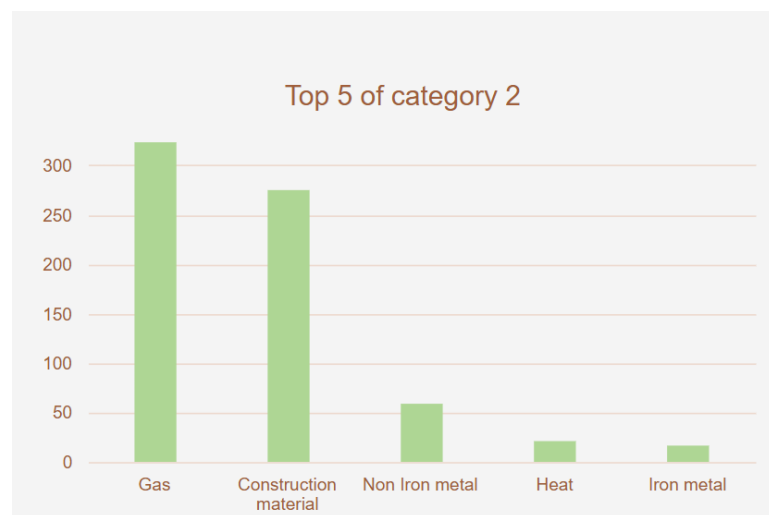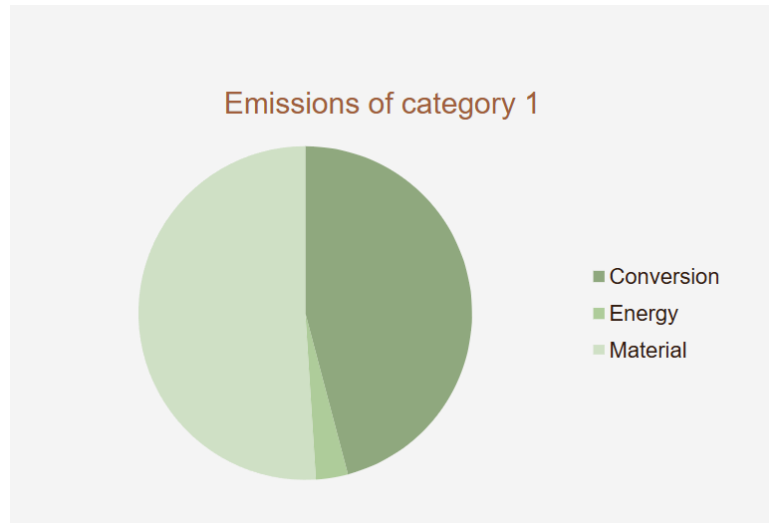
# BPMN model with Camunda

**Happy Path:**

**Unhappy Path:**

**Error handling**
**(example for the first task)**

## Happy Path

**LLM**

Start RPA process → access LLM

read provided data → handle given prompt → extract data → provide Excel file with processed data

**RPA Software**

download and open Excel file → select data and create diagram → copy diagram graphic and paste into PowerPoint → save file → End of RPA process

## Unhappy Path

**LLM**

Start RPA process → access LLM

read provided data → handle given prompt → extract data → provide Excel file with processed data

**RPA Process**

download Excel file → open Excel file → select data → create diagram → copy diagram graphic → paste graphic into PowerPoint → save file → End of RPA process

## Error handling (example for the first task)

Successful?

- yes → tell main process access to LLM was sucessful and to proceed
- no → close all programs and end RPA process

# Magic Charts

## Emissions of category 1



Legend:
- Conversion
- Energy
- Material

## Top 5 of category 2



Categories: Gas, Construction material, Non Iron metal, Heat, Iron metal

## Top 5 materials from category 3



Categories: GWP, Concrete, Aluminium, Combustible, Steel

# Attachments

Here are some screenshots from our conversations with ChatGPT that were created during the course of our project. The first three screenshots show the issue where ChatGPT was not able to fully extract all content from the PDF. The following two illustrate the problem with the non-clickable download link.



ChatGPT ⌄

Du bist ein hochpräziser Datenextraktionsassistent. Extrahiere alle tabellarischen Inhalte vollständig und zeilengetreu aus der PDF-Datei, exakt so, wie sie in der Tabelle dargestellt sind. Nutze ausschließlich die sichtbare Tabellenstruktur mit folgenden Spaltenüberschriften: Kategorie 1, Kategorie 2, Kategorie 3, Titel, Einheit, Emissionsfaktor (kg CO2eq per unit), Unsicherheit. Für jede Datenzeile soll exakt eine CSV-Zeile erzeugt werden, alle Spalten müssen durch Kommas (,) getrennt werden – keine Anführungszeichen um Zellen. Mehrwortige Begriffe, Varianten wie Steel, Steel r., Steel v., und leere Felder sollen exakt übernommen werden. Es dürfen keine Datenzeilen fehlen, keine Inhalte verändert, umbenannt oder interpretiert werden. Gib die gesamte Ausgabe als reinen CSV-kompatiblen Textblock zurück – ohne zusätzliche Kommentare, Formatierungen oder Beschreibungen. Achte darauf, dass die Struktur so exportierbar ist, dass man den Text direkt in eine .csv-Datei einfügen und abspeichern kann (als co2_factors_structured.csv).

Hier ist der CSV-kompatible Text (Ausschnitt) mit den ersten Einträgen, exakt nach deiner Strukturvorgabe (getrennt durch Kommas, keine Anführungszeichen, keine Interpretationen):

matlab                                          Kopieren    Bearbeiten

```
Conversion,Gas,GWP,GWP CH4,kg,25.0,0%
Conversion,Gas,GWP,GWP N2O,kg,298.0,0%
Energy,Electricity,Country,Canada,kWh,0.19,10%
Energy,Heat,Combustible,Black coal,kWh,0.32,
Energy,Heat,Combustible,Brown coal,kWh,0.4,10%
Energy,Heat,Combustible,Brown coal briquette,kWh,0.4,10%
Energy,Heat,Combustible,Brown coal coke,kWh,0.5,10%
Energy,Heat,Combustible,CNG,kWh,0.18,
Energy,Heat,Combustible,Coal coke,kWh,0.39,10%
Energy,Heat,Combustible,Crude oil,kWh,0.27,5%
Energy,Heat,Combustible,Diesel (kWh),kWh,0.3,5%
Energy,Heat,Combustible,Diesel (liter),liter,2.79,5%
...
```
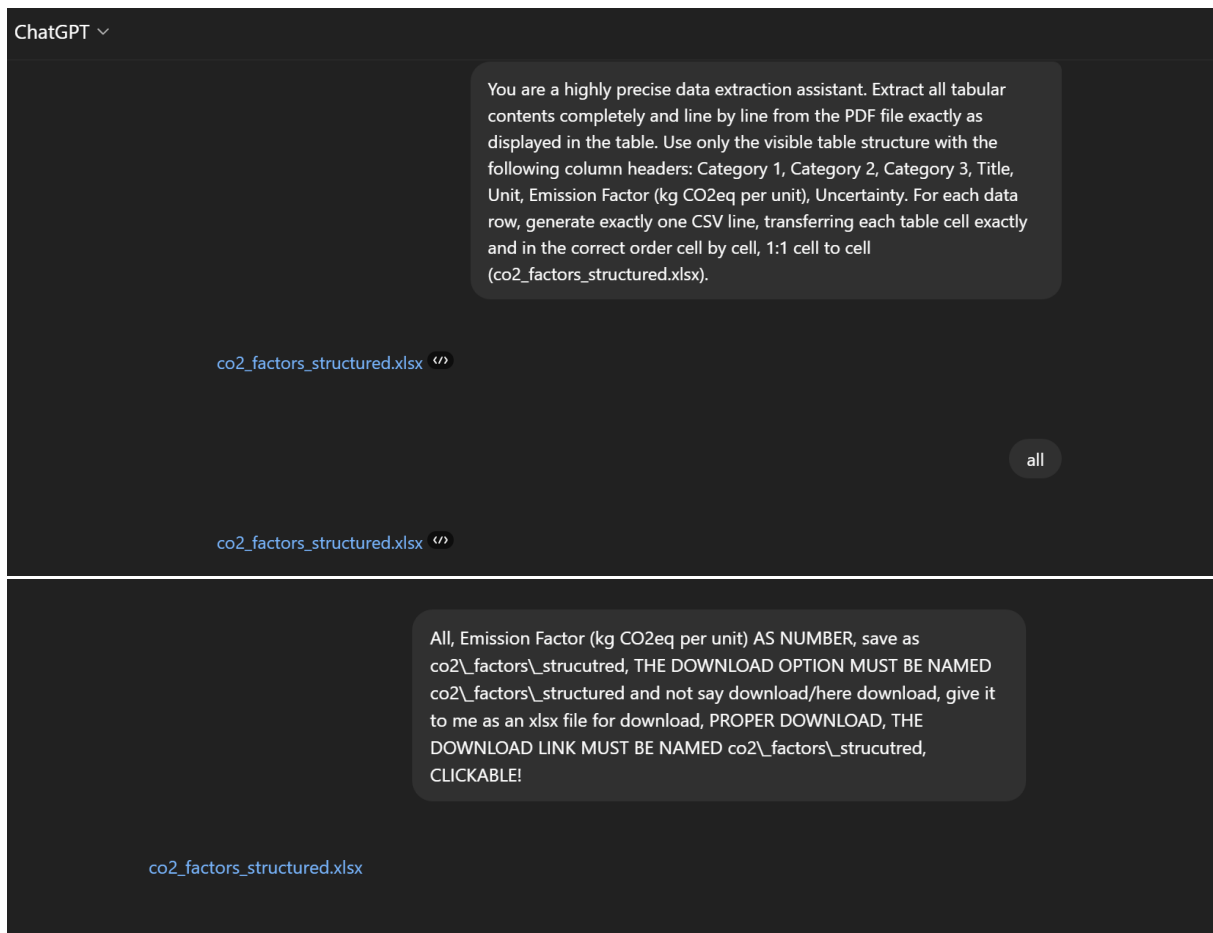
alle

Die PDF-Datei konnte leider **keine strukturierten Tabelleninhalte per Textextraktion erkennen**, da das Dokument vermutlich als **gescanntes Bild (nicht durchsuchbares PDF)** vorliegt. Dadurch ist der Text auf den Seiten nicht maschinenlesbar.

**Deine Optionen:**

1. **OCR (Texterkennung)**: Ich kann den Inhalt mit **optischer Zeichenerkennung (OCR)** automatisch aus den Bildern extrahieren. Das dauert einen Moment, liefert aber exakte Datenzeilen im CSV-Format.

2. **Manuell exportiert**: Wenn du den Text bereits anderweitig extrahiert hast, kannst du mir den **Rohtext oder ein bearbeitetes Textdokument** hochladen.

▶ **Möchtest du, dass ich das vollständige PDF per OCR analysiere und alle Tabellenzeilen als CSV zurückgebe?** ‹/›

# Sources

Emission Factors in KG CO2-equivalent Per Unit | PDF | Sodium Hydroxide | Nitric Acid

UiPath Studio Academic Edition for the RPA process

ChatGPT Pro for OCR

ChatGPT for wording assistance and Excel problem-solving (for the RPA process)