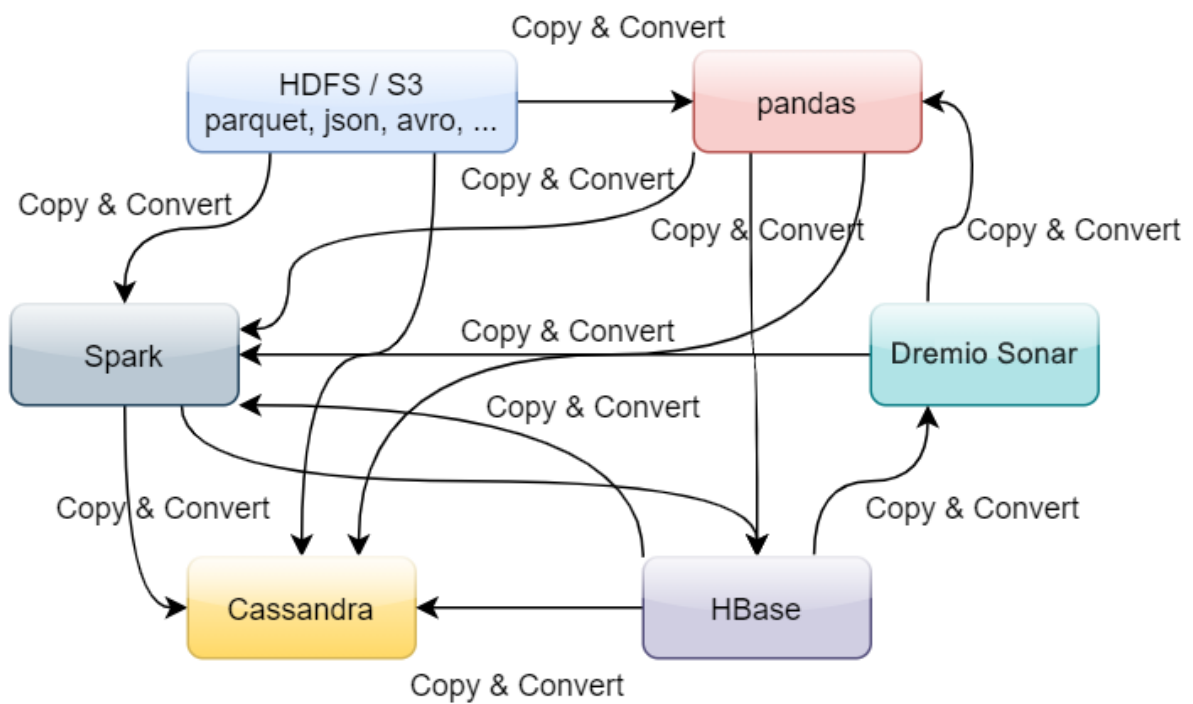
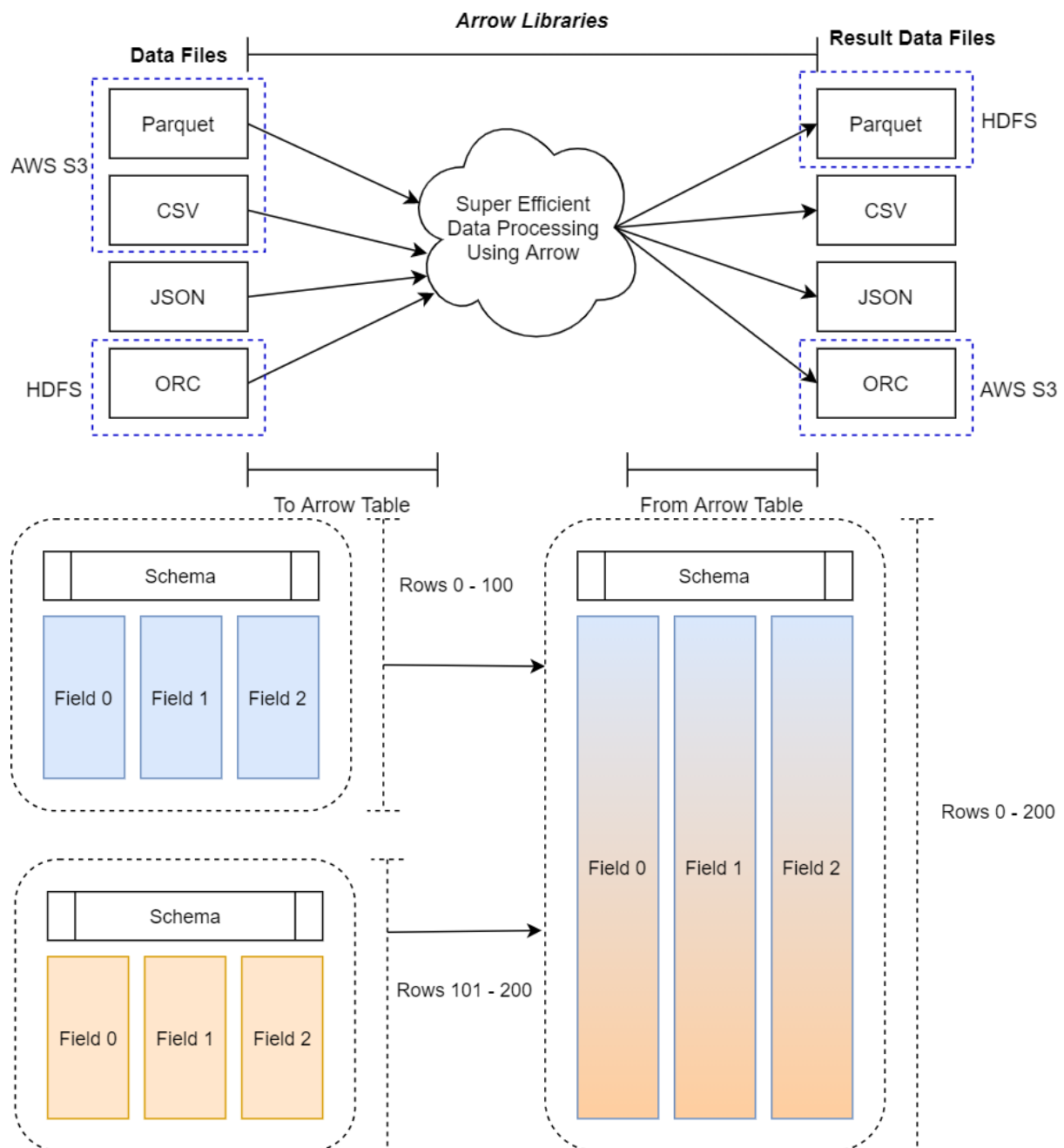
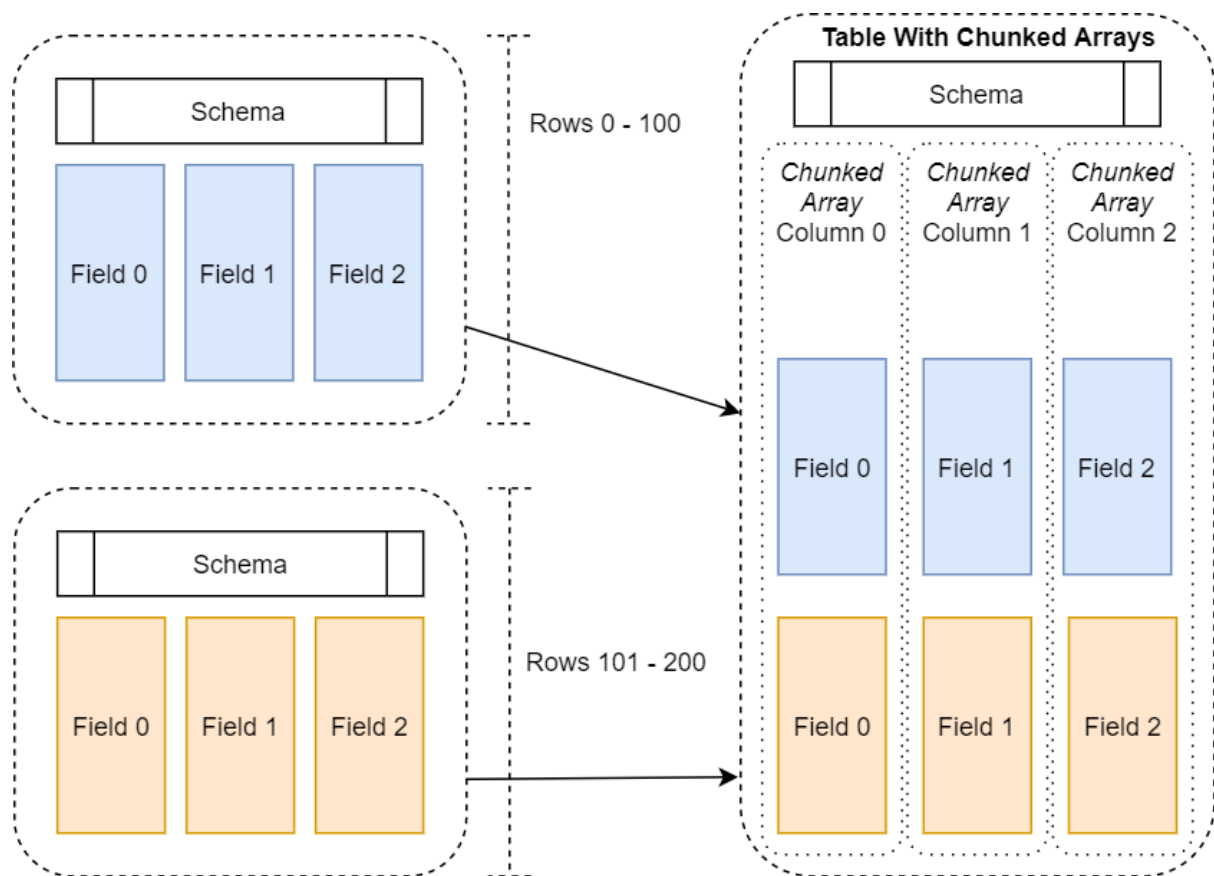


Chapter 1: Getting Started with Apache Arrow

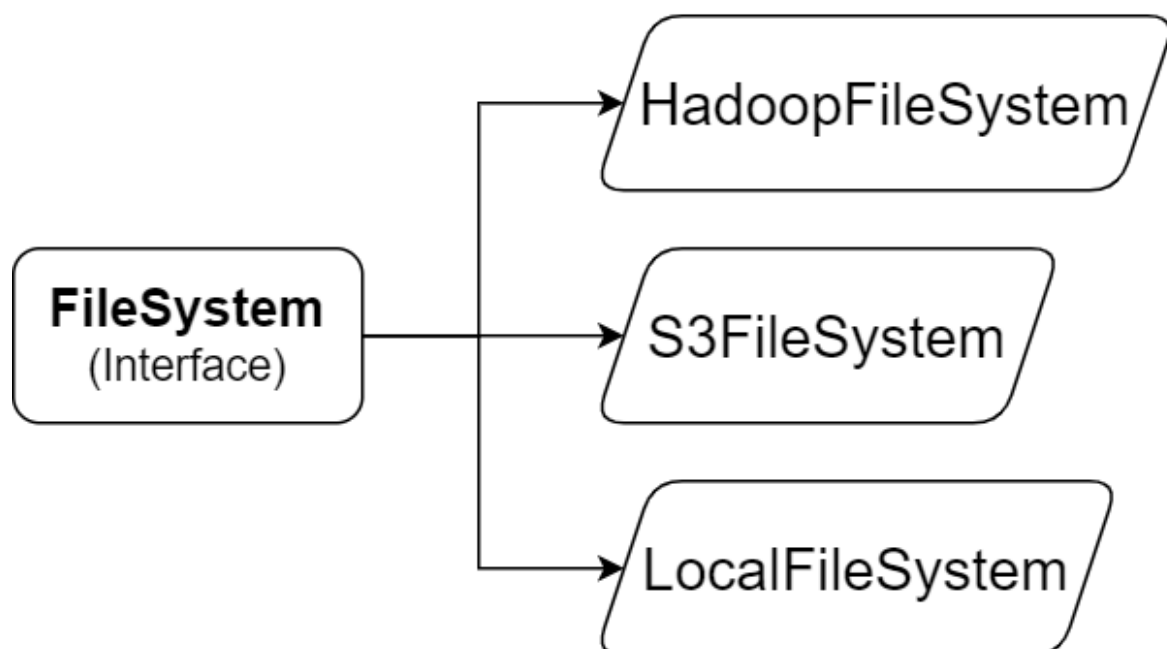


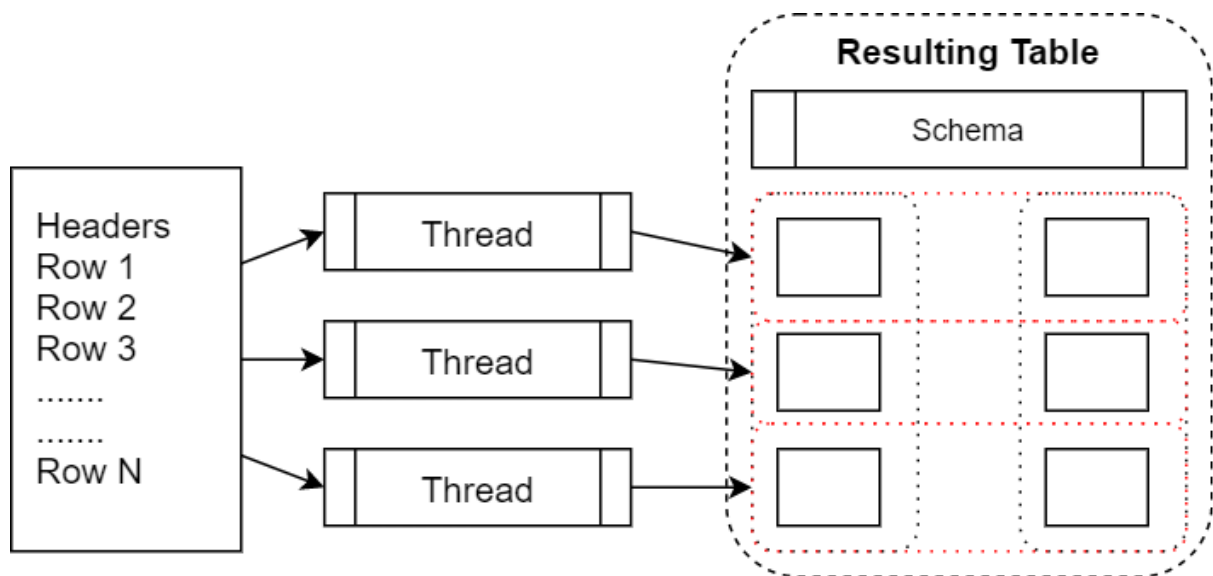
Chapter 2: Working with Key Arrow Specifications



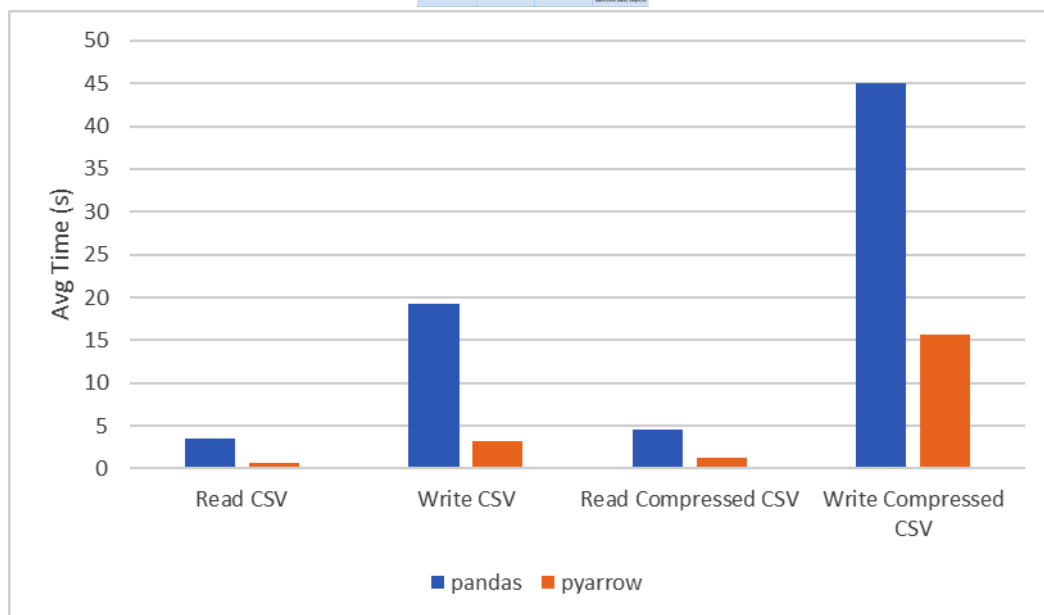


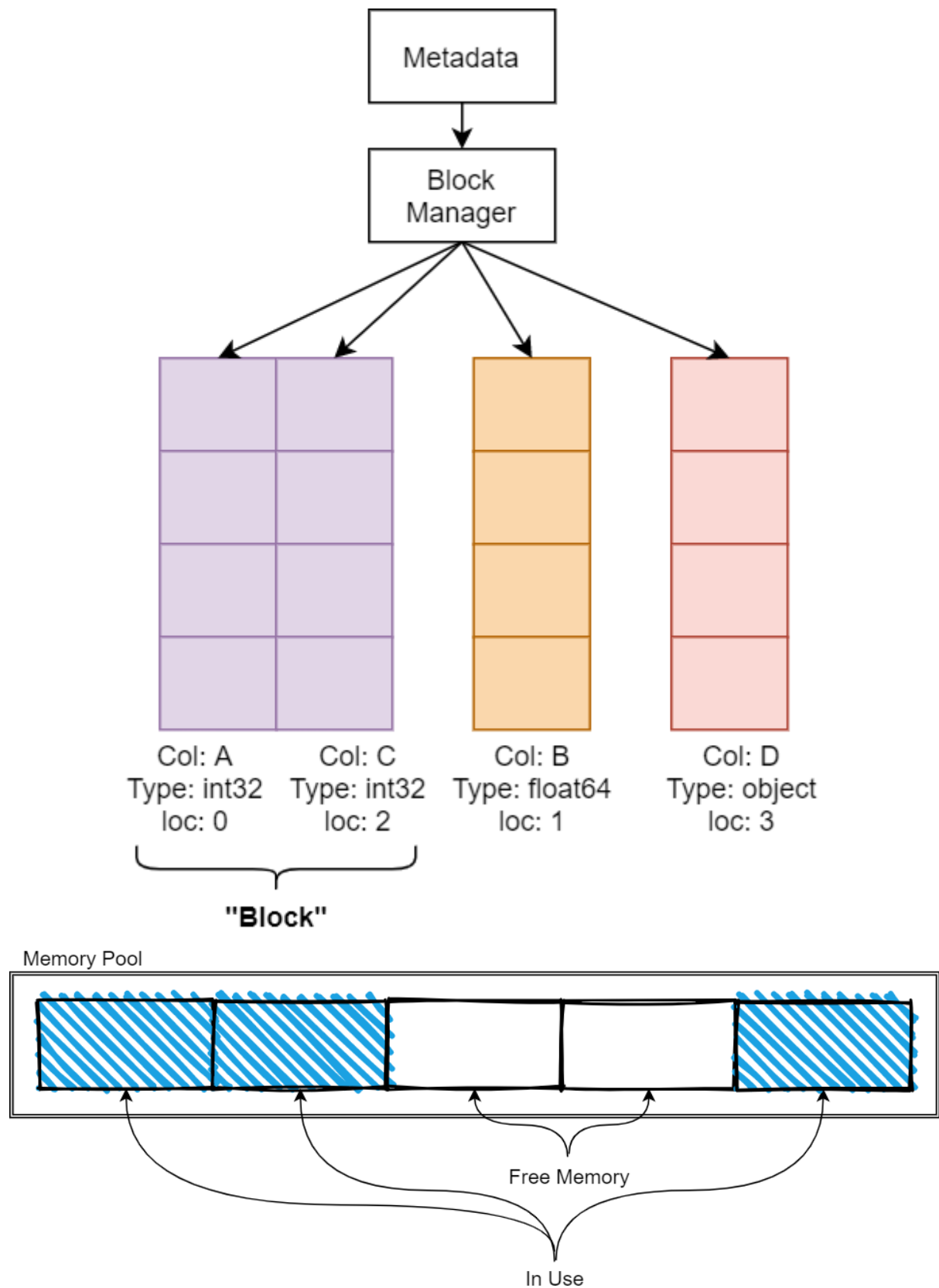
Concrete Implementations

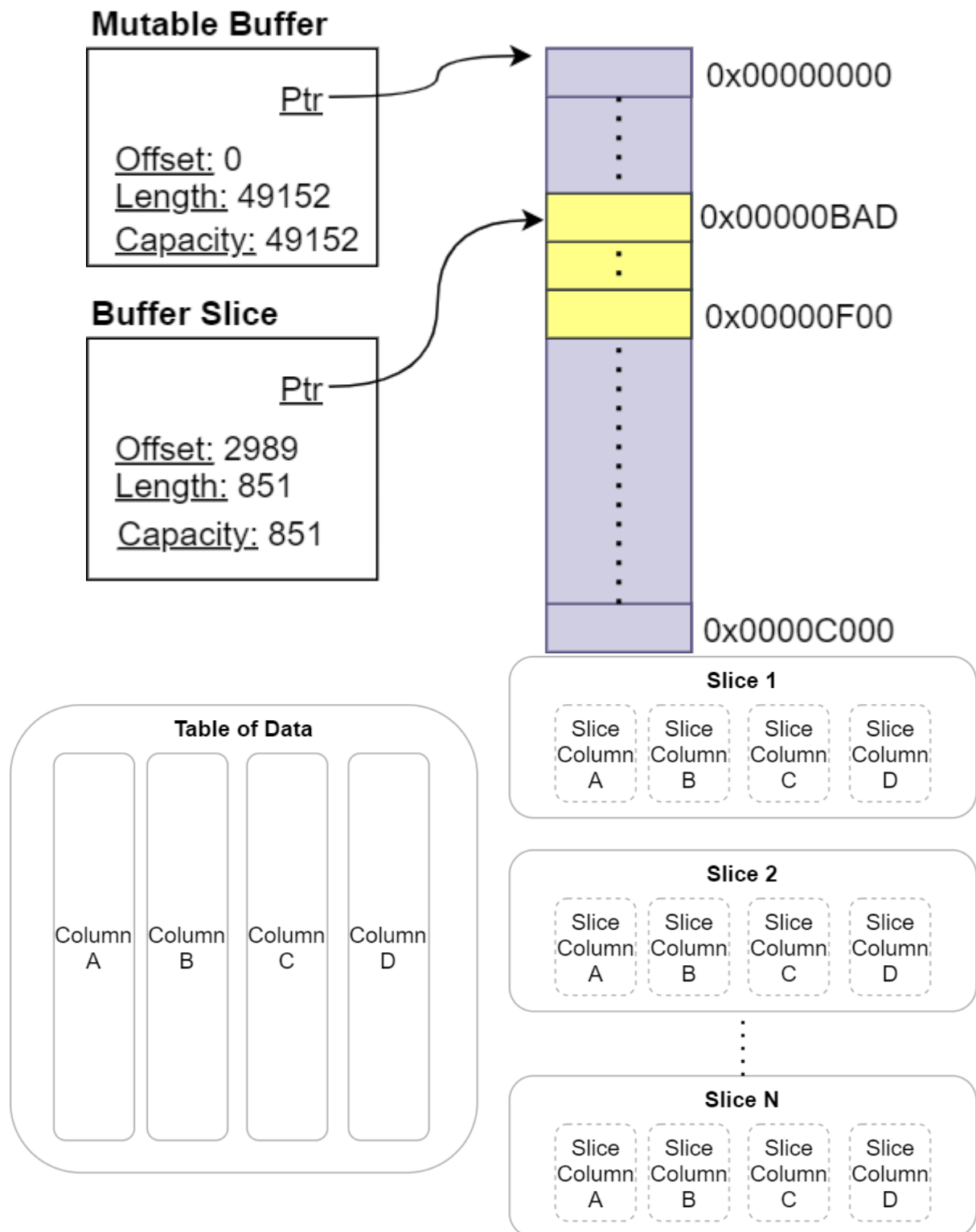


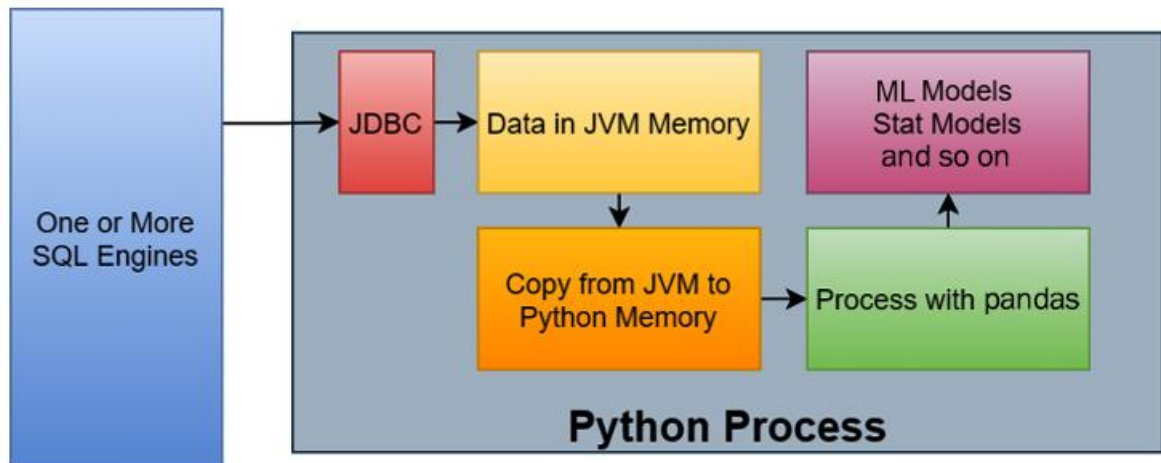
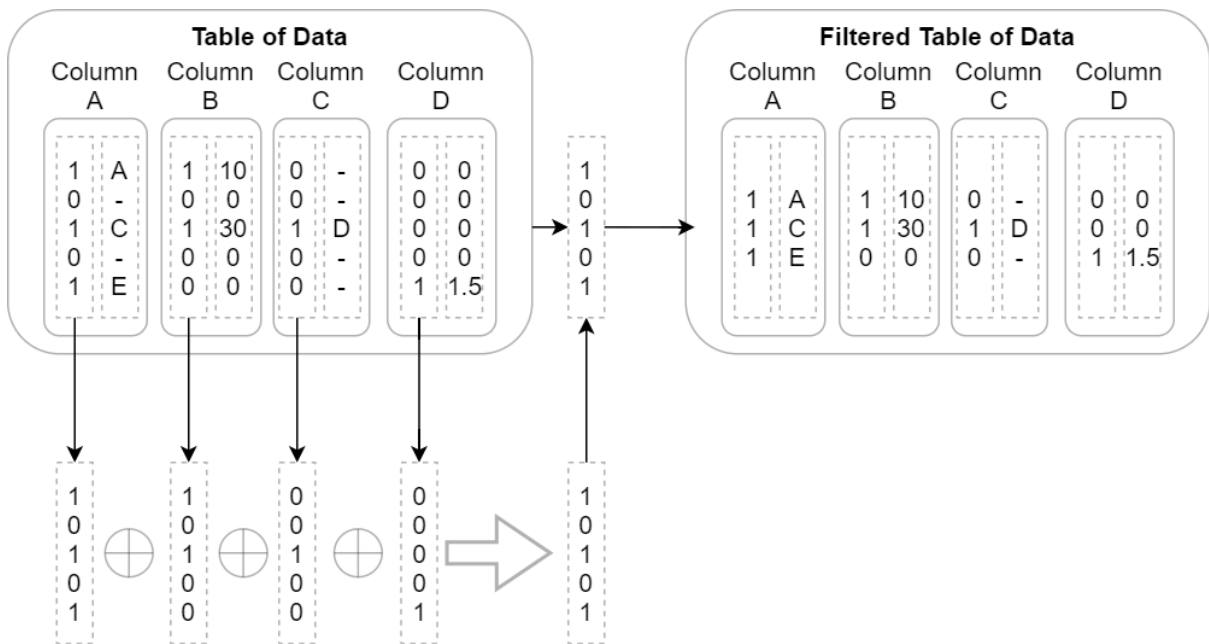


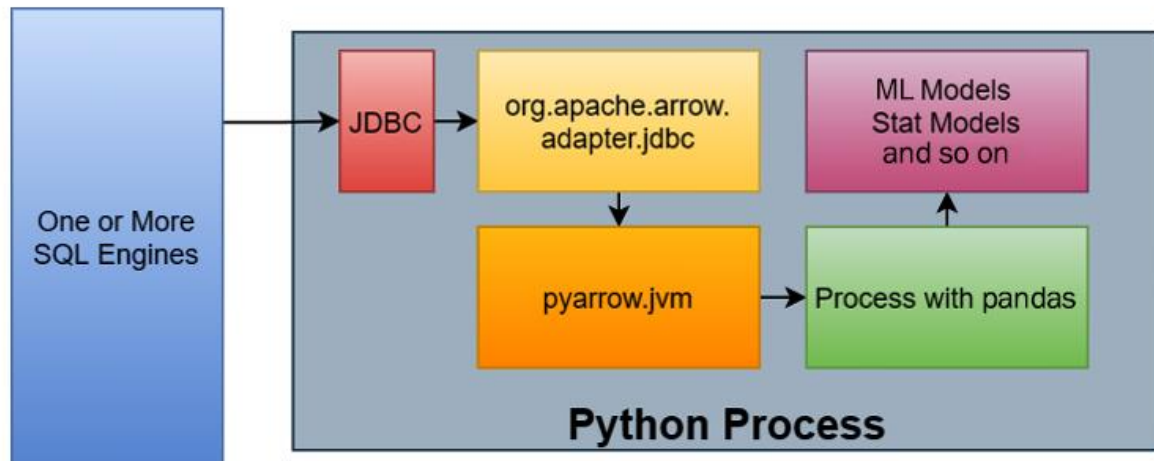
pandas -> Arrow		Arrow -> pandas	
pandas (source)	Arrow (destination)	Arrow (source)	pandas (destination)
bool	BOOL	BOOL	bool
int or uint [8,16,32,64]	int or uint [8,16,32,64]	BOOL with nulls	object with values: True, False, None
float32	FLOAT	int or uint [8,16,32,64]	int or uint [8,16,32,64]
float64	DOUBLE	int or uint [8,16,32,64] with nulls	float64
string[python]	STRING	FLOAT	float32
pandas.Categorical	CATEGORICAL	DOUBLE	float64
pandas.Timestamp	TIMESTAMP[us/nd]	STRING	str
datetime.date	DATE	CATEGORICAL	pandas.Categorical
datetime.time	TIME64	TIMESTAMP[us-7]	Pandas.Timestamp column identified by id
		DATE	object with datetime.date objects





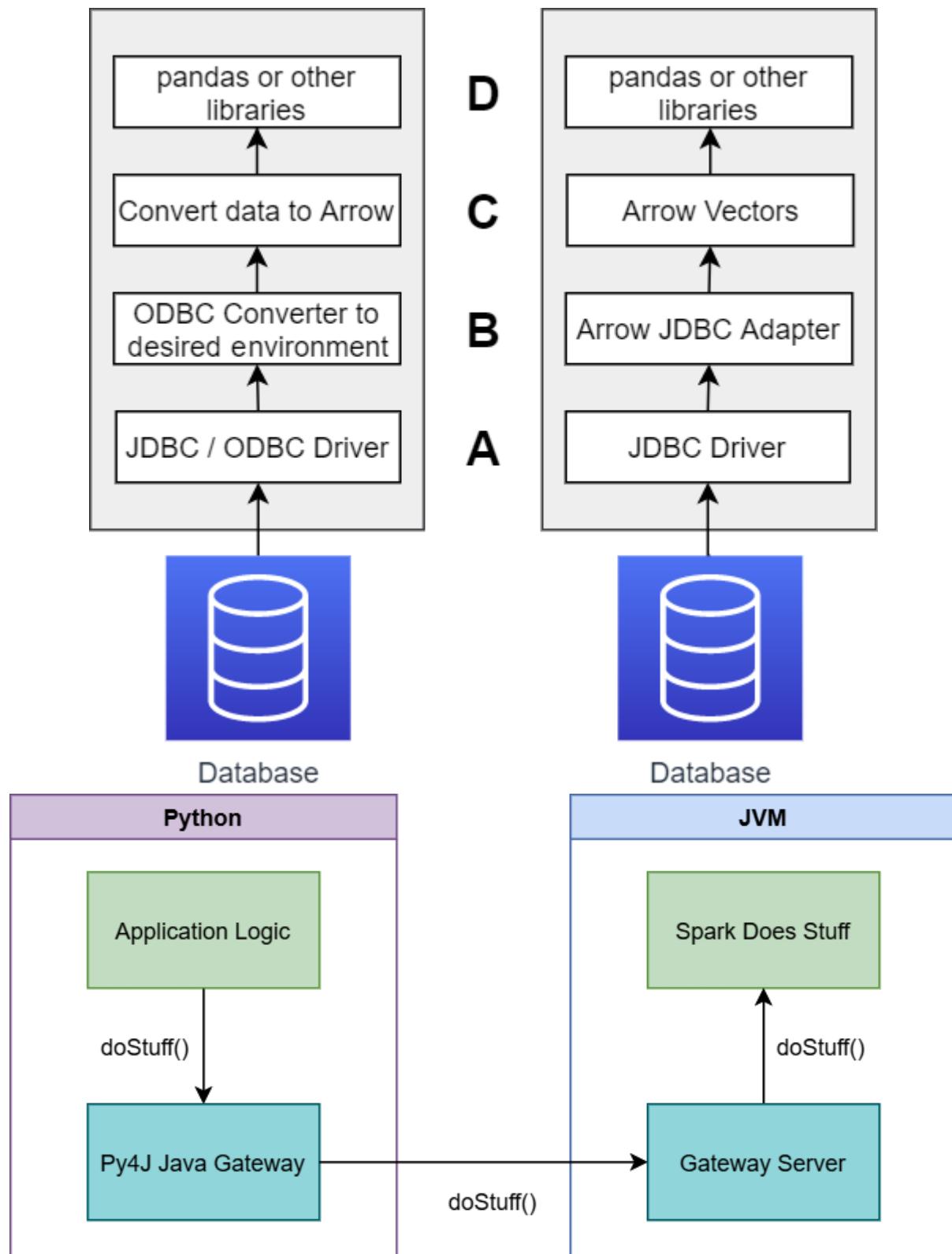


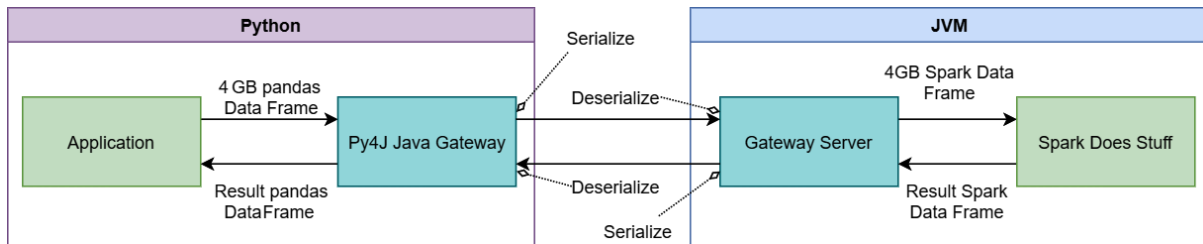




# Rows	Traditional (Copy)	Shared Memory (No Copy)	Speedup	% Improvement
950	274 ms	144 ms	1.90x	90.27%
10,000	1.29 s	175 ms	7.37x	637.14%
198,143	26.8 s	403 ms	66.50x	6550.12%
623,418	1 min 5 s	573 ms	113.44x	11243.80%

Chapter 3: Data Science with Apache Arrow





```

[I 2021-11-15 04:24:34.386 ServerApp] Writing Jupyter server cookie secret to /home/jovyan/.local/share/jupyter/runtime/jupyter_cookie_secret
[I 2021-11-15 04:24:34.605 ServerApp] nbclassic | extension was successfully linked.
[I 2021-11-15 04:24:34.637 ServerApp] nbclassic | extension was successfully loaded.
[I 2021-11-15 04:24:34.638 LabApp] JupyterLab extension loaded from /opt/conda/lib/python3.9/site-packages/jupyterlab
[I 2021-11-15 04:24:34.638 LabApp] JupyterLab application directory is /opt/conda/share/jupyter/lab
[I 2021-11-15 04:24:34.643 ServerApp] jupyterlab | extension was successfully loaded.
[I 2021-11-15 04:24:34.644 ServerApp] Serving notebooks from local directory: /home/jovyan
[I 2021-11-15 04:24:34.644 ServerApp] Jupyter Server 1.11.2 is running at:
[I 2021-11-15 04:24:34.644 ServerApp] http://2dd5d89e795a:8888/lab?token=a58056e05bf0e3b751c087689952ec5776085702347b1fe6
[I 2021-11-15 04:24:34.644 ServerApp] or http://127.0.0.1:8888/lab?token=a58056e05bf0e3b751c087689952ec5776085702347b1fe6
[I 2021-11-15 04:24:34.644 ServerApp] Use Control-C to stop this server and shut down all kernels (twice to skip confirmation).
[C 2021-11-15 04:24:34.650 ServerApp]

To access the server, open this file in a browser:
    file:///home/jovyan/.local/share/jupyter/runtime/jpserver-8-open.html
Or copy and paste one of these URLs:
    http://2dd5d89e795a:8888/lab?token=a58056e05bf0e3b751c087689952ec5776085702347b1fe6
    or http://127.0.0.1:8888/lab?token=a58056e05bf0e3b751c087689952ec5776085702347b1fe6

```

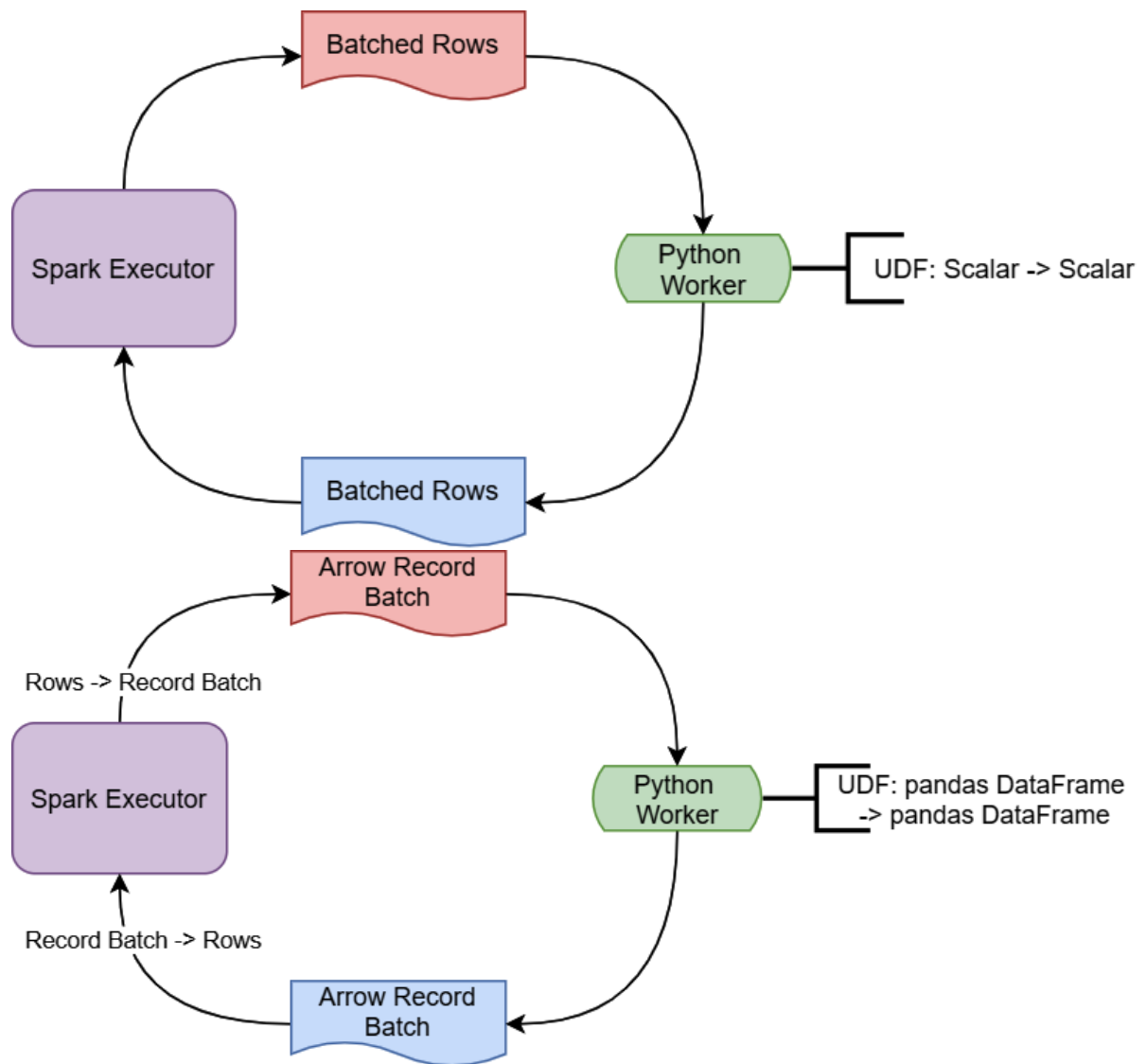
```

chapter3.ipynb
import pyspark as pyspark
from pyspark import SparkContext
conf = pyspark.SparkConf()
conf.set("spark.executor.memory", "8g")
conf.set("spark.driver.memory", "8g")
conf.set('spark.jars.packages', 'org.apache.hadoop:hadoop-aws:3.3.1')
sc = SparkContext(conf=conf)
sc._jsc.hadoopConfiguration().set('fs.s3a.aws.credentials.provider', 'org.apache.hadoop.fs.s3a.AnonymousAWSCredentialsProvider')
sc._jsc.hadoopConfiguration().set('fs.s3a.impl', 'org.apache.hadoop.fs.s3a.S3AFileSystem')

from pyspark.sql import SparkSession
spark = SparkSession(sc)

WARNING: An illegal reflective access operation has occurred
WARNING: Illegal reflective access by org.apache.spark.unsafe.Platform (file:/usr/local/spark-3.2.0-bin-hadoop3.2/jars/spark-unsafe-3.2.0.jar) to method java.lang.ProcessImpl::start()
WARNING: Please consider reporting this to the maintainers of org.apache.spark.unsafe.Platform
WARNING: Use --illegal-access=warn to enable warnings of further illegal reflective access operations
WARNING: All illegal access operations will be denied in a future release
:: loading settings :: url = jar:file:/usr/local/spark-3.2.0-bin-hadoop3.2/jars/ivy-2.5.0.jar!/org/apache/ivy/core/settings/ivysettings.xml
Ivy Default Cache set to: /home/jovyan/.ivy2/cache

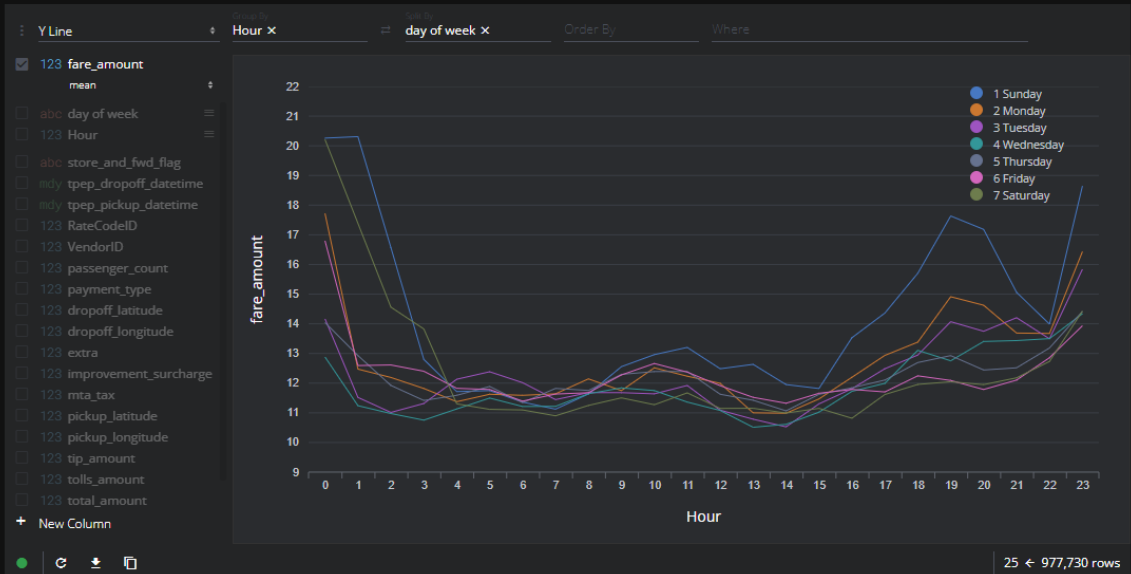
```

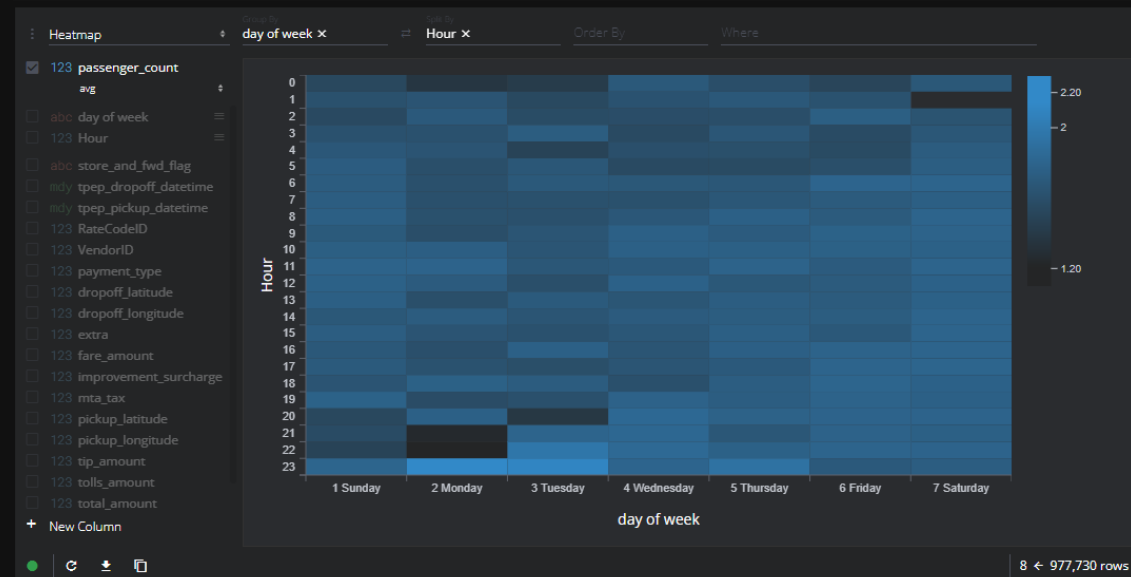
```
[8]: view = table.view(filter=[["tpep_pickup_datetime", "<", "2015-01-10"]])
display(view.num_rows())
```

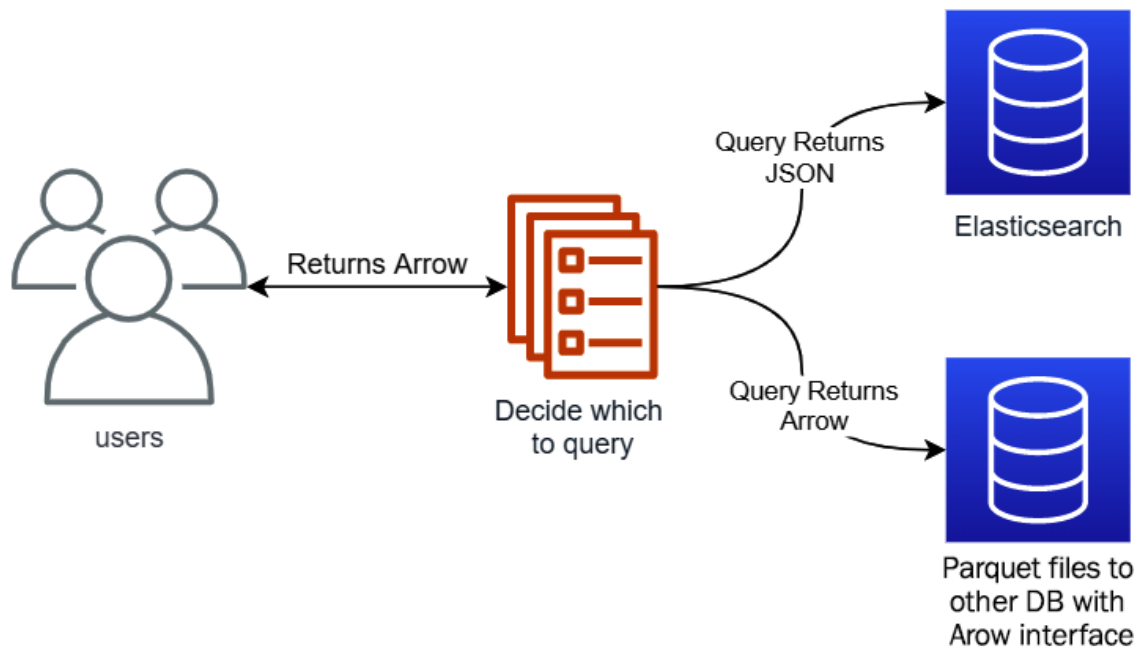
977730

```
[10]: widget = PerspectiveWidget(view.to_arrow())
widget
```

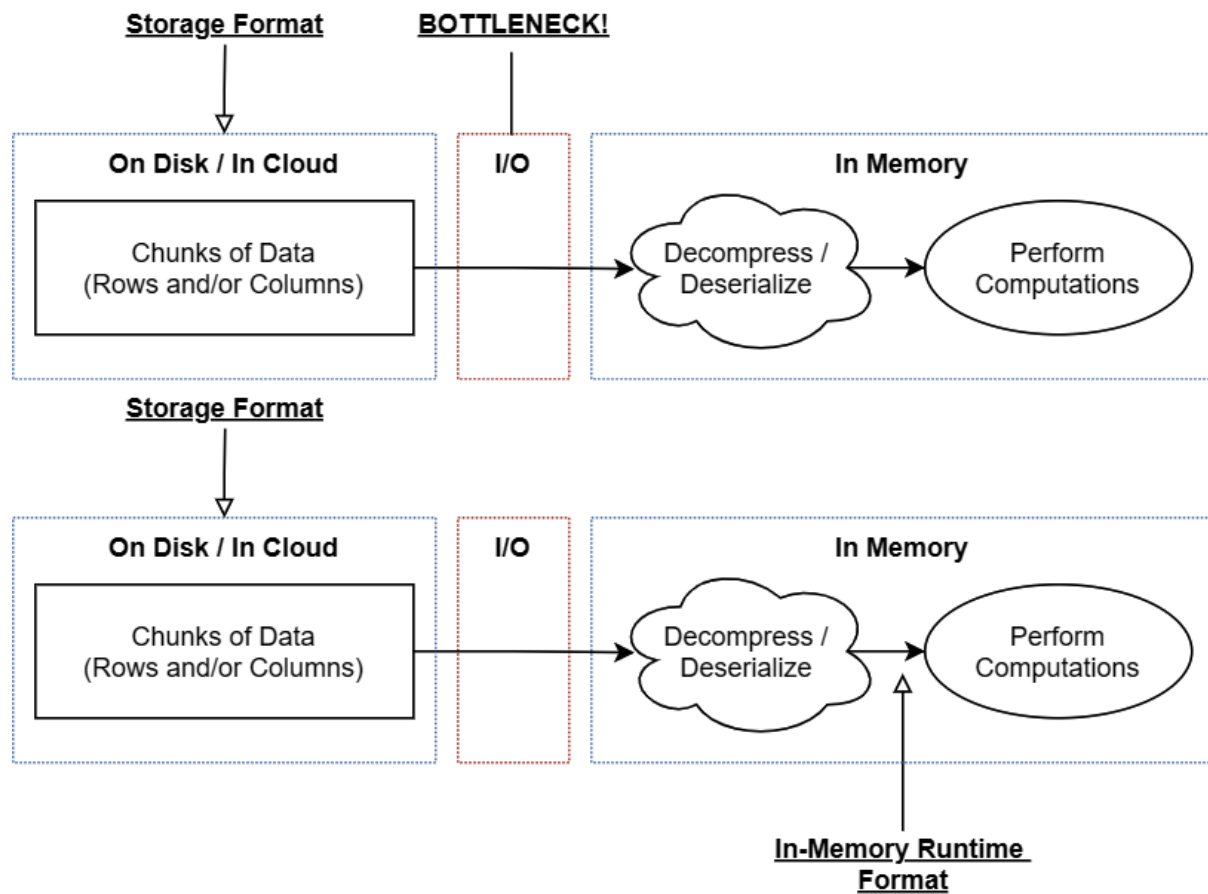


```
[10]: widget = PerspectiveWidget(view.to_arrow())
widget
```





Chapter 4: Format and Memory Handling



Stages											
Fetch	a	b	c	d							
Decode		a	b	c	d						
Execute			a	b	c	d					
Write Back				a	b	c	d				
	1	2	3	4	5	6	7	8	9	10	Clock Cycle

Stages											
Fetch	a	b	c	d	b	c	d				
Decode		a	b	c		b	c	d			
Execute			a	b			b	c	d		
Write Back				a				b	c	d	
	1	2	3	4	5	6	7	8	9	10	Clock Cycle

Bubble!

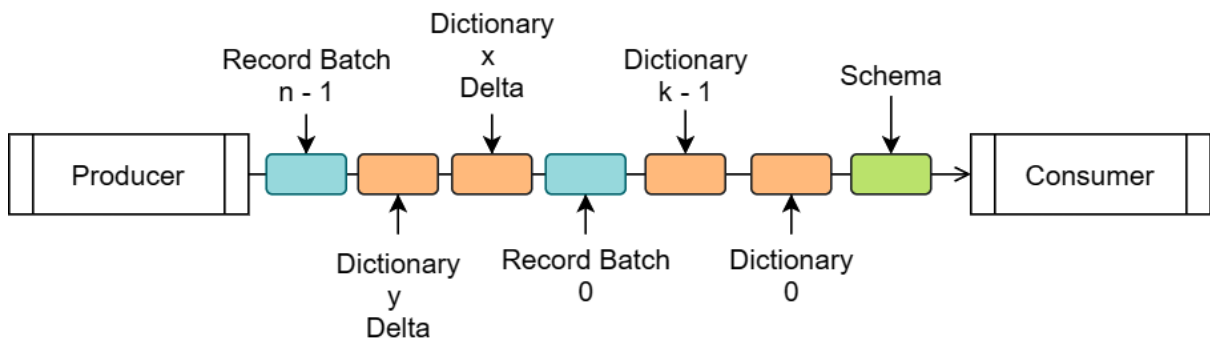
0xFFFFFFFF ← Continuation Indicator, 4 bytes

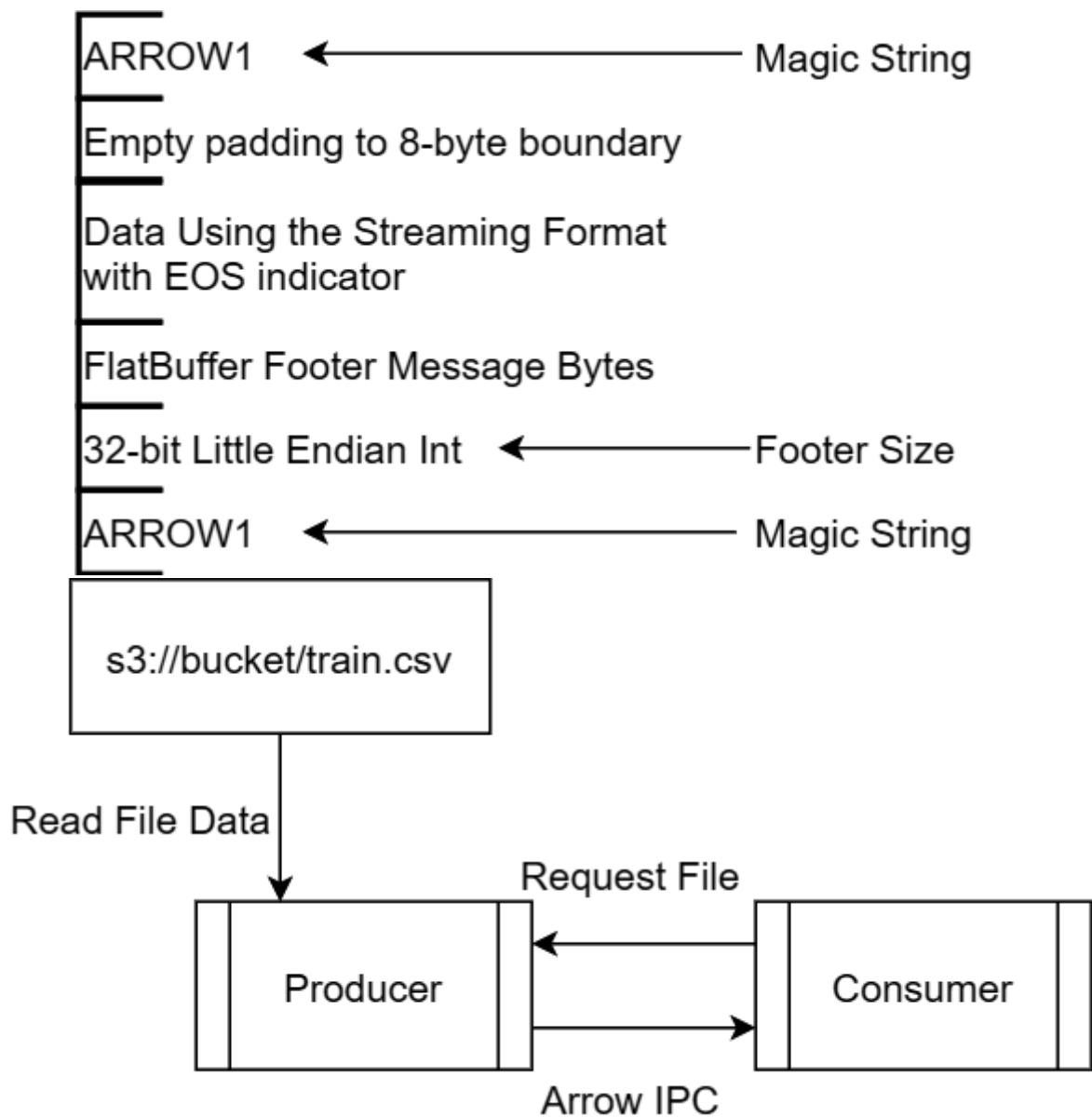
32-bit Little Endian Int ← Length of FlatBuffer Message

FlatBuffer Bytes ← Metadata

Padding to 8-bytes

Message Body Bytes

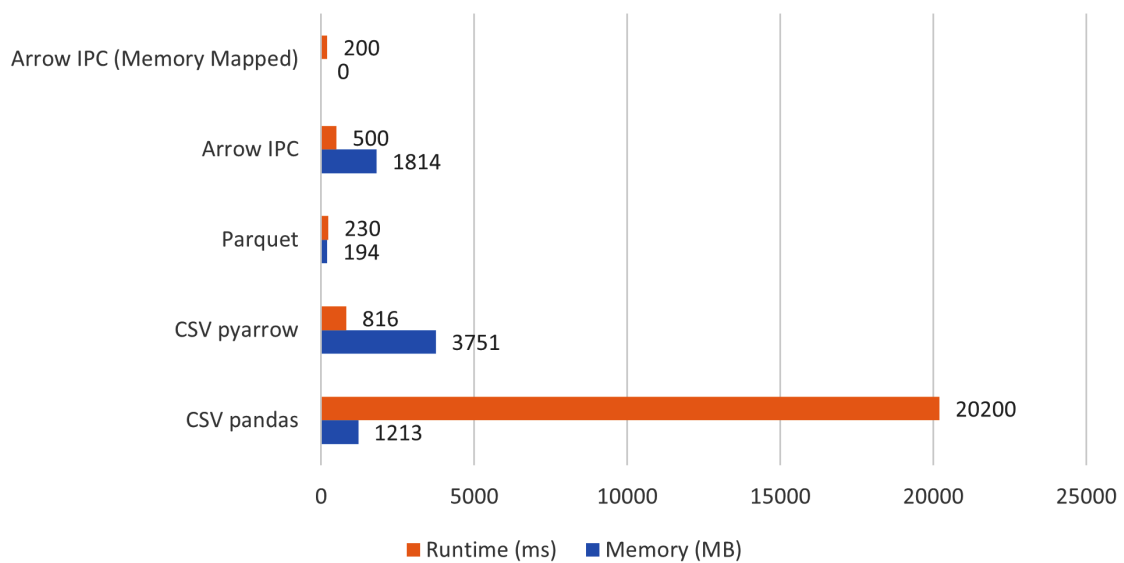


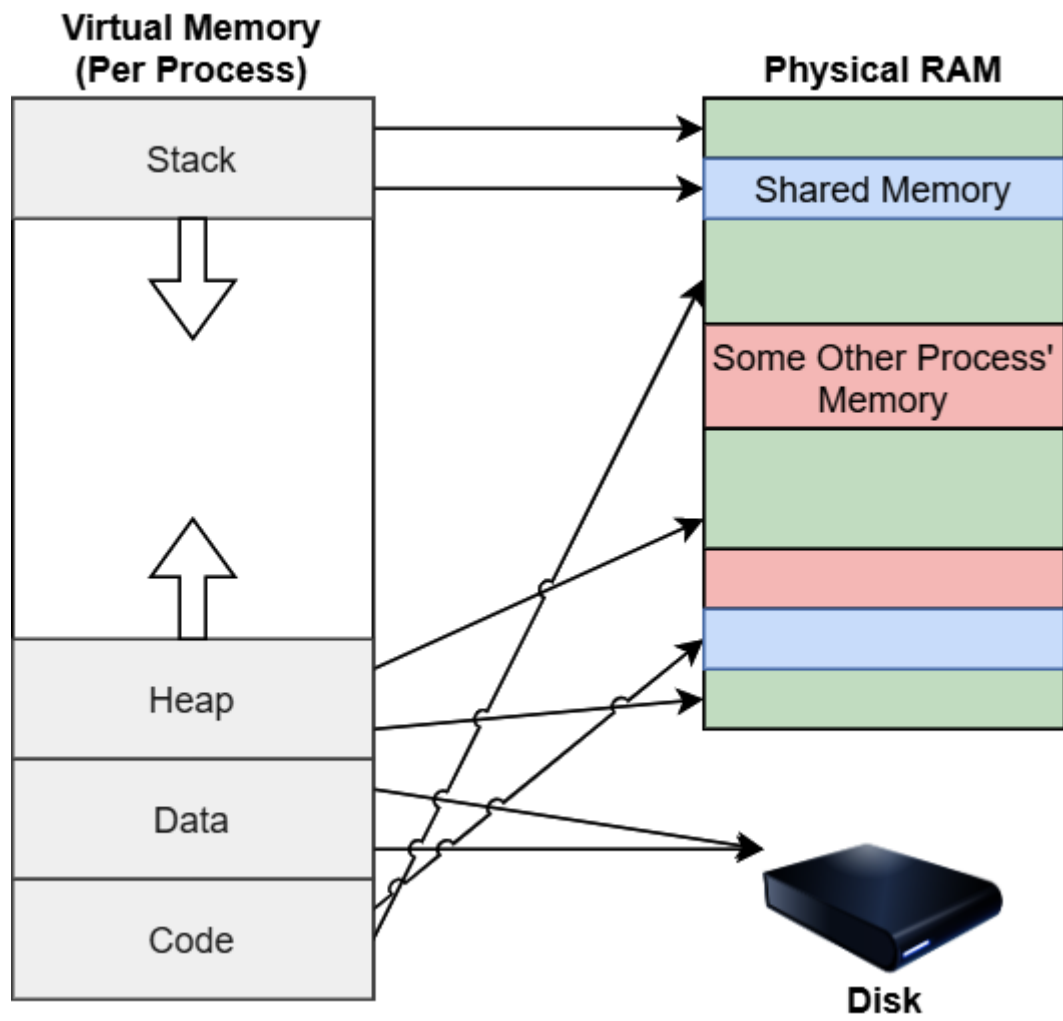


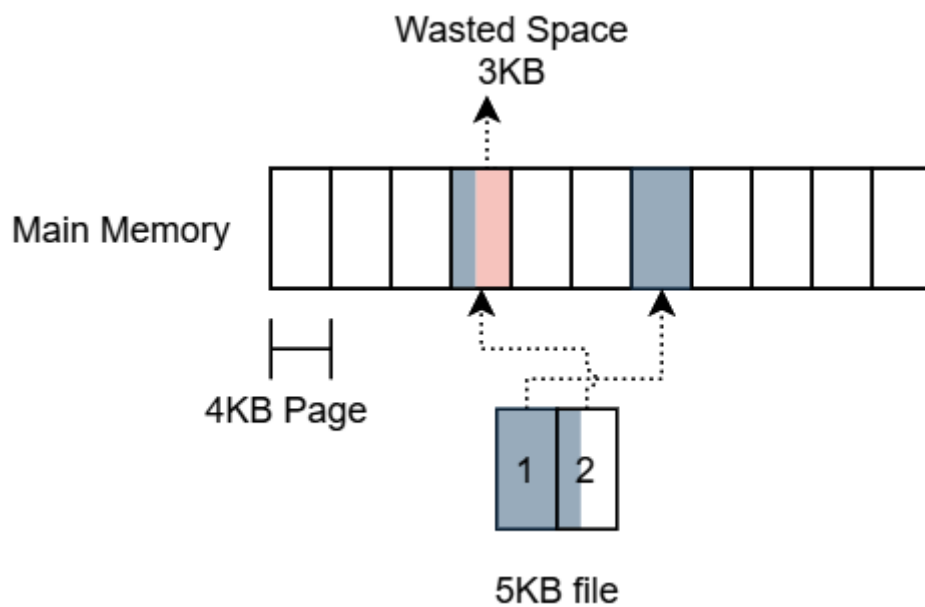
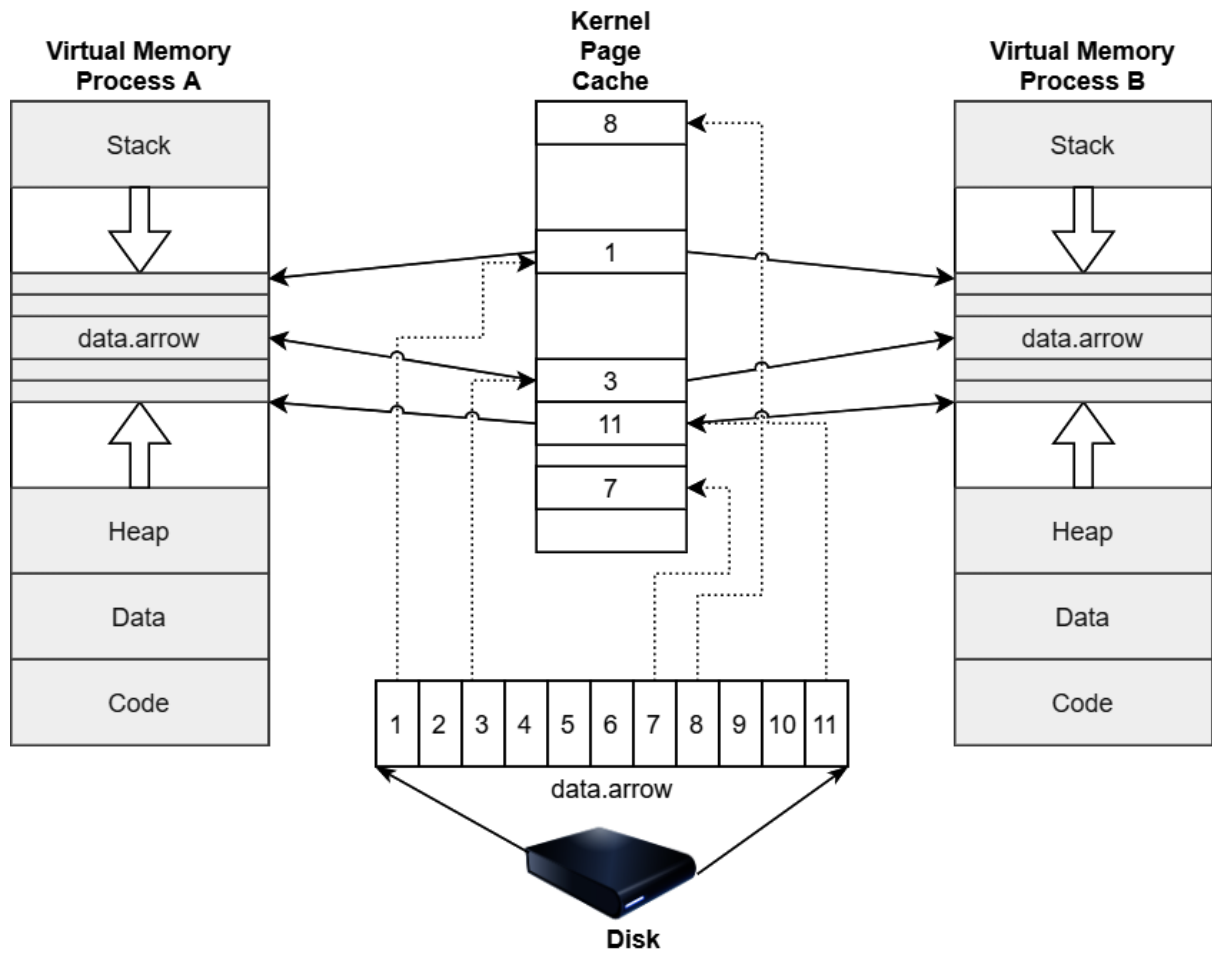
CSV	1.8 GB	<u>calc mean</u>
	pandas read_csv	20.2 s ± 400ms
	pyarrow.csv.read_csv	816 ms ± 27ms
Parquet	286 MB	<u>calc mean</u>
	pandas/pyarrow	230 ms

CSV	1.8 GB	<u>calc mean</u>
	pandas read_csv	20.2 s ± 400ms
	pyarrow.csv.read_csv	816 ms ± 27ms
Parquet	286 MB	<u>calc mean</u>
	pandas/pyarrow	230 ms
Arrow IPC	1.77 GB	<u>calc mean</u>
	pyarrow mmap	200 ms

Memory Usage and Runtime







Chapter 5: Crossing the Language Barrier with the Arrow C Data API

Arrow Type	Null	Boolean	Int8 Uint8	Int16 Uint16	Int32 Uint32	Int64 Uint64	Float16
Format String	n	B	c C	s S	i I	l L	e

Arrow Type	Float32	Float64	Binary	Large Binary	UTF-8 String	Large UTF-8 String
Format String	f	g	z	Z	u	U

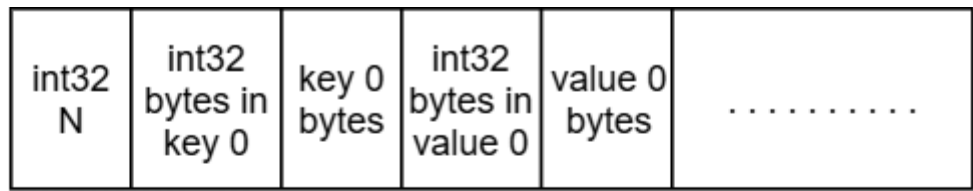
Arrow Type	Format String
Decimal128 [precision 19, scale 10]	d:19,10
Decimal Bit width = NNN [precision 19, scale 10]	d:19,10,NNN
Fixed-Width Binary [42 bytes]	w:42

Arrow Type	Date32 [days]	Date64 [milliseconds]	Time32 [seconds] [milliseconds]	Time64 [microseconds] [nanoseconds]
<i>Format String</i>	tdD	tdm	tts ttm	ttu ttn

Arrow Type	Duration [seconds] [milliseconds]	Duration [microseconds] [nanoseconds]	Interval [months] [day, time]	Interval [month, day, nanoseconds]
<i>Format String</i>	tDs tDm	tDu tDn	tiM tiD	tin

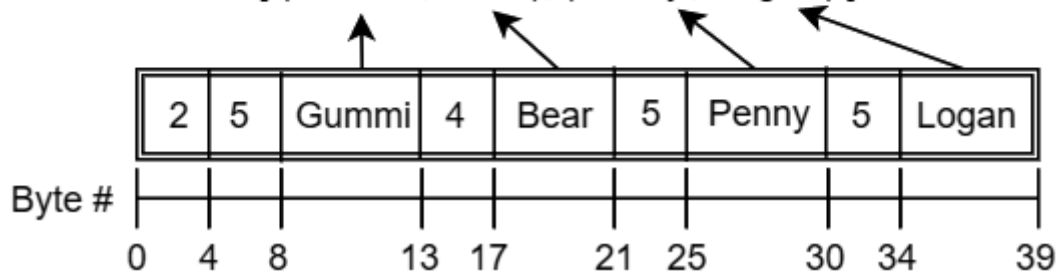
Arrow Type	Timestamp with Timezone "..." [seconds] [milliseconds]	Timestamp with Timezone "..." [microseconds] [nanoseconds]
<i>Format String</i>	tss:... tsm:...	tsu:... tsn:...

Arrow Type	List	Large List	Fixed-size List [123 items]	Struct	Map	Dense Union type- ids I,J,...	Sparse Union type- ids I,J,...
<i>Format String</i>	+l	+L	+w:123	+s	+m	+ud:I,J,...	+us:I,J,...



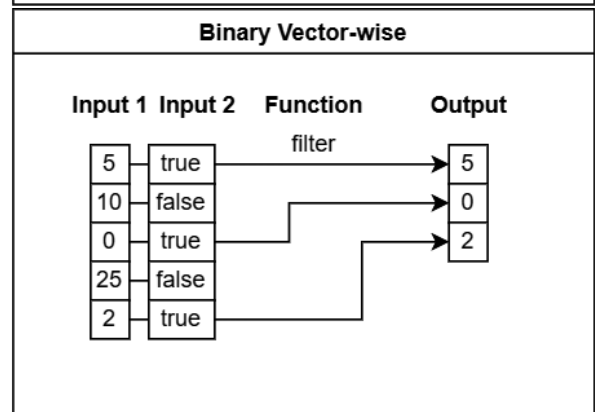
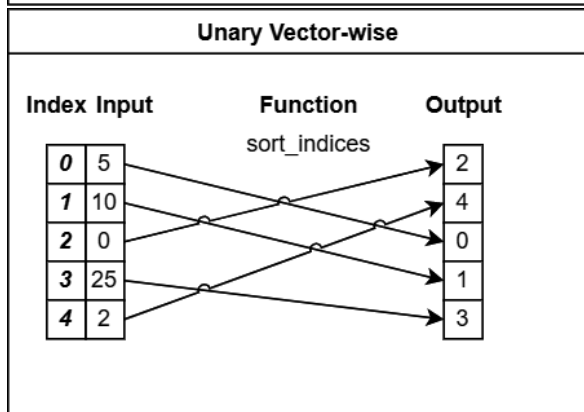
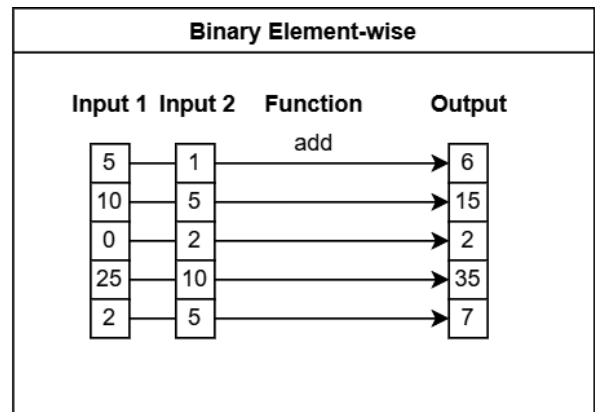
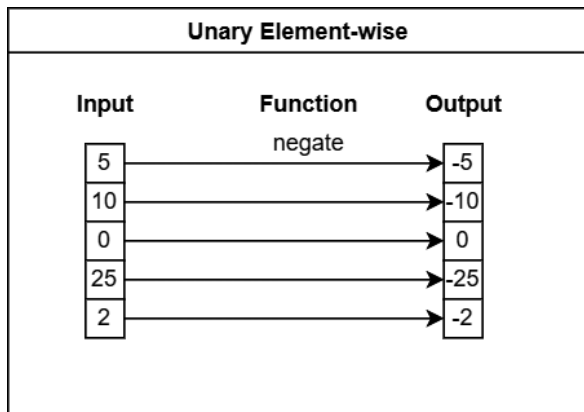
Repeat N times

[('Gummi', 'Bear'), ('Penny', 'Logan')]

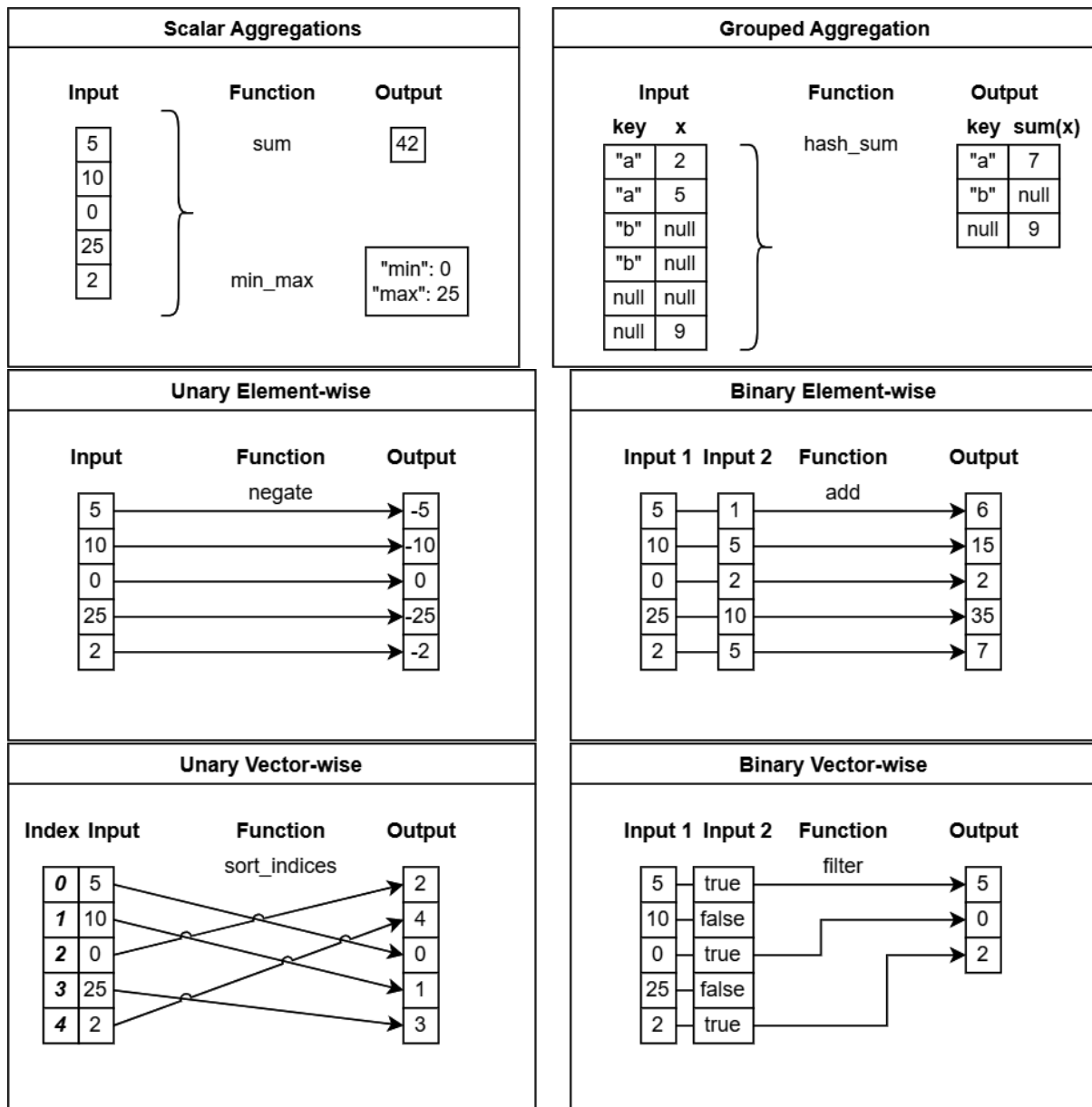


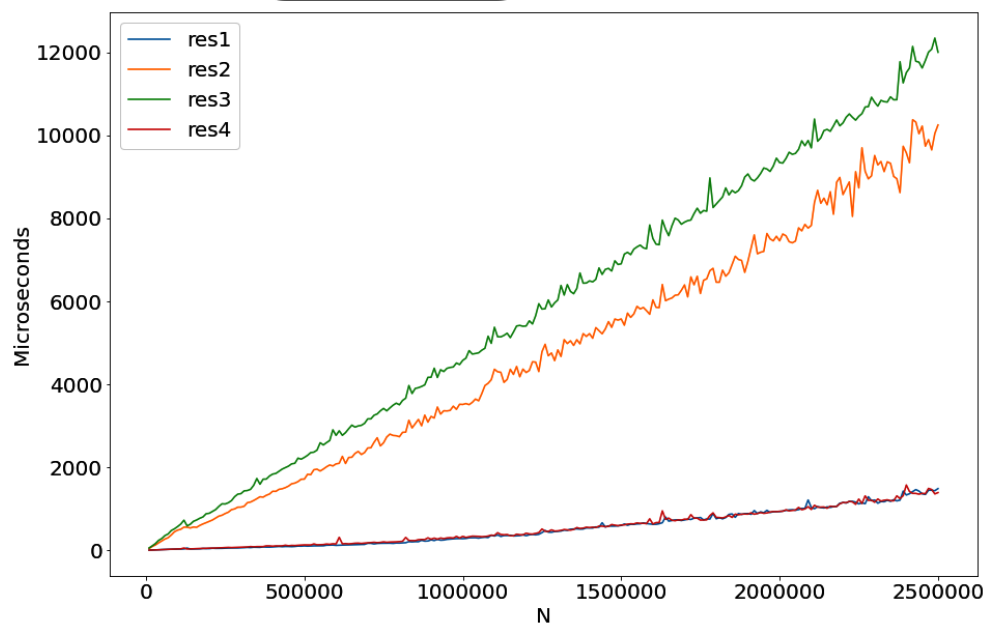
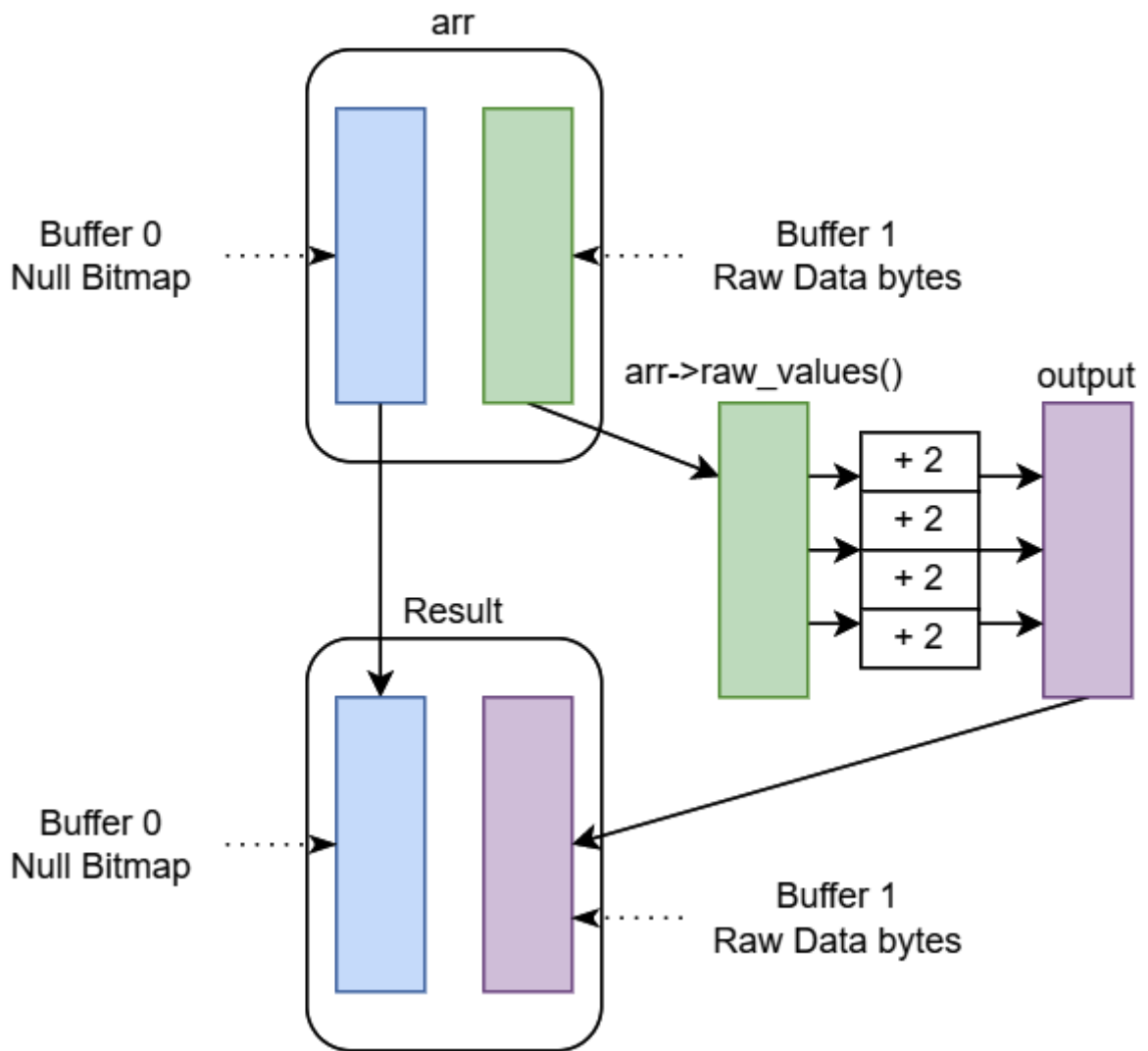
Scalar Aggregations		
Input	Function	Output
<div><div>5</div><div>10</div><div>0</div><div>25</div><div>2</div></div>	sum	<div>42</div>
	min_max	<div>"min": 0 "max": 25</div>

Grouped Aggregation		
Input	Function	Output
key x	hash_sum	key sum(x)
"a" 2	}	"a" 7
"a" 5		"b" null
"b" null		null 9
"b" null		
null null		
null 9		

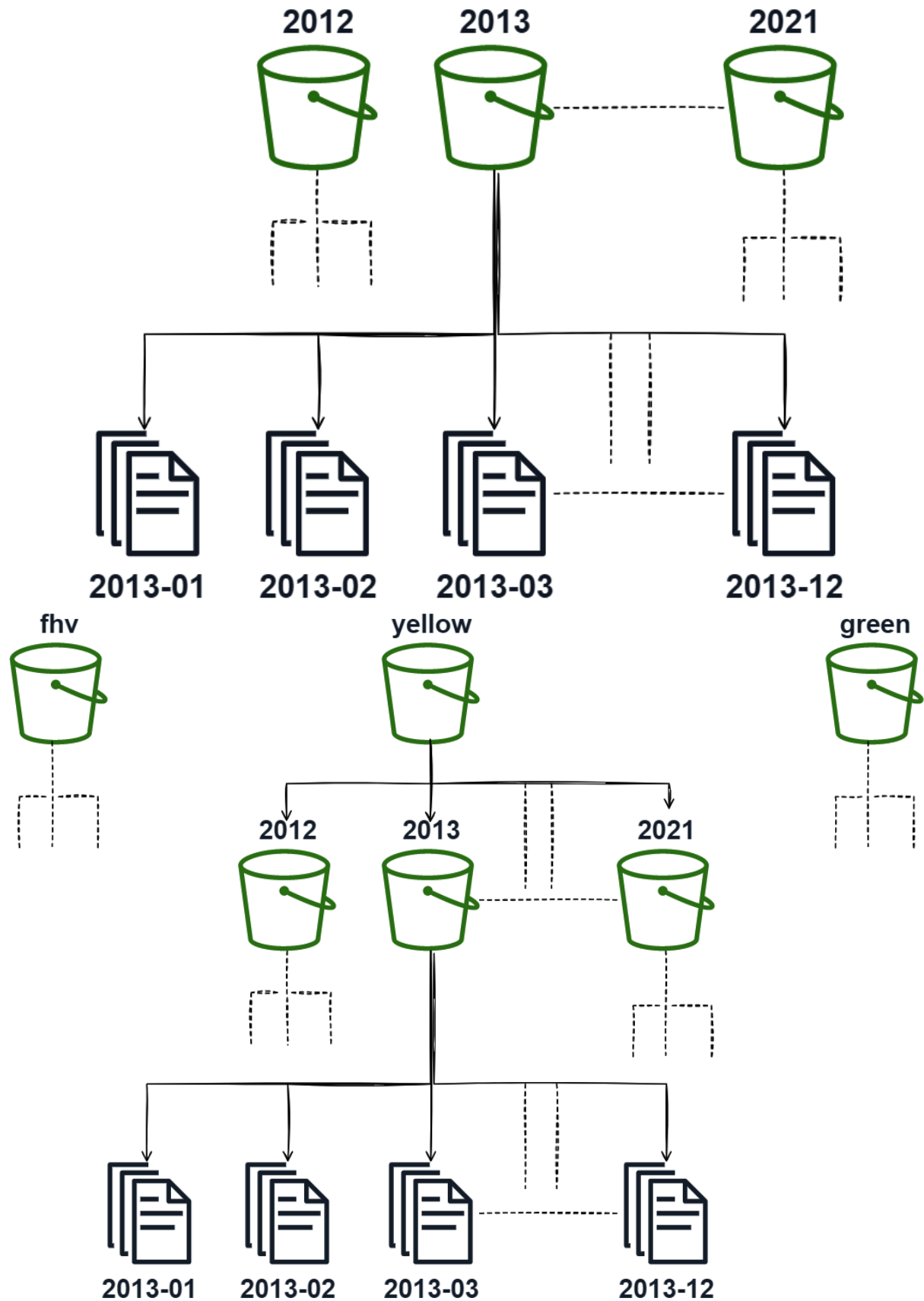


Chapter 6: Leveraging the Arrow Compute APIs





Chapter 7: Using the Arrow Datasets API



data1.parquet		
a	b	c
0	9	1
1	8	2
2	7	1
3	6	2
4	5	1

Input

data2.parquet		
a	b	c
5	4	2
6	3	1
7	2	2
8	1	1
9	0	2

Expression

Output

a	b	c
0	9	1
1	8	2
2	7	1
3	6	2
4	5	1

field_ref("b")

9
8
7
6
5

Input			Expression	Output
a	b	c		
0	9	1	call("less", [field_ref("a"), literal(4)])	true
1	8	2		true
2	7	1		true
3	6	2		true
4	5	1		false

ursa-labs-taxi-data



37 GB 1.54 billion rows

→ 2009/01/data.parquet
→ 2009/02/data.parquet
→ 2009/03/data.parquet
⋮
→ 2019/06/data.parquet

Per File:

read
passenger_count
split into batches

Step 1

Step 2

Per Batch:

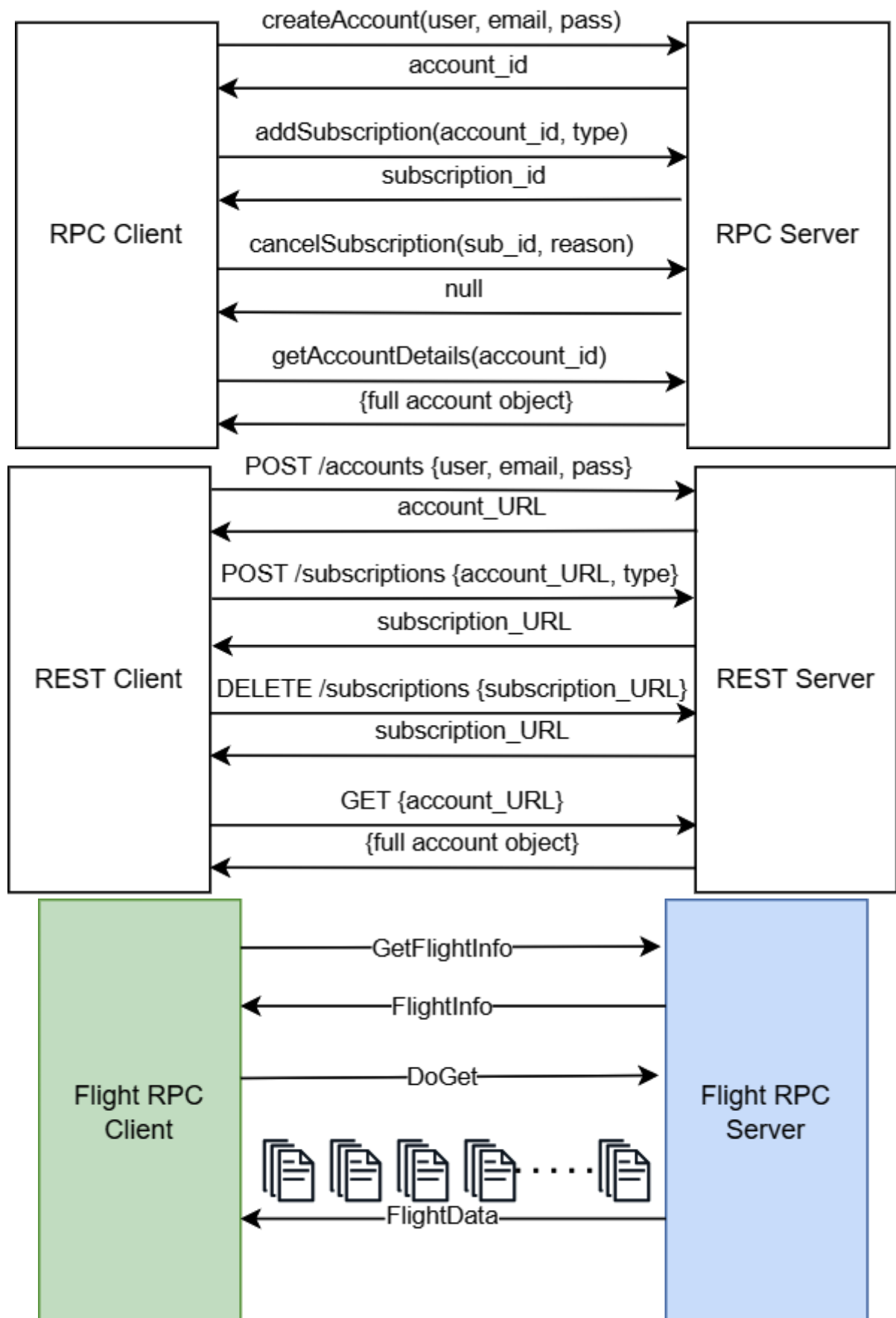
total_passengers += sum(passenger_count)
count += #rows

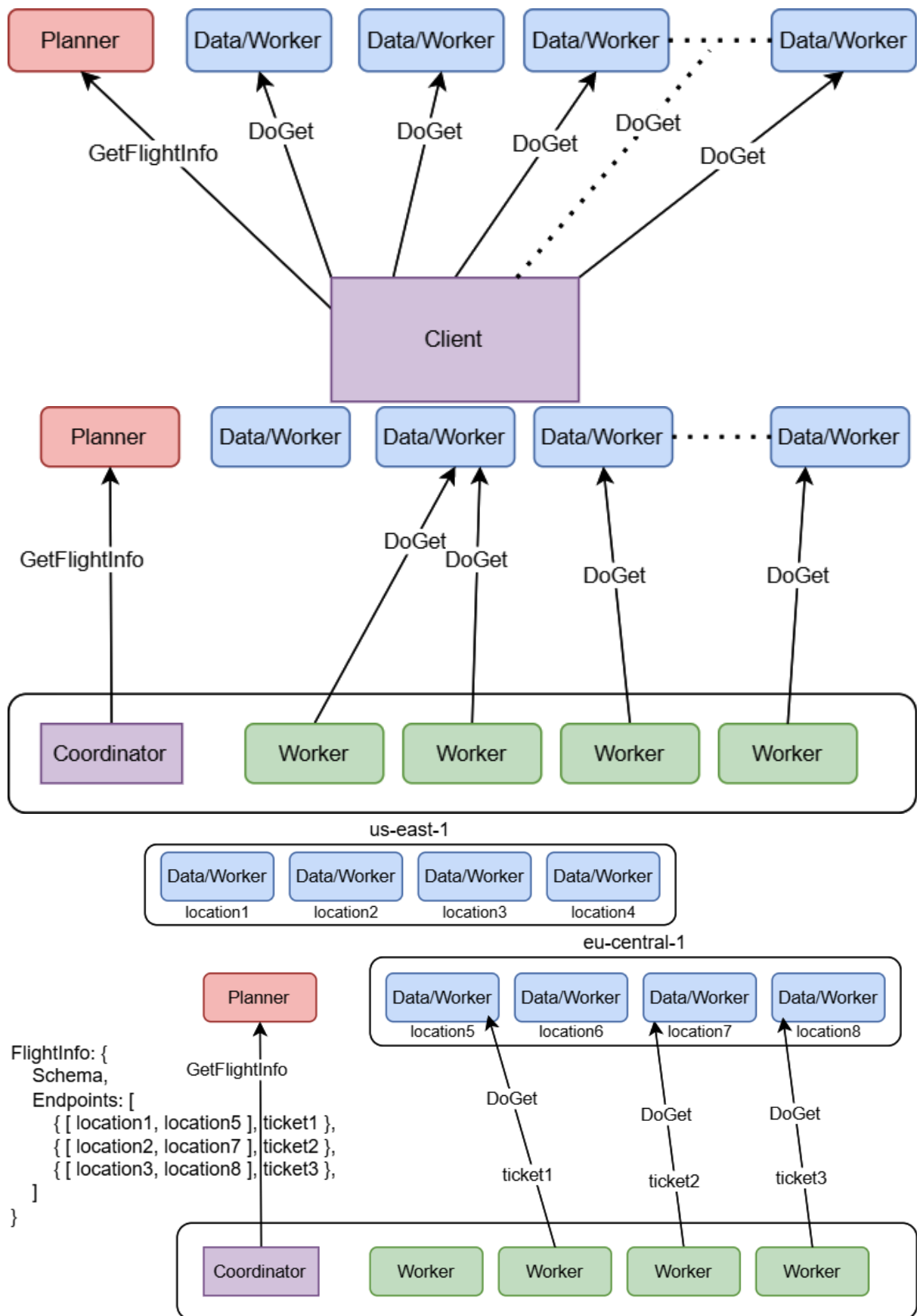
Step 3

total_passengers / count

Step 4

Chapter 8: Exploring Apache Arrow Flight RPC







Welcome to Dremio, please log in.

Username

dremio

Password

.....

Log In

[Privacy](#)

Data Lakes (0) »



You do not have any data lakes.

Add Sample Source

Add Data Lake

> External Sources (0) »





Name ▾

Action

1_0_0.parquet

1_1_0.parquet

1_2_0.parquet

1_3_0.parquet

1_4_0.parquet

1_5_0.parquet

1_5_1.parquet

 NYC-taxi-trips
Samples.samples.dre...

Data

Catalog

Reflections



Run

Preview

1 SELECT * FROM "NYC-taxi-trips"

6 fields

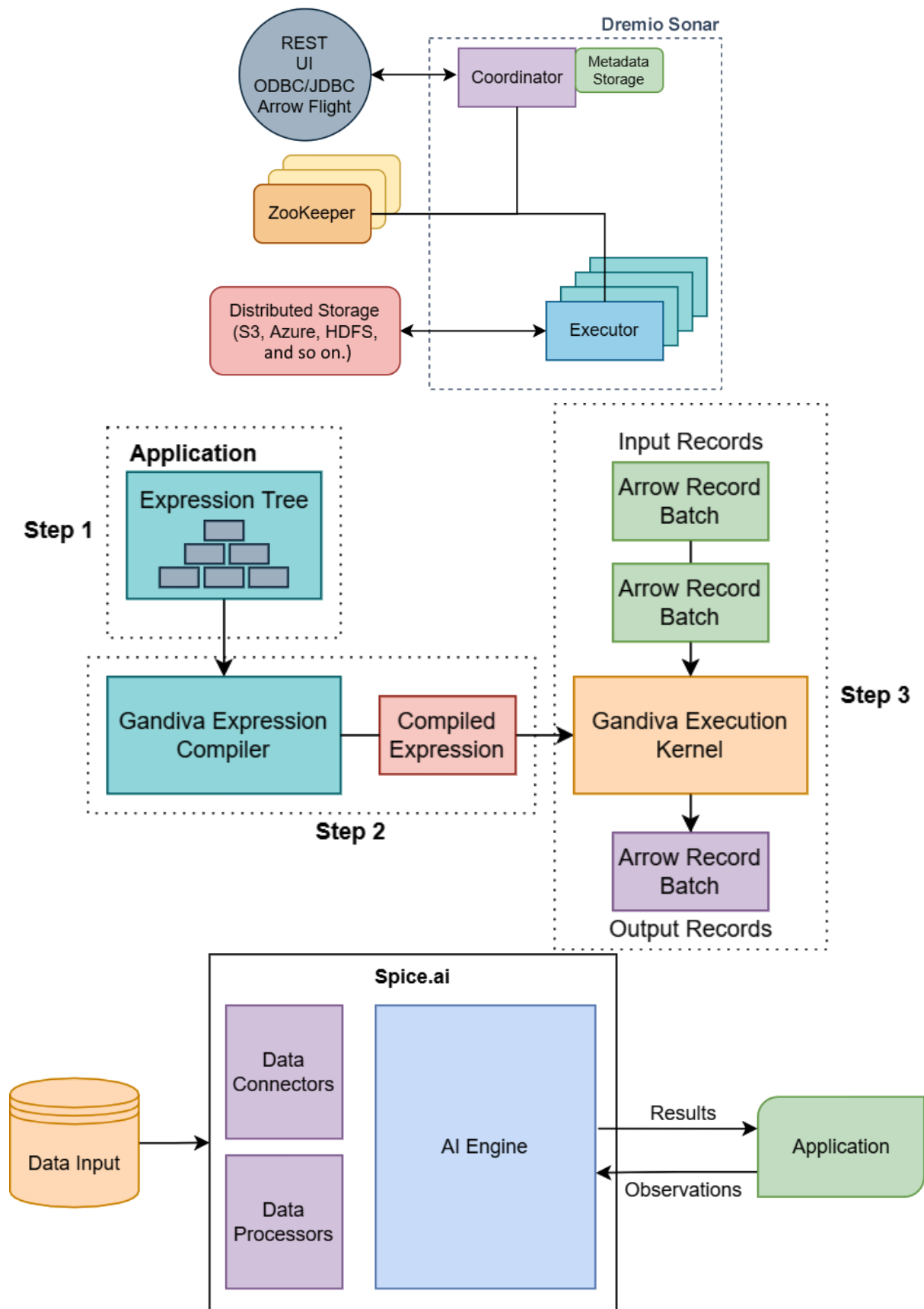
Add Field

Group By

Join

pickup_datetime	# passenger_count	## trip_distance
2013-05-31 07:53:00.000	1	
2013-05-31 08:23:00.000	1	
2013-05-31 08:15:00.000	1	
2013-05-30 22:09:00.000	1	
2013-05-30 22:04:00.000	2	
2013-05-30 21:56:00.000	1	
2013-05-30 21:52:00.000	1	
2013-05-30 22:03:00.000	6	
2013-05-30 22:10:00.000	5	

Chapter 9: Powered By Apache Arrow



Chapter 10: How to Leave Your Mark on Arrow

Create Issue

Configure Fields

All fields marked with an asterisk (*) are required

Project*

Apache Arrow (ARROW)

Issue Type*

Improvement



Apache Arrow / ARROW-7138

[Go][CI] Pre-install the go dependencies in the dockerfile get

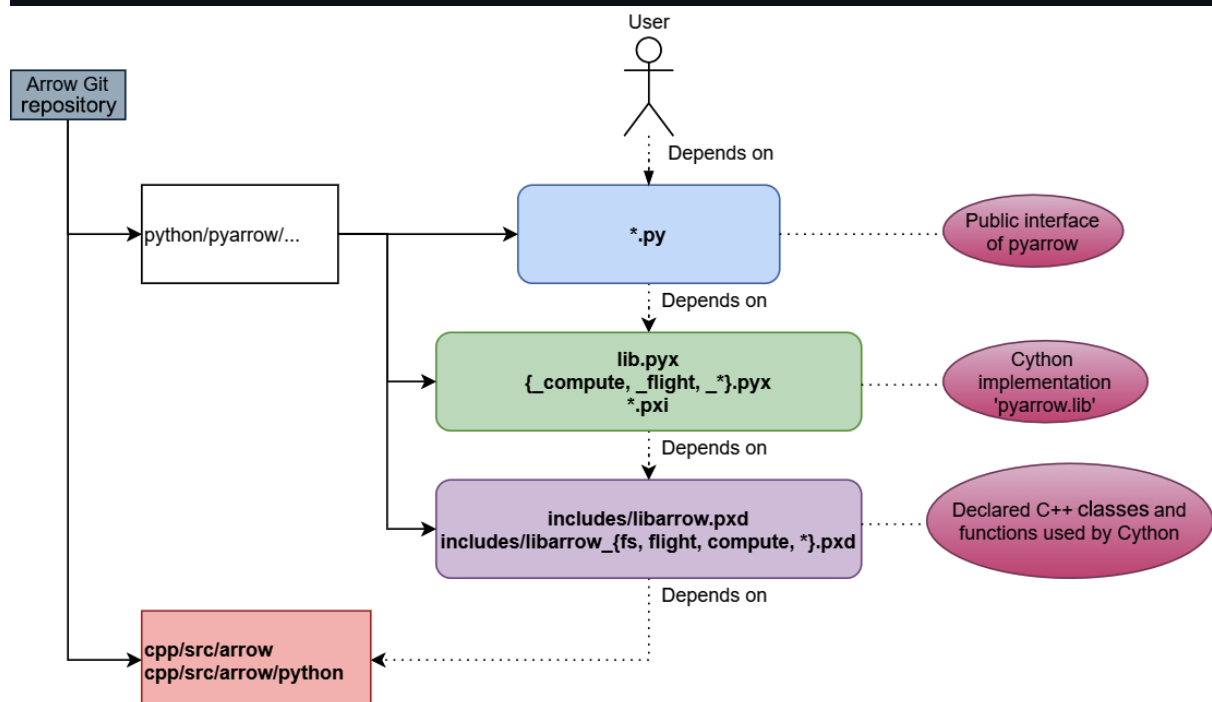
Edit Add comment Assign More Start Progress Resolve Issue Close Issue

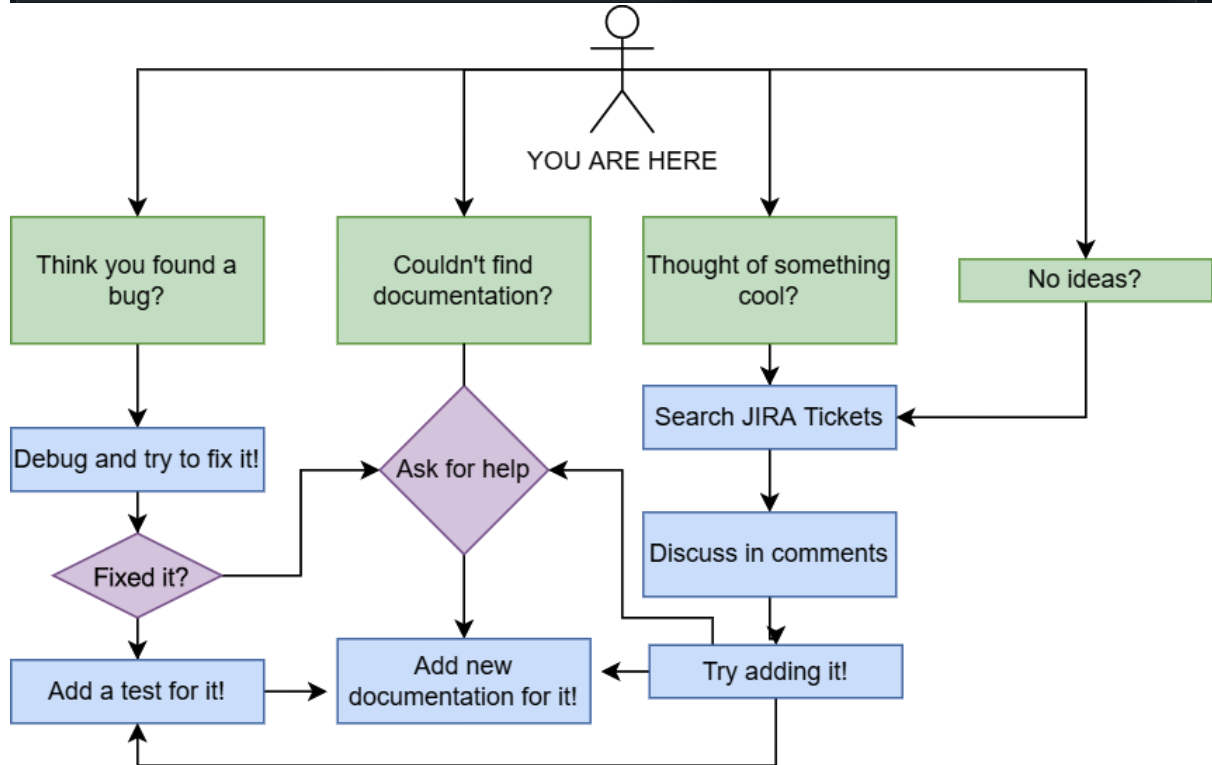
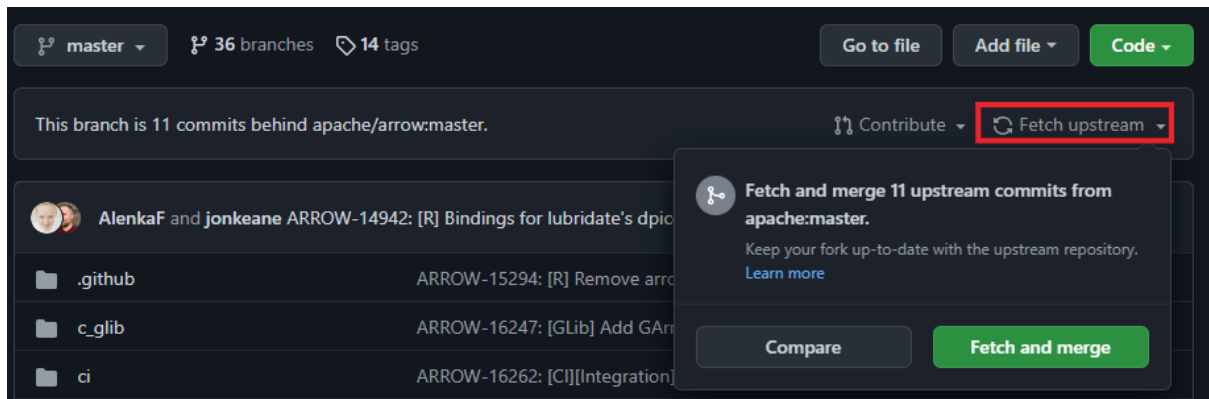
Details

Type:	Improvement	Status:	OPEN
Priority:	Major	Resolution:	Unresolved
Affects Version/s:	None	Fix Version/s:	None
Component/s:	Continuous Integration, Go		
Labels:	None		

apache / arrow Public

Error Unwatch 340 Fork 2.3k Star 9.2k





Chapter 11: Future Development and Plans

