

# ADVANCED TOPICS IN DATA ENGINEERING

*Entity Resolution Assignment*

*Supervisor: Mr. Giorgos Alexiou*

ΟΙΚΟΝΟΜΙΚΟ  
ΠΑΝΕΠΙΣΤΗΜΙΟ  
ΑΘΗΝΩΝ



ATHENS UNIVERSITY  
OF ECONOMICS  
AND BUSINESS

ΣΧΟΛΗ  
ΔΙΟΙΚΗΣΗΣ  
ΕΠΙΧΕΙΡΗΣΕΩΝ  
SCHOOL OF  
BUSINESS

ΤΜΗΜΑ  
ΔΙΟΙΚΗΤΙΚΗΣ  
ΕΠΙΣΤΗΜΗΣ &  
ΤΕΧΝΟΛΟΓΙΑΣ  
DEPARTMENT OF  
MANAGEMENT  
SCIENCE &  
TECHNOLOGY

Dimitris Matsanganis,  
f2822212

Thursday August 31<sup>st</sup>, 2023

## Abstract

This assignment focuses on the intricate domain of entity resolution and duplicate detection, utilizing advanced methodologies to augment data precision and integrity. Through a structured progression, each task contributes to a comprehensive comprehension of the methodologies' nuances.

### Task A: Token Blocking for Block Creation

Task A employs the Token Blocking approach, schema-agnostic in nature, to construct blocks represented as Key-Value (K-V) pairs. All attributes except the identifier (id) contribute to generating distinct Blocking Keys (BKs) from entity attribute values. The id column remains non-participatory in the block index creation process. Lowercasing string attributes during tokenization ensures precise alignment. A function is employed to visualize the BKs in a reader-friendly manner.

### Task B: Calculating Possible Comparisons

In Task B, we compute the full spectrum of necessary comparisons to rectify duplicates within the established blocks from Task A. The quantified number of comparisons offers insights into the computational intricacies of the entity resolution process.

### Task C: Meta-Blocking Graph with CBS Weighting Scheme

Task C introduces the Meta-Blocking graph, implementing the CBS (Common Block Scheme) Weighting Scheme. Edges with a weight below 2 are pruned, refining the block collection and eliminating unnecessary comparisons. The pruned collection serves as the foundation for recalculating the ultimate number of comparisons, paralleling Task B's methodology.

### Task D: Jaccard Similarity Function

Task D encompasses a custom function for computing Jaccard similarity, focusing on the "title" attribute. While not involving practical comparisons, this function serves as a gauge for attribute similarity.

This analysis is executed in Python, specifically version 3.10, using Jupyter Notebook. The assignment amalgamates theoretical understanding with hands-on application, fostering mastery of advanced entity resolution techniques and bolstering data analysis proficiencies.

# Contents

Abstract.....	2
Contents.....	3
Table of Figures.....	5
Assignment Description .....	6
Case Study Outline .....	7
I. Introduction .....	9
II. Data Preprocessing .....	10
Loading Data.....	10
Data Type Consistency .....	10
Handling Missing Values .....	10
Token Creation from Attributes .....	10
Cleaning Tokens .....	10
Concatenating Attributes .....	11
Final Data Preparation .....	11
Conclusion .....	11
III. Task A: Token Blocking for Block Creation.....	12
Creating Key-Value Pairs for Tokens .....	12
Cleaning the Key-Value Pairs.....	13
Displaying the Key-Value Pairs .....	13
Addressing Limitations in Jupyter Notebook and Solution .....	13
Conclusion .....	14
IV. Task B: Calculating Possible Comparisons .....	15
Methodology.....	15
Results & Findings .....	16
Conclusion .....	16
V. Task C: Meta-Blocking Graph with CBS Weighting Scheme .....	17
Initialization and Dictionary Construction .....	18
Entity Counter and Loop Control .....	18
Constructing Entity Pairs and Weights.....	18
Weight Aggregation .....	19
Why Aggregate Weights? .....	19
Iterating Through Entity Pairs.....	19
Summing the Weights .....	19
Outcome .....	19
Pruning Entity Pairs .....	20

Why Prune? .....	20
The Pruning Criterion .....	20
Python Dictionary Comprehension .....	20
Outcome .....	21
Displaying Filtered Entity Pairs .....	21
Calculating Total Comparisons After Pruning .....	22
Why Calculate Total Comparisons? .....	22
Summing the Weights .....	22
Displaying the Result .....	22
Interpreting Total Comparisons After Pruning .....	23
Preliminary Analysis .....	23
Sampling for Visualizations Prior and After Pruning .....	23
Visualizing the Meta-Blocking Graph Prior Pruning .....	23
Sort the Entity Pair Dictionary .....	24
Interpretation of the Meta-Blocking Graph Prior Pruning .....	25
Graphical Representation Post-Pruning with the Reiteration of the Established Procedures .....	25
Inspection of Pruned Entity Pairs .....	27
Visualization Insights .....	28
Conclusion .....	28
VI. Task D: Jaccard Similarity Function .....	30
Creation of Jaccard Similarity Function and Analysis .....	32
Using the Jaccard Similarity Function with Emphasis on Title .....	34
Jaccard Similarities for Titles Insights .....	35
Observation on Jaccard Similarity for Titles .....	35
Conclusion .....	36
VII. Final Analysis Conclusions .....	37
Concluding Thoughts .....	38
VIII. Deliverables .....	39

## Table of Figures

Figure 1: Snapshot of the First 10 Rows of the Final Preprocessed Dataset. ....	11
Figure 2: Displaying the Key-Value Pairs.....	13
Figure 3: IOPub Data Rate Exceeded Error. ....	14
Figure 4: Sample Representation of Entity Pairs and Weights .....	20
Figure 5: Illustrative Display of Filtered Entity Pairs. ....	21
<i>Figure 6: Displaying the result after pruning. ....</i>	<i>22</i>
Figure 7: Sort the Entity Pair Dictionary (1/2). ....	24
Figure 8: Sort the Entity Pair Dictionary (2/2). ....	24
Figure 9: Meta-Blocking Graph - Sample of 100 entities (Prior Pruning, Without Self Loops). .	25
Figure 10: Meta-Blocking Graph - Sample of 100 entities (Post Pruning, Without Self Loops). .	26
Figure 11: Inspection of Pruned Entity Pairs (1/2).....	27
Figure 12: Inspection of Pruned Entity Pairs (2/2).....	27
Figure 13: Representation of Dataset Highlighting Tokenized Attributes for Jaccard Similarity Analysis. ....	31
Figure 14: Jaccard Coefficient Analysis Across Attributes for Randomly Selected Entity Profile Pairs.....	32
Figure 15: Jaccard Coefficient Analysis Highlighting Entities with Shared Blocks, Referencing the Meta-Blocking Process from Task C. ....	33
Figure 16: Jaccard Similarity function with emphasis on the title attribute. ....	34
Figure 17: Comparison of two random entities using the Jaccard Similarity metric, focusing on the 'title' attribute. This visualization illustrates the degree of overlap or similarity between the titles of the entities. ....	34

## Assignment Description

### **Task A [30 points]**

Use the Token Blocking (not to be confused with Standard Blocking) method as a schema-agnostic approach to create blocks in the form of K-V (Key-value) pairs. The key for every entry will be each distinct Blocking Key (BK) derived from the entities' attribute values, except for the id column which should ONLY be used as a reference for the blocking index and should NOT be included in the blocking process (block index creation). Please note that Token Blocking, being schema-agnostic, allows the utilization of all attributes from every entity (except id) for creating the blocks. To ensure accurate matching, it is advised to transform every string to lowercase during the tokens' creation before inserting them into the index. At the end of the creation, use a function to pretty-print the index, displaying the Key-Value pairs in a clear and readable format.

By employing the Token Blocking method in a schema-agnostic manner, you will create blocks that capture the essence of the entities' attributes. Each block's key will represent a distinct Blocking Key (BK) derived from these attribute values, excluding the id column. This approach allows for comprehensive analysis by utilizing all available attributes from each entity. Lowercasing the attribute strings during token creation helps avoid mismatches. Finally, ensure the generated index is presented neatly using a function to pretty-print the Key-Value pairs.

### **Task B [25 points]**

Compute all the possible comparisons that shall be made to resolve the duplicates within the blocks that were created in Task A. After the computation, please print the final calculated number of comparisons.

### **Task C [30 points]**

Create a Meta-Blocking graph of the block collection (created in Task A) and using the CBS Weighting Scheme (i.e., Number of common blocks that entities in a specific comparison have in common) i) prune (delete) the edges that have weight  $< 2$  ii) re-calculate the final number of comparisons (like in step B) of the new block collection that will be created after the edge pruning.

### **Task D [15 points]**

Create a function that takes as input two entities and computes their Jaccard similarity based on the attribute title. You are not requested to perform any actual comparisons using this function.

## Case Study Outline

In this case study, we embark on a comprehensive exploration of advanced techniques for entity resolution and duplicate detection. In today's data-driven world, the importance of ensuring data accuracy and quality cannot be overstated. Duplicates lurking within datasets can taint analytical outcomes and introduce erroneous insights. The advanced methodologies uncovered in this case study furnish a robust framework to surmount these challenges and elevate data precision to a new level.

The systematic journey unfolds across four distinct tasks, each contributing to the development of a holistic understanding of entity resolution methodologies. These tasks guide us through the intricate landscape of data cleansing and entity consolidation, armed with techniques that stand at the forefront of data refinement. From the art of creating data blocks through Token Blocking, and the quantification of comparisons to efficiently resolve duplicates, to crafting Meta-Blocking graphs and gauging attribute similarity, this case study equips us with the precise tools essential for navigating complex data scenarios.

At the heart of this case study lies the integration of theoretical knowledge with practical implementation. Python 3.10, a powerful and versatile programming language, collaborates seamlessly with the dynamic environment of Jupyter Notebook. Together, they form the vessel for executing and demonstrating the advanced techniques outlined in this study. The ensuing sections of this case study unravel the intricate threads woven into the fabric of entity resolution and duplicate detection, inviting readers to immerse themselves in the richness of these methodologies.

### Assignment Structure:

#### **I. Introduction**

The foundation is laid with an exploration of data quality's paramount importance. Within the expansive realm of data management and analysis, the specter of duplicates threatens accuracy. This assignment's core objective revolves around the investigation of advanced methodologies that elevate data integrity through precise identification and resolution of duplicates.

#### **II. Data Preprocessing**

In the realm of accurate analysis, the importance of data preprocessing cannot be overstated. Our primary objective is to refine and transform the dataset, making it more amenable to token blocking and subsequent examination.

#### **III. Task A: Token Blocking for Block Creation**

The initial step introduces us to the Token Blocking method—an avant-garde schema-agnostic technique. By leveraging all attributes except the identifier (id), we construct Key-Value (K-V) pairs that encapsulate distinctive Blocking Keys (BKs). These BKs, derived from entity attribute values, become the building blocks for data consolidation. The id column's passive role in block index creation is emphasized, alongside the importance of lowercasing string attributes



for precise tokenization. The result is a structured index, elegantly visualized using a dedicated function.

#### **IV. Task B: Calculating Possible Comparisons**

In Task B, we tackle the computational intricacies entailed in resolving duplicates. The task entails a comprehensive calculation of all requisite comparisons within the established blocks from Task A. By quantifying the magnitude of computations required, we gain insight into the computational complexity inherent in precise entity resolution.

#### **V. Task C: Constructing a Meta-Blocking Graph**

Task C introduces us to the concept of Meta-Blocking, a powerful strategy that leverages the CBS (Common Block Scheme) Weighting Scheme. Assembling a graph based on the block collection from Task A, we prune edges with weights below 2—a strategic move to streamline the block collection and curtail superfluous comparisons. The pruned collection forms the bedrock for recalculating the definitive number of comparisons.

#### **VI. Task D: Jaccard Similarity Function**

Task D unveils a custom-crafted Jaccard similarity function, tailor-made for attribute comparison based on the "title" attribute. Though devoid of practical comparisons, the function serves as an invaluable tool for quantifying attribute similarity—a linchpin in entity resolution.

#### **VII. Conclusion**

As the case study culminates, the accomplishments of each task and the methodologies' intricacies converge. The assignment concludes with a reflective summary, shedding light on the multifaceted contributions of executed techniques to bolstered data quality, encapsulating the essence of accurate entity resolution and duplicate detection.

#### **VIII. Deliverables**

The assignments and analysis deliverables will be mentioned in this section.



## I. Introduction

In the intricate realm of **data management and analysis**, the **cornerstone of success** lies in unwavering **data quality** and unwavering **accuracy**. Amidst this pursuit, an **omnipresent challenge** emerges in the form of **duplicates** residing within datasets. These duplicates possess the **uncanny ability** to distort insights and introduce misinterpretations, demanding the application of nothing short of **robust and advanced techniques** for effective **entity resolution** and meticulous **duplicate detection**.

Within the **purview** of this assignment, the **profound importance** of **data quality** takes center stage. Duplicates, when left unchecked, can inflict **havoc** upon analytical endeavors, casting a **shadow of doubt** upon the **credibility** of outcomes. To address this **critical concern**, this assignment embarks on a journey to unravel a suite of **advanced methodologies** meticulously crafted to not only detect but to **precisely identify** and rectify duplicates.

Our **voyage of exploration** commences with **Task A**, where the innovative **Token Blocking method** is introduced. This **schema-agnostic approach** ingeniously creates **blocks** through the harmonious pairing of **Key-Value (K-V) pairs**. These pairs, stemming from distinctive **Blocking Keys (BKs)**, encapsulate the very **essence** of entity attributes. However, the **id column** stands distinct, reserved solely for referencing within the **blocking index**, thus serving as a marker but not actively participating in the block creation process. An important note to consider is the transformation of **string attributes** to **lowercase** during token creation, ensuring meticulous accuracy and circumventing **mismatches**. The culmination of this process is a **visually appealing** and **intelligible Key-Value pair representation**, presented through a dedicated **function**.

Advancing forward, **Task B** ventures into the domain of **computational complexity**. Here, we undertake the meticulous **quantification** of all **conceivable comparisons**. This calculation shines a **spotlight** on the intricate web of comparisons required for the **precise resolution of duplicates** within the blocks that were methodically constructed in **Task A**.

As our exploration unfolds, **Task C** unveils the realm of **Meta-Blocking**, introducing a thoughtfully constructed **graph**. This graph is birthed from the **block collection** forged in **Task A**, and it harnesses the potency of the **CBS (Common Block Scheme) Weighting Scheme**. Pruning edges with weights below the threshold of **2** stands as a **strategic decision**, streamlining the block collection and minimizing the **quantum of superfluous comparisons**. This conscientious curation forms the **bedrock** for recalculating the ultimate **number of comparisons**.

**Task D**, the final stride, delves into the creation of a function that mirrors the **Jaccard similarity concept**. This custom function, designed to compute similarity based on the attribute "**title**," emerges as an **indispensable tool** for gauging the **likeness of attributes**. Although this function does not perform **actual comparisons** within this assignment, it is poised to play a **pivotal role** in evaluating similarity across the **expanse of attributes**.

This assignment unfurls within the canvas of **Python 3.10** and the dynamic tapestry of **Jupyter Notebook**. United in purpose, these tools merge **theoretical understanding** with **hands-on application**. As we journey through these tasks, let's embrace the **opportunity to hone our skills**, cultivate a **profound understanding** of **advanced techniques** for entity resolution, and contribute to the realm of **data integrity** and **accuracy**. Through this **comprehensive voyage**, we shall refine our **expertise**, bridging the gap between **theory and practice**, and ultimately elevate our capability to wield **advanced methodologies** for meticulous **entity resolution** and **duplicate detection**.

## II. Data Preprocessing

In this section we are going to mention all the preprocessing procedures needed prior to initialize our analysis. **We need to point out that for a better understanding you can look up for our Jupyter notebook, where our analysis has taken place** (See [Deliverables](#)).

In other words, to ensure accurate analysis, we must first preprocess our data. This phase forms the bedrock of our subsequent analysis, ensuring that the dataset is in the right shape and format.

### Loading Data

To begin with, every analysis starts with data. For our case study, we'll be loading a dataset named **ER-Data.csv**. This dataset is a treasure trove of columns filled with entity attributes, which will be the core of our analysis.

### Data Type Consistency

In the initial phase, we focus on *converting* all relevant columns to **string type**. When dealing with multiple sources, it's common to encounter a mix of data types. **Ensuring** that all relevant columns have a consistent data type, specifically strings, is crucial. This uniformity guarantees that operations like concatenation proceed without any type-related errors, forming a *stable foundation* for the steps that follow.

### Handling Missing Values

Next, we address the ever-present challenge of *missing values*, often represented as NaN (Not a Number). These gaps, if overlooked, can disrupt the analysis, and leaving them untreated might lead to inconsistencies. By **replacing** these NaN values and the string 'nan' with 'None', we ensure uniformity across the dataset. This step underscores the importance of a clean dataset, which is vital for subsequent analysis.

### Token Creation from Attributes

Tokens, the atomic units of our dataset, form the core of our analysis. To create these tokens, we undertake the vital step of **converting all attribute values to lowercase and then splitting them into individual words** (as mentioned in the assignment's instructions). These tokens subsequently serve as the building blocks for the **blocking keys**, which will be instrumental in the entity resolution process. The transformation into tokens ensures that we have distinct, manageable units that *encapsulate the essence* of the information contained within the attributes.

### Cleaning Tokens

In our pursuit of perfection, we noticed that some tokens may contain extraneous characters, such as commas, especially after the split operation. These characters, though seemingly insignificant, can hinder our analysis. By **iteratively processing** each token and removing such characters, we ensure that our tokens are clean and truly represent unique words or attribute values. Furthermore, we remove the words **none** and **nan** to not be biased on the current analysis. This cleanliness is pivotal for the effectiveness of the subsequent analytical methods.

## Concatenating Attributes

With our tokens immaculately prepared, the next logical step is to *combine* them into a unified list for each entity. By concatenating the tokens derived from authors, venue, year, and title, we create a comprehensive **fingerprint** (unique identifier **id**) for each entity. This *amalgamation* ensures that we encapsulate all available information about an entity, significantly enhancing the chances of accurate matches in the subsequent stages.

## Final Data Preparation

As we edge towards concluding our preprocessing, we undertake two pivotal steps. First, we **discard** the original attribute columns, as they become redundant post token creation. Subsequently, we *reconfigure* the dataframe's index using the **id** column. These optimizations not only make our data leaner but also streamline its structure, making it more intuitive for subsequent analysis.

## Conclusion

The data preprocessing stage, often understated, is vital in any data analysis process. By meticulously following the aforementioned steps, we have transformed our dataset into an *optimal format*. This preprocessing lays the groundwork for the token blocking method and further analysis, ensuring that each entity in our dataset is represented by a comprehensive set of tokens derived from all its available attributes. This *structured and clean format* is the bedrock upon which the success of the subsequent stages of our analysis is built.

A sample figure of the first ten rows of our dataframe after the whole data preprocessing procedures follows below:

tokens	
id	
1	[qd, inc, san, diego, none, 11578, sorrento, v...
2	[as, argon, jg, hannaosh, phil., mag, none, in...
3	[gh, hansen, ll, wetterberg, h, sjä¶ strå¶ m, o,...
4	[tm, hammett, p, harmon, w, rhodes, see, none,...
5	[jr, cogdell, new, directions, for, teaching, ...
6	[wm, schmidt, to, none, the, zero, multiplic...
7	[ra, haats, none, 2002.0, predicitive, validity...
8	[jr, norris, j, deisenhofer, san, diego:, acad...
9	[f, bennour, res., rep., ceria, u., paris, non...
10	[i, borg, jc, lingoes, none, 1987.0, multidime...

Figure 1: Snapshot of the First 10 Rows of the Final Preprocessed Dataset.

With the completion of the meticulous data preprocessing phase, our dataset is now in an optimal state for in-depth analysis. We shall commence our investigative journey with *Task A*, laying the foundation for subsequent analytical endeavors.

### III. Task A: Token Blocking for Block Creation

In accordance with the task's directive within this assignment, we embark on the intriguing journey of **Task A** - a pursuit that dives deep into the realm of advanced entity resolution techniques. Employing the dataset presented in the **ER-Data.csv** file, this task harnesses the power of **Token Blocking** as a schema-agnostic methodology to create blocks, subsequently forming **Key-Value (K-V) pairs**.

The foundation of this task resides in the innate need for **data integrity** and the meticulous elimination of duplicates that permeate datasets. While the **Token Blocking** method shares the overarching objective of duplicate detection, it offers a distinct approach. Unlike standard blocking, Token Blocking isn't confined to specific attributes or schemas, thus encompassing the holistic range of attributes from each entity (excluding the id column). This approach fuels a **comprehensive analysis**, drawing insights from all available attributes.

Within this methodology, the creation of blocks transpires through the ingenious fusion of **Key-Value (K-V) pairs**. Each entry's key is ingeniously derived from a distinct **Blocking Key (BK)**, which in turn is an embodiment of the entity's attribute values. Notably, the id column plays a unique role as a mere reference point for the blocking index, abstaining from active participation in the block creation process. A crucial nuance within Token Blocking involves the **transformation of string attributes to lowercase** during tokenization. This tactical maneuver significantly contributes to ensuring accurate matching and precludes potential mismatches that might arise from case discrepancies.

Upon the completion of this process, we culminate with an index that magnificently captures the essence of each entity's attributes, showcased through the clarity of **Key-Value pairs**. These pairs, the embodiment of the schema-agnostic Token Blocking methodology, illuminate the unique attributes that each entity contributes to the data landscape. Through the thoughtful implementation of a dedicated function, the index is meticulously **pretty-printed**, allowing for clear and discernible observation of the Key-Value pairs.

As we embark on this Task A journey, we stand poised to unravel the intricacies of Token Blocking, where each Key-Value pair serves as a testament to our commitment to data accuracy and duplicate eradication. In parallel, we set the stage for the subsequent tasks, where these meticulously crafted blocks will play a pivotal role in the profound exploration of advanced techniques for entity resolution and duplicate detection.

#### Answer of Task A:

In this segment of our analysis, our primary focus lies in the Token Blocking method, an approach aimed at crafting blocks through Key-Value (K-V) pairs. The objective is clear: discern and record relationships between different entities based on shared tokens, facilitating efficient entity resolution.

#### Creating Key-Value Pairs for Tokens

Our analysis begins with the design of a strategic dictionary. Here, each unique token metamorphoses into a key, with corresponding values being the IDs of entities that harbor these tokens: Significance of Multi-entity Keys.

More precisely, **Significance of Multi-entity Keys** is one of the foundational pillars of this step was ensuring keys linked to multiple entities. By setting a cut-off at a minimum of 2 entities per key, we underscored the importance of each key in the blocking mechanism.

## Cleaning the Key-Value Pairs

Following the dictionary's inception, we pivot towards refining it **Pruning Redundancies**. Certain tokens, while present, don't contribute meaningfully to our analysis. Specifically, keys resonating with 'nan', 'none', and blank values were pruned to ensure the cleanliness and relevance of our constructed blocks.

## Displaying the Key-Value Pairs

With the dictionary structured and refined, our next step orbits around its presentation: **Showcasing Interrelationships**. We endeavored to present our Key-Value pairs in a manner that was both clear and insightful. Each token, juxtaposed with its associated entities, painted a vivid picture of interlinked relationships, all anchored by shared tokens.

```
Token: matchmaking  
Entities including it: [9841, 19732, 20516, 26686, 33229, 35380, 44786, 47939, 54039, 64713]
```

```
Token: instructor  
Entities including it: [9842, 21872, 29883, 39844, 47966, 54030]
```

```
Token: arthroscopic  
Entities including it: [9842, 23869, 38857, 50772]
```

Figure 2: Displaying the Key-Value Pairs.

## Addressing Limitations in Jupyter Notebook and Solution

During our intricate analysis within the Jupyter Notebook environment, we stumbled upon an unexpected challenge. As the volume of data we were processing and visualizing grew, we confronted an **error related to Jupyter Notebook's IOPub data rate limit**. The precise message was an "IOPub data rate exceeded" warning. This error acts as a safety net, designed to prevent Jupyter's server from overloading the client with vast outputs, which could potentially lead to crashes (See Figure 3).

The essence of this error is rooted in the default settings of Jupyter. Our expansive analysis, particularly during the display phases, was generating outputs that surpassed this default threshold. This disruption, while a minor setback, required swift troubleshooting to ensure the continuation of our workflow.

Our remedy was both immediate and effective. By invoking the command

**`jupyter notebook --NotebookApp.iopub_data_rate_limit=1e10,`**

we were able to elevate the IOPub data rate limit considerably. This adjustment paved the way for our data-intensive tasks to operate seamlessly, eliminating the risk of triggering the aforementioned error and **giving us the chance to output all the K-V pairs**. On top of that, we

have included a figure titled "IOPub Data Rate Exceeded Error" immediately following this section, offering a snapshot of the error we grappled with.

### Displaying the Key-Value Pairs

Let's display the resulting key-value pairs. This will show each token and the corresponding entities (represented by IDs) that contain it. Observing this structure gives us an **insight into how entities relate** based on shared tokens, which is the foundation of our blocking method.

```
In [17]: # Print the key-value pairs.
for key, values in kv_pairs.items():
    print('Token:', key)
    print('Entities including it:', values)
    print('\n')
```

IOPub data rate exceeded.  
The notebook server will temporarily stop sending output  
to the client in order to avoid crashing it.  
To change this limit, set the config variable  
`--NotebookApp.iopub\_data\_rate\_limit`.

```
Current values:
NotebookApp.iopub_data_rate_limit=1000000.0 (bytes/sec)
NotebookApp.rate_limit_window=3.0 (secs)
```

[49, 209, 423, 1003, 1206, 5235, 5400, 5911, 7500, 7944, 12006, 12342, 12342, 12651, 14519, 14848, 14868, 15373, 16120, 16125, 17474, 19364, 20106, 21500, 21627, 21627, 22381, 23416, 23813, 25784, 25919, 26845, 26976, 27153, 27576, 27896, 30168, 30347, 30946, 31340, 32608, 33502, 37980, 38960, 38960, 39339, 40641, 42522, 44572, 45424, 45424, 46165, 46395, 47866, 48692, 49160, 49519, 51370, 51611, 53412, 54760, 55285, 56448, 56637, 57985, 58608, 59646, 59759, 59816, 60481, 61185, 62063, 62063, 62803]

Token: injection

Entities including it: [40 768 810 1654 2021 2757 6055 7120 7500 8228 8501 8565 8570 0182 0086 12072 12621 1

Figure 3: IOPub Data Rate Exceeded Error.

## Conclusion

The culmination of **Task A** signifies a blend of methodical planning intertwined with meticulous execution. The Token Blocking method, augmented by our systematic approach, ensures that we're poised for success in the subsequent segments of our entity resolution process. Each stage, from token identification to visualization, has been engineered to elevate the accuracy and comprehensiveness of our blocks.



## IV. Task B: Calculating Possible Comparisons

In Task B, we embark on a crucial step of the entity resolution process – **Calculating Possible Comparisons**. This phase is integral to understanding the computational complexity inherent in resolving duplicates within the blocks created in Task A. By systematically quantifying the number of comparisons required, we gain insights into the intricacies of duplicate detection and entity consolidation.

**Duplicates** within datasets can obscure meaningful insights and hinder accurate analysis. The process of **entity resolution** entails identifying these duplicates and consolidating them to represent a single real-world entity. This task unveils the mechanics behind evaluating possible comparisons and lays the groundwork for efficient and effective duplicate resolution.

Our methodology involves examining each block crafted in Task A and computing the total number of comparisons needed to identify and merge duplicates within that block. The cumulative count of these comparisons offers a tangible measure of the computational effort required for robust entity resolution. By unveiling the inherent challenges and intricacies, this task underscores the significance of employing advanced techniques for data accuracy.

Let's dive into the code implementation and calculations to shed light on the magnitude of possible comparisons and pave the way for optimal duplicate detection and entity consolidation.

### Answer of Task B:

In the entity resolution process, understanding the computational complexity inherent in resolving duplicates within the blocks is crucial. **Task B** delves deep into this phase, emphasizing the significance of quantifying the number of comparisons required.

### Methodology

Our approach involves examining each block formed in *Task A* and tallying the total number of comparisons needed to identify and merge duplicates within that block. The combined count of these comparisons provides a tangible measure of the computational effort required for comprehensive entity resolution.

Here's a step-by-step breakdown of the methodology:

1. **Initialize** an empty list named ***comparisons\_per\_block*** to hold the number of comparisons per block.
2. **Iterate** through each block created in *Task A*.
3. For each block, **calculate** the number of entities, denoted as  $(n)$ .
4. **Compute** the number of comparisons for  $(n)$  entities using the formula:  $\frac{n(n-1)}{2}$ .
5. **Append** the computed number of comparisons to the ***comparisons\_per\_block*** list.
6. Finally, **sum up** the list to obtain the total number of comparisons.

The code that fulfills the above is provided below:

```
1. # Initialize an empty list to hold the number of comparisons per block
2. comparisons_per_block = []
3.
4. # Iterate through each block
5. for block_entities in kv_pairs.values():
```

```

6.     # Calculate the number of entities within the block
7.     num_entities = len(block_entities)
8.     # Calculate the number of comparisons for n entities within the block
9.     num_comparisons = num_entities * (num_entities - 1) / 2
10.    # Append the number of comparisons to the list
11.    comparisons_per_block.append(num_comparisons)
12.
13. # Sum up the list to obtain the total number of comparisons
14. total_comparisons = sum(comparisons_per_block)

```

## Results & Findings

After executing the prescribed methodology, we derived the total number of comparisons required to robustly resolve duplicates within the blocks of *Task A*. This calculated metric is pivotal in assessing the computational complexity tied to the entity resolution process. Such insights are vital for formulating strategies that augment duplicate detection accuracy and facilitate streamlined entity consolidation.

The total number of comparisons amounted to an impressive **2,648,668,047**. This staggering figure accentuates the intricate nature of duplicate detection. It underscores the imperative need for sophisticated techniques to deftly traverse the data landscape, ensuring utmost data integrity.

## Conclusion

Entity resolution is a complex task that necessitates meticulous attention to detail. Task B's emphasis on calculating possible comparisons offers a deep dive into the computational challenges involved in this process. *The revelation of over two billion comparisons (2,648,668,047) reiterates the sheer volume and intricacy of data we often grapple with.* Such insights fortify the rationale behind employing advanced techniques in entity resolution. By understanding the depth of the challenge, one can appreciate the importance of innovation and advanced methodologies in ensuring data accuracy and integrity in real-world applications. This task not only quantified the challenge but also shed light on the path forward, pushing the boundaries of what's possible in the realm of entity resolution and duplicate detection.





## V. Task C: Meta-Blocking Graph with CBS Weighting Scheme

*In the intricate tapestry of data management, **entity resolution** stands out as a cornerstone, ensuring the integrity and quality of our datasets. While **Task A** laid the foundation by introducing the Token Blocking technique, **Task C** beckons us into the realm of advanced optimization. Here, we unravel the transformative potential of **Meta-Blocking**.*

**Meta-Blocking** is not just another technique; it's an evolution in our approach to entity resolution. By constructing a graph, where each node is a beacon representing an entity, we begin to see the intricate web of relationships within our data. The edges of this graph aren't mere connectors; they are bridges of significance. These edges indicate that the connected entities coexist within the same block, hinting at potential similarities.

But, how do we measure the strength of these relationships? Enter the **CBS (Common Block Schemes) Weighting Scheme**. This ingenious mechanism breathes life into our graph by allocating weights to its edges. The weight of an edge, symbolizing the number of common blocks shared by two entities, is a testament to their potential similarity. It's a quantifiable measure, a score, that whispers tales of how intertwined two entities might be.

Yet, in the vast expanse of data, not all relationships are created equal. Some are fleeting, while others are profound. In our quest for efficiency, we introduce a pivotal step: **Edge Pruning**. By setting a threshold, we decide to prune edges with a weight less than 2. This deliberate act isn't about exclusion; it's about focus. By pruning, we choose to hone in on the most meaningful of comparisons, ensuring that our attention remains undivided.

Finally, the essence of **Task C** culminates in **Computational Refinement**. After our meticulous pruning process, we recompute the number of comparisons for our newly refined block collection. This is not just a mathematical exercise; it's a reflection of our journey towards unparalleled efficiency in entity resolution.

In the subsequent sections, we shall navigate this intricate journey of **Task C**. With Meta-Blocking as our compass, we will explore, refine, and optimize, ever eager to enhance the precision of our entity resolution process.

### Answer of Task C:

In this section, we delve into an advanced technique known as **Meta-Blocking**. Building upon the blocks generated in Task A, we aim to further optimize our blocking strategy by:

1. Constructing a graph where nodes represent entities and edges signify that these entities appear in the same block.
2. Leveraging the **CBS Weighting Scheme** to assign weights to the edges, with the weight indicating the number of common blocks shared by the two entities.
3. *Pruning* edges that have a weight of less than 2, ensuring a more streamlined comparison process.

Our overarching goal with this approach is to minimize both the number of duplicates and the volume of non-meaningful comparisons, thereby refining our blocking methodology. Let's embark on this task, by firstly initialize a dictionary.

## Initialization and Dictionary Construction

At the very onset of our Meta-Blocking journey, we lay the foundational groundwork by **initializing a dedicated dictionary**. This dictionary, aptly named **entity\_pair\_weights**, is conceptualized to be the linchpin holding together the relationships between entities and the weights they share. The weights, in essence, denote the number of common blocks that entities share, providing a quantifiable measure of their relationship.

## Entity Counter and Loop Control

The sheer volume of data mandates an **efficient** and **controlled** processing approach. To this end, we utilize an entity counter (**processed\_entities\_count**) to track processed entities, and a loop control mechanism (**loop\_control**) to halt processing upon reaching a predefined threshold. These tools work in tandem to ensure the Meta-Blocking process remains optimal and does not get overwhelmed by the data volume.

## Constructing Entity Pairs and Weights

The core of the Meta-Blocking process is to intelligently construct entity pairs and assign them appropriate weights. This step is pivotal as it shines a spotlight on potential duplicates within the dataset. The guiding principle is to traverse each block and pinpoint entities that share the same block. Their shared presence implies potential similarities. The significance of their relationship is then captured as a weight, derived from the *CBS (Common Blocks Scheme)*, which quantifies the blocks in which an entity pair co-exists.

The algorithm follows sequentially the series of following steps:

- 1. Loop through Blocks:** We initiate by navigating through each block. These blocks, clusters of entities, house potential duplicates and are the cornerstone of our Meta-Blocking process.
- 2. Process Each Entity in a Block:** Each entity is meticulously processed. To bolster efficiency, not all entities might be subjected to processing, a decision governed by the entity counter. Furthermore, singleton blocks, or blocks with a sole entity, are bypassed as they offer no potential for pairing.
- 3. Form Entity Pairs and Update Weights:** In this pivotal step, each entity within a block is juxtaposed with subsequent entities, forming pairs. These pairs, representative of potential duplicates, are then weighted based on their shared block occurrences.
- 4. Avoid Redundant Pairing:** Efficiency demands the avoidance of repetitiveness. Post pairing, entities are pruned from their respective blocks, ensuring they don't undergo redundant pairing in subsequent iterations.
- 5. Loop Termination:** The algorithm remains vigilant, monitoring the number of processed entities. Upon reaching a set threshold, it gracefully concludes its processing, emphasizing efficiency and precision.

## Weight Aggregation

Having identified pairs of entities that co-exist within common blocks, it becomes paramount to quantify the strength of their relationship. This strength, or weight, signifies the number of blocks in which each entity pair co-exists. The process of deriving this weight is known as **Weight Aggregation**, an integral step in the meta-blocking process. This weight, once calculated, serves as the very foundation of the meta-blocking graph, illuminating the connection between entities.

### *Why Aggregate Weights?*

During the meticulous construction of entity pairs, each occurrence of an entity pair within a block prompted the addition of a '1' to their associated list. This iterative addition resulted in lists that mirrored their frequency of co-existence. For instance, a list like **[1, 1, 1]** suggests that the corresponding entities co-existed across three distinct blocks. However, to derive meaningful insights and ensure optimal data representation, these lists need to be consolidated. Aggregating these weights paves the way for a clearer and more concise data representation, transforming lists such as **[1, 1, 1]** into a singular integer value, like 3, which stands as a testament to the entities' shared presence across three blocks.

### *Iterating Through Entity Pairs*

The crux of weight aggregation lies in iterating through each entity pair meticulously cataloged in the **entity\_pair\_weights** dictionary. Within this dictionary, each key symbolizes an entity pair, such as 'EntityA,EntityB', while the corresponding value captures their co-existence across blocks in the form of a list like **[1, 1, 1]**. This list represents the number of blocks in which the entities have been found together.

### *Summing the Weights*

The weight aggregation process culminates in the summation of these lists, encapsulating the frequency of co-existence of entity pairs into singular integer values. This transformation is executed by the code line: `entity_pair_weights[key] = sum(values)`. Here, each list, like **[1, 1, 1]**, undergoes summation, resulting in an integer, such as 3. This integer then replaces the original list within the **entity\_pair\_weights** dictionary, providing a succinct representation of the data.

### *Outcome*

Upon the successful culmination of the weight aggregation process, the **entity\_pair\_weights** dictionary undergoes a transformation. Each entity pair now boasts an associated integer weight, encapsulating the total number of blocks they share. This weight isn't just a number; it's the bedrock on which the meta-blocking graph is constructed, highlighting the relationship between entities. It is pivotal in determining the strength and significance of the relationship between entities in the ensuing meta-blocking graph.

Last but not least, for a comprehensive understanding, a sample representation of the transformed dictionary is provided below. This illustrative example offers a snapshot of the data post-aggregation, emphasizing the entity pairs and their corresponding weights (see *Figure 4*).

```
{'1,55360': 1, '1,852': 2, '1,923': 2, '1,2857': 2, '1,3057': 2, '1,3486': 1, '1,4378': 2, '1,4854': 1, '1,5589': 1, '1,6339': 1, '1,7038': 1, '1,8574': 1, '1,9368': 1, '1,10500': 1, '1,11596': 1, '1,15004': 2, '1,16358': 1, '1,17005': 1, '1,18337': 1, '1,21912': 1, '1,22216': 2, '1,22275': 2, '1,23987': 1, '1,24308': 2, '1,26475': 2, '1,27244': 1, '1,29327': 1, '1,30987': 2, '1,32596': 2, '1,36256': 1, '1,36411': 1, '1,38590': 1, '1,39000': 2, '1,40028': 1, '1,41393': 1, '1,42111': 1, '1,42685': 1, '1,43073': 2, '1,43918': 2, '1,44647': 1, '1,44908': 2, '1,45497': 2, '1,46763': 2, '1,47124': 2, '1,49406': 2, '1,50987': 2, '1,51988': 2, '1,52505': 2, '1,53149': 2, '1,56545': 1, '1,56805': 2, '1,57323': 2, '1,58857': 2, '1,59169': 1, '1,60213': 1, '1,61288': 1, '1,61397': 1, '1,62978': 1, '1,65238': 1, '923,923': 1, '923,2857': 2, '923,3057': 2, '923,3486': 1, '923,4378': 2, '923,4854': 1, '923,5589': 1, '923,6339': 1,
```

Figure 4: Sample Representation of Entity Pairs and Weights

## Pruning Entity Pairs

Building on the momentum from the previous steps, we now embark on a mission to refine our dataset by homing in on entity pairs that exhibit a strong likelihood of being genuine matches. This refinement is achieved through a process known as **pruning**, which involves filtering out entity pairs based on the weights they possess. Here, the weight symbolizes the number of blocks where a particular entity pair coexists. Intuitively, a larger weight suggests frequent co-appearances, pointing towards potential similarities or even duplicates among entities.

### Why Prune?

Pruning isn't just a process; it's a strategic optimization step. The rationale behind pruning is twofold. Firstly, by meticulously removing entity pairs that carry lower weights, we drastically slash the number of comparisons awaiting us in the subsequent stages. This not only infuses speed into the process but also fortifies its accuracy. Secondly, pruning minimizes false positives, reducing the risk of mistakenly flagging non-matching entities as potential matches, ensuring precision.

### The Pruning Criterion

Every pruning process requires a guiding light, a criterion that steers its decisions. In our endeavor, this guiding light is a weight threshold. The principle is straightforward: entity pairs burdened with a weight of 2 or less are deemed to have fragile associations and are thus excised. The logic, encapsulated in the code snippet *if value > 2*, suggests a discerning approach: if entities coalesce in just one or two blocks, their association might merely be serendipitous. Conversely, if their camaraderie spans three or more blocks, they are held in higher regard, earning a spot in the subsequent stages of the meta-blocking process.

### Python Dictionary Comprehension

The beauty of the Python language shines through during the pruning process. The task, which might seem daunting, is elegantly executed using Python's dictionary comprehension. This prowess allows us to effortlessly traverse the initial **entity\_pair\_weights** dictionary and chisel out a refined dictionary, aptly named **pruned\_entity\_pair\_weights**. The line of code, *pruned\_entity\_pair\_weights = {key: value for key, value in entity\_pair\_weights.items() if value > 2}*, is a testament to Python's efficiency. It ensures that only those pairs, whose weight surpasses 2, find a place in the pruned dictionary.

## Outcome

Upon the culmination of the pruning process, what emerges is a streamlined dictionary - **pruned\_entity\_pair\_weights**. This dictionary, a subset of the original, proudly showcases entity pairs that exude confidence in their association. Each pair, backed by a weight, indicates the number of blocks they harmoniously coexist in. As we venture ahead, these pairs, symbolizing strong associations, will be at the heart of our comparison or matching endeavors.

## Displaying Filtered Entity Pairs

With a pruned set of entity pairs in our arsenal, it becomes imperative to visualize a subset of this data. The goal? To glean insights into the refined data and ensure transparency in the process. The subsequent code segments are designed to showcase the first few entity pairs, precisely five in our case, that boast a weight denoting a shared presence across two or more blocks.

To be more precise the displaying filtered entity pairs procedure steps follows below:

- 1. Initialization:** At the outset, a counter springs into action. Its sole purpose? To keep a vigilant eye on the number of displayed entity pairs, ensuring the output remains within manageable bounds, capped at five in our scenario.
- 2. Iterate and Display:** With the counter in place, the next step embarks on a journey through the **pruned\_entity\_pair\_weights** dictionary. This dictionary is a repository of pruned entity pairs, each vouching for their strong associations. As we navigate through the dictionary, for each entity pair:
  - a. The respective entities (Nodes) and the count of their shared blocks (Number of Common Blocks) are displayed.
  - b. The counter gets incremented.
  - c. Post the display of five pairs, the loop gracefully concludes its task.

For a holistic understanding, an illustrative snapshot of the displayed entity pairs is provided, showcasing the entities and their shared blocks count (**all have 2 or 3 so fulfill the demands**).

```
In [65]: for key, values in pruned_entity_pair_weights.items():
         if counter < 10:
             print('Nodes:', key, '\n' 'Number of Common Blocks:', values, '\n')
             counter += 1

Nodes: 53149,15004
Number of Common Blocks: 2

Nodes: 58857,15004
Number of Common Blocks: 2

Nodes: 8,17069
Number of Common Blocks: 2

Nodes: 8,24512
Number of Common Blocks: 2

Nodes: 8,25955
Number of Common Blocks: 2

Nodes: 8,29089
Number of Common Blocks: 3
```

Figure 5: Illustrative Display of Filtered Entity Pairs.

## Calculating Total Comparisons After Pruning

Emerging from the rigors of pruning, our next endeavor is to discern the impact of our efforts. The beacon guiding this quest? The total number of entity comparisons post-pruning. This metric not only offers a lens into the efficiency gains secured through the meta-blocking process but also charts the course for subsequent endeavors.

### Why Calculate Total Comparisons?

At the heart of the entity resolution process lies the essence of comparisons. The foundational blocking mechanism herded entities into distinct blocks, where each block echoed with potential matches. But without a refined lens, the sheer volume of comparisons can be daunting, a veritable challenge for any system. Herein lies the beauty of the meta-blocking technique, especially the act of pruning. It doesn't just reduce numbers; it curates them, ensuring only the most promising candidates make the cut. Determining the total number of comparisons post this refinement illuminates two crucial aspects:

1. The prowess and effectiveness of the meta-blocking strategy.
2. The computational trajectory of the upcoming matching processes.

### Summing the Weights

Emerging from the rigors of pruning, our next endeavor is to discern the impact of our efforts. The beacon guiding this quest? The total number of entity comparisons post-pruning. This metric not only offers a lens into the efficiency gains secured through the meta-blocking process but also charts the course for subsequent endeavors.

Aggregating these weights, a task made seamless by Python's `sum()` function, equips us with the total comparisons required in the aftermath of pruning.

### Displaying the Result

Transparency is key. Thus, post the aggregation, the resultant number is presented, offering a crystal-clear perspective on the computational efficiency achieved: **1127 for the sampled 100 entities.**

```
In [38]: print(f'The number of total comparisons after pruning is: {total_comparisons_after_pruning}')  
The number of total comparisons after pruning is: 1127
```

### Interpreting Total Comparisons After Pruning

Given the result:

**The number of total comparisons is: 1127**

Figure 6: Displaying the result after pruning.



## Interpreting Total Comparisons After Pruning

Numbers, while quantitative, often tell a qualitative tale. The figure 1127, emerging from the shadows of the meta-blocking process, carries with it a narrative of refinement. It signifies that after the meticulous pruning, hinged on the CBS (Common Blocks Scheme) weighting scheme, a total of 1127 per 100 entities, entity-to-entity comparisons beckon us.

## Preliminary Analysis

In our preliminary implementation, we closely examined a subset of **100 entities**. Our objective was to determine how many pairs from the entire set of **66,879 entities**, as highlighted in Task A, were present within our initial dataset. Our findings revealed that the number of final comparisons retained amounted to 1127 for the sampled 100 entities so we can **estimate** about **753,726 total comparisons** ( $1127 \times 66879/100$ ).

However, it's crucial to highlight that our solution was based on a **significantly limited sample size**. Consequently, the real number of final comparisons might differ when analyzed on the complete dataset. Regardless, what stands out from our initial assessment is the **substantial reduction** in potential comparisons that could lead to actual matches.

## Sampling for Visualizations Prior and After Pruning

Given the substantial size and intricacy of the complete graph, visualizing it in its entirety can be computationally challenging and might require significant time. To circumvent this challenge, and to provide a clearer illustration, we opted to work with a subset of our data.

We sampled a representative subset comprising **100 entities** from our block collection. This strategic sampling ensured that we could generate a more manageable visualization while still capturing the fundamental relationships that exist between the entities.

For readers interested in visualizing the complete dataset, which consists of *66,879 entities*, we have provided the necessary code in the associated Jupyter notebook. As a note, executing this code will demand substantial computational resources and may take an extended period (the commented section on the Jupyter notebook).

## Visualizing the Meta-Blocking Graph Prior Pruning

The visualization procedure elucidates relationships between a subset of entities based on shared blocks. By employing a seed for randomness, the process ensures that the visual representation remains consistent across multiple runs, although complete reproducibility might not always be guaranteed across various environments.

Entities within the dataset are paired when they share common blocks. The more blocks two entities have in common, the stronger their connection, which is quantified by assigning weights to these entity pairs. This weight system effectively highlights the closeness or frequency of interaction between two entities in the dataset.

Once these weighted relationships are established, a graph is constructed with entities as nodes and their shared-block connections as edges. For clarity in representation, self-looping entities (entities connected to themselves) are excluded, and the layout employs a spring mechanism, which arranges nodes based on their weighted relationships, ensuring a visually intuitive and meaningful display. To further enhance legibility, the visualization slightly fades the connections, allowing for a clearer view of the relationships between entities.

### Sort the Entity Pair Dictionary

### Sort the Entity Pair Dictionary

```
Out[43]: {'10229,29089': 3,
          '10466,29089': 3,
          '10615,29089': 3,
          '10724,29089': 3,
          '10916,29089': 3,
          '11122,29089': 3,
          '11379,29089': 3,
          '11496,29089': 3,
          '11662,29089': 3,
          '11844,29089': 3,
          '11891,29089': 3,
          '12193,29089': 3,
          '12533,29089': 3,
          '12689,29089': 3,
          '13103,29089': 3,
          '13317,29089': 3,
          '13561,29089': 3,
          '13607,29089': 3,
          '13723,29089': 3}
```

Figure 7: Sort the Entity Pair Dictionary (1/2).

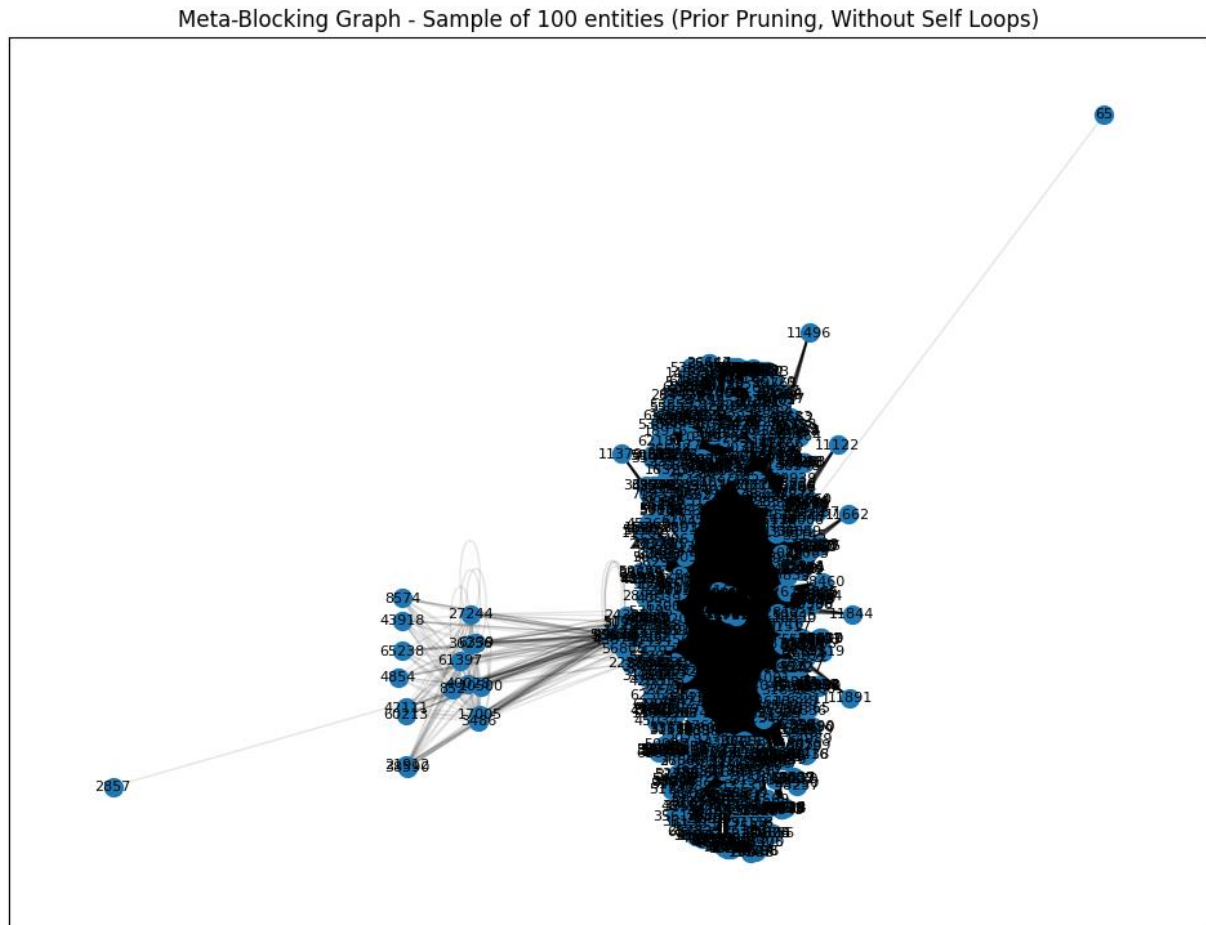
'10229,51988': 1,  
'10229,52114': 1,  
'10229,52164': 1,  
'10229,52287': 1,  
'10229,52449': 1,  
'10229,52505': 1,  
'10229,52536': 1,  
'10229,52542': 1,  
'10229,52739': 1,  
'10229,52812': 1,  
'10229,52974': 1,  
'10229,52993': 1,  
'10229,53131': 1,  
'10229,53149': 1,  
'10229,53195': 1,  
'10229,53305': 1,  
'10229,53365': 1,  
'10229,53469': 1,  
...}

Figure 8: Sort the Entity Pair Dictionary (2/2).



### Interpretation of the Meta-Blocking Graph Prior Pruning

In the presented visualization below, we observe a "Meta-Blocking Graph" that encapsulates a subset of 100 entities. Prior to the visualization, the graph underwent a pruning process, ensuring the elimination of self-loops. Within this graphical representation, individual nodes denote distinct entities. The connecting edges, on the other hand, symbolize the relationships or associations between these entities. These associations are delineated by the **entity\_pair\_weights** dictionary. To enhance clarity, especially in regions with overlapping edges, we opted for a reduced edge opacity.



The nodes, indicative of individual entities, are distinctly marked and labeled. For the sake of visual clarity, especially in regions of overlap, we've moderated the opacity of the connecting edges. Let's plot the Meta-Blocking Graph - Sample of 100 entities (Post Pruning, Without Self Loops) below and then interpret the results.

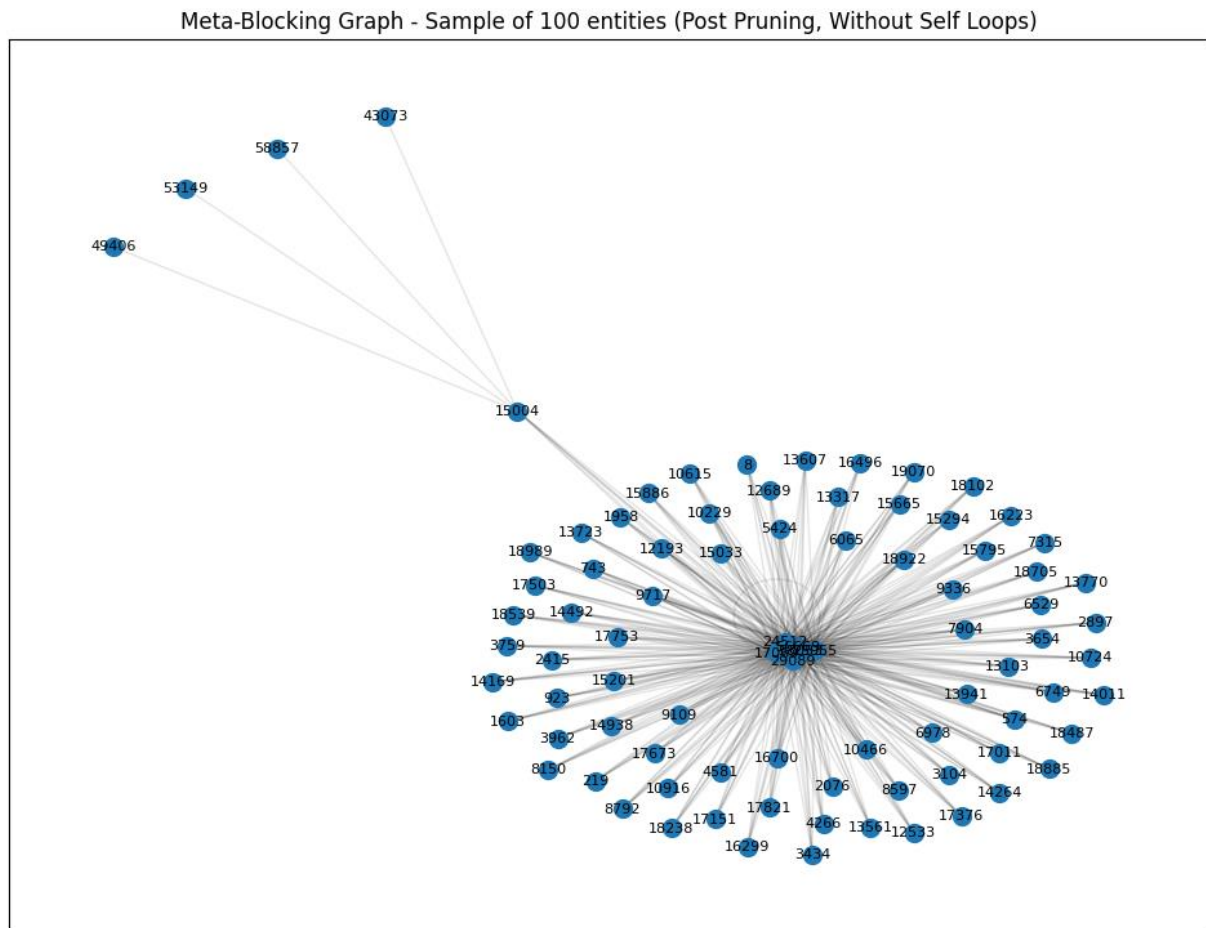


Figure 10: Meta-Blocking Graph - Sample of 100 entities (Post Pruning, Without Self Loops).

## Inspection of Pruned Entity Pairs

Post visualization, we present a sorted dictionary of the pruned entity pairs. This dictionary, *pruned\_entity\_pair\_weights\_sorted*, arranges the entity pairs in descending order of their associated weights. Through this arrangement, we aim to facilitate a seamless and efficient inspection process, enabling the reader to discern patterns or anomalies with ease. Moreover, through this procedure we aim to **validate that there are no pairs with weight less than 2** at the bottom of the dictionary.

### Inspection of Pruned Entity Pairs

Post visualization, we present a sorted dictionary of the pruned entity pairs. This dictionary, *pruned\_entity\_pair\_weights\_sorted*, arranges the entity pairs in descending order of their associated weights. Through this arrangement, we aim to facilitate a seamless and efficient inspection process, enabling the reader to discern patterns or anomalies with ease.

Moreover, through this procedure we aim to **validate that there are no pairs with weight less than 2** at the bottom of the dictionary.

```

In [45]: # Sorted to validate only weights >=2 are remained on the pruned dictionary.
pruned_entity_pair_weights_sorted = dict(sorted(pruned_entity_pair_weights.items(), key=lambda item: (-item[1], item[0])))

# Display the dictionary for validating purposes.
pruned_entity_pair_weights_sorted

Out[45]: {'10229,29089': 3,
          '10466,29089': 3,
          '10615,29089': 3,
          '10724,29089': 3,
          '10916,29089': 3,
          '11122,29089': 3,
          '11379,29089': 3,
          '11496,29089': 3,
          '11662,29089': 3,
          '11844,29089': 3,
          '11891,29089': 3,
          '12193,29089': 3,
          '12533,29089': 3,
          '12689,29089': 3,
          '13103,29089': 3,
          '13317,29089': 3,
          '13561,29089': 3,
          '13607,29089': 3,
          '13723,29089': 3,
          '13779,29089': 3,
          '13844,29089': 3,
          '13891,29089': 3,
          '14193,29089': 3,
          '14533,29089': 3,
          '14689,29089': 3,
          '15103,29089': 3,
          '15317,29089': 3,
          '15561,29089': 3,
          '15607,29089': 3,
          '15723,29089': 3,
          '15779,29089': 3,
          '15844,29089': 3,
          '15891,29089': 3,
          '16193,29089': 3,
          '16533,29089': 3,
          '16689,29089': 3,
          '17103,29089': 3,
          '17317,29089': 3,
          '17561,29089': 3,
          '17607,29089': 3,
          '17723,29089': 3,
          '17779,29089': 3,
          '17844,29089': 3,
          '17891,29089': 3,
          '18193,29089': 3,
          '18533,29089': 3,
          '18689,29089': 3,
          '19103,29089': 3,
          '19317,29089': 3,
          '19561,29089': 3,
          '19607,29089': 3,
          '19723,29089': 3,
          '19779,29089': 3,
          '19844,29089': 3,
          '19891,29089': 3,
          '20193,29089': 3,
          '20533,29089': 3,
          '20689,29089': 3,
          '21103,29089': 3,
          '21317,29089': 3,
          '21561,29089': 3,
          '21607,29089': 3,
          '21723,29089': 3,
          '21779,29089': 3,
          '21844,29089': 3,
          '21891,29089': 3,
          '22193,29089': 3,
          '22533,29089': 3,
          '22689,29089': 3,
          '23103,29089': 3,
          '23317,29089': 3,
          '23561,29089': 3,
          '23607,29089': 3,
          '23723,29089': 3,
          '23779,29089': 3,
          '23844,29089': 3,
          '23891,29089': 3,
          '24193,29089': 3,
          '24533,29089': 3,
          '24689,29089': 3,
          '25103,29089': 3,
          '25317,29089': 3,
          '25561,29089': 3,
          '25607,29089': 3,
          '25723,29089': 3,
          '25779,29089': 3,
          '25844,29089': 3,
          '25891,29089': 3,
          '26193,29089': 3,
          '26533,29089': 3,
          '26689,29089': 3,
          '27103,29089': 3,
          '27317,29089': 3,
          '27561,29089': 3,
          '27607,29089': 3,
          '27723,29089': 3,
          '27779,29089': 3,
          '27844,29089': 3,
          '27891,29089': 3,
          '28193,29089': 3,
          '28533,29089': 3,
          '28689,29089': 3,
          '29103,29089': 3,
          '29317,29089': 3,
          '29561,29089': 3,
          '29607,29089': 3,
          '29723,29089': 3,
          '29779,29089': 3,
          '29844,29089': 3,
          '29891,29089': 3,
          '30193,29089': 3,
          '30533,29089': 3,
          '30689,29089': 3,
          '31103,29089': 3,
          '31317,29089': 3,
          '31561,29089': 3,
          '31607,29089': 3,
          '31723,29089': 3,
          '31779,29089': 3,
          '31844,29089': 3,
          '31891,29089': 3,
          '32193,29089': 3,
          '32533,29089': 3,
          '32689,29089': 3,
          '33103,29089': 3,
          '33317,29089': 3,
          '33561,29089': 3,
          '33607,29089': 3,
          '33723,29089': 3,
          '33779,29089': 3,
          '33844,29089': 3,
          '33891,29089': 3,
          '34193,29089': 3,
          '34533,29089': 3,
          '34689,29089': 3,
          '35103,29089': 3,
          '35317,29089': 3,
          '35561,29089': 3,
          '35607,29089': 3,
          '35723,29089': 3,
          '35779,29089': 3,
          '35844,29089': 3,
          '35891,29089': 3,
          '36193,29089': 3,
          '36533,29089': 3,
          '36689,29089': 3,
          '37103,29089': 3,
          '37317,29089': 3,
          '37561,29089': 3,
          '37607,29089': 3,
          '37723,29089': 3,
          '37779,29089': 3,
          '37844,29089': 3,
          '37891,29089': 3,
          '38193,29089': 3,
          '38533,29089': 3,
          '38689,29089': 3,
          '39103,29089': 3,
          '39317,29089': 3,
          '39561,29089': 3,
          '39607,29089': 3,
          '39723,29089': 3,
          '39779,29089': 3,
          '39844,29089': 3,
          '39891,29089': 3,
          '40193,29089': 3,
          '40533,29089': 3,
          '40689,29089': 3,
          '41103,29089': 3,
          '41317,29089': 3,
          '41561,29089': 3,
          '41607,29089': 3,
          '41723,29089': 3,
          '41779,29089': 3,
          '41844,29089': 3,
          '41891,29089': 3,
          '42193,29089': 3,
          '42533,29089': 3,
          '42689,29089': 3,
          '43103,29089': 3,
          '43317,29089': 3,
          '43561,29089': 3,
          '43607,29089': 3,
          '43723,29089': 3,
          '43779,29089': 3,
          '43844,29089': 3,
          '43891,29089': 3,
          '44193,29089': 3,
          '44533,29089': 3,
          '44689,29089': 3,
          '45103,29089': 3,
          '45317,29089': 3,
          '45561,29089': 3,
          '45607,29089': 3,
          '45723,29089': 3,
          '45779,29089': 3,
          '45844,29089': 3,
          '45891,29089': 3,
          '46193,29089': 3,
          '46533,29089': 3,
          '46689,29089': 3,
          '47103,29089': 3,
          '47317,29089': 3,
          '47561,29089': 3,
          '47607,29089': 3,
          '47723,29089': 3,
          '47779,29089': 3,
          '47844,29089': 3,
          '47891,29089': 3,
          '48193,29089': 3,
          '48533,29089': 3,
          '48689,29089': 3,
          '49103,29089': 3,
          '49317,29089': 3,
          '49561,29089': 3,
          '49607,29089': 3,
          '49723,29089': 3,
          '49779,29089': 3,
          '49844,29089': 3,
          '49891,29089': 3,
          '50193,29089': 3,
          '50533,29089': 3,
          '50689,29089': 3,
          '51103,29089': 3,
          '51317,29089': 3,
          '51561,29089': 3,
          '51607,29089': 3,
          '51723,29089': 3,
          '51779,29089': 3,
          '51844,29089': 3,
          '51891,29089': 3,
          '52193,29089': 3,
          '52533,29089': 3,
          '52689,29089': 3,
          '53103,29089': 3,
          '53317,29089': 3,
          '53561,29089': 3,
          '53607,29089': 3,
          '53723,29089': 3,
          '53779,29089': 3,
          '53844,29089': 3,
          '53891,29089': 3,
          '54193,29089': 3,
          '54533,29089': 3,
          '54689,29089': 3,
          '55103,29089': 3,
          '55317,29089': 3,
          '55561,29089': 3,
          '55607,29089': 3,
          '55723,29089': 3,
          '55779,29089': 3,
          '55844,29089': 3,
          '55891,29089': 3,
          '56193,29089': 3,
          '56533,29089': 3,
          '56689,29089': 3,
          '57103,29089': 3,
          '57317,29089': 3,
          '57561,29089': 3,
          '57607,29089': 3,
          '57723,29089': 3,
          '57779,29089': 3,
          '57844,29089': 3,
          '57891,29089': 3,
          '58193,29089': 3,
          '58533,29089': 3,
          '58689,29089': 3,
          '59103,29089': 3,
          '59317,29089': 3,
          '59561,29089': 3,
          '59607,29089': 3,
          '59723,29089': 3,
          '59779,29089': 3,
          '59844,29089': 3,
          '59891,29089': 3,
          '60193,29089': 3,
          '60533,29089': 3,
          '60689,29089': 3,
          '61103,29089': 3,
          '61317,29089': 3,
          '61561,29089': 3,
          '61607,29089': 3,
          '61723,29089': 3,
          '61779,29089': 3,
          '61844,29089': 3,
          '61891,29089': 3,
          '62193,29089': 3,
          '62533,29089': 3,
          '62689,29089': 3,
          '63103,29089': 3,
          '63317,29089': 3,
          '63561,29089': 3,
          '63607,29089': 3,
          '63723,29089': 3,
          '63779,29089': 3,
          '63844,29089': 3,
          '63891,29089': 3,
          '64193,29089': 3,
          '64533,29089': 3,
          '64689,29089': 3,
          '65103,29089': 3,
          '65317,29089': 3,
          '65561,29089': 3,
          '65607,29089': 3,
          '65723,29089': 3,
          '65779,29089': 3,
          '65844,29089': 3,
          '65891,29089': 3,
          '66193,29089': 3,
          '66533,29089': 3,
          '66689,29089': 3,
          '67103,29089': 3,
          '67317,29089': 3,
          '67561,29089': 3,
          '67607,29089': 3,
          '67723,29089': 3,
          '67779,29089': 3,
          '67844,29089': 3,
          '67891,29089': 3,
          '68193,29089': 3,
          '68533,29089': 3,
          '68689,29089': 3,
          '69103,29089': 3,
          '69317,29089': 3,
          '69561,29089': 3,
          '69607,29089': 3,
          '69723,29089': 3,
          '69779,29089': 3,
          '69844,29089': 3,
          '69891,29089': 3,
          '70193,29089': 3,
          '70533,29089': 3,
          '70689,29089': 3,
          '71103,29089': 3,
          '71317,29089': 3,
          '71561,29089': 3,
          '71607,29089': 3,
          '71723,29089': 3,
          '71779,29089': 3,
          '71844,29089': 3,
          '71891,29089': 3,
          '72193,29089': 3,
          '72533,29089': 3,
          '72689,29089': 3,
          '73103,29089': 3,
          '73317,29089': 3,
          '73561,29089': 3,
          '73607,29089': 3,
          '73723,29089': 3,
          '73779,29089': 3,
          '73844,29089': 3,
          '73891,29089': 3,
          '74193,29089': 3,
          '74533,29089': 3,
          '74689,29089': 3,
          '75103,29089': 3,
          '75317,29089': 3,
          '75561,29089': 3,
          '75607,29089': 3,
          '75723,29089': 3,
          '75779,29089': 3,
          '75844,29089': 3,
          '75891,29089': 3,
          '76193,29089': 3,
          '76533,29089': 3,
          '76689,29089': 3,
          '77103,29089': 3,
          '77317,29089': 3,
          '77561,29089': 3,
          '77607,29089': 3,
          '77723,29089': 3,
          '77779,29089': 3,
          '77844,29089': 3,
          '77891,29089': 3,
          '78193,29089': 3,
          '78533,29089': 3,
          '78689,29089': 3,
          '79103,29089': 3,
          '79317,29089': 3,
          '79561,29089': 3,
          '79607,29089': 3,
          '79723,29089': 3,
          '79779,29089': 3,
          '79844,29089': 3,
          '79891,29089': 3,
          '80193,29089': 3,
          '80533,29089': 3,
          '80689,29089': 3,
          '81103,29089': 3,
          '81317,29089': 3,
          '81561,29089': 3,
          '81607,29089': 3,
          '81723,29089': 3,
          '81779,29089': 3,
          '81844,29089': 3,
          '81891,29089': 3,
          '82193,29089': 3,
          '82533,29089': 3,
          '82689,29089': 3,
          '83103,29089': 3,
          '83317,29089': 3,
          '83561,29089': 3,
          '83607,29089': 3,
          '83723,29089': 3,
          '83779,29089': 3,
          '83844,29089': 3,
          '83891,29089': 3,
          '84193,29089': 3,
          '84533,29089': 3,
          '84689,29089': 3,
          '85103,29089': 3,
          '85317,29089': 3,
          '85561,29089': 3,
          '85607,29089': 3,
          '85723,29089': 3,
          '85779,29089': 3,
          '85844,29089': 3,
          '85891,29089': 3,
          '86193,29089': 3,
          '86533,29089': 3,
          '86689,29089': 3,
          '87103,29089': 3,
          '87317,29089': 3,
          '87561,29089': 3,
          '87607,29089': 3,
          '87723,29089': 3,
          '87779,29089': 3,
          '87844,29089': 3,
          '87891,29089': 3,
          '88193,29089': 3,
          '88533,29089': 3,
          '88689,29089': 3,
          '89103,29089': 3,
          '89317,29089': 3,
          '89561,29089': 3,
          '89607,29089': 3,
          '89723,29089': 3,
          '89779,29089': 3,
          '89844,29089': 3,
          '89891,29089': 3,
          '90193,29089': 3,
          '90533,29089': 3,
          '90689,29089': 3,
          '91103,29089': 3,
          '91317,29089': 3,
          '91561,29089': 3,
          '91607,29089': 3,
          '91723,29089': 3,
          '91779,29089': 3,
          '91844,29089': 3,
          '91891,29089': 3,
          '92193,29089': 3,
          '92533,29089': 3,
          '92689,29089': 3,
          '93103,29089': 3,
          '93317,29089': 3,
          '93561,29089': 3,
          '93607,29089': 3,
          '93723,29089': 3,
          '93779,29089': 3,
          '93844,29089': 3,
          '93891,29089': 3,
          '94193,29089': 3,
          '94533,29089': 3,
          '94689,29089': 3,
          '95103,29089': 3,
          '95317,29089': 3,
          '95561,29089': 3,
          '95607,29089': 3,
          '95723,29089': 3,
          '95779,29089': 3,
          '95844,29089': 3,
          '95891,29089': 3,
          '96193,29089': 3,
          '96533,29089': 3,
          '96689,29089': 3,
          '97103,29089': 3,
          '97317,29089': 3,
          '97561,29089': 3,
          '97607,29089': 3,
          '97723,29089': 3,
          '97779,29089': 3,
          '97844,29089': 3,
          '97891,29089': 3,
          '98193,29089': 3,
          '98533,29089': 3,
          '98689,29089': 3,
          '99103,29089': 3,
          '99317,29089': 3,
          '99561,29089': 3,
          '99607,29089': 3,
          '99723,29089': 3,
          '99779,29089': 3,
          '99844,29089': 3,
          '99891,29089': 3,
          '100193,29089': 3,
          '100533,29089': 3,
          '100689,29089': 3,
          '101103,29089': 3,
          '101317,29089': 3,
          '101561,29089': 3,
          '101607,29089': 3,
          '101723,29089': 3,
          '101779,29089': 3,
          '101844,29089': 3,
          '101891,29089': 3,
          '102193,29089': 3,
          '102533,29089': 3,
          '102689,29089': 3,
          '103103,29089': 3,
          '103317,29089': 3,
          '103561,29089': 3,
          '103607,29089': 3,
          '103723,29089': 3,
          '103779,29089': 3,
          '103844,29089': 3,
          '103891,29089': 3,
          '104193,29089': 3,
          '104533,29089': 3,
          '104689,29089': 3,
          '105103,29089': 3,
          '105317,29089': 3,
          '105561,29089': 3,
          '105607,29089': 3,
          '105723,29089': 3,
          '105779,29089': 3,
          '105844,29089': 3,
          '105891,29089': 3,
          '106193,29089': 3,
          '106533,29089': 3,
          '106689,29089': 3,
          '107103,29089': 3,
          '107317,29089': 3,
          '107561,29089': 3,
          '107607,29089': 3,
          '107723,29089': 3,
          '107779,29089': 3,
          '107844,29089': 3,
          '107891,29089': 3,
          '108193,29089': 3,
          '108533,29089': 3,
          '108689,29089': 3,
          '109103,29089': 3,
          '109317,29089': 3,
          '109561,29089': 3,
          '109607,29089': 3,
          '109723,29089': 3,
          '109779,29089': 3,
          '109844,29089': 3,
          '109891,29089': 3,
          '110193,29089': 3,
          '110533,29089': 3,
          '110689,29089': 3,
          '111103,29089': 3,
          '111317,29089': 3,
          '111561,29089': 3,
          '111607,29089': 3,
          '111723,29089': 3,
          '111779,29089': 3,
          '111844,29089': 3,
          '111891,29089': 3,
          '112193,29089': 3,
          '112533,29089': 3,
          '112689,29089': 3,
          '113103,29089': 3,
          '113317,29089': 3,
          '113561,29089': 3,
          '113607,29089': 3,
          '113723,29089': 3,
          '113779,29089': 3,
          '113844,29089': 3,
          '113891,29089': 3,
          '114193,29089': 3,
          '114533,29089': 3,
          '114689,29089': 3,
          '115103,29089': 3,
          '115317,29089': 3,
          '115561,29089': 3,
          '115607,29089': 3,
          '115723,29089': 3,
          '115779,29089': 3,
          '115844,29089': 3,
          '115891,29089': 3,
          '116193,29089': 3,
          '116533,29089': 3,
          '116689,29089': 3,
          '117103,29089': 3,
          '117317,29089': 3,
          '117561,29089': 3,
          '117607,29089': 3,
          '117723,29089': 3,
          '117779,29089': 3,
          '117844,29089': 3,
          '117891,29089': 3,
          '118193,29089': 3,
          '118533,29089': 3,
          '118689,29089': 3,
          '119103,29089': 3,
          '119317,29089': 3,
          '119561,29089': 3,
          '119607,29089': 3,
          '119723,29089': 3,
          '119779,29089': 3,
          '119844,29089': 3,
          '119891,29089': 3,
          '120193,29089': 3,
          '120533,29089': 3,
          '120689,29089': 3,
          '121103,29089': 3,
          '121317,29089': 3,
          '121561,29089': 3,
          '121607,29089': 3,
          '121723,29089': 3,
          '121779,29089': 3,
          '121844,29089': 3,
          '121891,29089': 3,
          '122193,29089': 3,
          '122533,29089': 3,
          '122689,29089': 3,
          '123103,29089': 3,
          '123317,29089': 3,
          '123561,29089': 3,
          '123607,29089': 3,
          '123723,29089': 3,
          '123779,29089': 3,
          '123844,29089': 3,
          '123891,29089': 3,
          '124193,29089': 3,
          '124533,29089': 3,
          '124689,29089': 3,
          '125103,29089': 3,
          '125317,29089': 3,
          '125561,29089': 3,
          '125607,29089': 3,
          '125723,29089': 3,
          '125779,29089': 3,
          '125844,29089': 3,
          '125891,29089': 3,
          '126193,29089': 3,
          '126533,29089': 3,
          '126689,29089': 
```

## Visualization Insights

After a carefully look to both image and to the resulting dictionaries we can ensure **that only edges that have weight < 2** are kept after the pruning procedure.

By comparing the 2 inputted images - since the other 2 will be different per execution due to sampling - we can highlight the following points:

- About the **1/3 of the nodes have been removed** (initially there were 100, while after the pruning there are 69).
- Most of the edges have been removed something clearly visually noticed since the second graph is way less dense than the first one.
- Nodes like the **65** and **2857** have been removed since their edges weight does not fulfill the prune criterion.
- The **six nodes** in the middle of the cycle part of the graph maintain their spot to the post prune graph but lots of their edges were removed (the ones with weight equal to 1).
- On the other hand nodes like the **49406**, **53149**, **58857**, and **43073** (on the top left part) may not having the highest betweenness in the graph but their edges - validated through the dictionary print also - fulfill the given criterion. For more information you can see the documentation, but for example the node **49406** had lots of edges with weight 1 - see this depicted at the documentation - but only one edge with weight equal to 2 so kept after the prune. The previous mentioned edge connects the **49406** with the node **15004** - the one that can be noticed in the above graph.

## Conclusion

Task C wasn't just a task; it was a journey, a voyage through the intricate corridors of entity resolution. Commencing with the solid foundation of Task A and propelled by the revelations of Task B, Task C soared into the realms of advanced techniques like meta-blocking and pruning.

The initial phase of Task B, with its revelation of 2.6 billion comparisons, was akin to a spotlight, highlighting the vastness and intricacies of our data landscape. Such a mammoth figure underscored the challenges that lurked in the entity resolution domain. It also amplified the clarion call for further refinement.

Answering this call was Task C, armed with the Meta-Blocking graph. The graph, in its essence, was a tapestry, weaving together entities based on their shared blocks. The CBS weighting scheme added another dimension, infusing the tapestry with depth by attributing weights to each thread, each connection.

But what truly set Task C apart was the art of pruning. Through discerning choices, edges with weights less than two were gracefully excised, reducing computational demands and accentuating potential matches. This pruning phase was a testament to efficiency, a strategy to sharpen focus and enhance accuracy in the subsequent phases of entity resolution.

Re-evaluating the number of comparisons post this pruning, juxtaposed against the initial 2.6 billion, offered a tangible metric of our accomplishments. But beyond numbers, it painted a narrative of effectiveness, of the meta-blocking strategy's transformative potential.

In its entirety, Task C was both a challenge and a revelation. While it highlighted the intricacies of handling vast datasets, it also showcased the potential of innovative techniques in navigating this vastness. By integrating advanced methodologies, Task C not only ensured computational efficiency but also championed the cause of data accuracy and integrity.

## VI. Task D: Jaccard Similarity Function

In the final and intricate stage of the **Token Blocking method**, our objective takes a nuanced turn. We are venturing into the realm of assessing the likeness between different entity profiles. This might seem simple at first glance, but the intricacies lie in the details. Each entity, represented by its attributes and values, carries unique information. However, due to various reasons like data entry errors, different data sources, or alternate representations, multiple profiles might refer to the same real-world entity. This is where the challenge of *entity resolution* comes into play.

Entity resolution is pivotal in the vast domain of data management and analytics. The process of determining similarities between entities aids in deducing if two entities, despite potential differences in their data representations, refer to an identical real-world entity. This isn't just about finding exact matches but involves gauging the degree of similarity between entities. If two entities are found to have a high similarity score, especially using robust measures like the Jaccard Similarity, it becomes a compelling indication of their potential identity. Such a revelation is not just an academic exercise. Identifying that multiple records refer to the same real-world entity can have profound implications. It can assist businesses in avoiding redundancy, ensuring data accuracy, providing consistent customer experiences, and even in areas like fraud detection. Thus, the broader goal of data de-duplication and ensuring data accuracy becomes paramount in various data-driven decision processes.

### Steps for Actual Comparison Calculations

**Attribute Selection:** Selecting the right attributes is foundational. These attributes will be used to calculate the similarity between the two entity profiles. Their selection plays a pivotal role as it can significantly influence the accuracy and reliability of the resulting similarity score.

**Selecting an Appropriate Similarity Metric:** The nature of the attribute values dictates the choice of the similarity metric. For instance, for string values, metrics like the *Heinz distance* might be apt. For sets, the *Jaccard Similarity* is a popular choice, and for vectors, metrics like the *cosine similarity* might be relevant.

**Structuring the Problem:** Before jumping into crafting the function to compute similarities, it's crucial to frame the problem correctly. Decisions like whether the function should take in the entire dataset or just individual entities for comparison need to be made. Additionally, the format in which data will be provided to the function is also a significant consideration.

### Our Chosen Approach

In this specific exercise, we are centering our attention on a subset of the initial dataset, which results after the tokenization process of each column or attribute. Post-tokenization, each attribute essentially transforms into a set of tokens or words. Instead of encompassing all attributes in our similarity calculations, we're narrowing our focus to base our calculations predominantly on one attribute. The metric that has been zeroed in for this task is the renowned **Jaccard Similarity**.

The Jaccard Similarity offers a mathematical approach to gauge the similarity between two sets. For two sets, A and B, it's defined as:

$$J(A, B) = |A \cap B| / |A \cup B|$$

Here,  $|A \cap B|$  signifies the size of the intersection of the sets A and B. In other words, it's the count of elements common to both A and B. On the other hand,  $|A \cup B|$  represents the size of the union of the sets, i.e., the total count of distinct elements present in either A or B or both.

The resulting value of  $J(A, B)$  will always lie between 0 and 1. A value of 1 suggests that the sets A and B are identical, containing the exact same elements. Conversely, a value of 0 indicates that the sets A and B share no common elements.

### Answer of Task D:

Within this illustrative scenario, the similarity computation's domain narrows down to a particular segment of the original dataset. This segment is distinguished by its **tokenized attributes/columns**. By converting each attribute into a distinct set of tokens, the methodology emphasizes assessing similarity based on individual attributes. The **Jaccard Similarity** stands out as the preferred metric, adeptly gauging the overlap between token sets.

By adopting this streamlined approach, the **Token Blocking method** amplifies the precision in deriving similarity scores, shedding light on the likely overlap between potential entity profile matches. This methodological refinement encapsulates the essence of the method's concluding phase. It champions a detail-oriented strategy, anchored by deliberate attribute choice, metric synchronization, and thoughtful organization, all aiming to pinpoint genuine matches amidst the intricate web of entity profiles.

Thus, it's evident that the dataset utilized in computing the **Jaccard Similarity Coefficient** showcases token lists, distinctly maintained for each attribute. **These lists are meticulously cleansed of superfluous comma characters at the token level.** You can check the first row under the field venue "san, diego", in the initial data there is an extra comma after the word diego - "san, diego,".

Therefore our procedure works efficiently and since we delve deeper into our discussion on similarity analysis, it's pivotal to visualize the data structure we're working with. Now follows *Figure 13*, which offers a clear representation of the dataset, emphasizing the tokenized nature of the attributes tailored for Jaccard Similarity computations.

```
In [41]: # Overview of the Task D's dataset.
data_taskD.head(10)
```

id	authors	venue	year	title
0 1	[qd, inc]	[san, diego]	[none]	[11578, sorrento, valley, road]
1 2	[as, argon, jg, hannoosh]	[phil, mag]	[none]	[initiation, of, crazes, in, polystyrene]
2 3	[gh, hansen, ll, wetterberg, h, sjöström, o,...]	[the, histochemical, journal]	[1992]	[immunogold, labelling, is, a, quantitative, m...
3 4	[tm, hammett, p, harmon, w, rhodes]	[see]	[none]	[the, burden, of, infectious, disease, among, ...]
4 5	[jr, cogdell]	[new, directions, for, teaching, and, learning]	[1995]	[the, role, of, faculty, advising, in, science...
5 6	[wm, schmidt]	[to]	[none]	[the, zero, multiplicity, of, linear, recurren...
6 7	[ra, haats]	[none]	[2002]	[predictive, validity, of, kindergarten, scree...
7 8	[jr, norris, j, deisenhofer]	[san, diego:, academic]	[none]	[the, photosynthetic, reaction, center]
8 9	[f, bennour]	[res., rep., ceria, u., paris]	[none]	[f., diene, av, ndiaye, y., hachage, lin@aire...
9 10	[i, borg, jc, lingoes]	[none]	[1987]	[multidimensional, similarity, structure, anal...

Figure 13: Representation of Dataset Highlighting Tokenized Attributes for Jaccard Similarity Analysis.



## Creation of Jaccard Similarity Function and Analysis

In our ongoing exploration of entity profiles, we consistently employ the **Jaccard coefficient** as a crucial metric, aiming to understand the nuanced similarities between pairs of entities. What sets our approach apart is the meticulous attention we pay to each attribute independently. Consequently, during each pair's evaluation, we immerse ourselves in a comprehensive analysis, encompassing as many as four distinct **Jaccard Coefficients**. The bespoke function we've architected is primed to accept not only the unique IDs of entities from our dataset but also the entirety of the dataset itself, replete with the multifaceted attributes that constitute these entity profiles. Each attribute, in turn, is delineated as a list of tokens. Notably, in situations where there exists a conspicuous absence of data for a given attribute, a placeholder token, aptly named "none", is invoked to signify this void.

Furthermore, it's imperative to mention that whenever either of the entities in a given pair manifests a missing data point for a specific attribute, our approach prudently omits the calculation of the **Jaccard coefficient** for that attribute. Instead of leaving this oversight unaddressed, our system is designed to articulate a clarifying message, thereby ensuring transparency in our methodology.

We've encapsulated our findings in *Figure 14 and 15*, a visual representation that meticulously elucidates the outcomes of our Jaccard Coefficient calculations, emphasizing the distinctions across diverse attributes for a trio of entity profile pairs. The rationale behind our pair selection is twofold: the first three pairs are serendipitously chosen (*Figure 14: Jaccard Coefficient Analysis Across Attributes for Randomly Selected Entity Profile Pairs*), whereas the three following pairs are purposefully selected owing to the entities' shared blocks, a critical linkage to the meta-blocking process we executed in *Task C - some with 3 common blocks and not 2* (*Figure 15: Jaccard Coefficient Analysis Highlighting Entities with Shared Blocks, Referencing the Meta-Blocking Process from Task C*).

### Random Entities

We will provide two examples of three random entities:

```
In [41]: # An example of two random entities.
compute_jaccard_similarity(1, 2, data_taskD)

The Jaccard Similarity between entities 1 and 2 based on attribute authors is: 0.00
The Jaccard Similarity between entities 1 and 2 based on attribute venue is: 0.00
The Jaccard Similarity between entities 1 and 2 based on attribute year is: 0.00
The Jaccard Similarity between entities 1 and 2 based on attribute title is: 0.06

In [42]: # An example of two random entities.
compute_jaccard_similarity(66877, 66878, data_taskD)

The Jaccard Similarity between entities 66877 and 66878 based on attribute authors is: 0.00
The Jaccard Similarity between entities 66877 and 66878 based on attribute venue is: 0.00
The Jaccard Similarity between entities 66877 and 66878 based on attribute year is: 0.00
The Jaccard Similarity between entities 66877 and 66878 based on attribute title is: 0.00

In [43]: # An example of two random entities.
compute_jaccard_similarity(1565, 112, data_taskD) # Same year!

The Jaccard Similarity between entities 1565 and 112 based on attribute authors is: 0.00
The Jaccard Similarity between entities 1565 and 112 based on attribute venue is: 0.00
The Jaccard Similarity between entities 1565 and 112 based on attribute year is: 1.00
The Jaccard Similarity between entities 1565 and 112 based on attribute title is: 0.00

In [44]: # Same attribute check that the functions outputs 1 for all fields and work right as it does!
compute_jaccard_similarity(1, 1, data_taskD)

The Jaccard Similarity between entities 1 and 1 based on attribute authors is: 1.00
The Jaccard Similarity between entities 1 and 1 based on attribute venue is: 1.00
The Jaccard Similarity between entities 1 and 1 based on attribute year is: 1.00
The Jaccard Similarity between entities 1 and 1 based on attribute title is: 1.00
```

Figure 14: Jaccard Coefficient Analysis Across Attributes for Randomly Selected Entity Profile Pairs.



#### Entities that coincide inside the same blocks the most (3 times)

We will provide two examples of two entities that coincide inside the same blocks 3 times:

```

In [45]: # An example of two entities that coincide inside the same blocks 3 times.
compute_jaccard_similarity(1, 29089, data_taskD)

The Jaccard Similarity between entities 1 and 29089 based on attribute authors is: 0.00
The Jaccard Similarity between entities 1 and 29089 based on attribute venue is: 0.00
The Jaccard Similarity between entities 1 and 29089 based on attribute year is: 0.00
The Jaccard Similarity between entities 1 and 29089 based on attribute title is: 0.09

In [46]: # An example of two entities that coincide inside the same blocks 3 times.
compute_jaccard_similarity(352, 29089, data_taskD)

The Jaccard Similarity between entities 352 and 29089 based on attribute authors is: 0.00
The Jaccard Similarity between entities 352 and 29089 based on attribute venue is: 0.00
The Jaccard Similarity between entities 352 and 29089 based on attribute year is: 0.00
The Jaccard Similarity between entities 352 and 29089 based on attribute title is: 0.08

In [47]: # An example of two entities that coincide inside the same blocks 3 times.
compute_jaccard_similarity(496, 29089, data_taskD)

The Jaccard Similarity between entities 496 and 29089 based on attribute authors is: 0.00
The Jaccard Similarity between entities 496 and 29089 based on attribute venue is: 0.00
The Jaccard Similarity between entities 496 and 29089 based on attribute year is: 0.00
The Jaccard Similarity between entities 496 and 29089 based on attribute title is: 0.00
  
```

Figure 15: Jaccard Coefficient Analysis Highlighting Entities with Shared Blocks, Referencing the Meta-Blocking Process from Task C.

Within the framework of our analytical setup, the articulation of entity profile attributes is strategically framed as either lists or sets of tokens. This deliberate configuration naturally gravitated us towards adopting the Jaccard similarity measure as our chosen metric. A pivotal crossroads we encountered pertained to the criteria for determining similarity. On one hand, we contemplated a probabilistic trajectory, wherein a Jaccard similarity surpassing a predetermined threshold would be deemed sufficient. Conversely, a deterministic paradigm beckoned, where an unequivocal Jaccard similarity of 1 would be necessitated for the amalgamation of entities.

With respect to our function's structural design, we've adopted an expansive approach. Armed with the reference IDs of entity profiles and the encompassing dataset, our function embarks on a deep dive into Jaccard similarity computations that are attribute-centric. The culmination of this intricate process hinges on a comprehensive decision, one that synergistically amalgamates the collective similarity spread across all individual attribute domains.

Our initial illustrative examples, characterized by three randomly selected entity pairs (Figure 8), unravel a discerning revelation: despite rigorous evaluations, the entities under consideration do not resonate in harmony and remain distinct. The "year" attribute, in particular, emerges as a challenging element due to its pervasive missing values or how high is the probability to match (see Figure 8), prompting us to judiciously possible exclude it from our comparative analysis.

In a later example (Figure 9), our choice of entities was deliberate. We paired entities (three pairs), acutely aware of the shared blocks they possessed, a connection traced back to Task C. Yet, the findings were unanticipated. Post a detailed scrutiny of their attributes and subsequent Jaccard similarity computations, it became unequivocally clear that these entities weren't identical. This revelation beckons a potential re-evaluation of our foundational strategy. Given the potency of the Token Block method within the confines of this dataset, it's plausible that the threshold defining shared blocks should be elevated. By making this strategic shift, we foresee a substantial decline in candidate comparisons, without compromising on precision, ensuring that true entity matches aren't inadvertently missed.

## Using the Jaccard Similarity Function with Emphasis on Title

In the broader scope of our analysis, while we considered various attributes of the entities, we recognize the significance of the 'title' attribute. Titles, in many datasets, capture the essence and primary information of entities. Therefore, when determining the similarity between entities, it's imperative to understand how their titles align.

Codewise, we modify our prior created function in order to return the titles of the selected entities along with their Jaccard similarity score.

```
In [57]: def compute_jaccard_similarity_title(entity_1, entity_2, data_frame):

    # Extracting titles for the entities.
    title_1 = data_frame['title'].iloc[entity_1]
    title_2 = data_frame['title'].iloc[entity_2]

    # Check for missing values.
    if title_1 == "none" or title_2 == "none":
        print("For the attribute 'title', there are missing values; hence, Jaccard coefficient cannot be calculated.")
        return

    # Compute Jaccard Similarity for the titles.
    intersection = len(set(title_1) & set(title_2))
    union = len(set(title_1) | set(title_2))
    jaccard_coefficient = intersection / union

    # Print the results.
    print(f"The Jaccard Similarity between entities {entity_1} and {entity_2} based on the title attribute is: {jaccard_coefficient}")
    print(f"Title of entity {entity_1}: {title_1}")
    print(f"Title of entity {entity_2}: {title_2}")
```

Figure 16: Jaccard Similarity function with emphasis on the title attribute.

The updated function is designed to iterate over each attribute, while it provides a specialized focus on the **title** attribute. The rationale behind this emphasis is the unique and significant information that titles often encapsulate about an entity. By computing the Jaccard Similarity specifically for titles, we can discern how similar two entities are in terms of their content or thematic subject matter.

To provide a clear understanding of the Jaccard Similarity function's application on the title attribute, we will compare two random entities on the following examples:

```
In [58]: import random

# Random Example 1.
# Select two random entities from the dataset.
entity_index_1 = random.randint(0, len(data_taskD) - 1)
entity_index_2 = random.randint(0, len(data_taskD) - 1)

# Comparing two random entities and their titles.
compute_jaccard_similarity_title(entity_index_1, entity_index_2, data_taskD)

The Jaccard Similarity between entities 55012 and 3243 based on the title attribute is: 0.00
Title of entity 55012: ['intellectual', 'capital', 'and', 'wealth', 'creation:', 'a', 'case', 'study', 'of', 'total', 'qualit', 'y', 'management', 'at', 'a', 'large', 'cargo']
Title of entity 3243: ['molecules', 'in', 'carbon', 'nanotubes.']

In [59]: # Random Example 2.

# Select two random entities from the dataset.
entity_index_1 = random.randint(0, len(data_taskD) - 1)
entity_index_2 = random.randint(0, len(data_taskD) - 1)

# Comparing two random entities and their titles.
compute_jaccard_similarity_title(entity_index_1, entity_index_2, data_taskD)

The Jaccard Similarity between entities 39505 and 6553 based on the title attribute is: 0.00
Title of entity 39505: ['a', '3d', 'neural', 'network', 'for', 'business', 'forecasting']
Title of entity 6553: ['data', 'grid', 'management', 'systems']

In [60]: # Random Example 3.

# Select two random entities from the dataset.
entity_index_1 = random.randint(0, len(data_taskD) - 1)
entity_index_2 = random.randint(0, len(data_taskD) - 1)

# Comparing two random entities and their titles.
compute_jaccard_similarity_title(entity_index_1, entity_index_2, data_taskD)

The Jaccard Similarity between entities 957 and 41316 based on the title attribute is: 0.12
Title of entity 957: ['application', 'of', 'hla', 'based', 'solutions', 'for', 'modeling', 'and', 'control', 'of', 'chemical', 'plants']
Title of entity 41316: ['method', 'for', 'the', 'synthesis', 'of', 'interactive', 'system', 'specifications.']
```

Figure 17: Comparison of two random entities using the Jaccard Similarity metric, focusing on the 'title' attribute. This visualization illustrates the degree of overlap or similarity between the titles of the entities.

## Jaccard Similarities for Titles Insights

Here are some key insights derived from the previous figure's comparisons:

**Vast Diversity in Titles:** In the first two examples, the Jaccard Similarity is 0.00, indicating that there's no word overlap in the titles of the compared entities. This showcases the vast diversity in titles within the dataset.

**Rare Overlaps:** The third example has a Jaccard Similarity of 0.12, which means there's a slight overlap in the words of the titles. However, even in this case, the overlap is minimal, suggesting that identical words in titles don't necessarily indicate a thematic similarity between entities.

**Importance of Common Words:** The word "of" appears in both the titles of the third example. While such common words can increase the Jaccard Similarity, they might not contribute meaningfully to the overall thematic similarity. Hence, caution is needed when interpreting results based solely on this metric.

These insights underline the importance of the **title** attribute in differentiating entities and also highlight the potential challenges in relying solely on Jaccard Similarity for understanding thematic overlap (***disclaimer our trial examples were random but as many as we did, we got similar results with low Jaccard score and similarity based mainly on articles***).

## Observation on Jaccard Similarity for Titles

Analyzing the Jaccard Similarity on the 'title' attribute of various entity pairs reveals a telling narrative. The bulk of our comparisons yield similarity percentages that are strikingly low, often veering close to 0. This observation underscores that the dataset's titles are predominantly unique, showcasing minimal overlap in terms of their individual words.

Yet, there are anomalies. In certain comparisons, we observe marginally higher similarity values. Diving deeper into these instances, we identify that the uptick in similarity often stems from common words, especially articles like "the", "an", and "a". Such frequently occurring words, while increasing the mathematical similarity, might not lend significant thematic overlap between the titles.

In summary, while the Jaccard Similarity metric is instrumental in quantifying word overlap in titles, it's imperative to discern the true nature of this overlap. Over-reliance on the metric, without considering the context of common or generic words, can misguide the interpretation. This becomes especially salient when leveraging title-based similarity for intricate tasks like entity resolution, where understanding the true thematic resemblance is paramount.

## Conclusion

**Task D** delved deep into the intricacies of the **Jaccard Similarity** metric within the landscape of entity resolution. In our specific setup, entities' attribute values manifested as sets or lists of tokens. Such a structure naturally steered us towards the **Jaccard similarity**, a tool that effectively measures the similarity between sets, thereby aiding in determining if entities mirror the same real-world object.

The early examples served as an eye-opener, revealing that certain entity profiles, despite being juxtaposed for comparison, did not resonate in similarity. Some attributes, notably "year", posed challenges, predominantly due to their riddled missing values or easily to get matched, compelling us to exclude them from the similarity evaluations.

Moreover, a poignant observation surfaced when two entity profiles, already known to share a considerable number of blocks from **Task C**, were analyzed. Their **Jaccard similarity** did not corroborate them as identical entities. Such revelations hint towards a potential recalibration: within the context of this dataset, the **token block method** might be more efficacious with a stringent threshold. By amplifying this threshold, we stand a chance to prune candidate comparisons without eroding efficiency. The goal is to strike a balance, ensuring that in our pursuit to reduce comparisons, we don't inadvertently inflate false negatives - those elusive identical entity pairs that slip detection.

Navigating through Task D was akin to embarking on an enlightening journey, one that underscored the meticulous strategy requisite for adept entity resolution. Through a tapestry of theoretical exposition, hands-on examples, and astute observations, the task illuminated the core mechanics and applications of the **Jaccard Similarity metric**.

In retrospect, the emphasis on the 'title' attribute has elucidated the subtleties and challenges inherent in entity resolution. As we've seen, while metrics like Jaccard Similarity offer quantifiable insights, their interpretation, especially when analyzing titles, demands a nuanced approach. As future work unfolds, refining the thresholds and understanding the intrinsic semantics of titles could pave the way for more accurate and discerning entity resolutions.

## VII. Final Analysis Conclusions

The challenge of entity resolution stands prominently at the forefront of data integration endeavors. Ensuring data accuracy and functionality in a centralized data warehouse is non-negotiable, with entity resolution playing a pivotal role in achieving this feat. The ultimate aim remains clear-cut: a unified data repository where entities, along with their linked information, are uniquely identified and unambiguously represented.

To sum up, our Journey Through Numbers and Insights:

**Task A - Token Blocking:** Our expedition commenced with **Token Blocking**. Here, distinctive Blocking Keys (BKs) served as the backbone, facilitating the creation of cohesive blocks. This method, unbiased by schema, set the cornerstone for the subsequent meticulous entity resolution stages.

**Task B - The Magnitude of Comparisons:** Transitioning to Task B, we were met with the enormity of the computational challenge ahead. A staggering 2,648,668,047 comparisons were computed, underscoring the intricacies embedded in the task of duplicate detection.

**Task C - Meta-Blocking, CBS Weighting, and Pruning:** Our next stride led us to the construction of the Meta-Blocking graph using the **CBS Weighting Scheme**. In our focused study on a subset of just 100 entities, aiming to discern the pairs from the overarching set of 66,879 entities (as unveiled in Task A), we noted an interesting outcome. From the billions of comparisons in Task B, the number dwindled to a mere **1127 final comparisons** in this subset study. Yet, a word of caution is in order. Our analysis, confined to a notably small sample size, might not reflect the exact numbers for the entire dataset (100 entries). Nevertheless, the substantial reduction in comparisons, even within this microcosm, cannot be overlooked. It attests to the potential efficiencies we can harness, ensuring a balance between thoroughness and computational feasibility.

**Task D - Jaccard Similarity:** Our final task was emblematic of the entire process. We designed a custom function to gauge the Jaccard Similarity between entities, focusing specifically on the attribute "title". This measure, which quantifies the similarity between sets, exemplified the nuanced approach required for effective entity resolution. However, the analysis underscored the intricacies of entity resolution using the 'title' attribute, revealing that while Jaccard Similarity provides measurable insights, its nuanced interpretation is crucial for enhanced accuracy in future resolutions.

## Concluding Thoughts

As we wrap up our exploration into advanced techniques for entity resolution and duplicate detection, we stand enriched by our findings and experiences. We've traversed a landscape marked by complex methodologies, state-of-the-art techniques, and intricate challenges. The lessons imbibed, the insights garnered, and the methodologies explored will undoubtedly serve as guiding beacons in our future endeavors in the realm of data management. In our data-dominant era, entity resolution isn't merely a computational task; it's a commitment to upholding the sanctity and quality of data.

To be more precise, the process of **entity resolution** stands as an indispensable cornerstone in the vast landscape of data management. In every scenario, especially before the crucial step of data integration, this process emerges as both vital and computationally intensive. As we advance into an era dominated by colossal volumes of data, ensuring the integrity and quality of data warehouses becomes paramount.

Throughout the evolution of data science, myriad methodologies have been pioneered to address this challenge. The overarching goal has always been twofold: First, to maximize **Efficiency** by accurately detecting genuine matches between entity profiles that correspond to identical real-world objects. Second, to optimize **Quality** by ensuring these matches are pinpointed with minimal computational overhead, thereby reducing the number of actual comparisons made.

Among the plethora of techniques at our disposal, the **Token Blocking method** shines brightly. This innovative approach not only serves the dual purpose of enhancing efficiency and quality but also strikes a harmonious balance between them. By judiciously reducing the computational burden while maintaining the accuracy of entity matches, Token Blocking exemplifies a nuanced strategy in the ever-evolving challenge of entity resolution.

Regarding the final task, building on the exploration of entity resolution with a focus on the **title** attribute, it's evident that while quantitative metrics like Jaccard Similarity are invaluable, their true potential is unlocked when applied with thoughtful interpretation, ensuring both precision and depth in discerning entity relationships.

## VIII. Deliverables

To ensure a thorough understanding of the process and the methodologies employed in our entity resolution tasks, the following deliverables are provided:

1. **Documentation Report:** This document provides an in-depth description of the assignments tasks, elucidating the methodology, findings, and implications. It serves as a comprehensive guide for readers to understand the significance of the entity resolution process.
2. **Analysis Jupyter Notebook:** Alongside this documentation, we are providing a Jupyter Notebook that contains the actual code implementation and analysis. The notebook, titled "**advanced\_techniques\_for\_entity\_resolution\_and\_duplicate\_detection.ipynb**", is developed using *Python 3.10*. It offers a hands-on view of the techniques described in this documentation, allowing readers to delve into the actual coding and analytical practices employed during the process.

These deliverables are designed to offer a holistic view of the entity resolution tasks, combining both theoretical insights from the documentation and practical application from the Jupyter Notebook. By leveraging both, readers can gain a profound grasp of the intricate subject matter at hand. **Therefore, please have a look on both deliverables files.**