# MTH 765P Mini-project

Thanh Trung Vu

January 15, 2024

## 1 Introduction

This dataset comprises a comprehensive collection of Solid State Drives (SSDs). It has been gathered by enthusiasts and serves as a research tool for consumers contemplating purchases.

The analysis of this dataset can be beneficial for those seeking detailed information to guide their SSD purchasing decisions, especially when the market is saturated with various SSD options. This variety makes it very difficult to choose the right one as each individual's needs and preferences vary. Using this dataset, one can for example compare the latest models with older generations, which might reveal that older versions are still adequate, or perhaps even better in certain respects.

By analysing this dataset, we can learn some interesting information, such as the annual trends in SSD production, prices, DRAM cache sizes among others.

The bulk of this project lies in the process of collecting data by web scraping and then converting HTML files to JSON for storage. I will also briefly describe how the data is processed before it is used in some simple analysis and visualisation.

## 2 Data Collection

The process of collecting and storing data looks as follows:



*Process of collecting and storing data*

This can be summarised into two main steps:

- Step 1: **Retrieving URLs for SSD Data Pages**

  This step involves identifying and collecting the web addresses (URLs) of the data pages for various SSDs. This is where specific links to each product's detailed information are gathered.

- Step 2: **Obtaining HTML Documents of the Data Pages**

  In this step, the actual HTML documents of the SSD data pages, accessible through the URLs obtained in Step 1, are downloaded. These documents contain the detailed information and specifications of each SSD, which will be further processed and analyzed.

## 2.1 Retrieving URLs for SSD Data Pages



*The SSD database page (the "root" for scraping)*

I began the process by accessing the SSD database at **TechPowerUp** (*https://www.techpowerup.com/ssd-specs/*). Above is a screenshot of the SSD database page. It's important to note that the screenshot was taken after the data collection was completed, which may result in some discrepancies between the displayed information and the data used in this project.

This page contains a list of SSDs, from which I retrieved the URL for each drive's data page. Each of the small buttons with a drive capacity value contains the link to a data page. However, the page only shows at most 100 drives at a time. To obtain all the URLs, I had to make several queries, each targeting the drives produced by a specific manufacturer. This means that I first had to compile a list of URLs for each SSD brand.

Acer (27)    ADATA (90)    Addlink (30)    AMD (2)    Aorus (26)    Apacer (13)    Apple (2)    Asgard (7)    Asura (4)    ASUS (4)    ASUS ROG (1)    Biostar (5)    CFD Gaming (3)    Colorful (2)    Corsair (54)    Crucial (112)    DapuStor (18)    Dera (8)    Digifast (7)    Digma (3)    Drevo (2)    Enmotus (2)    Fanxiang (18)    Galax (12)    Gigabyte (9)    Goodram (4)    Hikvision (24)    HP (28)    HyperX (6)    Inland (14)    Intel (27)    IRDM (19)    Kingspec (16)    Kingston (87)    Kioxia (31)    Klevv (8)    Kodak (4)    Lexar (46)    Longsys (6)    Magix (4)    Mancer (6)    Mega Electronics Fastro (3)    Memblaze (7)    Micron (43)    MiWhole (9)    MLD (1)    Movespeed (3)    MSI (29)    Mushkin (16)    Neo Forza (6)    Netac (41)    Nextorage (5)    OCZ (9)    Patriot (47)    Phison (3)    Pichau (5)    Plextor (23)    PNY (30)    Ramaxel (1)    Redragon (2)    Reletech (28)    RZX (1)    Sabrent (37)    Samsung (159)    SanDisk (4)    Seagate (63)    Silicon Power (32)    SK Hynix (25)    Solidigm (16)    Sony (4)    SSTC (6)    Synology (2)    Teamgroup (109)    TLET (1)    Toshiba (5)    Transcend (14)    Union Memory (9)    Valve (2)    VisionTek (4)    Western Digital (83)    Wicgtyp (3)    Xbox Series X (1)    XPG (107)    Xraydisk (1)    Zadak (17)    Zhitai (8)    Zircon (10)

*The manufacturer list in the Adanvanced Filters*

To get the URLs of all SSD manufacturers, I accessed the Advanced Filters feature on the page. Here I found the full list of brands as shown above. You might notice a count next to each brand, representing the number of data pages, not the actual number of drives. The number of drives is smaller than the number of data pages since a drive can have multiple versions with different capacities. This discrepancy arises because a single drive model can have multiple versions, each with different storage capacities. For instance, the brand Samsung is listed with 159 drives, but this number includes all the different capacity versions of each model. However, the actual number of distinct drive models is fewer than 100, allowing us to easily obtain all URLs from the database page after applying the filter.

Below is a simplified excerpt from the database page where I found the URLs of the manufacturers.

```
1    <html lang="en">
2      <head><title>SSD Database | TechPowerUp</title></head>
3      <body class="body ssddb"><h1 class="pagetitle">SSD Specs Database</h1>
4      <div class="filter-list">
5        <div class="filter-list-item" data-title="Manufacturer">
6          <div class="filter-list-item-header"><span>Manufacturer</span></div>
7          <div class="filter-list-item-entries">
8            <a href="/ssd-specs/filter/?mfgr=Acer"
9              >Acer <span class="filter-list-item-entry-count">(27)</span></a>
10           <a href="/ssd-specs/filter/?mfgr=ADATA"
11             >ADATA <span class="filter-list-item-entry-count">(90)</span></a>
12           <a href="/ssd-specs/filter/?mfgr=Addlink"
13             >Addlink <span class="filter-list-item-entry-count">(30)</span></a>
14           ...
15         </div>
16       </div>
17     </div>
18     </body>
19   </html>
```

*A simplied excerpt from the HTML document of the SSD database page*

Within this HTML document, I searched for the `<div>` element with attribute `data-title="Manufacturer"`. This element represents the filtering-by-brand functionality, which contains the full list of manufacturers. Then I simply extracted the URLs from the attribute `href` in all the `<a>` elements. For example, the URL for the brand Acer is shown here as */ssd-specs/filter/?mfgr=Acer*. Note that this has to be joined with the home page *https://www.techpowerup.com* at the front to get the full working URL.

After filtering by brand, the next step is to extract the URLs of the individual data pages for the SSD drives from that particular brand.

## 2.2 Obtaining HTML Documents of the Data Pages

```
1    <html lang="en">
2      <head><title>SSD Database | TechPowerUp</title></head>
3      <body class="body ssddb"><h1 class="pagetitle">SSD Specs Database</h1>
4      <table class="drives-desktop-table">
5        <tr>
6          <td>
7            <div class="drive-title">
8              <a href="/ssd-specs/acer-fa100-1-tb.d333" class="drive-name"
9                >Acer FA100 (Micron B27B)</a>
10             <div class="drive-capacities">
11               <a
12                 class="drive-capacity"
13                 href="/ssd-specs/acer-fa100-512-gb.d334"
14                 style="white-space: nowrap"
15                 >512 GB</a>
16               <a
17                 class="drive-capacity"
18                 href="/ssd-specs/acer-fa100-1-tb.d333"
19                 style="white-space: nowrap"
20                 >1 TB</a>
21               ...
22             </div>
23           </div>
24         </td>
```

```
25          </tr>
26          ...
27        </table>
28        </body>
29      </html>
```

*A simplied excerpt from the HTML document of the database page after filtered by brand*

This simplied excerpt is taken from the HTML document from the URL *https://www.techpowerup.com/ssd-specs/filter/?mfgr=Acer*. This is basically the same as the database page, except that the SSD drive list now only contains SSD drives manufactured by Acer.

To get the URLs for the data pages for all the individual Acer drives (each with all of its variants with different capacities), I searched for all the `<a>` elements with attribute `class="drive-capacity"` and extracted the URL from the `href` attribute.

After gathering all the URLs, it was simply the matter of making a HTTP request to each URL to get all the HTML elements. This process took the largest amount of time by far. This is because I could not make too many requests in a short time. If I made too many requests, the site would return `HTTP 429` (Too Many Requests) errors. Therefore, I got to include a delay time between 60 and 90 seconds before each making a request. Since there are well over one thousand URLs to access, this step took multiple hours, including retries when something went wrong.

In the next section, I will describe the structure of the HTML files downloaded from the data pages and how I extracted and stored these data.

# 3    Data Storing

## 3.1    HTML Structure

```
1   <html lang="en">
2     <head>
3       <title>Acer FA100 1 TB Specs | TechPowerUp SSD Database</title>
4       <meta property="og:url"
5       content="https://www.techpowerup.com/ssd-specs/acer-fa100-1-tb.d333"/>
6     </head>
7
8     <body class="body ssddb-details">
9       <h1 class="drivename">Acer FA100 1 TB (Micron B27B)</h1>
10
11      <div class="sectioncontainer">
12        <section class="details">
13          <h1>Solid-State-Drive</h1>
14          <table>
15            <tr> <th>Capacity:</th> <td>1 TB (1024 GB)</td> </tr>
16            <tr> <th>Production:</th> <td>Active</td> </tr>
17            <tr> <th>Released:</th> <td>Apr 30th, 2021</td> </tr>
18            <tr> <th>Price at Launch:</th> <td>105 USD</td> </tr>
19            <tr> <th>Market:</th> <td>Consumer</td> </tr>
20            ...
21          </table>
22        </section>
23
24        <section class="details">
25          <h1>Physical</h1>
26          <table>
27            <tr> <th>Power Draw:</th> <td>
28                1.00 W (Idle)<br />
29                2.0 W (Avg)<br />
```

```
30          3.7 W (Max)<br />
31        </td> </tr>
32      ...
33    </table>
34  </section>

36  <section class="details">
37    <h1>NAND Flash</h1>
38    <table>
39      <tr> <th>Type:</th> <td>TLC</td> </tr>
40      ...
41    </table>
42  </section>

44  <section class="details">
45    <h1>DRAM Cache</h1>
46    <table>
47      <tr> <th>Type:</th> <td>None</td> </tr>
48      ...
49    </table>
50  </section>

52  <section class="details">
53    <h1>Performance</h1>
54    <table>
55      <tr> <th>Sequential Read:</th> <td>3,300 MB/s</td> </tr>
56      <tr> <th>Sequential Write:</th> <td>2,700 MB/s</td> </tr>
57      <tr> <th>Random Read:</th> <td>325,000 IOPS</td> </tr>
58      <tr> <th>Random Write:</th> <td>293,000 IOPS</td> </tr>
59      <tr> <th>Warranty:</th> <td>5 Years</td> </tr>
60      ...
61    </table>
62  </section>

64    ...

66  </div>
67  </body>
68  </html>
```

*A simplified excerpt from the HTML document of a data page*

This simplied excerpt is taken from the HTML document from the URL *https://www.techpowerup.com/ssd-specs/acer-fa100-1-tb.d333*. This is the data page for a version of the drive Acer FA100. From the name `Acer FA100 1 TB (Micron B27B)`, this version has 1 TB of capacity and uses Micron B27B NAND Flash chips.

To distinguish between different SSD drives, I used the ending of a drive's URL as its ID. For example, this drive was taken from the URL *https://www.techpowerup.com/ssd-specs/acer-fa100-1-tb.d333* so I used `acer-fa100-1-tb.d333` as its ID.

I also wanted each drive to have a name. For this I take the text within the `<h1>` element with attribute `class="drivename"`. In this case the drive's name would be `Acer FA100 1 TB (Micron B27B)`.

The rest of the data are found within the `<section>` elements with attribute `class="details"`. Each of these elements contains a `<table>` element. The name of the table is found in the inner `<h1>` element. In this example, the tables are `Solid-State-Drive`, `Physical`, `NAND Flash`, `DRAM Cache` and `Performance`. Each table contains a group of attributes, each attribute is stored in a row represented by the `<tr>` tag. From each row, I extracted the name of the attribute's name from the `<th>` element and the attribute's value from the `<td>` element.

## 3.2 Storing as a JSON file

The attributes extracted from a HTML document were then temporarily stored in memory in a large dictionary object, where each element is an SSD drive with the key being its ID. I repeated this process for all the remaining HTML documents. Then I wrote the dictionary into a JSON file.
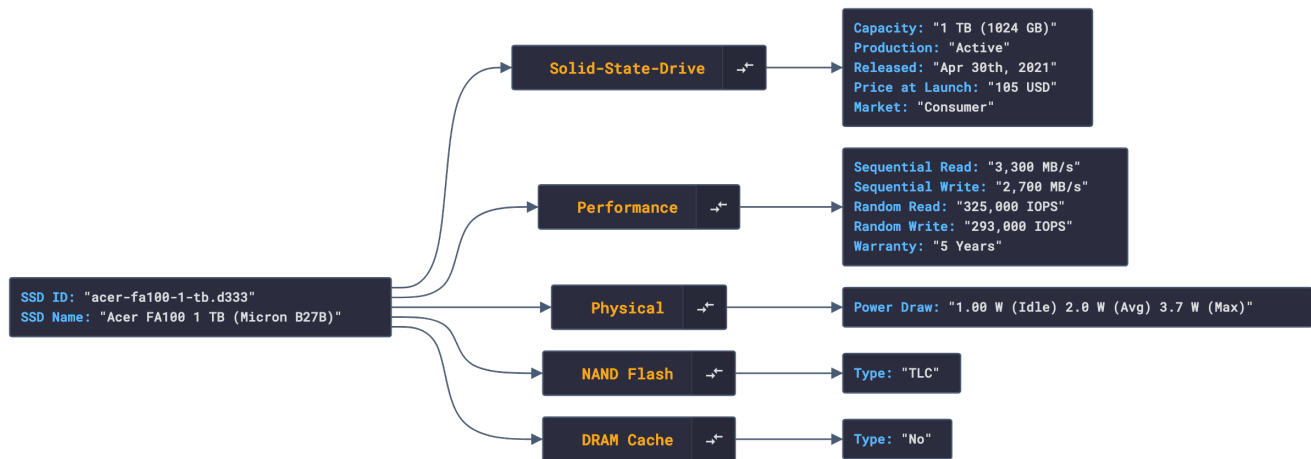
The resulting JSON file looks like this. Similar to all the examples above, this is only a simplified excerpt from the real data file. This section shows the resulting entry for the drive `Acer FA100 1 TB (Micron B27B)` which was discussed above.

```
1   {
2     "acer-fa100-1-tb.d333": {
3       "SSD ID": "acer-fa100-1-tb.d333",
4       "SSD Name": "Acer FA100 1 TB (Micron B27B)",
5       "Solid-State-Drive": {
6         "Capacity": "1 TB (1024 GB)",
7         "Production": "Active",
8         "Released": "Apr 30th, 2021",
9         "Price at Launch": "105 USD",
10        "Market": "Consumer",
11        ...
12      },
13      "Performance": {
14        "Sequential Read": "3,300 MB/s",
15        "Sequential Write": "2,700 MB/s",
16        "Random Read": "325,000 IOPS",
17        "Random Write": "293,000 IOPS",
18        "Warranty": "5 Years",
19        ...
20      },
21      "Physical": {
22        "Power Draw": "1.00 W (Idle) 2.0 W (Avg) 3.7 W (Max)",
23        ...
24      },
25      "NAND Flash": {
26        "Type": "TLC",
27        ...
28      },
29      "DRAM Cache": {
30        "Type": "No",
31        ...
32      },
33      ...
34    },
35    ...
36  }
```

*An excerpt of the JSON file that stores attributes of all SSDs*

For this step, I chose JSON as the data storage method for its simplicity as well as readability. It is easy to import the data from JSON into a dictionary which can then be flattened, a processing step which I will discuss in the next section, and loaded into a DataFrame. If I chose XML, I would end up with a longer file that is less readable. Alternatively, if I chose CSV, I would have to do the flattening first, the take the union of all the columns as there are attributes that are only present in some drives, and then write to the CSV file. Yet other alternatives include database solutions such as SQL and NoSQL which might be more efficient for large datasets but will increase the complexity of the code. For this project, since the amount of data is relatively small, I decided to keep the data storing simple to lower the chances of errors.

## 3.3 A brief desciption of the data



*Tree graph view produced by [https://omute.net/editor](https://omute.net/editor)*

To get a better feel of the data, I visualised the data of an SSD drive as a tree graph using *[https://omute.net/editor](https://omute.net/editor)*. The data of the SSD drive shown here contains five groups of attributes corresponding the five tables mentioned above. Each group contains a set of attributes and their corresponding values. For example, here the group `Solid-State-Drive` contains five attributes: `Capacity`, `Production`, `Released`, `Price` and `Market`. The `Solid-State-Drive` group contains some general information about the drive.

Next we have the group `Performance`. This group contains attributes about different aspects of a drive's performance, including the read and write speeds. It also contains the length of warranty given by the manufacturer.

The group `Physical` contains more information about the drive. An example is the attribute `Power Draw` which shows how the drive's power consumption. It includes the idle or minimum value, the maximum value and the average value of Watts that the drive draws

The group `NAND Flash` shows information about the memory chips used on the drive, such as the type of chips (SLC, MLC, TLC, QLC and so on...). This represents the number of layers or bits that can be stored in a memory cell. More layers mean more capacity per cell, allowing for higher-capacity drives at the same size. In exchange, this lowers the performance and durability of the drive.

The group `DRAM Cache` tells us whether a drive has dedicated memory cache to help improving performance and how much cache capacity there is if any. Drives without DRAM Cache usually have lower prices in exchange for the performance drop.

Please note that these are only the attribute groups and the attributes that I selected for processing and analysis. There are more attribute groups (such as `Controller` and `Features`) and also more attributes in these selected groups (such as `Endurance` and `MTBF` in the group `Performance`) which I will not use in this project but can be of interest for some others.

## 4 Data Preparation

Before processing, I first had to import the data from the JSON file into a DataFrame. Even though `pandas` has a convenient function `read_json`, I could not use it to directly import all the data because of the multi-level structure of the data in the JSON file. The attributes are still divided into groups. To solve this, I first read the JSON file into a dictionary using `json.load`. Then I flattened each of the dictionary by combining the group name with the attribute name. This gave me a dictionary that looks like this

```
1  {
2      "acer-fa100-1-tb.d333": {
3          "SSD ID": "acer-fa100-1-tb.d333",
```

```
 4          "SSD Name": "Acer FA100 1 TB (Micron B27B)",
 5          "Solid-State-Drive - Capacity": "1 TB (1024 GB)",
 6          "Solid-State-Drive - Market": "Consumer",
 7          "Solid-State-Drive - Price at Launch": "105 USD",
 8          "Solid-State-Drive - Production": "Active",
 9          "Solid-State-Drive - Released": "Apr 30th, 2021",
10          "Performance - Random Read": "325,000 IOPS",
11          "Performance - Random Write": "293,000 IOPS",
12          "Performance - Sequential Read": "3,300 MB/s",
13          "Performance - Sequential Write": "2,700 MB/s",
14          "Performance - Warranty": "5 Years",
15          "NAND Flash - Type": "TLC",
16          "DRAM Cache - Type": "No",
17          "Physical - Power Draw": "1.00 W (Idle) 2.0 W (Avg) 3.7 W (Max)",
18          ...
19       },
20       ...
21    }
```

*The dictionary after flattening*

The flattened dictionary was then loaded into a DataFrame which then contained over 100 columns.

If we were to use all the attributes, there would potentially be a huge amount of processing that needed to be done. I decided to keep only the columns as shown in the dictionary above.

Some of these columns can be left as text while others need to be converted to other data types. For instance, columns about price at launch, read/write speeds and warranty need to be converted to numeric data types. To extract the number, I simply split the string using whitespace as the separator and took converted the first token to a number. I added the unit to the column name, for example from `Performance - Warranty` to `Performance - Warranty (Years)`. Some values such as the read/write speeds also contain commas as thousands separators. To parse this correctly, I set a suitable `locale` configuration and used `locale.atof` for number conversion.

The column `Solid-State-Drive - Released` contains information about when the drive is released by the manufacturer. Usually this would be converted to a datetime data type. However, the value in this column can be a full date like `Apr 30th, 2021`, or just the month `Apr 2021`, or just the year `2021`. Because of the inconsistency, I decided to keep only the year number and store it as an integer.

Another column with some consistency is the capacity where the capacity could be shown in GB unit only, like `512 GB`, or in both TB and GB units, like `1 TB (1024 GB)`. I kept only the values in GB.

The trickiest column was the `Physical - Power Draw` attribute. Each value has three numbers for Idle, Avg and Max power consumption. I separated these to 3 separate columns. More specifically, I replaced `"Power Draw":"1.00 W (Idle) 2.0 W (Avg) 3.7 W (Max)"` with `"Idle Power Draw (W)":"1.00"`, `"Avg Power Draw (W)":"2.0"` and `"Max Power Draw (W)":"3.7"`

The resulting DataFrame looks like this

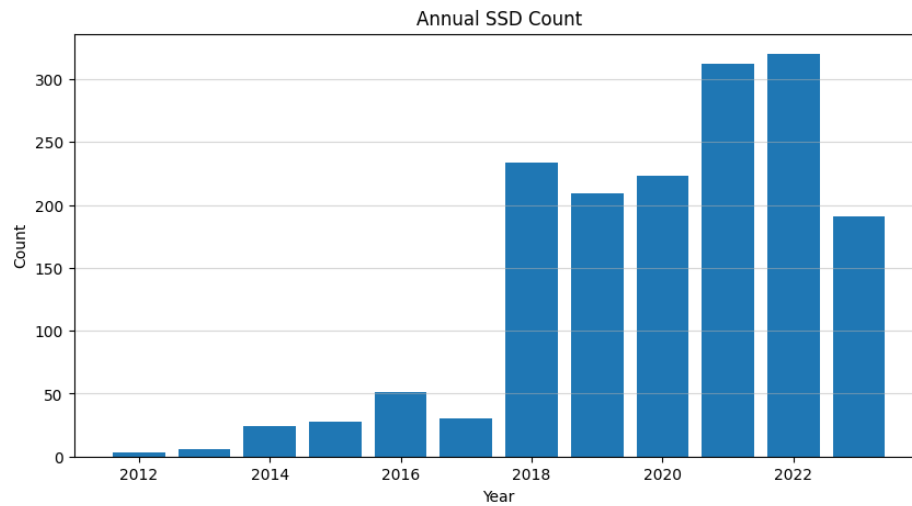| | SSD Name | Solid-State-Drive - Capacity (GB) | Solid-State-Drive - Market | Solid-State-Drive - Price at Launch (USD) | Solid-State-Drive - Production | Solid-State-Drive - Released Year | Performance - Random Read (IOPS) | Performance - Random Write (IOPS) | Performance - Sequential Read (MB/s) | Performance - Sequential Write (MB/s) | Performance - Warranty (Years) | NAND Flash - Type | DRAM Cache - Capacity (MB) | DRAM Cache - Type | Physical - Idle Power Draw (W) | Physical - Average Power Draw (W) | Physical - Maximum Power Draw (W) |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| acer-fa100-1-tb.d333 | Acer FA100 1 TB (Micron B27B) | 1024 | Consumer | 105.0 | Active | 2021 | 325000 | 293000 | 3300 | 2700 | 5 | TLC | <NA> | No | 1.0 | 2.0 | 3.7 |

*The DataFrame after processing*

In the next section I will make some simple visualisations to highlight a few characteristics of this dataset.
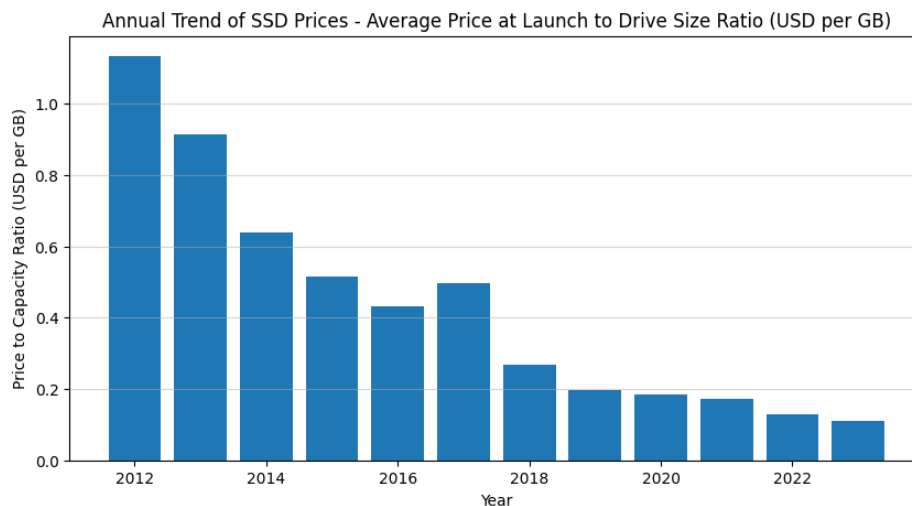
# 5 Data Analysis

## 5.1 SSD Production



*Annual SSD Production*

By plotting the number of SSDs for each year, we can see that there was a sudden jump in SSD production in 2018. It seems that there were not many SSDs before 2018. This would make sense as SSDs were much more expensive back then and Hard Disk Drives (HDDs) were very common. This is assuming that the data is more or less complete, i.e. there were not many missing SSD data for the years before 2018.
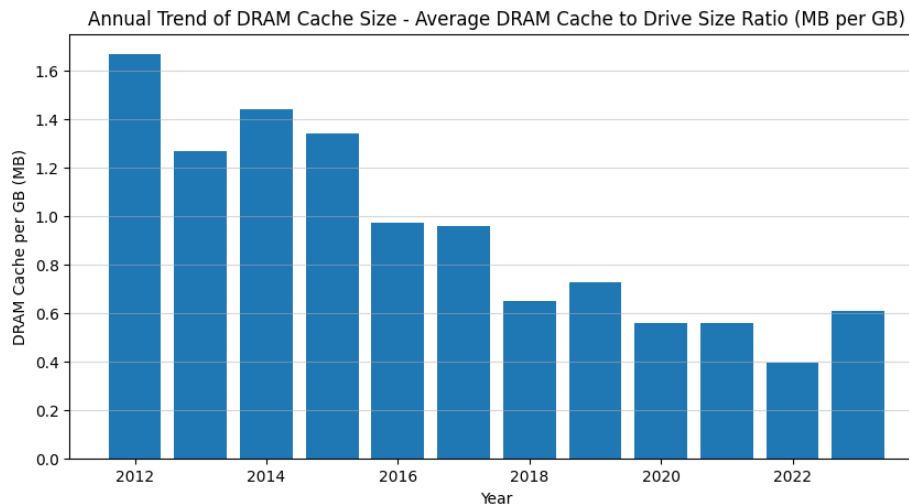
## 5.2 Pricing



*Annual Average "Price to Capacity" Ratio*

A simple way to get a sense of value is to divide the price by the capacity of the drive. This consistently decreasing pattern indicates that drives are either becoming cheaper for the same capacity, or larger drives are being produced at a similar or lower cost due to advancements in production and technology.

While the decrease in price is beneficial, it raises the question: does this reduction come with with a cost in other aspects?

## 5.3 DRAM Cache

Annual Trend of DRAM Cache Size - Average DRAM Cache to Drive Size Ratio (MB per GB)

*Annual Average "DRAM Cache Size to Capacity" Ratio*

One noticeable trend is the decreasing amount of DRAM Cache, which suggests a potential regression. DRAM Cache is beneficial for performance, so its reduction could represent a trade-off to decrease prices. Consequently, consumers might be inclined to seek older generation SSDs with DRAM Cache, rather than newer models without it. However, this trend might also indicate that increasing DRAM Cache size no longer significantly enhances speed performance.

# 6   Libraries

Some important libraries used for this project are:

```
1   import requests # send HTTP requests to get the HTML code of a web page
2   from bs4 import BeautifulSoup # parse HTML files
3   import time    # time.sleep to pause between requests to the server
4   import random # randomize the time between requests
5   import json # read and write json files
6   import glob # get the list of file names in a directory
7   import pandas as pd
8   import matplotlib.pyplot as plt
9   import re       # regular expression for finding numbers in strings
10  import locale   # convert string to numbers with thousand separators
```

# 7   Usage

I submitted three Jupyter notebooks as well as the accompanying data for demonstration purposes. These were designed to be run in order.

The first notebook, titled `step_1_collect_html_files.ipynb`, is designed for downloading the HTML files that store the data of the SSD drives. Please note that parts of the code might not work correctly, as it is primarily for demonstration purposes. Additionally, it makes several assumptions about the website's structure, which could change in the future.

The second notebook, `step_2_convert_html_to_json.ipynb`, is for extracting data from the HTML files and writing it to a JSON file. Since the JSON file is also provided, it is not necessary to run this notebook again to extract the data.

The last notebook, presumably titled `step_3_process_and_visualize_data.ipynb`, is for importing data from the JSON file into DataFrames for processing and visualisation.

# 8    Conclusion

The biggest challenge in this project was perhaps the process of collecting the data by scraping it directly from the website. It was difficult and time-consuming to ensure that the data was correctly scraped.

The second challenge involved storing the data. I opted to convert the individual HTML files into individual JSON files, which, while not ideal for large datasets, provided an easy starting point.

In the end, I obtained a collection of SSD data on which I did a preliminary processing and created a few visualisations showing some annual trends.

Anyone interested in reproducing this work can follow these steps to obtain more current data and conduct their own analysis.