

## Grouping 10,000 CSV records

comparing PowerShell and MySQL

Mauricé Ricardo Bärisch

Term Paper in course type Computer Science

Supervisor: Prof. Dr. Stefan Schiffner

Submission Date: May 16, 2024

I confirm that this thesis type is my own we	Disclaimer ork and I have documented all sources and material used.
Munich, May 16, 2024	Author

## **Abstract**

Performing aggregation and visualization on large datasets has become a major part in Information Technology. Such datasets are often exported by a third party software or service provider, and arrive in a transmittable format like JSON or CSV. This term paper compares the grouping operation on large CSV files between PowerShell's Group-Object command and MySQL's GROUP BY statement in regards of their interface, underlying algorithm and RAM usage. On two identical machines limited in RAM, we simulate how PowerShell's and MySQL's grouping operations perform. As a result, we formulate best practices for aggregating large CSV files.

Keywords: Database, Aggregation, Grouping, CSV

## **Contents**

1	Intr	oduction	1
	1.1	Motivation	1
	1.2	Main contributions	1
2	Bacl	ground	2
	2.1	Comma-Separated Values (CSV)	2
	2.2		2
		2.2.1 Relational Schema	2
		2.2.2 Selection and Projection	3
	2.3	Grouping	4
		2.3.1 Split into groups	4
			4
3	Sim	ulation	5
	3.1		5
	3.2		5
	3.3		5
4	Con	clusion	5
-	4.1		5
	4.2	±	5
5	Gen	eral Addenda	6
•	5.1		6
6	Figu	res	6
	6.1	Example 1	6
	6.2	•	6
Gl	lossar	y	8
p:	bliogi	anhy	10
וע	UHUZI	արույ	··

## 1 Introduction

Nowadays, many processes that handle large amounts of data (not to confuse with Big Data<sup>1</sup>) already exist. Still, they must be maintained, verified and optimized to ensure both quality and efficiency.

#### 1.1 Motivation

At q.beyond, we had to deal with just that: large monthly CSV exports (each file contains the data of a month) by a third-party service provider that our process was operating on. For verification, we collected the exports from a whole year, grouped the records by a column and exported each group as an individual Excel file. Each file contained two sheets: one with an aggregated monthly overview and the other with all the individual records of that particular group. In 2.3, we will explain why those two sheets require grouping, but only the latter requires every particular record to be stored (in memory or storage).

We ended up writing a PowerShell<sup>2</sup> Script, which, in our first draft, crashed the Azure Virtual Desktop (AVD) it was running on due to insufficient memory.

#### 1.2 Main contributions

Based on this experience, we can define the purpose of this paper:

- **Problem**: We need to aggregate CSV data bigger than the available memory
- Objectives:
  - 1. Find a solution that can group large CSV data
  - 2. Proof by tests and simulation that the solution works
- Questions:
  - 1. Why does the Group-Object command not work on limited memory?
  - 2. What are requirements for a grouping algorithm in order to run on low memory?

As an alternative solution, we will consider MySQL, a Database Management System (DBMS).

<sup>&</sup>lt;sup>1</sup>Data with high variety, volume and velocity. Cannot be processed by conventional data processing software.

<sup>&</sup>lt;sup>2</sup>A scripting language coming with Windows. Used for administrative and automation tasks.

## 2 Background

### 2.1 Comma-Separated Values (CSV)

Specified in Shafranovich (2005), CSV is a file format for storing a table in plain text. It has the following characteristics:

- 1. each row starts on a new line
- 2. the values for each column are separated by a delimiter (usually a comma)
- 3. there may be a header line those values specify the names of the columns
- 4. whitespaces other than line breaks for 1. are ignored, unless part of a value

Our CSV export is structured in exactly this way. Let it be<sup>3</sup>:

```
tid,
                   amount,
          pname,
                                comment
2 1,
          Alice,
                   -1,
                                daily expenses
          Bob,
                   2,
                                pocket money
3 2,
4 3,
          Alice,
                  1,
                                friend
5 4,
          Carl,
                   0,
                                dummy
                  2.99,
                                card game
6 5,
          Dan,
```

Listing 1: CSV export

#### 2.2 Relational Model

Introduced by E. F. Codd (1970), the relational model is a mathematical description of data management, structuring data in tuples and tuples in relations. A relational algebra is provided to act as a simple mathematical query language. Industrial query languages like SQL later implemented the features described in the relational model.

As further explained by Arenas et al. (2022), a relation is just a pair of a relation schema and a set of tuples (the rows of the relation).

#### 2.2.1 Relational Schema

Arenas defines the relation schema based on three characteristics:

- relation name: the name of the table
- relation attributes: a list of uniquely identifiable attributes/columns
- relation arity: the number of attributes/columns

<sup>&</sup>lt;sup>3</sup>In this paper, we will use a simplification and abstraction of the original data mentioned in 1.1. The "..." indicates that it is indeed a large CSV with many more records.

Schweikardt et al. (2010) defines the relation schema more formally as a database-wide mapping from relation names to either their respective arity (unnamed perspective) or their respective attributes (named perspective). Then, a row in a relation is either defined as a tuple of values  $a = (a_1, \ldots, a_n)$  (unnamed perspective) or a set of pairs  $a = \{(c_1, a_1), \ldots, (c_n, a_n)\}$ , where each pair consists of the column name  $c_n$  and the value  $a_n$  assigned to it (named perspective).

In this paper, we will use relations as described in the named perspective. For that, we must ensure that each column got a unique name. Luckily, the header from our CSV export statisfies this condition (see listing 1). Let its schema be:

```
Transaction [ tid, pname, amount, comment ]
```

Note that we do not assign any types to the attributes, as the CSV header row only consists of attribute names. This conforms with the understanding of the relational model from Schweikardt, who just assumes that the values for our columns origin from an infinite set containing any values for any columns.

The visual representation of a relation is a table. Below figure shows the corresponding relation to the first 5 records of the CSV export from listing 1.

tid	pname	amount	comment
1	Alice	-1	daily expense
2	Bob	2	pocket money
3	Alice	1	friend
4	Carl	0	dummy
5	Dan	2.99	card game

Table 1: Transaction

#### 2.2.2 Selection and Projection

In chapter 12.1, Halpin & Morgan (2008) introduces relational algebra. It includes six comparison operators  $\theta \in \{=,<>,<,>,\leq,\geq\}$  as well as the three logical operators  $\{\land,\lor,\neg\}$ . Furthermore, some language-specific operations are defined. We will focus on the selection and projection.

Given a relation R, the selection operator  $\sigma_c(R)$  selects all rows from that relation that match the condition c (similar to the **where** clause in SQL). Such condition can consist of any valid expression using comparison operators, logical operators, attributes and constants.

For example, the selection  $\sigma_{\text{pname}=\text{Alice}}$  (Transaction) returns the following relation:

tid	pname	amount	comment
1	Alice	-1	daily expense
3	Alice	1	friend

Table 2: Selection of *pname=Alice* on Transaction

Similarly, the projection operator  $\pi_{a1,...,a_n}(R)$  projects all columns from a relation R stated in the projection list  $a1,...,a_n$ , and ignores all others. Consider the query  $\pi_{\text{pname}}(\text{Transaction})$ :

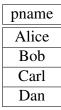


Table 3: Projection of pname on Transaction

Note that the projection removes fully duplicate rows, hence why Alice only exists once in the resulting relation.

As the visualizations suggests, both the selection and projection return a relation.

### 2.3 Grouping

Since a projection removes duplicated rows, projecting only the grouping column (the column whose distinct values should be the name of our groups), results in a relation containing all group names (see table 3).

#### 2.3.1 Split into groups

We can use these group names to split the relation into distinct sub-relations consisting only of the elements with the same group name:

$$G := \{ \sigma_{\text{pname}=x}(\text{Transaction}) \mid (\text{pname}, x) \in \pi_{\text{pname}}(\text{Transaction}) \}$$

with *G* being a set of (sub)relations. We will refer to this kind of grouping as **split into groups**. An algorithm which splits a relation into groups hence must put every row into its respective group. In other words: a split into groups requires all rows to be loaded in memory (or storage).

Apart from this observation, our definition of G does not further specify the algorithm a database could use to perform the split. Usually, databases simply sort the relation on the grouping attribute, such that rows belonging to the same group are next to each other.

#### 2.3.2 Aggregate on groups

Oftentimes, we do not require all the details of all the rows in a group. Instead, we only want a summary of our grouped data, also known as aggregation. On p. 48, while listing features missing in traditional relational algebra, Arenas et al. (2022) find a great explanation for aggregation and grouping:

An extremely common feature of SQL queries is the use of aggregation and grouping. Aggregation allows numerical functions to be applied to entire columns, for example, to find the total salary of all employees in a company. Grouping allows such columns to be split according to a value of some attribute (...)

They also formally define and analyze aggregation as an extension to relational algebra, but we will

- 3 Simulation
- 3.1 Methods
- 3.2 Implementation
- 3.3 Results
- 4 Conclusion
- 4.1 Comparison
- 4.2 Future Work

## 5 General Addenda

If there are several additions you want to add, but they do not fit into the thesis itself, they belong here.

## 5.1 Detailed Addition

Even sections are possible, but usually only used for several elements in, e.g. tables, images, etc.

- 6 Figures
- 6.1 Example 1
- 6.2 Example 2

# **List of Figures**

## **List of Tables**

1	Transaction	3
2	Selection of <i>pname=Alice</i> on Transaction	3
3	Projection of <i>pname</i> on Transaction	4

## Glossary

**AVD** Azure Virtual Desktop. 1

**DBMS** Database Management System. 1, 8

**SQL** The Structured Query Language is a popular query language based on the relational model. It is used by many known relational DBMS. 3

L	is	ti	n	gs
	-~		:	_~

1	CSV export	 . <b>.</b>	 	 		 						 	2

## References

- ARENAS, MARCELO; PABLO BARCEL'O; LEONID LIBKIN; WIM MARTENS; und ANDREAS PIERIS. 2022. *Database theory querying data*. Preliminary Ed. San Francisco, CA, USA: Santiago Paris Bayreuth Edinburgh.
- CODD, E. F. 1970. A relational model of data for large shared data banks. *Commun. ACM* 13.377–387. URL https://doi.org/10.1145/362384.362685.
- HALPIN, TERRY, und TONY MORGAN. 2008. *Information modeling and relational databases*. 2 Ed. San Francisco, CA, USA: Morgan Kaufmann Publishers Inc.
- SCHWEIKARDT, NICOLE; THOMAS SCHWENTICK; und LUC SEGOUFIN. 2010. *Database theory:* query languages, 19. 2 Ed. Chapman & Hall/CRC.
- SHAFRANOVICH, YAKOV. 2005. Common Format and MIME Type for Comma-Separated Values (CSV) Files. RFC 4180. URL https://www.rfc-editor.org/info/rfc4180.