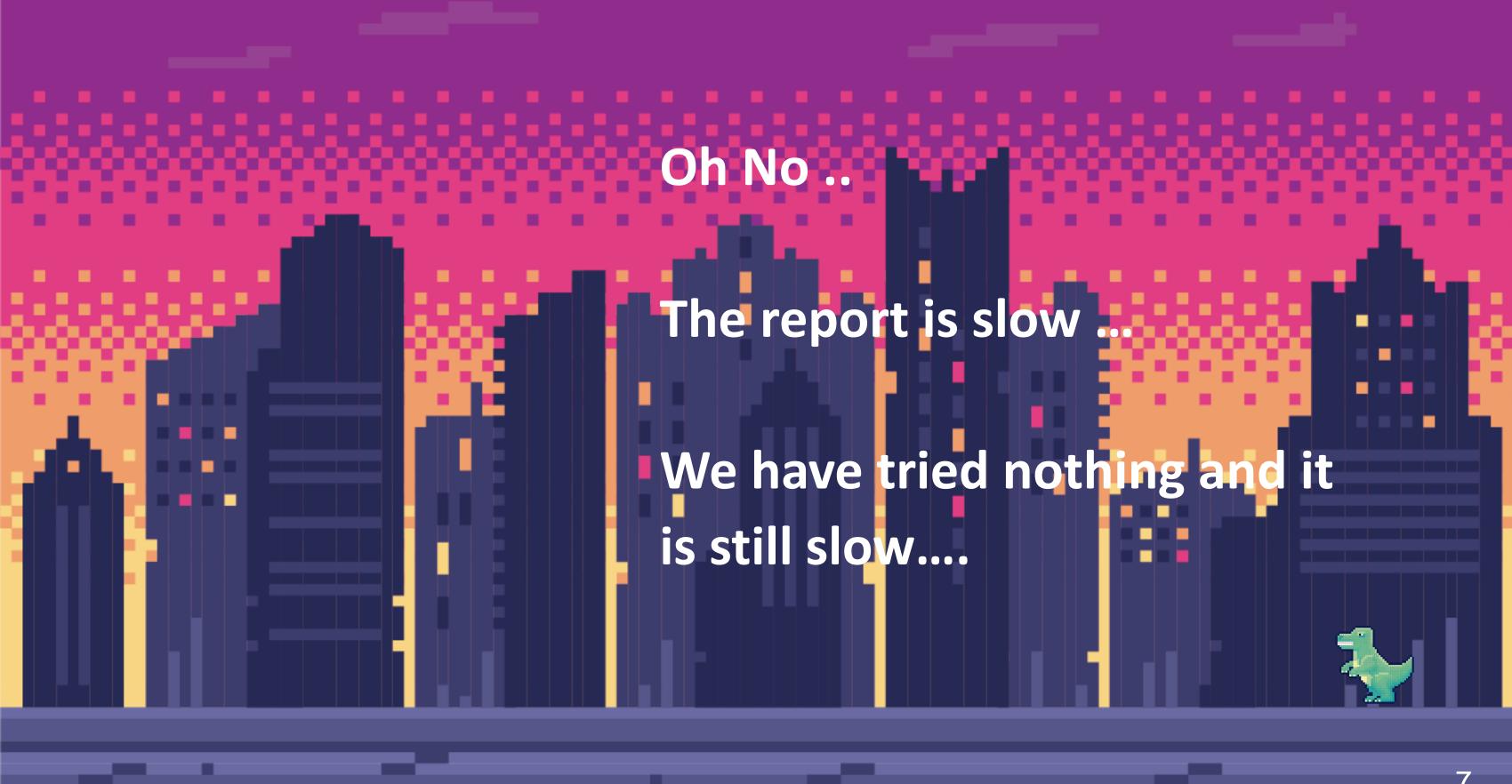


Reporting on 1.4B rows with Serverless SQL Pools

Life is good





Introduce yourself Say hello to the team.



Stijn Wynants Senior FastTrack Engineer



Mark Pryce-Maher Senior Program Manager



Benni De Jagere Program Manager

What is Serverless SQL Pools

- Our Unicorn
- Polaris
- On-Demand
- Pay per TB

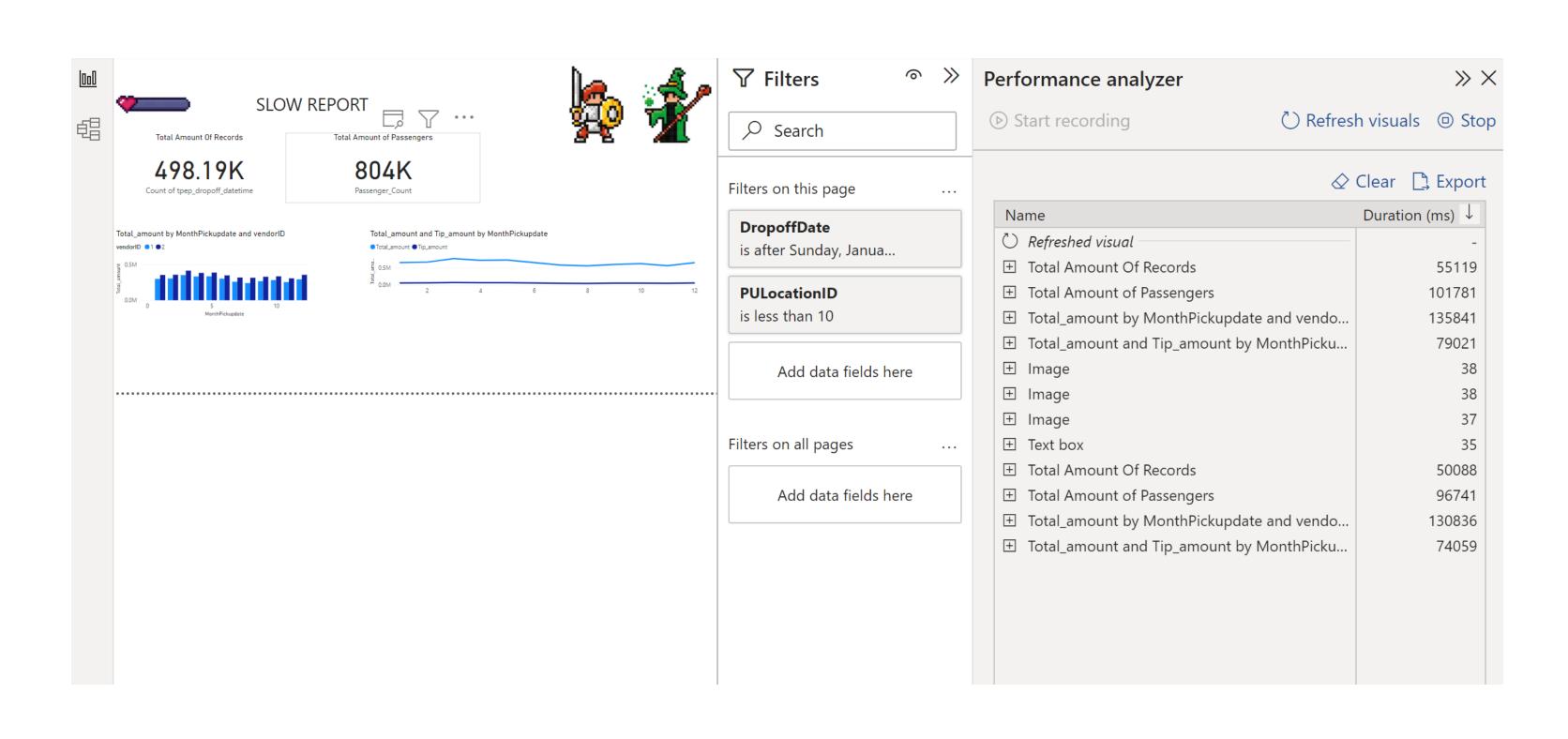




Let's look at the original

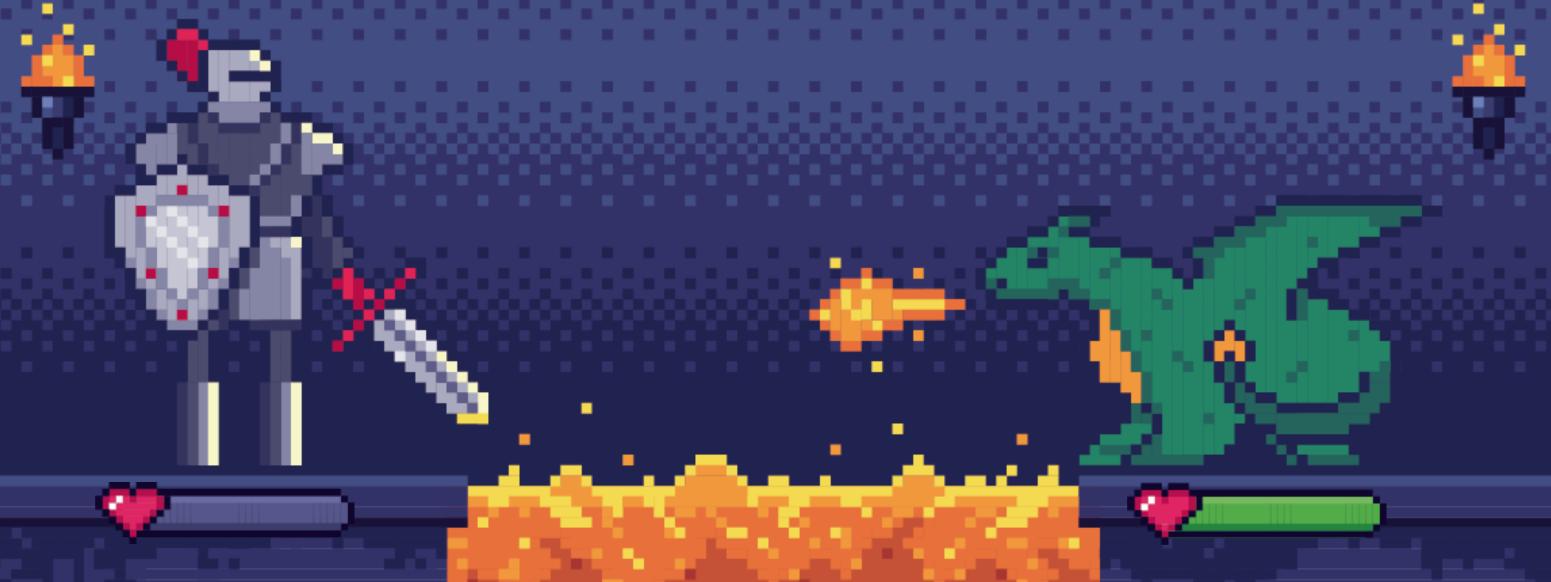
- CSV List
- Any Calculation will trigger a full scan of the CSV Files
- Can be used for import, direct query will be slow as the volume is large
- Optimizations: Statistics, Partitioning





Let's start Optimizing...

Parquet Files.



Parquet Files







- Spark
- Mapping Data Flow
- Copy Data Task

Compresses data in a columnar format

Contains metadata about the files

Built out of Rowgroups & Column Chunks

Serverless can eliminate column chunks, rowgroups



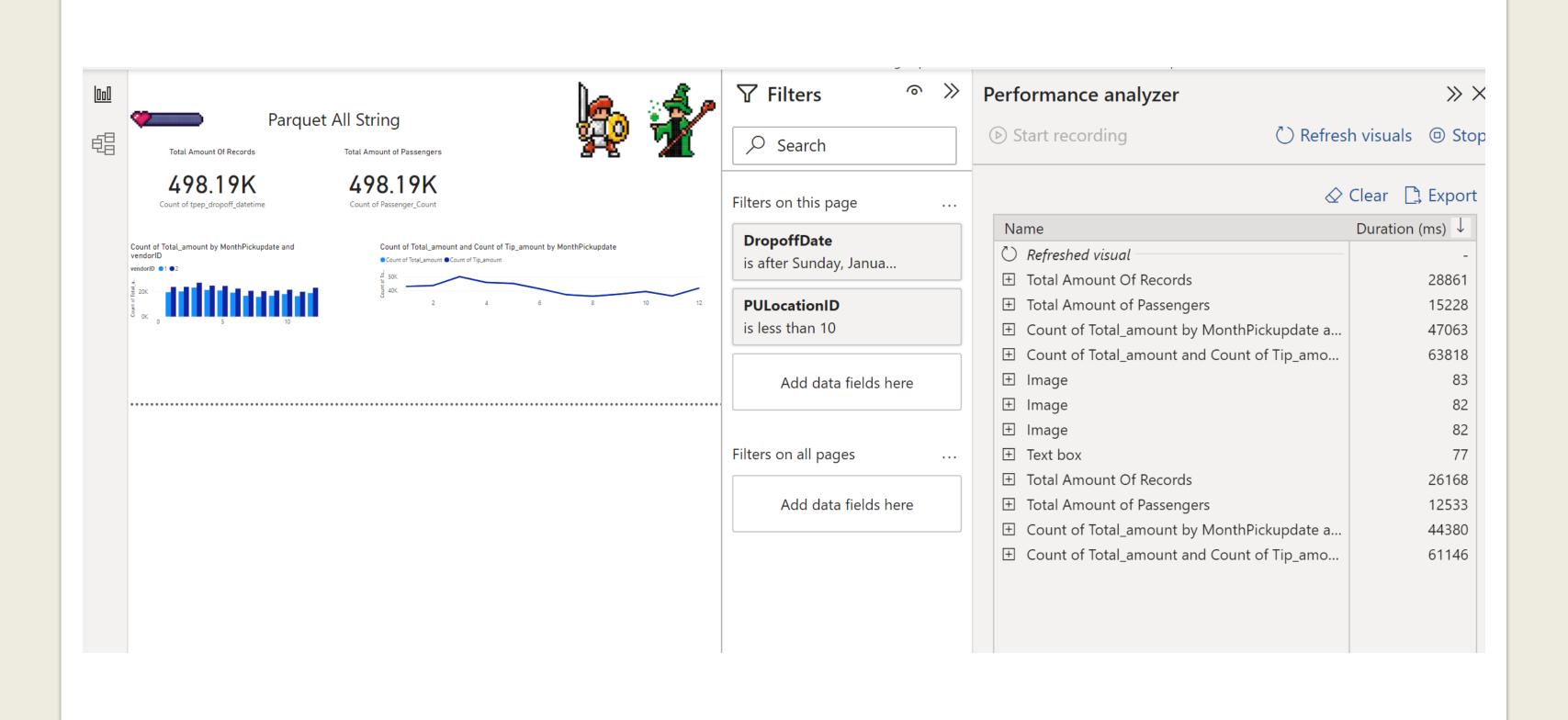
Header Data Body Footer



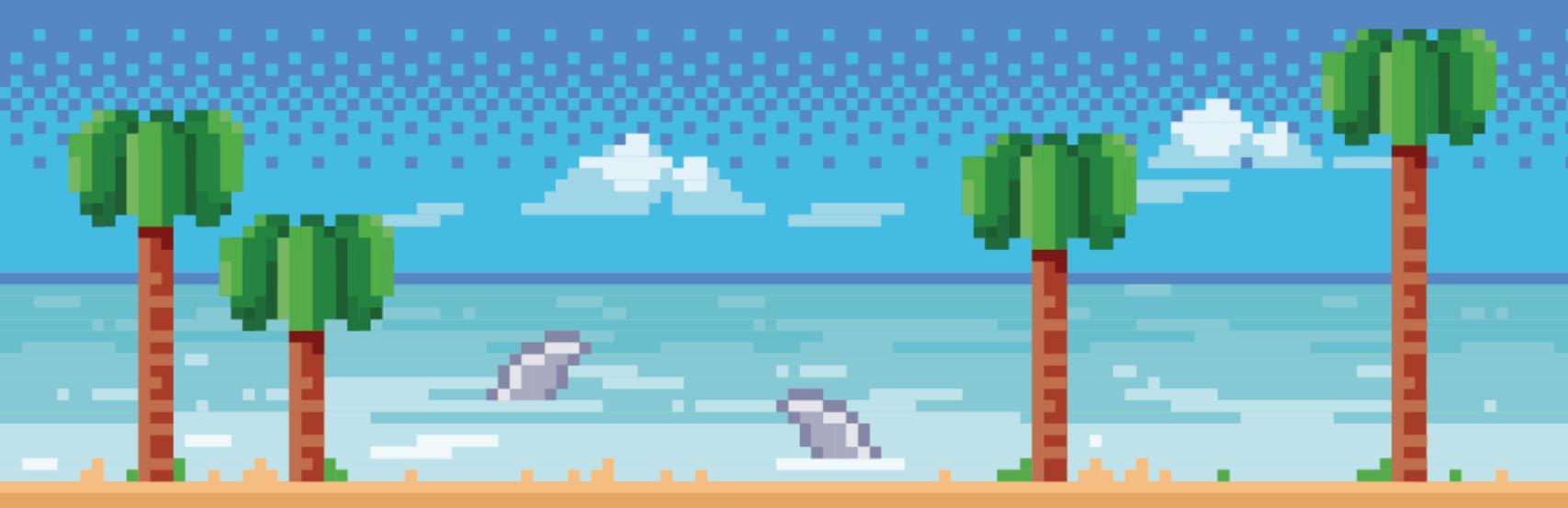
Data Body	
Data Block 1	
Row Group 1	
Columnar Chunk	Columnar Chunk
Row Group 2	
Columnar Chunk	Columnar Chunk
Data Block 2	
Row Group 1	
Columnar Chunk	Columnar Chunk
Row Group 2	
Columnar Chunk	Columnar Chunk



Let's start Optimizing... DEMO



Well that is already a big difference, but can we do more?



Let's keep on Optimizing...

Data Types.



Data Types









Use the smallest possible data types; smallint

If character field is fixed width use char,nchar or nvarchar,varchar

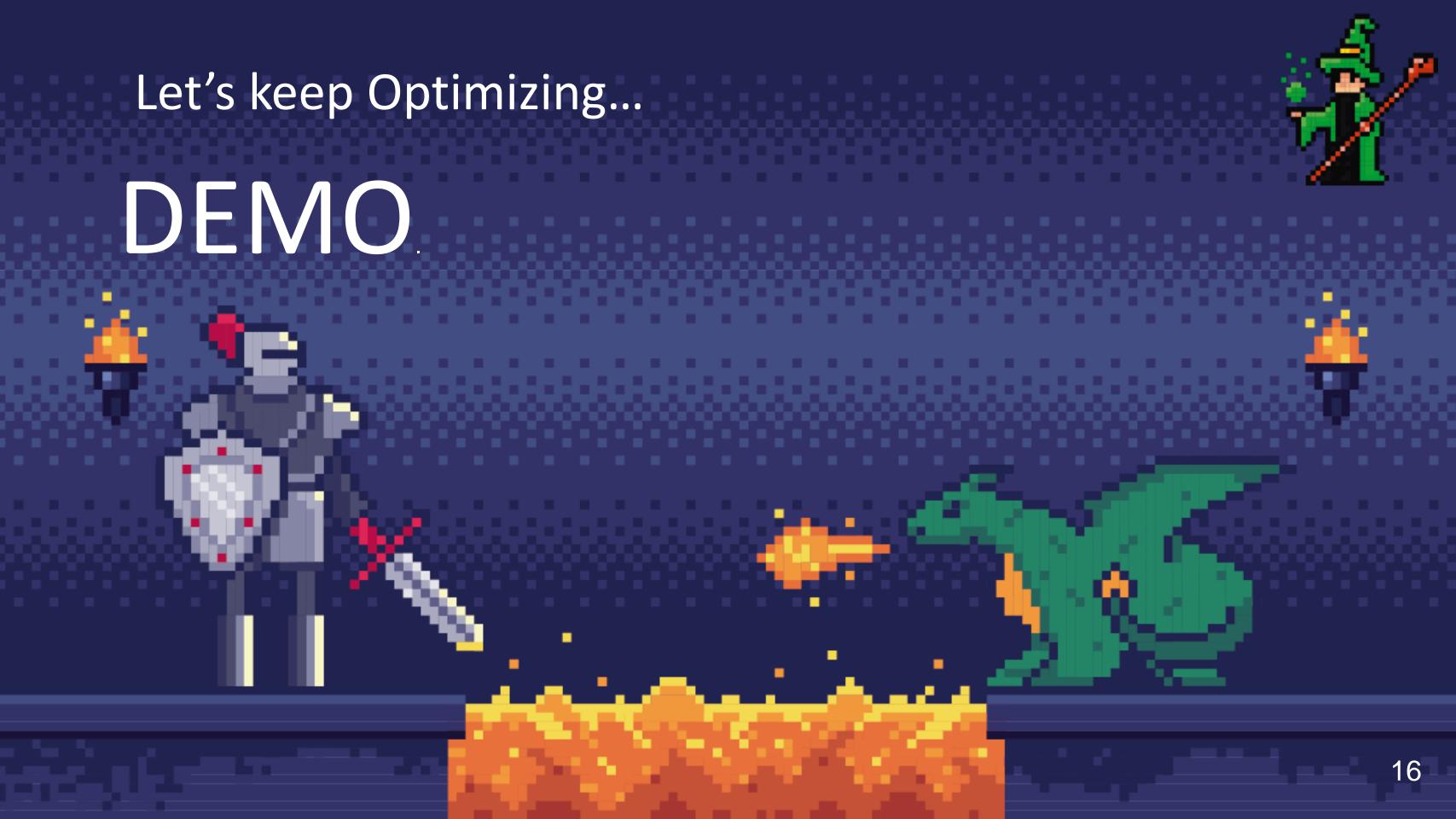
Use varchar and char over nvarchar, ncahr

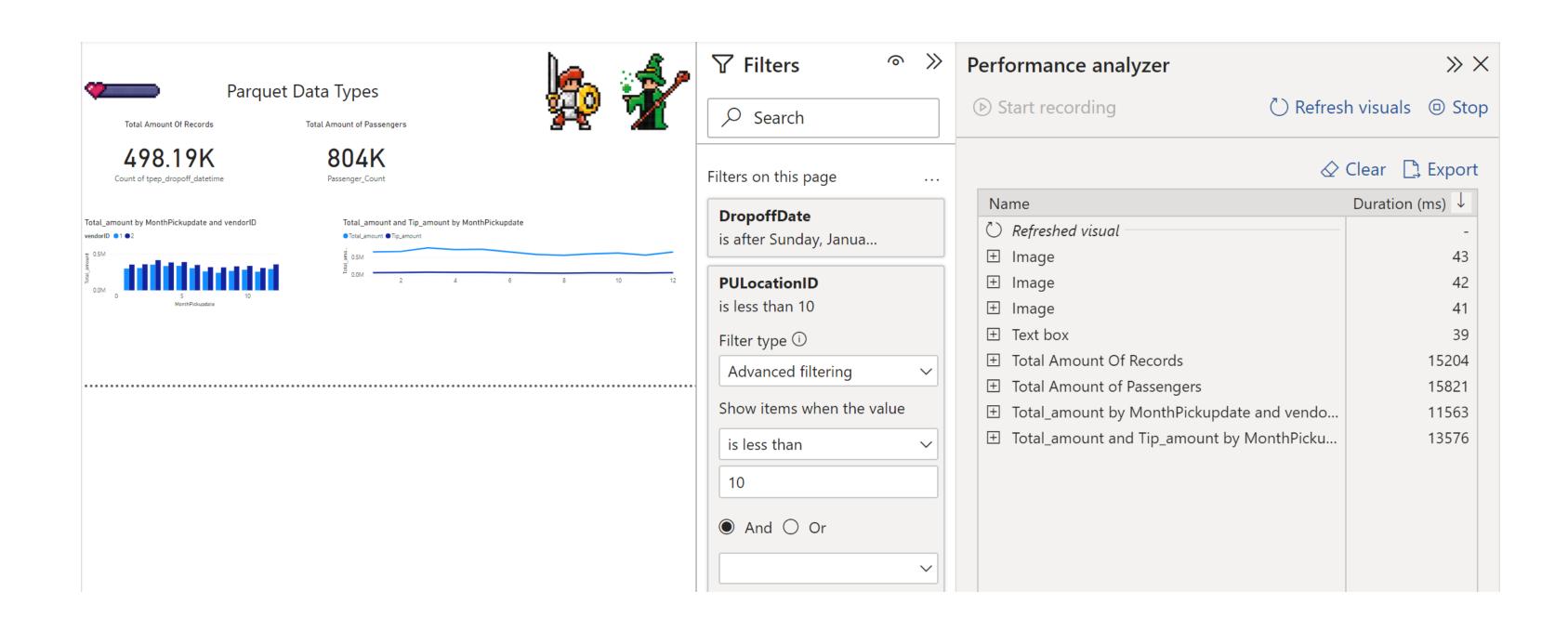
Use nvarchar for UTF-16 files

Use integers for SORT, JOIN, GROUP BY

If using schema interence, check inferred datatypes – override if possible.

Remember: It's SQL
Server – datatypes are
important for ALL SQL
Servers









Partitioning





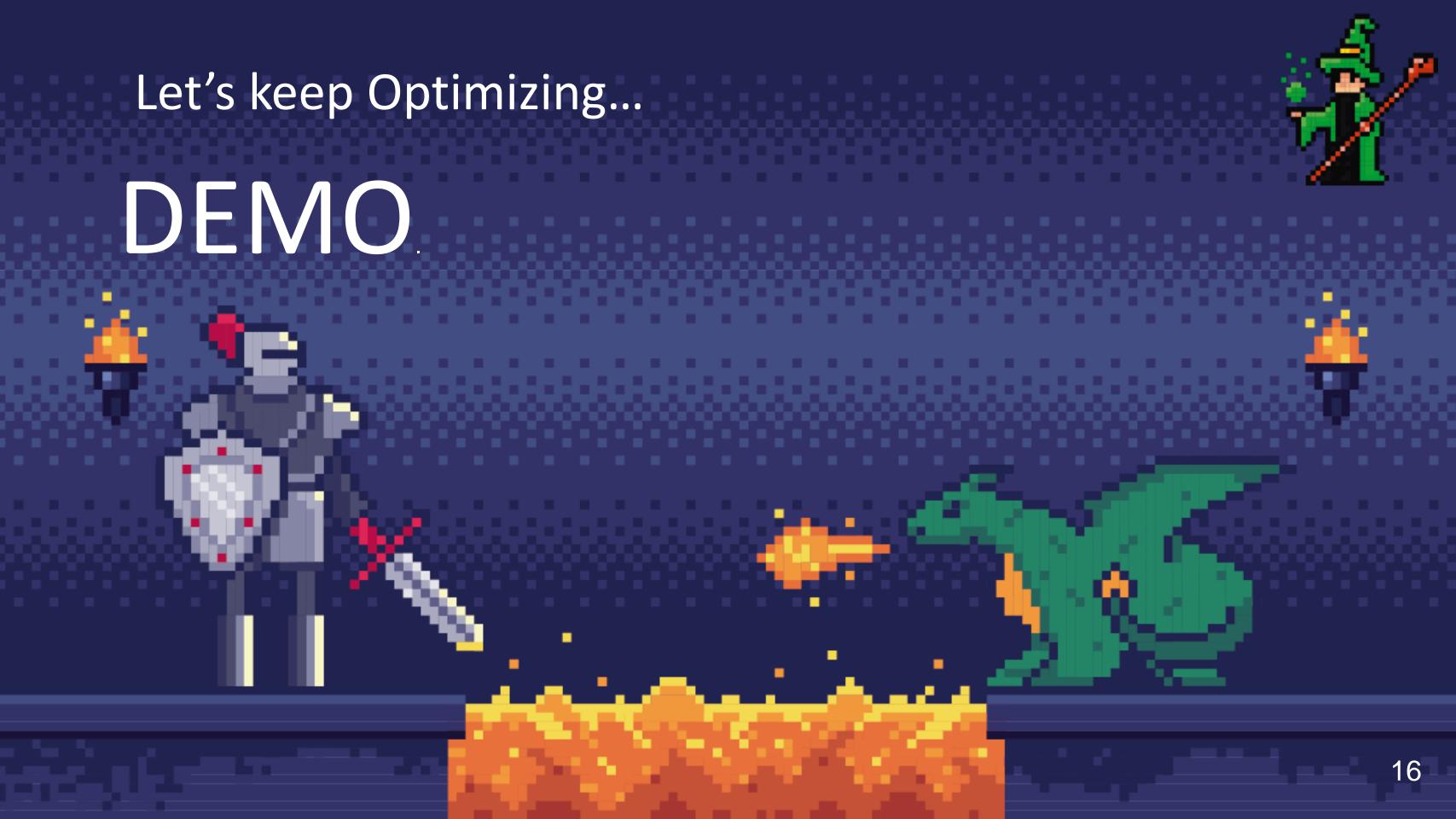


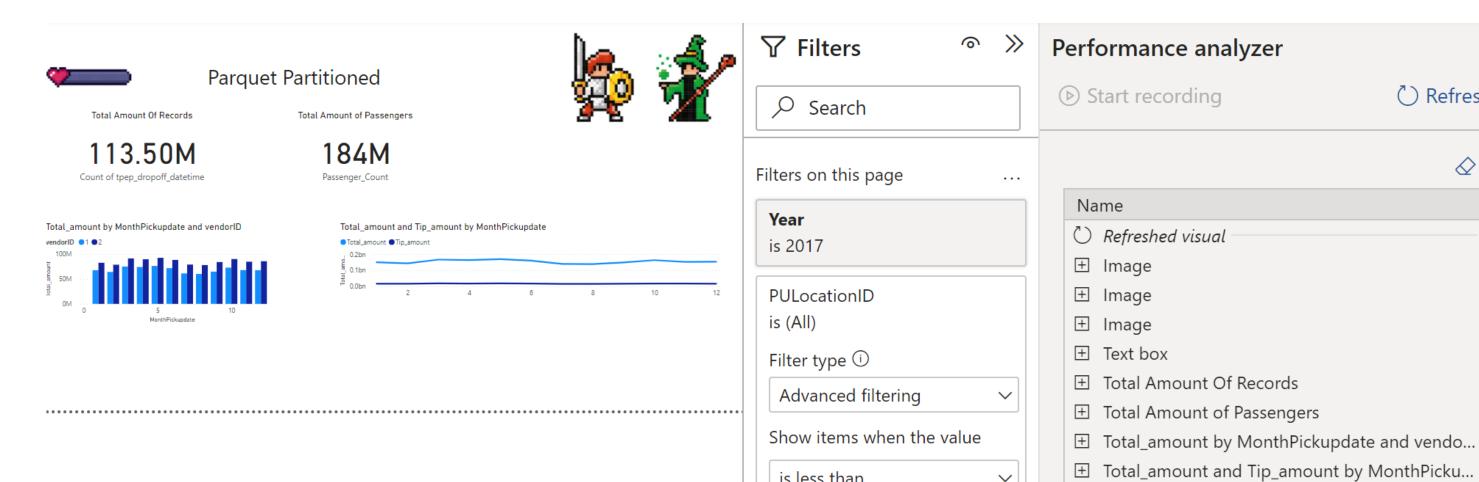


Spark SQL to write data to different folders

Serverless will read out the folder structure first and skip data

Openrowset filepath() is your best friend!

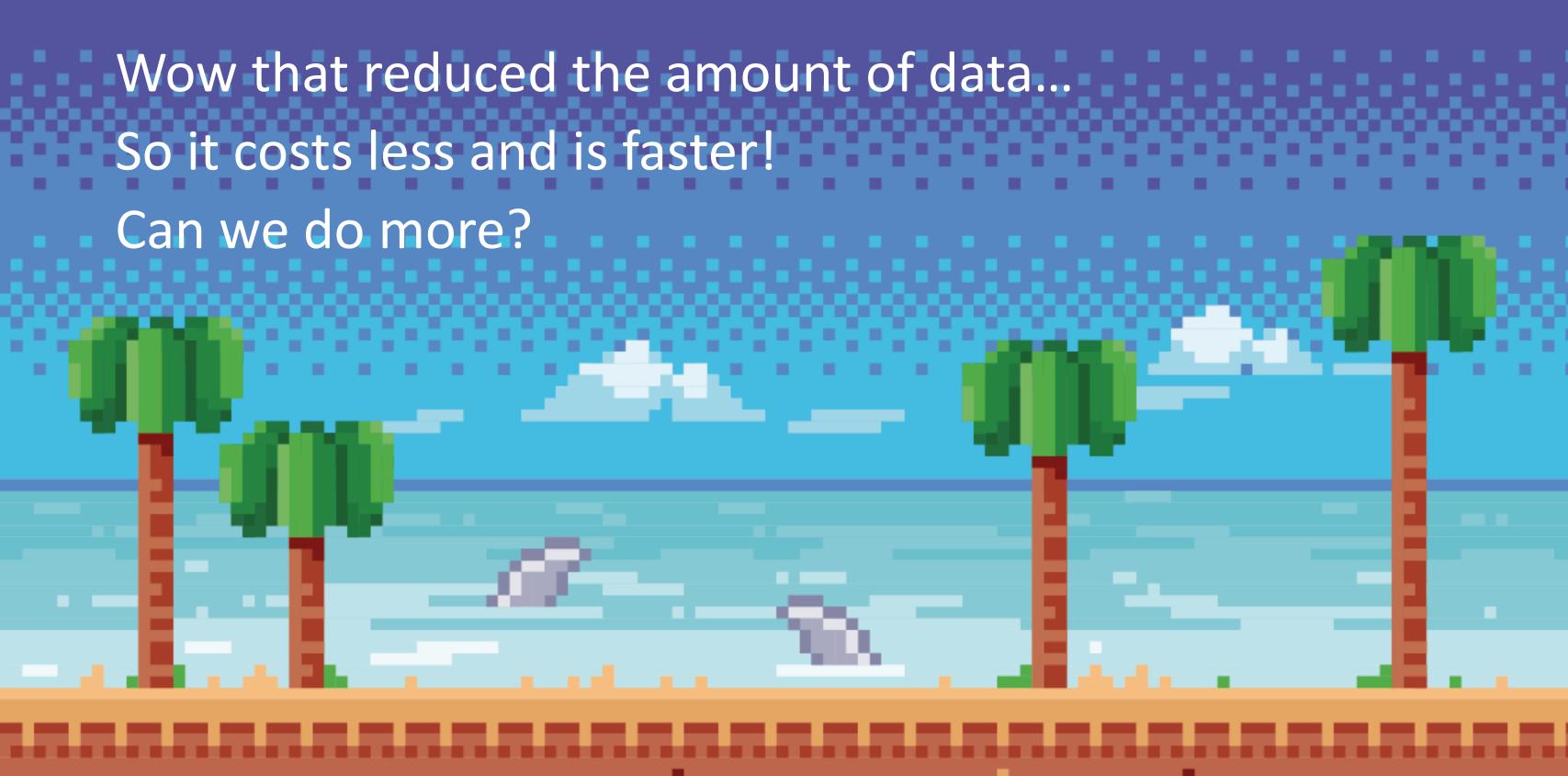




is less than

 $\gg \times$

Duration (ms) ↓





Column/Segment Elimination









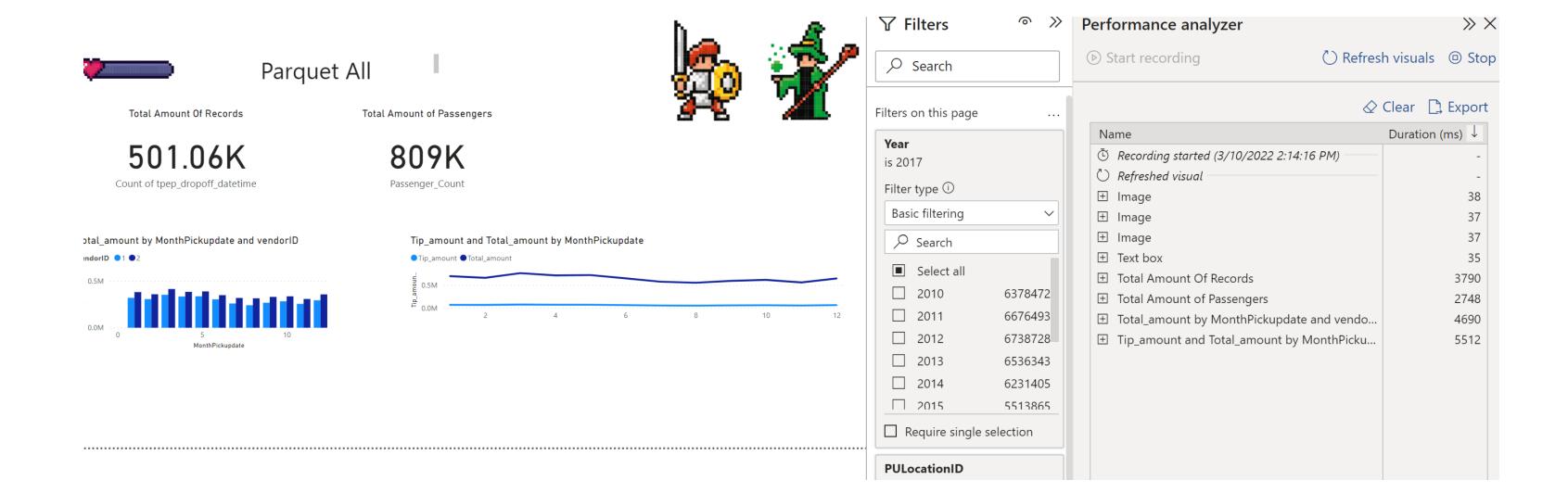
Spark SQL to write your data to your files in a specific order

Reduce column chunks

Serverless will read the file metadata and skip rowgroups & column chunks

Openrowset filepath() is your best friend!

Let's keep Optimizing... DEMO.



The Results



We reduced the report response time from 2 minutes and 10 seconds to 5 seconds !!!!

We are happy ...



As it always is, the dragon returns!



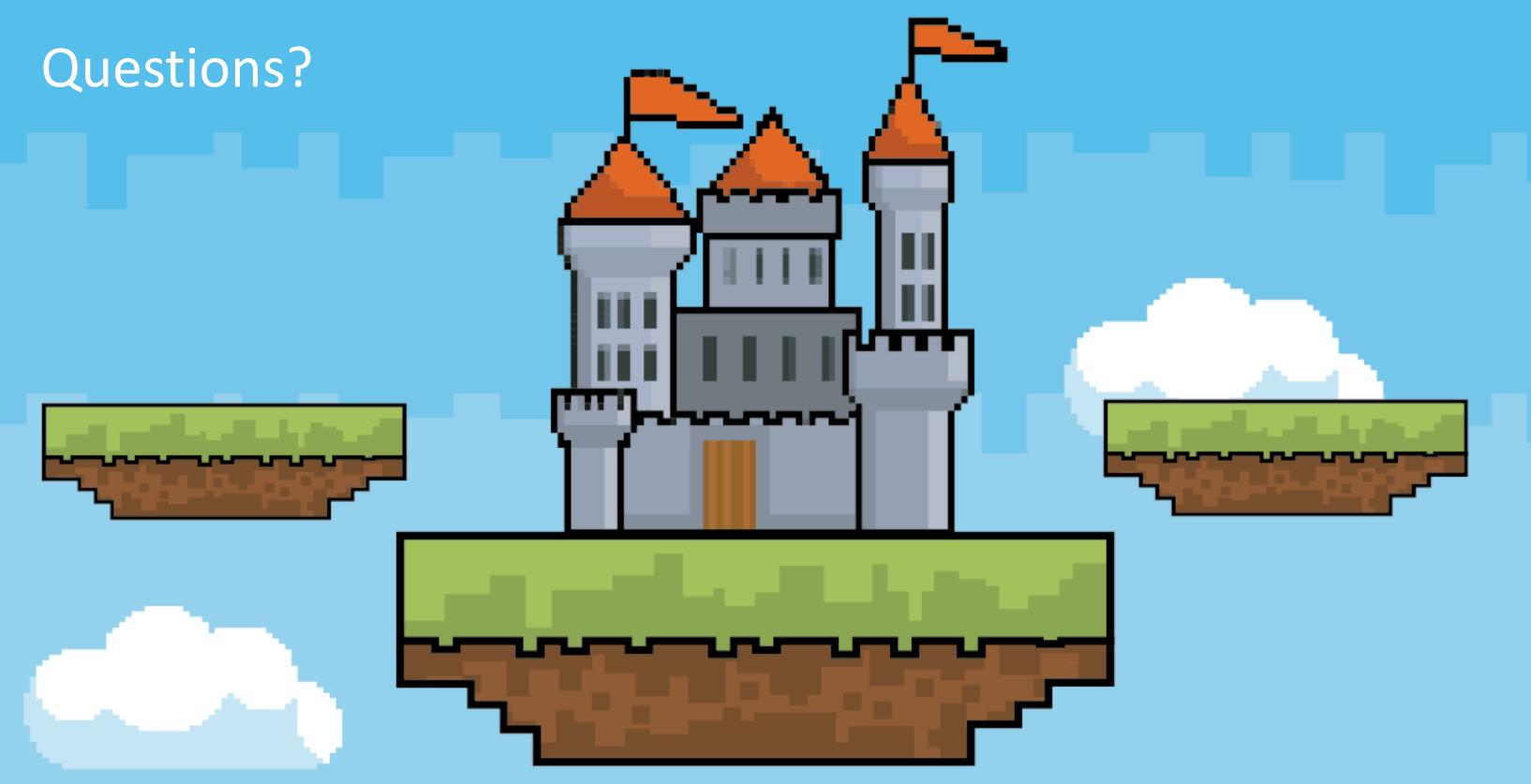
Recap!

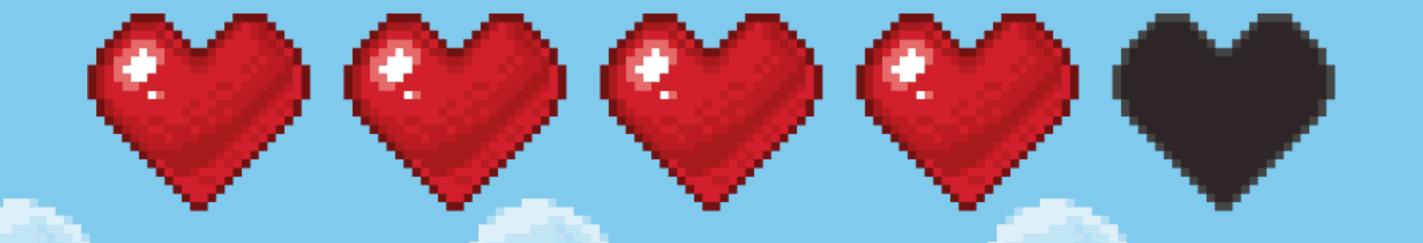
- 1. Use Parquet instead of CSV (if you can)
- 2. Check the data types.
- 3. Partition and sort your data (files).
- 4. Only select the data/columns you need.
- 5. Serverless is amazing and cost effective... ©











GRACE OUER

Please give us your Feedback...



Credits

Presentation template by Slides for Education

Illustrations and infographics by Freepik Photos created by Unsplash



