**Parquet**
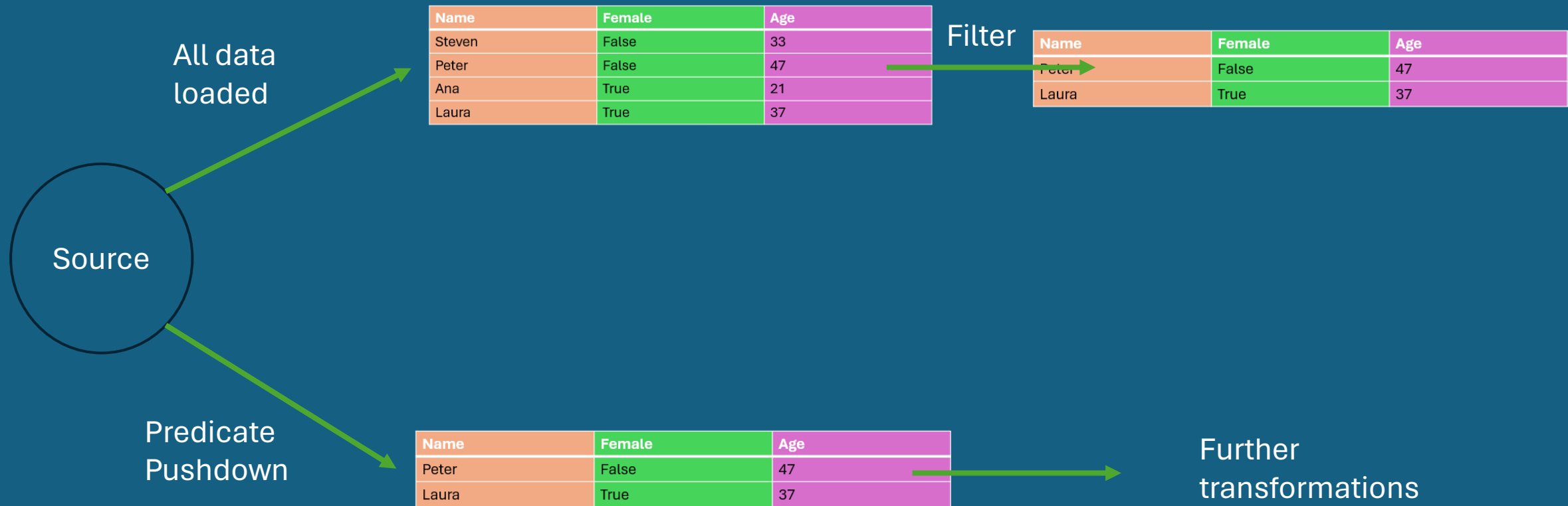
- Row Groups are a logical division on row level of a parquet defaulting to 128 MB

- Column part relates to column chunk of row groups

- Pages are invisible units where the encoding and compression happens

- Footer containing file metadata which can be used for predicate pushdown:
  - File level: num rows/ columns, schema
  - Row group: num rows/ columns
  - Column level: min, max, null count, distinct values, page indexes etc.

Diagram text:

File
Magic Number (4 bytes): "PAR1"
Row group 0
Column a
Page 0
Page header (ThriftCompactProtocol)
Repetition levels
Definition levels
values
Page 1
Column b
Row group 1

Footer
FileMetaData (ThriftCompactProtocol)
- Version (of the format)
- Schema
- extra key/value pairs
Row group 0 meta data:
Column a meta data:
- type / path / encodings / codec
- num values
- offset of first data page
- offset of first index page
- compressed/uncompressed size
- extra key/value pairs
column "b" meta data
Row group 1 meta data
Footer length (4 bytes)
Magic Number (4 bytes): "PAR1"

# Predicate and Aggregate Pushdown

All data
loaded

| Name | Female | Age |
|------|--------|-----|
| Steven | False | 33 |
| Peter | False | 47 |
| Ana | True | 21 |
| Laura | True | 37 |

Filter

| Name | Female | Age |
|------|--------|-----|
| Peter | False | 47 |
| Laura | True | 37 |

Source

Predicate
Pushdown

| Name | Female | Age |
|------|--------|-----|
| Peter | False | 47 |
| Laura | True | 37 |

Further
transformations

# Predicate and Aggregate Pushdown

- **Predicate Pushdown is an optimization technique filtering data at the source and often relies on statistics**

- **Benefits:**
  - Less I/O meaning less data to load
  - Less memory usage
  - Faster queries

- **Parquet supports Predicate Pushdown using statistics saved in meta data footer**

- **Since Spark 3.1.0 also possible on Avro, CSV, JSON**

- **Parquet supports since Spark 3.3.0 also Pushed Aggregates**