Looking under the hood of the:

# Parquet
# Format

André Kamman

MVP
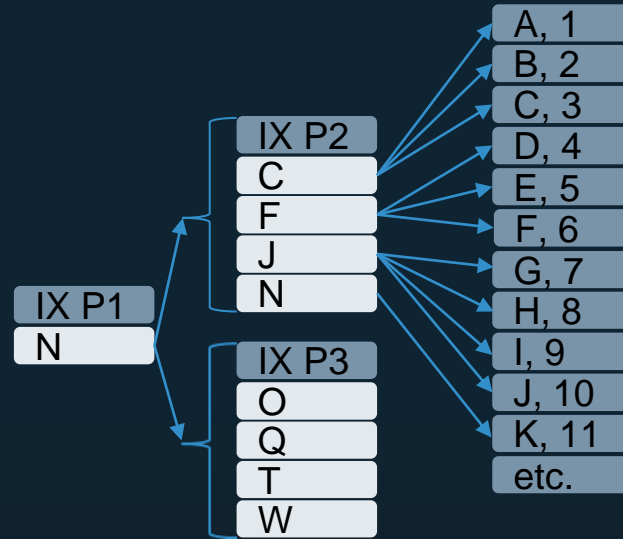Microsoft®
Most Valuable
Professional

@andrekamman

andrekamman@gmail.com

**https://evals.datagrillen.com/**

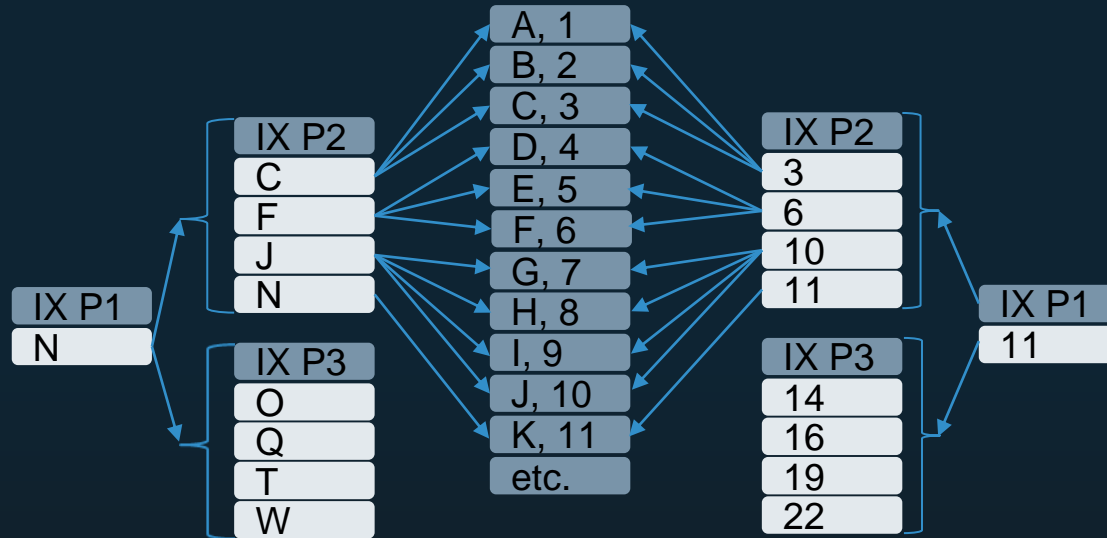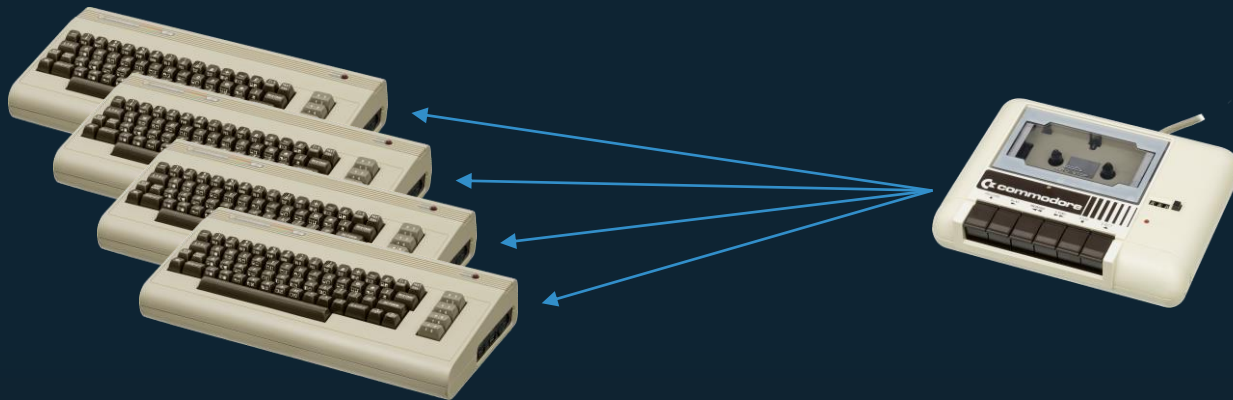# B-Tree

# B-Tree

Spark + Databricks

bits  Python

Schedule ▾

⊕ ⊘ democluster ▾

Cmd 1

```python
spark.read.format("csv") \
  .option("header", "true") \
  .option("inferSchema", "true") \
  .load("/mnt/lake/bits/videogames.csv") \
  .display()

```

▶ (3) Spark Jobs

Table    Data Profile

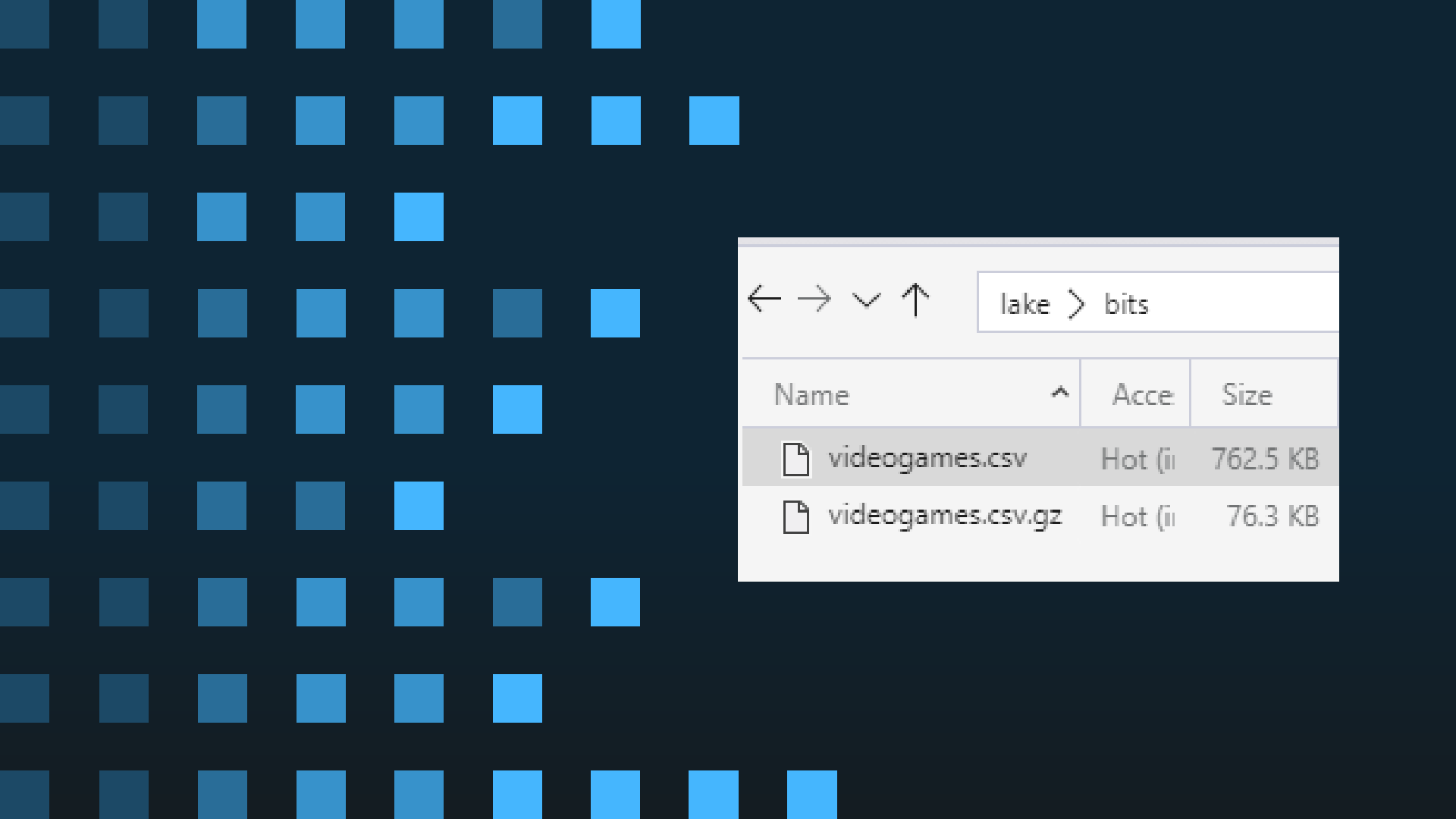| | Console ▲ | Title |
|---|---|---|
| 1 | Nintendo DS | Super Mario 64 DS |
| 2 | Sony PSP | Lumines: Puzzle Fusion |

```sql
%sql
select
  count(*) as Nr
  ,Genre
  ,Console
from videogames
group by Console, Genre
order by count(*) desc
```

▸ (2) Spark Jobs

**Table**   Data Profile

| | Nr | Genre | Console |
|---|---|---|---|
| 1 | 178 | Action | X360 |
| 2 | 137 | Action | PlayStation 3 |
| 3 | 123 | Action | Nintendo DS |
| 4 | 119 | Action | Nintendo Wii |
| 5 | 103 | Action | Sony PSP |
| 6 | 77 | Sports | X360 |
| 7 | 58 | Sports | PlayStation 3 |

lake > bits

| Name | ^ | Acce | Size |
|------|---|------|------|
| videogames.csv | | Hot (i | 762.5 KB |
| videogames.csv.gz | | Hot (i | 76.3 KB |

# Possible Issues

- No schema (only column names)
- Have to uncompress on one node
- Have to read everything every time
- Only flat data, no nesting possible

# Possible Issues

- No schema (only column names)
- Have to uncompress on one node
- **Have to read everything every time**
- Only flat data, no nesting possible

bits  [ Python ]

⛁  ⊘ democluster  ▾     📄 ▾   ✎ ▾   🖼 ▾   ▶   ◈ ▾

```sql
1  %sql
2  select
3      count(*) as nr
4      ,Console
5  from videogames
6  where Genre = "Action"
7  group by Console          ⬅
8  order by count(*) desc
```
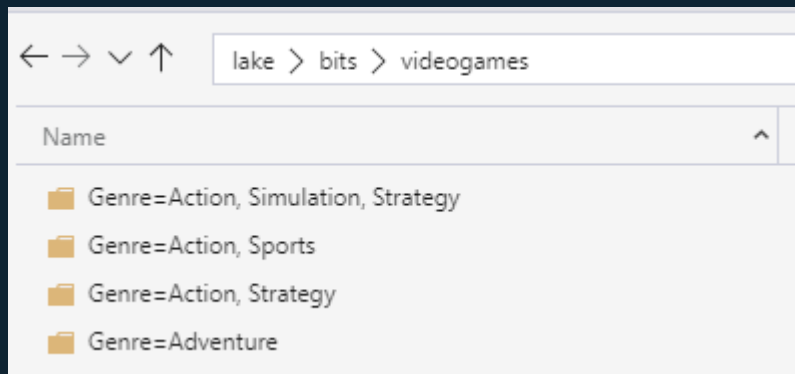
**Submitted Time:** 2022/03/07 17:18:54
**Duration:** 0.5 s
**Succeeded Jobs:** 46 47

☐ Expand all the details in the query plan

▾ (2) Spark Jobs

    ▸ Job 46   View (Stages: 1/1)

    ▸ Job 47   View (Stages: 1/1, 1 skipped)

Scan csv +details

Stages: 48.0

| | |
|---|---|
| file sorting by size time | 0 ms |
| number of files read | 1 |
| filesystem read data size | 27.7 KiB |
| filesystem read data size (sampled) | 55.4 KiB |
| filesystem read time (sampled) | 26 ms |
| metadata time | 2 ms |
| size of files read | 27.7 KiB |
| number of partitions read | 1 |
| rows output | 660 |

**Table**   Data Profile

| | nr ▲ | Console ▲ |
|---|---|---|
| 1 | 178 | X360 |
| 2 | 137 | PlayStation 3 |
| 3 | 123 | Nintendo DS |
| 4 | 119 | Nintendo Wii |
| 5 | 103 | Sony PSP |

# Possible Issues

- No schema (only column names)
- Have to uncompress on one node
- Have to read everything every time
- **Only flat data, no nesting possible**

JSON

# Columns instead of Rows

|        | Col A | Col B | Col C | Col D |
|--------|-------|-------|-------|-------|
| Row 0  | A0    | B0    | C0    | D0    |
| Row 1  | A1    | B1    | C1    | D1    |
| Row 2  | A2    | B2    | C2    | D2    |
| Row 3  | A3    | B3    | C3    | D3    |
| Row 4  | A4    | B4    | C4    | D4    |

|        | Col A | Col B | Col C | Col D |
|--------|-------|-------|-------|-------|
| Row 0  | A0    | B0    | C0    | D0    |
| Row 1  | A1    | B1    | C1    | D1    |
| Row 2  | A2    | B2    | C2    | D2    |
| Row 3  | A3    | B3    | C3    | D3    |
| Row 4  | A4    | B4    | C4    | D4    |

A0 A1 A2 A3 A4 B1 B2 B3
B4 C0 C1 C2 C3 C4 D0 D1
D2 D3 D4

| | Col A | Col B | Col C | Col D |
|---|---|---|---|---|
| Row 0 | A0 | B0 | C0 | D0 |
| Row 1 | A1 | B1 | C1 | D1 |
| Row 2 | A2 | B2 | C2 | D2 |
| Row 3 | A3 | B3 | C3 | D3 |
| Row 4 | A4 | B4 | C4 | D4 |

# Apache Parquet

Born out of a cooperation between Twitter & Cloudera
Open source storage format: PAX – Partition Attributes Across

Data Page Layouts for Relational Databases
on Deep Memory Hierarchies
Anastassia Ailamaki
David J. DeWitt
Mark D. Hill
The VLDB Journal 11 (2002)

PAR1

-

-

-

-

-

PAR1

**PAR1**

-

-

-
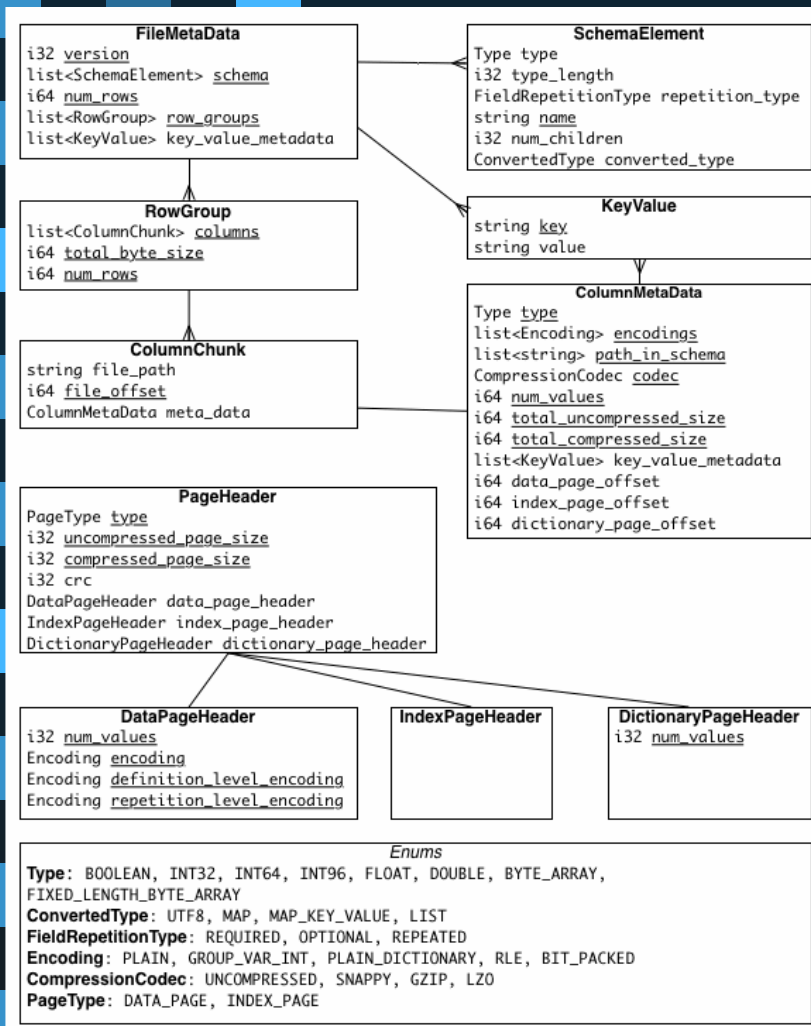
-

**Footer Length (4 bytes)**
**PAR1**

PAR1

-

-

-

Footer (metadata)
Footer Length (4 bytes)
PAR1

Thrift Compact Protocol

Row group 0

## Column Chuck – Column A

Page 0 (Dictionary Page)

Page 1 (Data Page)

## Page 1

### Metadata
Min / Max / Count
Encoding and Compression Codec info

### Repetition levels
(only used for nested columns)

### Definition levels
(only used for nullable columns)

### Encoded Data

# 3 Types of Metadata

- File metadata
- Column (chunk) metadata
- Page header metadata

# 2 Locations for metadata
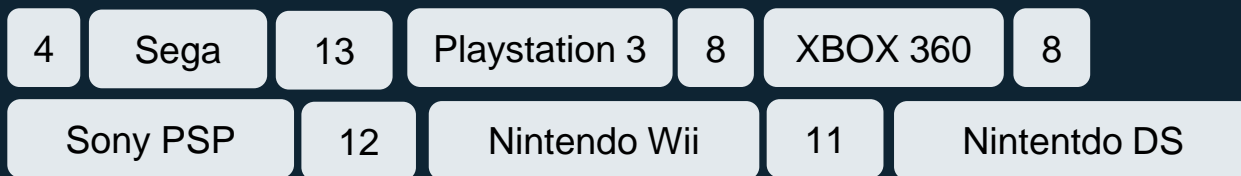
Page Header
File Footer

# Advantages So Far:

- Metadata at the end, one pass write.
- Column locations known, can grab only the columns needed.
- Min / Max values known, can skip pages, or even whole files.
- (Parquet table is usually a bunch of files and not just one)
- Directory based partitioning used to skip files, don't even have to parse the metadata.
- Data portion of pages are compressed, allowing compression without having to resort to a single node processor.

# Encoding

# Plain

- Fixed-width: back-to-back
- Variable Length: prefixed

| 4 | Sega | 13 | Playstation 3 | 8 | XBOX 360 | 8 |

| Sony PSP | 12 | Nintendo Wii | 11 | Nintentdo DS |

# RLE_Dictionary

- Run-length encoding
- Bit packing
- Dictionary Compression

| | |
|---|---|
| Sony PSP | |
| XBOX 360 | |
| SEGA | |
| Nintendo Wii | |
| Nintendo Wii | |
| Nintendo Wii | |
| Nintendo DS | |
| Nintendo DS | |

| | |
|---|---|
| 0 | Sony PSP |
| 1 | XBOX 360 |
| 2 | SEGA |
| 3 | Nintendo Wii |
| 4 | Nintendo DS |

| 0 | 1 | 2 | 3 | 3 |
|---|---|---|---|---|
| 3 | 4 | 4 | | |

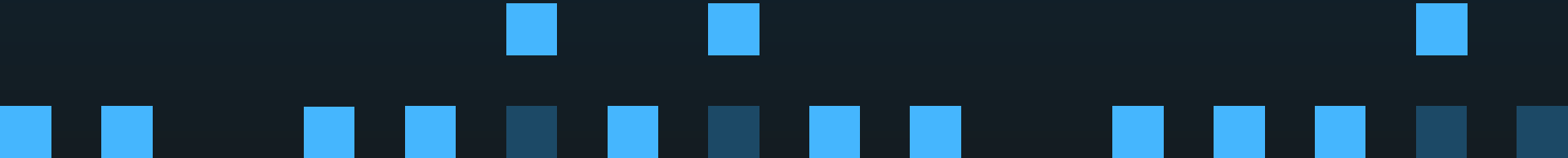| | |
|---|---|
| 0 | Sony PSP |
| 1 | XBOX 360 |
| 2 | SEGA |
| 3 | Nintendo Wii |
| 4 | Nintendo DS |

| 0 | 1 | 2 | 3,3 | 2,4 |
|---|---|---|---|---|

# Delta_Length_Byte_Array

V2

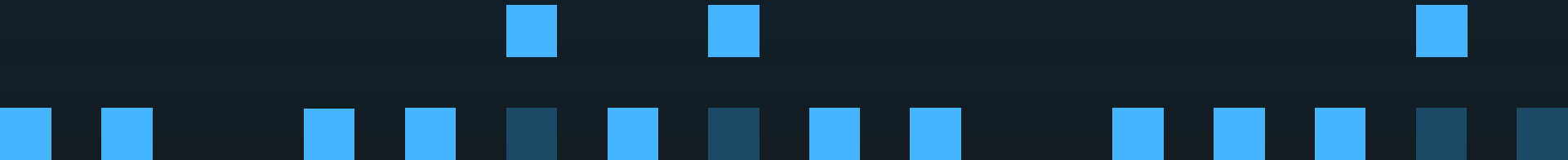| 4,13,8,8,12,11 | SegaPlaystation 3XBOX 360Sony PSPNintendo WiiNintendo DS |

Delta_Binary_Packed

V2

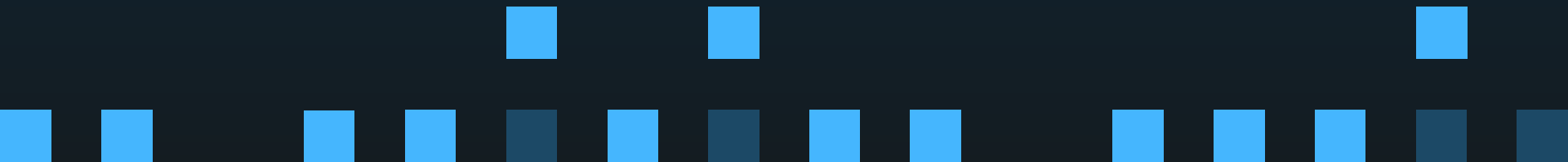| 1 | 2 | 3 | 4 | 5 |
| --- | --- | --- | --- | --- |
| | 1 | 1 | 1 | 1 |
| 1 | 0 | 0 | 0 | 0 |

Value Count : 5
First Value     : 1
Minimum delta: 1
Bitwidth : Not Needed!!

Adapted from "Decoding billions of integers per second through vectorization (D. Lemire, L.Boytsov)

# Random Learnings

Dictionary pages are always PLAIN encoded
Delta Encoding with signed integers is called "ZigZag Encoding"
Row-Groups are 128MB, pages are 1MB by default (In spark I guess)

# Is V2 The Future?

Maybe, a bit like IPv6 perhaps?
Check out Databricks Delta though!