# Lightweight Correlation-Aware Table Compression

**Mihail Stoian**, Alexander van Renen, Jan Kobiolka, Ping-Lin Kuo, Josif Grabocka, Andreas Kipf

University of Technology Nuremberg

UTN

```
Compression
import pandas as pd
import virtual

# Read your data.
df = pd.read_csv('file.csv')

# Your operations.
df = ...

# Virtualize + save to Parquet.
virtual
  .to_format(df, 'file.parquet')
```

Compression

```python
import pandas as pd
import virtual

# Read your data.
df = pd.read_csv('file.csv')

# Your operations.
df = ...

# Virtualize + save to Parquet.
virtual
  .to_format(df, 'file.parquet')
```

Query

```python
import virtual

virtual.query('''
  select avg(price)
  from read_parquet(
    "file.parquet"
  ) where year >= 2024''',
  engine = 'duckdb'
)
```

Columnar Encoding Schemes

- Frame-of-Reference (FOR), Run-Length-Encoding (RLE) etc.
- Pretty lightweight $\Rightarrow$ Fast decompression ✓
- File sizes: Could be better.. 🙄

### Columnar Encoding Schemes

- Frame-of-Reference (FOR), Run-Length-Encoding (RLE) etc.
- Pretty lightweight ⇒ Fast decompression ✓
- File sizes: Could be better.. 🙄

### Open File Formats

- Recent surge: Apache Parquet, ORC, etc.
- Research prototypes: BtrBlocks, FastLanes.

### Columnar Encoding Schemes

- Frame-of-Reference (FOR), Run-Length-Encoding (RLE) etc.
- Pretty lightweight ⇒ Fast decompression ✓
- File sizes: Could be better.. 🙄

### Open File Formats

- Recent surge: Apache Parquet, ORC, etc.
- Research prototypes: BtrBlocks, FastLanes.
- *Still using the standard encoding schemes.*
  ⇒ They have reached a plateau.

| Property Total | Burglary | Larceny | Motor Vehicle Theft |
|:---:|:---:|:---:|:---:|
| 5583 | 1884 | 3264 | 435 |
| 6368 | 1988 | 3878 | 502 |
| 6641 | 2246 | 3858 | 537 |

Property Total = Burglary + Larceny + Motor Vehicle Theft

| Property Total | Burglary | Larceny | Motor Vehicle Theft |
|:---:|:---:|:---:|:---:|
| 5583 | 1884 | 3264 | 435 |
| 6368 | 1988 | 3878 | 502 |
| 6641 | 2246 | 3858 | 537 |

Property Total = Burglary + Larceny + Motor Vehicle Theft

| Property Total_offset | Burglary | Larceny | Motor Vehicle Theft |
|:---:|:---:|:---:|:---:|
| 0 | 1884 | 3264 | 435 |
| 0 | 1988 | 3878 | 502 |
| 0 | 2246 | 3858 | 537 |

Requirements

(a) Make target column redundant $\Rightarrow$ Zero storage fingerprint.

(b) Touch as few columns as possible $\Rightarrow$ Fast table scans.

(c) Allow multiple functions $\Rightarrow$ Even better compression.

### Sparse Linear Regression

- Linear regression $\Rightarrow$ linear functions.

Sparse Linear Regression

- Linear regression $\Rightarrow$ linear functions.
- However: We have to enforce (b).
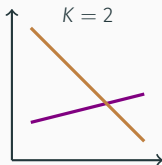
### Sparse Linear Regression

- Linear regression $\Rightarrow$ linear functions.
- However: We have to enforce (b).
  $\Rightarrow$ Optimize via $\mathcal{L}_0$ penalty.
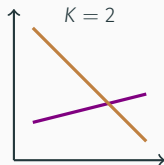
### Sparse Linear Regression

- Linear regression $\Rightarrow$ linear functions.
- However: We have to enforce (b).
  $\Rightarrow$ Optimize via $\mathcal{L}_0$ penalty.

### *K*-Regression

- Train *multiple* linear (sparse) regressors.

### Sparse Linear Regression

- Linear regression $\Rightarrow$ linear functions.
- However: We have to enforce (b).
  $\Rightarrow$ Optimize via $\mathcal{L}_0$ penalty.

### $K$-Regression

- Train *multiple* linear (sparse) regressors.



- Encode which regression we select via an auxiliary column.

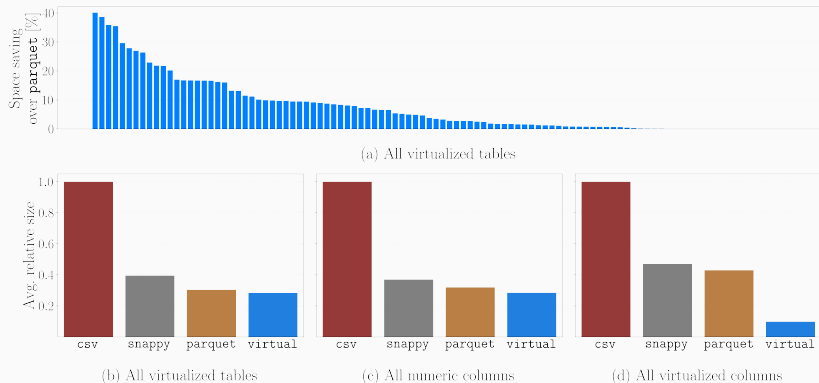(a) All virtualized tables

(b) All virtualized tables

(c) All numeric columns

(d) All virtualized columns

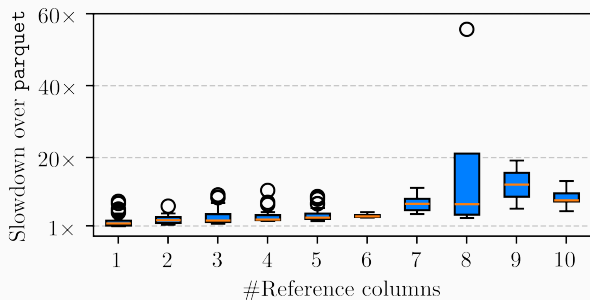**Figure 1:** Comparison to Parquet+Snappy (`parquet`) on 103 `data.gov` tables

Figure 2: Linear column scan slowdown

Lightweight Correlation-Aware Compression

· Learn multiple sparse linear regressors.

Lightweight Correlation-Aware Compression

- Learn multiple sparse linear regressors.
- Exploit them in compression and query execution.

github.com/utndatasystems/virtual