

Duale Hochschule Baden-Württemberg Mannheim

Projektberichtarbeit
Intel Core I7 Prozessoren

Studiengang Informatik
Studienrichtung Angewandte Informatik

Verfasser(in):	Moritz Werr, Phil Richter, Max Stege
Matrikelnummer:	5401527
Matrikelnummer:	4164342
Matrikelnummer:	<Ihre Martikelnummer>
Kurs:	TINF22AI2
Dozent:	Dr. Frank Schulz
Bearbeitungszeitraum:	01.10.2024 – 01.12.2024

Ehrenwörtliche Erklärung

Ich versichere hiermit, dass ich die vorliegende Arbeit mit dem Titel "*Intel Core I7 Prozessoren*" selbstständig verfasst und keine anderen als die angegebenen Quellen und Hilfsmittel benutzt habe. Ich versichere zudem, dass die eingereichte elektronische Fassung mit der gedruckten Fassung übereinstimmt.

Quelltextverzeichnis

2.1	<code>clean.sh</code>	4
3.1	<code>database.R</code>	6
5.1	load data from the database	13
5.2	Funktion <code>extraxct_generation</code>	14
5.3	Code für den Scatterplot in Figur 5.1	15

1 Einleitung

Dieses Kapitel enthält die Einleitung mit ihren verschiedenen Abschnitten/Sections und Unterabschnitten.

1.1 Beispiel Abschnitt: \LaTeX -Installation

Zur Verwendung von \LaTeX -Installation einer Distribution z. B. TeXLive, MikTeX etc. sowie eines Editors z. B. TeXStudio, TeXnicCenter etc. notwendig.

Installieren Sie zunächst die Distribution und anschließend den Editor. Beim ersten Start des Editors öffnet sich ein Konfigurationsassistent, der zunächst nach dem Pfad der installierten Distribution fragt.

Nach der Installation können Einstellungen z. B. für einen PostScript-Viewer gemacht werden. Dieser Schritt kann ohne Weiteres übersprungen werden. Entscheidend sind die Einstellungen für den pdf-Viewer.

Jetzt kann \LaTeX verwendet werden. Um die Ausgabe eines Dokumentes zu erzeugen, muss das Dokument kompiliert werden (Ausgabe > Aktives Dokument > Erstellen und betrachten).

1.1.1 Beispiel Unterabschnitt: Aufbau eines \LaTeX -Dokuments

Ein \LaTeX -Dokument besteht in der Regel aus folgenden Komponenten:

- Präambel
- Titelseite
- Textteil

1.1.2 Beispiel Unterabschnitt auf zweiter Ebene: Präambel

In der Präambel werden global die Einstellungen für das gesamte Dokument definiert. Hierbei können z. B. die Seitenränder, der Zeilenabstand oder auch die Sprache für die Silbentrennung festgelegt werden. In der ersten Zeile eines jeden Dokumentes wird dabei immer die zu verwendende Klasse festgelegt. Standardmäßig kann hier die Artikel-Klasse gewählt werden:

```
\documentclass[12pt,titlepage]{article}
```

In den eckigen Klammern wird dabei u.a. die Standardschriftgröße für das gesamte Dokument festgelegt.

Außerdem werden in der Präambel die für das Dokument benötigten Pakete festgelegt. Gebräuchlich sind vor allem folgende Pakete:

- `\usepackage[ngerman]{babel}`
- `\usepackage[latin1]{inputenc}`
- `\usepackage{color}`
- `\usepackage[a4paper]{geometry}`
- `\usepackage{amssymb}`
- `\usepackage{amsthm}`
- `\usepackage{graphicx}`

Im vorliegenden Fall werden die Pakete in der Konfigurationsdatei `config.tex` festgelegt, deren Inhalt durch `\input{config}` in das Hauptdokument `master.tex` inkludiert wird.

Beispiel Unterabschnitt auf zweiter Ebene: Titelseite

Nachdem die Dokumenten-Klasse und die zu verwendenden Pakete festgelegt worden sind, folgt die Titelseite. Da die Titelseite bereits Teil des eigentlichen Dokuments ist, muss ihr unbedingt der Befehl `\begin{document}` vorausgehen. Am Ende des Dokuments sollte der Befehl `\end{document}` gesetzt werden. Alles was nach diesem Befehl steht, wird vom Compiler nicht mehr beachtet.

1.2 Noch ein Beispiel-Abschnitt

Der Textteil beinhaltet nun den eigentlichen Text des Dokuments.

2 Vorbereitung der Daten (Moritz Werr)

Intel bietet auf ihrer Webseite die Möglichkeit unterschiedliche Prozessoren miteinander zu vergleichen. Es werden alle Desktop-Prozessoren der Intel Core I7-Reihe an Prozessoren ausgewählt, damit es eine große Anzahl an Datensätzen gibt und auch die Spannweite zwischen gewissen Werten groß genug ist, dass Unterschiede leichter zu erkennen sind. Z.B. ist die Kernanzahl pro Prozessor seit der 4. Generation stark gestiegen, wohingegen die Anzahl an Prozessorkernen bei Intel Core I3-Modellen nicht so stark angestiegen sind. Die Daten werden direkt von Intels Webseite bezogen. Da Intel keine öffentliche API anbietet, um deren CPU-Spezifikationen zu vergleichen und zu exportieren, muss die CSV-Datei selber heruntergeladen werden. Bevor die Daten aus der CSV ausgelesen werden können, müssen diese noch etwas aufbereitet werden. Wie in Abbildung 2.1 von der Herstellerseite zu sehen ist, sind hier leider die Prozessormodelle die Spalten anstatt der Zeilen.

Bevor die CSV bereinigen können, muss also die Matrix der CSV noch einmal transponiert werden. Sonst wären die einzelnen Prozessoren die Spalten in der Datenbank, in der die bereinigte Daten der CSV abgespeichert werden. Weiterhin sind nicht alle Spalten in der CSV interessant für die Auswertung der Daten. Viele der Einträge in der CSV listen nur Features oder Prozessorerweiterungen auf, die nicht einfach zu vergleichen sind oder sich nur schwierig visualisieren lassen. Weiterhin sind manche Einträge bei älteren Prozessoren komplett leer, weswegen die Datenmenge in diesen Spalten sehr gering ist. Daher werden nur ausgewählte Spalten gespeichert, die für die Fragestellung interessant sind und auch genügend Datenpunkte haben.

56 Intel® Produkte vergleichen

Vergleich exportieren

Intel® Core™ i7-8700K Processor

Intel® Core™ i7-8700 Processor

Intel® Core™ i7-8700T Processor

Intel® Core™ i7-12700 Processor

Intel® Core™ i7-12700F Processor

Intel® Core™ i7-12700K Processor

Hauptdaten

Produktsortiment	Intel® Core™ i7 Prozessoren der 8. Generation	Intel® Core™ i7 Prozessoren der 8. Generation	Intel® Core™ i7 Prozessoren der 8. Generation	Intel® Core™ i7 Prozessoren der 12. Generation	Intel® Core™ i7 Prozessoren der 12. Generation	Intel® Core™ i7 Prozessoren der 12. Generation
Vertikales Segment	Desktop	Desktop	Desktop	Desktop	Desktop	Desktop
Prozessornummer	i7-8700K	i7-8700	i7-8700T	i7-12700	i7-12700F	i7-12700K
Lithographie	14 nm	14 nm	14 nm	Intel 7	Intel 7	Intel 7
Empfohlener Kundenpreis	\$359.00	\$303.00	\$303.00	\$373.00-\$383.00	\$345.00-\$355.00	\$450.00

CPU-Spezifikationen

Anzahl der Kerne	6	6	6	12	12	12
------------------	---	---	---	----	----	----

Figure 2.1: Webseite zum Produktvergleich von Intel-Prozessoren
Quelle: [5]

Aufgrund der mächtigen Standardwerkzeugen eines POSIX-kompatiblen Betriebssystems wie Linux, wurde `awk` verwendet, um die Spalten in Zeilen zu konvertieren bzw. die Zeilen zu Spalten zu konvertieren. `awk` ist ein Werkzeug zur Erkennung und Verarbeitung von Textmustern [2]. `awk` arbeitet nicht nur mit Zeilen, wie z.B. `grep` oder `sed`, sondern kann auch in einzelnen Spalten bzw. Feldern filtern [2]. Normalerweise sind die Standardtrennzeichen von Spalten für `awk` Leerzeichen und Tabulatoren. Es ist aber möglich mit dem `-F` Flag anzugeben, dass ein anderes Trennzeichen verwendet wird [2]. In unserem Fall wird das Trennzeichen das Komma sein, weil wir mit einer CSV arbeiten. In `awk` wird häufig mit regulären Ausdrücken gefiltert, aber da unsere Datei eine statische CSV-Datei ist, kann direkt nach Spaltennamen sortiert werden. Das entsprechende Skript sieht folgendermaßen aus:

```
1 #!/bin/bash
2
3 awk -F ',' '
4 NR==5 {
5 NR==7 {
6 NR==8 {
7 NR==12 {
8 NR==13 {
9 NR==14 {
10 NR==15 {
11 NR==16 {
12 NR==17 {
13 NR==18 {
14 NR==19 {
```

```

15 NR==22  ||
16 NR==23  ||
17 NR==24  ||
18 NR==25  ||
19 NR==26  ||
20 NR==27  ||
21 NR==28  ||
22 NR==29  ||
23 NR==30  ||
24 NR==31  ||
25 NR==36  ||
26 NR==42  ||
27 NR==56  ||
28 NR==57  ||
29 NR==58  ||
30 NR==59  ||
31 NR==60  ||
32 NR==61  ||
33 NR==62  ||
34 NR==82  ||
35 NR==90
36 {print $0}' "Intel_UPE_ComparisonChart_2024_11_04_i7.csv" |
37
38 awk '
39 {
40     for (i=1; i<=NF; i++)
41         a[i] = a[i] ? a[i] "," $i : $i
42     }
43     END {
44         for (i=1; i in a; i++)
45             print a[i]
46     }
47 ' FS=, OFS=, > clean.csv

```

Quelltext 2.1: clean.sh

Das Skript besteht aus zwei Teilen: Im ersten Teil werden die entsprechenden Spalten mit `awk` nach Spaltennummer ausgefiltert, welche vorher als interessant deklariert wurden, und in den zweiten Teil gepiped. Im zweiten Teil werden die jeweiligen Spalten in einer `for`-Schleife in neue Zeilen geladen. Nachdem die gesamte Datei eingelesen wurde, kann diese wieder transponiert ausgegeben werden. Nun kann mit dem eigentlichen Säubern der Daten begonnen werden.

3 Säubern der Daten (Moritz Werr)

Nun ist die CSV vorbereitet für die eigentliche Säuberung. Die Daten liegen in den entsprechenden Spalten vor und alle uninteressanten Daten wurden entfernt. Bevor die Daten jedoch in einer Datenbank in Tabellenform gespeichert werden können, sollten noch entsprechende Spalten bereinigt werden, dass nur noch Zahlenwerte ohne Einheiten dort aufzufinden sind. Weiterhin beinhalten viele Spaltennamen immer noch Sonderzeichen und Leerzeichen, welche es schwierig macht die Spalten zu lesen.

Um diese Probleme zu adressieren, wurde R zur Bereinigung der Daten ausgewählt. R ist eine Programmiersprache für statistische Berechnungen [19]. Es bietet mächtige Werkzeuge zur Auswertung und Speicherung von Daten [19]. Damit ist es eine exzellente Wahl um dieses Problem zu lösen. Es bietet Möglichkeiten die Spaltennamen automatisch mit der Funktion `clean_names()` aus der Bibliothek `janitor` umzubenennen [6], indem es problematische Zeichen, wie z.B. Leerzeichen und Sonderzeichen, erkennt und durch unproblematische Zeichen ersetzt. Weiterhin werden auch alle Großbuchstaben durch Kleinbuchstaben ersetzt. Darüber hinaus ist R eine dynamisch typisierte Programmiersprache [4]. Datentypen werden je nach Kontext automatisiert erkannt. Wenn die entsprechenden Spalten von allen nicht-numerischen Zeichen gesäubert werden, lassen sich diese problemlos in Zahlen umwandeln. Die entsprechenden Einheiten werden nach der Spaltensäuberung stattdessen noch in den Spaltennamen hinzugefügt.

Wenn alle Spalten und Spaltennamen gesäubert bzw. angepasst wurden, können die Daten in einer Datenbank gespeichert werden.

Mit den vorher genannten Überlegungen im Hinterkopf ist folgendes Skript entstanden:

```
1 #!/usr/bin/Rscript
2
3 library(RPostgres)
4
5 intel <- read.csv("clean.csv")           #Reading
6 intel <- janitor::clean_names(intel)      #Cleaning the column
      names
7
8 #cleaning and mutating columns to numbers
9 intel$max_turbo_taktfrequenz <- as.numeric(gsub(" GHz", "",
      intel$max_turbo_taktfrequenz))
10 intel$lithographie <- gsub("Intel 7", "10 nm",
      intel$lithographie)
```

```

11 intel$lithographie <- as.numeric(gsub(" nm", "",
    intel$lithographie))
12 intel$intel_turbo_boost_technik_2_0_taktfrequenz <- as.
    numeric(gsub(" GHz", "",
    intel$intel_turbo_boost_technik_2_0_taktfrequenz))
13 intel$grundtaktfrequenz_des_prozessors <- gsub(" \\| ", ".",
    intel$grundtaktfrequenz_des_prozessors)
14 intel$grundtaktfrequenz_des_prozessors <- as.numeric(gsub("
    GHz", "", intel$grundtaktfrequenz_des_prozessors))
15 intel$cache <- gsub(" MB Intel Smart Cache", "", intel$cache)
16 intel$cache <- as.numeric(gsub(" MB", "", intel$cache))
17 intel$bus_taktfrequenz <- as.numeric(gsub(" GT/s", "",
    intel$bus_taktfrequenz))
18 intel$verlustleistung_tdp <- as.numeric(gsub(" W", "",
    intel$verlustleistung_tdp))
19 intel$intel_turbo_boost_max_technology_3_0_frequency <- gsub
    (" \\| ", ".",
    intel$intel_turbo_boost_max_technology_3_0_frequency)
20 intel$intel_turbo_boost_max_technology_3_0_frequency <- as.
    numeric(gsub(" GHz", "",
    intel$intel_turbo_boost_max_technology_3_0_frequency))
21 intel$single_p_core_turbo_frequency <- as.numeric(gsub(" GHz
    ", "", intel$single_p_core_turbo_frequency))
22 intel$single_e_core_turbo_frequency <- as.numeric(gsub(" GHz
    ", "", intel$single_e_core_turbo_frequency))
23 intel$e_core_base_frequency <- gsub(" GHz", "",
    intel$e_core_base_frequency)
24 intel$e_core_base_frequency <- as.numeric(gsub("900 MHz",
    "0.9", intel$e_core_base_frequency))
25 intel$total_l2_cache <- as.numeric(gsub(" MB", "",
    intel$total_l2_cache))
26 intel$processor_base_power <- as.numeric(gsub(" W", "",
    intel$processor_base_power))
27 intel$maximum_turbo_power <- as.numeric(gsub(" W", "",
    intel$maximum_turbo_power))
28 intel$grundtaktfrequenz_der_grafik <- as.numeric(gsub(" MHz",
    "", intel$grundtaktfrequenz_der_grafik))
29 intel$max_dynamische_grafikfrequenz <- as.numeric(gsub(" GHz
    ", "", intel$max_dynamische_grafikfrequenz))

```

```

30 intel$max_videospeicher_der_grafik <- as.numeric(gsub(" GB",
    "", intel$max_videospeicher_der_grafik))
31 intel$x4k_unterstutzung <- gsub("Hz", "",
    intel$x4k_unterstutzung)
32 intel$x4k_unterstutzung <- as.numeric(gsub("Yes \\| at ",
    "", intel$x4k_unterstutzung))
33
34
35 #Putting the units back into the column names
36 colnames(intel)[colnames(intel) == "max_turbo_taktfrequenz"]
    <- "max_turbo_taktfrequenz_GHz"
37 colnames(intel)[colnames(intel) == "lithographie"] <- "
    litographie_nm"
38 colnames(intel)[colnames(intel) == "
    intel_turbo_boost_technik_2_0_taktfrequenz"] <- "
    intel_turbo_boost_technik_2_0_taktfrequenz_GHz"
39 colnames(intel)[colnames(intel) == "
    grundtaktfrequenz_des_prozessors"] <- "
    grundtaktfrequenz_des_prozessors_GHz"
40 colnames(intel)[colnames(intel) == "cache"] <- "cache_MB"
41 colnames(intel)[colnames(intel) == "bus_taktfrequenz"] <- "
    bus_taktfrequenz_GT_per_s"
42 colnames(intel)[colnames(intel) == "verlustleistung_tdp"] <-
    "verlustleistung_tdp_W"
43 colnames(intel)[colnames(intel) == "
    intel_turbo_boost_max_technology_3_0_frequency"] <- "
    intel_turbo_boost_max_technology_3_0_frequency_GHz"
44 colnames(intel)[colnames(intel) == "
    single_p_core_turbo_frequency"] <- "
    single_p_core_turbo_frequency_GHz"
45 colnames(intel)[colnames(intel) == "
    single_e_core_turbo_frequency"] <- "
    single_e_core_turbo_frequency_GHz"
46 colnames(intel)[colnames(intel) == "e_core_base_frequency"]
    <- "e_core_base_frequency_GHz"
47 colnames(intel)[colnames(intel) == "total_l2_cache"] <- "
    total_l2_cache_MB"
48 colnames(intel)[colnames(intel) == "processor_base_power"] <-
    "processor_base_power_W"

```

```

49 colnames(intel)[colnames(intel) == "maximum_turbo_power"] <-
    "maximum_turbo_power_W"
50 colnames(intel)[colnames(intel) == "
    grundtaktfrequenz_der_grafik"] <- "
    grundtaktfrequenz_der_grafik_MHz"
51 colnames(intel)[colnames(intel) == "
    max_dynamische_grafikfrequenz"] <- "
    max_dynamische_grafikfrequenz_GHz"
52 colnames(intel)[colnames(intel) == "
    max_videospeicher_der_grafik"] <- "
    max_videospeicher_der_grafik_GB"
53 colnames(intel)[colnames(intel) == "x4k_unterstutzung"] <- "
    x4k_unterstutzung_at"
54
55
56 con <- dbConnect(
57   RPostgres::Postgres(),
58   dbname = "bda",
59   host = "localhost" ,
60   port = 5432,
61   user = "bda",
62   password = "bda",
63 )
64
65 dbWriteTable(con, "intel", intel, row.names = FALSE,
66   overwrite = TRUE)
67 dbDisconnect(con)

```

Quelltext 3.1: database.R

Als erstes wird im Skript die transponierte CSV eingelesen und die Spaltennamen mit Hilfe der *janitor* Library gesäubert. Danach kommt ein großer Block an Funktionen, die die einzelnen Spalten von Einheiten und komischer Formatierung säubern. Danach kommt ein großer Block zur Umbenennung der Spaltennamen. Im letzten Schritt wird der Dataframe in einer existierenden Datenbank gespeichert. Ein großer Vorteil an der Funktion ist, dass die Tabelle innerhalb der Datenbank vorher nicht existieren muss. R erstellt von selber eine neue Tabelle mit entsprechenden Datentypen basierend auf den Datentypen innerhalb des Dataframes. Daher ist es wichtig die Datentypen innerhalb des Dataframes schon zu Nummern zu konvertieren.

4 Speichern der Daten (Moritz Werr)

Aus dem vorherigen Kapitel ist bereits hervorgegangen, dass die Daten in einer Datenbank gespeichert werden. In diesem Fall ist die Wahl auf PostgreSQL gefallen. PostgreSQL ist ein Open-Source, relationales Datenbankmanagementsystem [17]. Es ist komplett ACID-kompatibel [17] und speichert Daten in geordneter Form als Tabellen ab [1]. Da die Datenquelle bereits als CSV geliefert wird, bietet es sich an diese in einer entsprechenden Datenbank zu speichern. CSV-Dateien speichern bereits Daten in Tabellenform, weswegen die Wahl direkt auf eine relationale SQL-Datenbank gefallen ist.

Es gibt viele Alternativen zu PostgreSQL, diese sind aber teils kostenpflichtig oder bieten die gleichen bzw. weniger Features wie PostgreSQL. Proprietäre Datenbankmanagementsysteme wie Oracle RDBMS oder Microsoft SQL-Server sind für große, Enterprise-Skala Deployments konzipiert und entsprechend kommerzielle Lizenzen [8][11]. Kostenlose Optionen sind daher bevorzugt. MariaDB bzw. MySQL und SQLite wären Open-Source Alternativen zu PostgreSQL, welche auch gut geeignet sind[21][10]. SQLite jedoch bietet keine Serverfunktionalität weswegen nur lokal gearbeitet werden kann[21]. MariaDB/MySQL haben keinen nennenswerten Funktionalitätsunterschied, aber wurden nicht gewählt, weil mit PostgreSQL mehr Erfahrung im Team existiert[10][16].

Nachdem das Skript aus dem vorherigen Kapitel durchgelaufen ist, ist das Datenbankschema in Abbildung 4.1 zu sehen.

Tabelle »public.intel«				
Spalte	Typ	Sortierfolge	NULL erlaubt?	Vorgabewert
produktsortiment	text			
prozessornummer	text			
litographie_nm	double precision			
empfohlener_kundenpreis	text			
anzahl_der_kerne	integer			
gesamte_threads	integer			
max_turbo_taktfrequenz_GHz	double precision			
intel_turbo_boost_technik_2_0_taktfrequenz_GHz	double precision			
grundtaktfrequenz_des_prozessors_GHz	double precision			
cache_MB	double precision			
bus_taktfrequenz_GT_per_s	double precision			
verlustleistung_tdp_W	double precision			
x_of_performance_cores	integer			
x_of_efficiency_cores	integer			
intel_turbo_boost_max_technology_3_0_frequency_GHz	double precision			
single_p_core_turbo_frequency_GHz	double precision			
single_e_core_turbo_frequency_GHz	double precision			
p_core_base_frequency	text			
e_core_base_frequency_GHz	double precision			
total_l2_cache_MB	double precision			
processor_base_power_W	double precision			
maximum_turbo_power_W	double precision			
einfuhrungsdatum	text			
expected_discontinuance	text			
grundtaktfrequenz_der_grafik_MHz	double precision			
max_dynamische_grafikfrequenz_GHz	double precision			
max_videospeicher_der_grafik_GB	double precision			
x4k_unterstutzung_at	double precision			
max_auflosung_hdmi	text			
max_auflosung_dp	text			
max_auflosung_e_dp_integrierter_flachbildschirm	text			
pci_express_konfigurationen	text			
thermische_spezifikation_des_kuhlers	text			

Figure 4.1: Datenbankschema erstellt durch das R-Skript

5 Visualisieren der Daten (Phil Richter)

Nach der Bereinigung und dem Speichern der Daten in einer Datenbank, wie in den vorherigen Kapiteln beschrieben, können die Daten visualisiert werden. Dafür wird als Grundlage die Programmiersprache Python verwendet. Diese ist vorallem im Bereich der Datenanalyse und -visualisierung sehr weit verbreitet [9, 15]. Eine alternative Möglichkeit wäre die Verwendung von R, welches ebenfalls sehr weit verbreitet ist und speziell für die Datenanalyse entwickelt wurde [19]. Zudem wurde R bereits für die Säuberung der Daten, wie im Kapitel 3 beschrieben, verwendet. Da im Team allerdings mehr Erfahrung mit Python vorhanden war, wurde sich bei der Visualisierung für die Verwendung von Python entschieden.

Für das Laden und Visualisieren der Daten werden folgende Packages verwendet:

Package	Version	Beschreibung
dotenv	1.0.1	Lädt Umgebungsvariablen aus einer <i>.env</i> Datei
sqlalchemy	2.0.36	Ermöglicht Zugriff auf die Datenbank, sowie Datenbankabfragen
pandas	2.2.3	Ermöglicht die Manipulation, Analyse und Verarbeitung von Daten
matplotlib	3.9.2	Erstellt aus gegebenen Daten anpassbare Diagramme
seaborn	0.13.2	Basiert auf <i>matplotlib</i> und wird ebenfalls zur Datenvisualisierung verwendet

Table 5.1: Auflistung aller verwendeten Packages, sowie ihrer Versionen

5.1 Laden der Daten (Phil Richter)

Im ersten Schritt der Visualisieren werden die Daten aus der Datenbank geladen. Dafür werden die Datenbankverbindungsinformationen mit Hilfe von *dotenv* [18] aus einer lokal definierten *.env* Datei geladen. Es werden in diesem Fall Umgebungsvariablen verwendet, damit die Verbindungsinformationen nicht im Code hart codiert sind und somit nicht versehentlich veröffentlicht werden. Nach dem Laden der Variablen wird mit *sqlalchemy* [20] eine Verbindung zur Datenbank aufgebaut und die Daten werden mit einer SQL-Query abgefragt [12]. Die abgefragten Daten werden in einem *pandas* [13] DataFrame gespeichert. Ein DataFrame ist dabei eine zweidimensionale Datenstruktur, wie eine Tabelle [14]. Das DataFrame wird verwendet, um die Daten vor der Visualisierung manipulieren zu können, ohne die Originaldaten zu verändern. Im folgendem Python-Code werden die Daten wie beschrieben geladen:

```

1      import os
2      import re
3      import numpy as np
4      import pandas as pd
5      import seaborn as sn
6      import matplotlib.pyplot as plt
7      from sqlalchemy import create_engine
8      from dotenv import load_dotenv
9
10     load_dotenv() # load environment variables
11
12     # get the environment variables
13     b_user = os.getenv("DB_USER")
14     db_password = os.getenv("DB_PASSWORD")
15     db_name = os.getenv("DB_NAME")
16     db_host = os.getenv("DB_HOST")
17     db_port = os.getenv("DB_PORT")
18
19     db_url = f"postgresql://{db_user}:{db_password}@{db_host}
20             :{db_port}/{db_name}" # create the db url
21     engine = create_engine(db_url)
22
23     query = "SELECT * FROM intel;" # query to get the data
24     from the database
25     df = pd.read_sql(query, engine) # store the data in a
26     dataframe

```

Quelltext 5.1: load data from the database

5.2 Bereinigen der geladenen Daten (Phil Richter)

Auch wenn die Daten, wie im Kapitel 3 beschrieben, bereits vor dem Speichern in die Datenbank bereinigt werden, bleiben kleine Unreinheiten in den Daten bestehen. Daher werden bestimmte Spalten im DataFrame nochmals angepasst.

Ein Beispiel dafür ist die Spalte *produktsortiment*. Diese enthält die jeweilige Generation des Prozessors, welche im Datensatz in drei verschiedenen Schreibweisen vorkommt:

- Die Generation als Zahl, z.B. *4. Generation*
- Die Generation als Zahl ausgeschrieben, z.B. *vierte Generation*

- Die Generation als englische Zahl, z.B. *4th Generation*

Um die Prozessoren für die Visualisierung nach ihrer Generation Gruppieren/Sortieren zu können, wird die Spalte *produktsortiment* in eine neue Spalte *generation* umgewandelt, wobei die Schreibweise vereinheitlicht wird. Dies wird mit Hilfe der Funktion *extract_generation* gemacht. Der Code dazu sieht wie folgt aus:

```

1      number_words = {
2          "erste": 1, "zweite": 2, "dritte": 3, "vierte": 4, "
          fuenfte": 5,
3          "sechste": 6, "siebte": 7, "achte": 8, "neunte": 9, "
          zehnte": 10
4      }
5
6      def extract_generation(gen_text):
7          # match the numeric generation
8          match_numeric = re.search(r'(\d+)(?:\.|th)?\s?Gen',
          gen_text, re.IGNORECASE)
9          if match_numeric:
10             # return the numeric generation
11             return int(match_numeric.group(1))
12     for word, num in number_words.items():
13         # check if the word is in the text
14         if word in gen_text.lower():
15             return num # return the number
16     return None
17
18     # extract the generation from the produktsortiment column
19     df['generation'] = df['produktsortiment'].apply(
        extract_generation)

```

Quelltext 5.2: Funktion *extraxct_generation*

Im Code wird zuerst ein Dictionary *number_words* definiert, welches die Zahlenwörter auf die entsprechende Zahl abbildet. Anschließend wird die Funktion *extract_generation* definiert, welche einen String als Parameter erwartet und die Prozessor-Generation aus dem Text extrahiert. Dabei wird zuerst mit Hilfe eines regulären Ausdrucks, falls vorhanden, die numerische Generation extrahiert. Ein regulärer Ausdruck ist dabei eine Zeichenkette, die ein bestimmtes Suchmuster für ein gegebenen String definiert. Der in der Funktion verwendete reguläre Ausdruck `"(\d+)(?:\.|th)?\s?Gen"` sucht nach einer Zahl gefolgt von einem Punkt oder *th*, sowie einem Leerzeichen und den Buchstaben *Gen*. Falls dabei ein übereinstimmender Text gefunden wird, wird dieser zurückgegeben. Anderenfalls wird überprüft, ob eines

der im Dictionary definierten Wörter im Text vorkommt. Falls ja, wird die entsprechende Zahl zurückgegeben. Abschließend wird die Funktion auf die Spalte *produktsortiment* angewendet und das Ergebnis in der neuen Spalte *generation* gespeichert, welche nur noch die Generation als Zahl enthält. Die neue Spalte kann wie jede andere Spalte im DataFrame zur Visualisierung verwendet werden.

5.3 Erstellen von Graphen (Phil Richter)

Nachdem die Daten geladen und bereinigt wurden, können diese nun visualisiert werden. Dafür wird das Package *matplotlib*, sowie *seaborn* verwendet. Diese sind sehr weit verbreitet und bieten eine Vielzahl an Möglichkeiten, um Daten zu visualisieren. Mit ihnen können verschiedene Diagramme wie Histogramme, Scatterplots, Balkendiagramme und viele mehr erstellt werden. Im folgenden ist ein Beispiel für ein Scatterplot und dessen Code dargestellt:

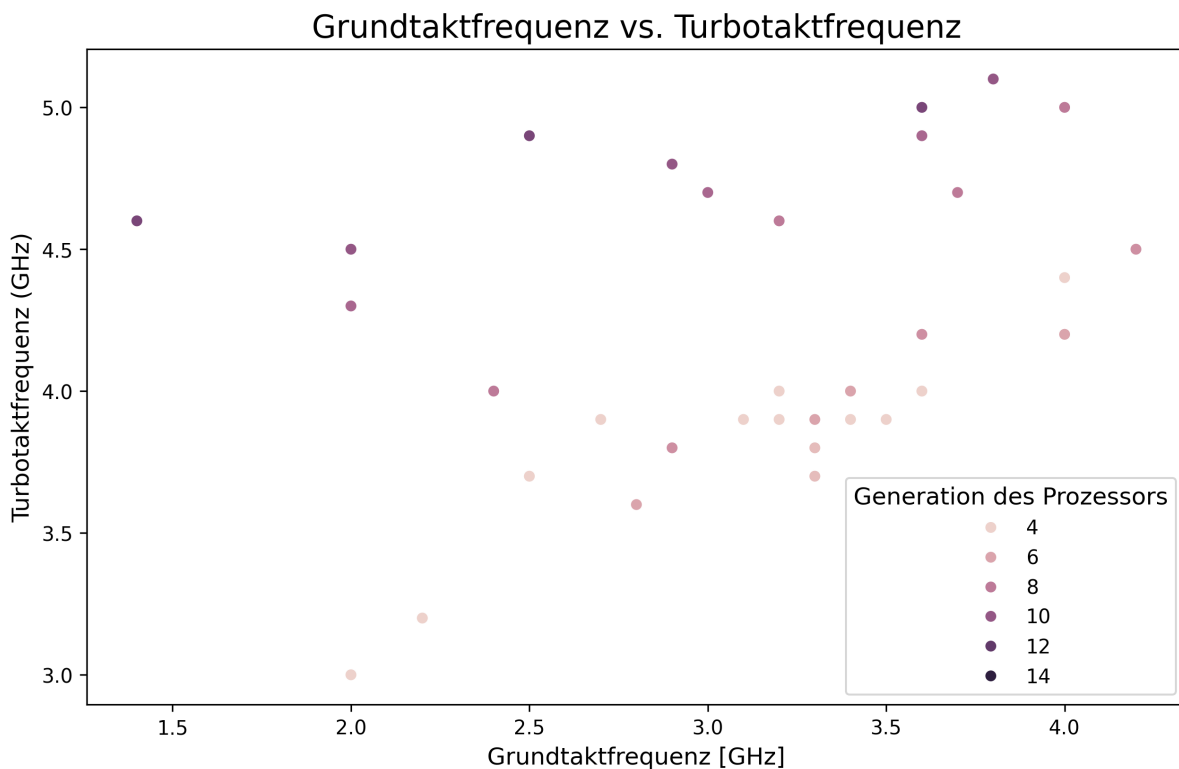


Figure 5.1: Scatterplot erstellt mit *matplotlib* und *seaborn*

```

1 plt.figure(figsize=(10, 6))
2 sn.scatterplot(data=df, x='
    grundtaktfrequenz_des_prozessors_GHz', y='
    max_turbo_taktfrequenz_GHz', hue='generation')
3 plt.title("Grundtaktfrequenz vs. Turbotaktfrequenz",
    fontsize=16)

```

```
4     plt.legend(title='Generation des Prozessors ', fontsize
           =10, title_fontsize=12)
5     plt.xlabel("Grundtaktfrequenz [GHz]", fontsize=12)
6     plt.ylabel("Turbotaktfrequenz (GHz)", fontsize=12)
7     plt.show()
```

Quelltext 5.3: Code für den Scatterplot in Figur 5.1

Im Code wird zuerst ein neues Diagramm mit der Funktion *plt.figure()* erstellt. Dabei wird die Größe des Diagramms auf 10x6 Zoll festgelegt. In Zeile zwei wird mit der Funktion *sn.scatterplot()* bereits der Scatterplot erstellt. Dabei wird als Datenquelle das im Kapitel 5.1 erstellte DataFrame *df* verwendet. Die x-Achse wird mit der Spalte *grundtaktfrequenz_des_prozessors_GHz* und die y-Achse mit der Spalte *max_turbo_taktfrequenz_GHz* belegt. Die Farbe des Punktes repräsentiert die Generation des Prozessors, welche in der im Kapitel 5.2 erstellten Spalte *generation* gespeichert ist. In den Spalten drei bis sechs werden der Titel, die Legende, sowie die Beschriftungen der x- und y-Achse festgelegt. Mit der Funktion *plt.show()* in Spalte sieben wird das Diagramm angezeigt.

Die restlichen Diagramme, die in dieser Arbeit erstellt wurden, sind nach einem ähnlichen Schema wie eben beschrieben erstellt worden. Dabei wurden unterschiedliche Diagrammtypen sowie unterschiedliche Spalten des DataFrames verwendet. Alle Diagramme sind im Ordner *Datenanalyse* in der Datei *visualization.ipynb* zu finden.

6 Zusammenfassung

Dieses Kapitel enthält die Zusammenfassung der Arbeit mit Fazit und Ausblick.

6.1 Fazit

...

6.2 Ausblick

...

A Beispiel-Anhang: Testanhang

Anhänge werden am Ende Ihrer Arbeit vor dem Literaturverzeichnis und dem Index eingefügt.

A.1 Abschnitt im Anhang

Blabla

A.2 Noch ein Abschnitt im Anhang

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Ut purus elit, vestibulum ut, placerat ac, adipiscing vitae, felis. Curabitur dictum gravida mauris. Nam arcu libero, nonummy eget, consectetur id, vulputate a, magna. Donec vehicula augue eu neque. Pellentesque habitant morbi tristique senectus et netus et malesuada fames ac turpis egestas. Mauris ut leo. Cras viverra metus rhoncus sem. Nulla et lectus vestibulum urna fringilla ultrices. Phasellus eu tellus sit amet tortor gravida placerat. Integer sapien est, iaculis in, pretium quis, viverra ac, nunc. Praesent eget sem vel leo ultrices bibendum. Aenean faucibus. Morbi dolor nulla, malesuada eu, pulvinar at, mollis ac, nulla. Curabitur auctor semper nulla. Donec varius orci eget risus. Duis nibh mi, congue eu, accumsan eleifend, sagittis quis, diam. Duis eget orci sit amet orci dignissim rutrum.

Nam dui ligula, fringilla a, euismod sodales, sollicitudin vel, wisi. Morbi auctor lorem non justo. Nam lacus libero, pretium at, lobortis vitae, ultricies et, tellus. Donec aliquet, tortor sed accumsan bibendum, erat ligula aliquet magna, vitae ornare odio metus a mi. Morbi ac orci et nisl hendrerit mollis. Suspendisse ut massa. Cras nec ante. Pellentesque a nulla. Cum sociis natoque penatibus et magnis dis parturient montes, nascetur ridiculus mus. Aliquam tincidunt urna. Nulla ullamcorper vestibulum turpis. Pellentesque cursus luctus mauris.

Nulla malesuada porttitor diam. Donec felis erat, congue non, volutpat at, tincidunt tristique, libero. Vivamus viverra fermentum felis. Donec nonummy pellentesque ante. Phasellus adipiscing semper elit. Proin fermentum massa ac quam. Sed diam turpis, molestie vitae, placerat a, molestie nec, leo. Maecenas lacinia. Nam ipsum ligula, eleifend at, accumsan nec, suscipit a, ipsum. Morbi blandit ligula feugiat magna. Nunc eleifend consequat lorem. Sed lacinia nulla vitae enim. Pellentesque tincidunt purus vel magna. Integer non enim. Praesent euismod nunc eu purus. Donec bibendum quam in tellus. Nullam cursus pulvinar lectus. Donec et mi. Nam vulputate metus eu enim. Vestibulum pellentesque felis eu massa.

Quisque ullamcorper placerat ipsum. Cras nibh. Morbi vel justo vitae lacus tincidunt ultrices. Lorem ipsum dolor sit amet, consectetur adipiscing elit. In hac habitasse platea dictumst. Integer tempus convallis augue. Etiam facilisis. Nunc elementum fermentum wisi. Aenean placerat. Ut imperdiet, enim sed gravida sollicitudin, felis odio placerat quam, ac pulvinar elit purus eget enim. Nunc vitae tortor. Proin tempus nibh sit amet nisl. Vivamus quis tortor vitae risus porta vehicula.

Fusce mauris. Vestibulum luctus nibh at lectus. Sed bibendum, nulla a faucibus semper, leo velit ultricies tellus, ac venenatis arcu wisi vel nisl. Vestibulum diam. Aliquam pellentesque, augue quis sagittis posuere, turpis lacus congue quam, in hendrerit risus eros eget felis. Maecenas eget erat in sapien mattis porttitor. Vestibulum porttitor. Nulla facilisi. Sed a turpis eu lacus commodo facilisis. Morbi fringilla, wisi in dignissim interdum, justo lectus sagittis dui, et vehicula libero dui cursus dui. Mauris tempor ligula sed lacus. Duis cursus enim ut augue. Cras ac magna. Cras nulla. Nulla egestas. Curabitur a leo. Quisque egestas wisi eget nunc. Nam feugiat lacus vel est. Curabitur consectetur.

Suspendisse vel felis. Ut lorem lorem, interdum eu, tincidunt sit amet, laoreet vitae, arcu. Aenean faucibus pede eu ante. Praesent enim elit, rutrum at, molestie non, nonummy vel, nisl. Ut lectus eros, malesuada sit amet, fermentum eu, sodales cursus, magna. Donec eu purus. Quisque vehicula, urna sed ultricies auctor, pede lorem egestas dui, et convallis elit erat sed nulla. Donec luctus. Curabitur et nunc. Aliquam dolor odio, commodo pretium, ultricies non, pharetra in, velit. Integer arcu est, nonummy in, fermentum faucibus, egestas vel, odio.

Sed commodo posuere pede. Mauris ut est. Ut quis purus. Sed ac odio. Sed vehicula hendrerit sem. Duis non odio. Morbi ut dui. Sed accumsan risus eget odio. In hac habitasse platea dictumst. Pellentesque non elit. Fusce sed justo eu urna porta tincidunt. Mauris felis odio, sollicitudin sed, volutpat a, ornare ac, erat. Morbi quis dolor. Donec pellentesque, erat ac sagittis semper, nunc dui lobortis purus, quis congue purus metus ultricies tellus. Proin et quam. Class aptent taciti sociosqu ad litora torquent per conubia nostra, per inceptos hymenaeos. Praesent sapien turpis, fermentum vel, eleifend faucibus, vehicula eu, lacus.

B Beispiel-Anhang: Noch ein Testanhang

nochmal: lipsum ...

Bibliography

- [1] 5.1. *Table Basics*. PostgreSQL Documentation. Nov. 21, 2024. URL: <https://www.postgresql.org/docs/17/ddl-basics.html> (visited on 11/30/2024).
- [2] *awk(1): pattern scanning/processing - Linux man page*. URL: <https://linux.die.net/man/1/awk> (visited on 11/28/2024).
- [3] Min Chen et al. *Big Data: Related Technologies, Challenges and Future Prospects*. SpringerBriefs in Computer Science. Cham: Springer International Publishing, 2014. ISBN: 978-3-319-06244-0 978-3-319-06245-7. DOI: 10.1007/978-3-319-06245-7. URL: <https://link.springer.com/10.1007/978-3-319-06245-7> (visited on 11/13/2024).
- [4] ewendorf. *The Algorithmic Expert's Tool: An introduction to the Python programming language*. ML Conference. Apr. 24, 2019. URL: <https://mlconference.ai/blog/introduction-to-the-r-programming-language/> (visited on 11/30/2024).
- [5] i7. *Intel Core i7 Generation 4 bis 14*. i7. Nov. 4, 2024. URL: <https://www.intel.de/content/www/de/de/products/compare.html?productIds=212279,212280,212047,212048,212251,199325,199335,199314,199316,199318,191048,191792,193738,190885,186604,148263,140642,126684,97129,93339,88200,88195,87718,88040,80807,80808,80809,80814,80806,77656,76642,75121,75122,75123,75124,75125,236781,236794,236854,236783,236789,230492,230490,230491,230500,230489,134596,134591,134592,134594,134595,129948,126686,97122,97128,88196> (visited on 11/04/2024).
- [6] *janitor package - RDocumentation*. URL: <https://www.rdocumentation.org/packages/janitor/versions/2.2.0> (visited on 11/28/2024).
- [7] Mohsen Marjani et al. "Big IoT Data Analytics: Architecture, Opportunities, and Open Research Challenges". In: *IEEE Access* 5 (2017). Conference Name: IEEE Access, pp. 5247–5261. ISSN: 2169-3536. DOI: 10.1109/ACCESS.2017.2689040. URL: <https://ieeexplore.ieee.org/abstract/document/7888916> (visited on 11/13/2024).
- [8] *Microsoft SQL Server System Properties*. URL: <https://db-engines.com/en/system/Microsoft+SQL+Server> (visited on 11/30/2024).
- [9] *Most wanted data science skills U.S. 2019*. Statista. URL: <https://www.statista.com/statistics/1016247/united-states-wanted-data-science-skills/> (visited on 12/01/2024).
- [10] *MySQL System Properties*. URL: <https://db-engines.com/en/system/MySQL> (visited on 11/30/2024).
- [11] *Oracle System Properties*. URL: <https://db-engines.com/en/system/Oracle> (visited on 11/30/2024).
- [12] *ORM Quick Start — SQLAlchemy 2.0 Documentation*. URL: <https://docs.sqlalchemy.org/en/20/orm/quickstart.html> (visited on 12/01/2024).
- [13] *pandas - Python Data Analysis Library*. URL: <https://pandas.pydata.org/> (visited on 12/01/2024).

- [14] *pandas.DataFrame* — *pandas 2.2.3 documentation*. URL: <https://pandas.pydata.org/docs/reference/api/pandas.DataFrame.html#pandas.DataFrame> (visited on 12/01/2024).
- [15] *Popular technologies in the data science tech stack 2024*. Statista. URL: <https://www.statista.com/statistics/1292394/popular-technologies-in-the-data-science-tech-stack/> (visited on 12/01/2024).
- [16] *PostgreSQL System Properties*. URL: <https://db-engines.com/en/system/PostgreSQL> (visited on 11/30/2024).
- [17] *PostgreSQL: About*. URL: <https://www.postgresql.org/about/> (visited on 11/30/2024).
- [18] *python-dotenv: Read key-value pairs from a .env file and set them as environment variables*. Version 1.0.1. URL: <https://github.com/theskumar/python-dotenv> (visited on 12/01/2024).
- [19] *R: What is R?* URL: <https://www.r-project.org/about.html> (visited on 11/28/2024).
- [20] *SQLAlchemy*. URL: <https://www.sqlalchemy.org> (visited on 12/01/2024).
- [21] *SQLite System Properties*. URL: <https://db-engines.com/en/system/SQLite> (visited on 11/30/2024).