

Maximizing Apache Druid performance

Beyond the Basics

September 2020

Gian Merlino, Cofounder & CTO, Imply



VIRTUAL
DRUID SUMMIT

Who am I?



Gian Merlino



Committer & PMC chair at



Cofounder at



(we're hiring!)



Druid in the wild



VIRTUAL
DRUID SUMMIT

Druid powers analytic applications

VIRTUAL
DRUID SUMMIT



looker



Metabase



Superset

Druid in the wild



How Netflix uses Druid for Real-time Insights to Ensure a High-Quality Experience



Netflix Technology Blog

Mar 3 · 9 min read



By [Ben Sykes](#)

100+ billion rows/day

1+ trillion rows, 1+ year retained

100s of servers

sub-second to few seconds query latency

mix of streaming and batch ingest

Interactive Analytics at MoPub: Querying Terabytes of Data in Seconds

By [Aaron Rolett](#) and [Shravana Krishnamurthy](#)

Wednesday, 3 July 2019



How Druid enables analytics at Airbnb

Realtime and batch analytics at Airbnb and the role Druid plays in our analytics system architecture



Pala

[Follow](#)

Nov 13, 2018 · 10 min read

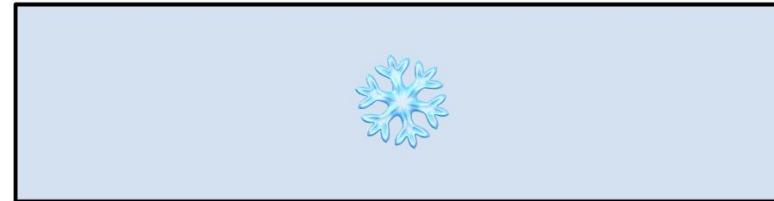


By [Pala Muthiah](#) and [Jinyang Li](#)

Hot vs. cold



- ⌚ 0.1–3s query
- /fa water/ fresh data
- /fa trophy/ high concurrency
- /fa bicycle/ highly interactive (OLAP)



- ⌚ slow queries are ok
- /fa water/ less fresh data is ok
- /fa trophy/ low concurrency
- /fa bicycle/ reporting / planning



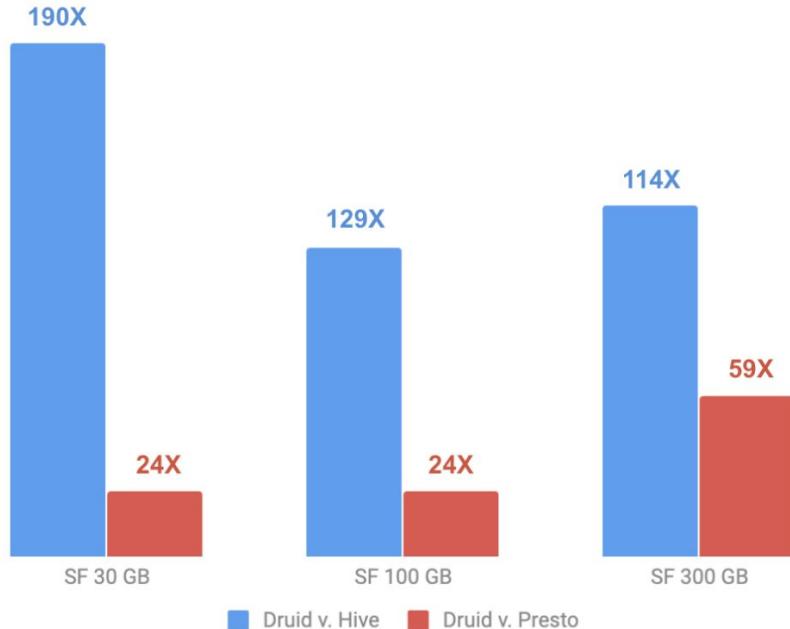
Designed for performance



VIRTUAL
DRUID SUMMIT

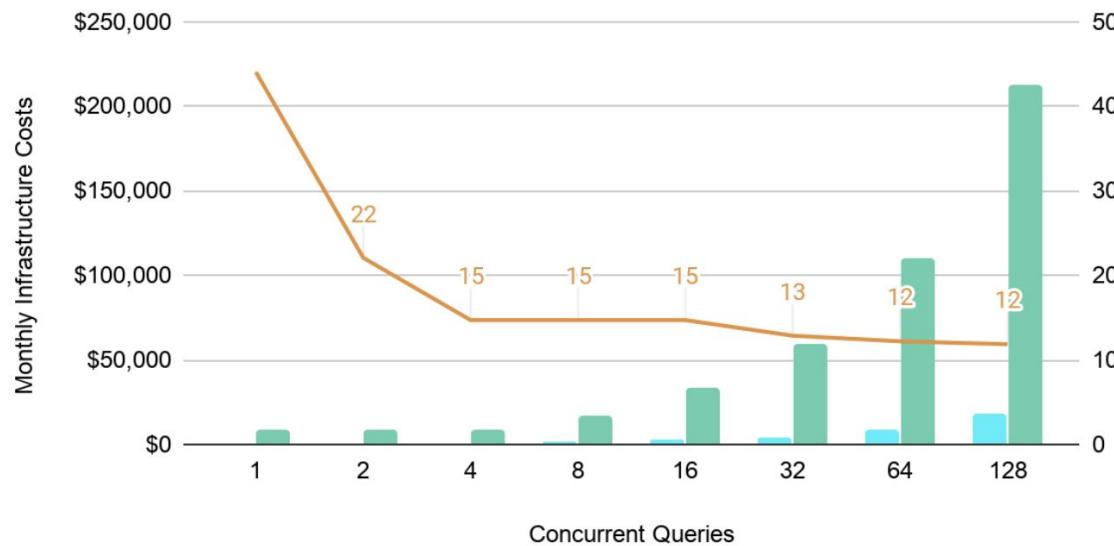
Druid is designed for performance

vs. leading open source SQL engines

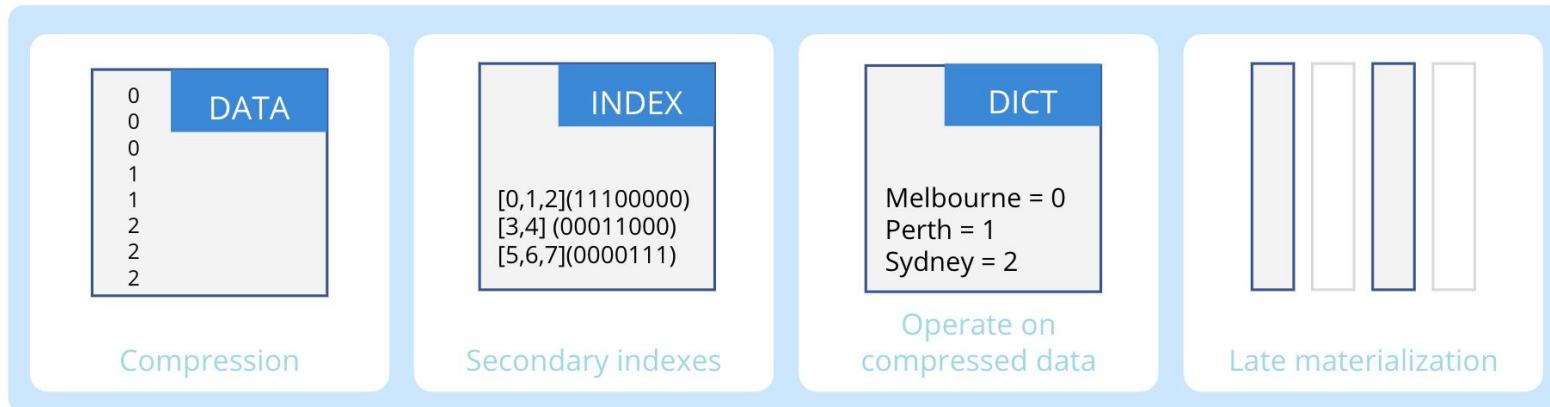


Druid is designed for performance

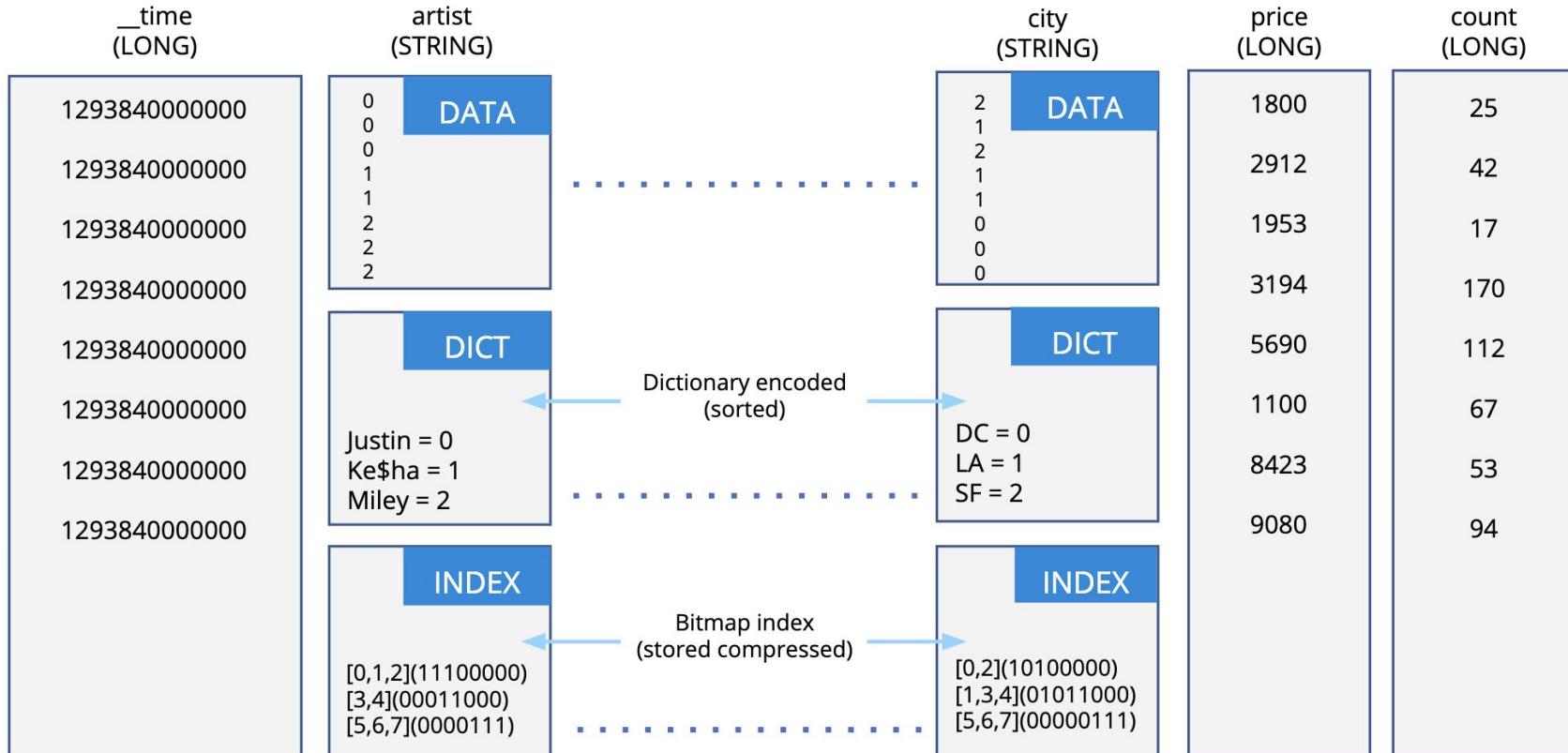
vs. a leading cloud-based data warehouse

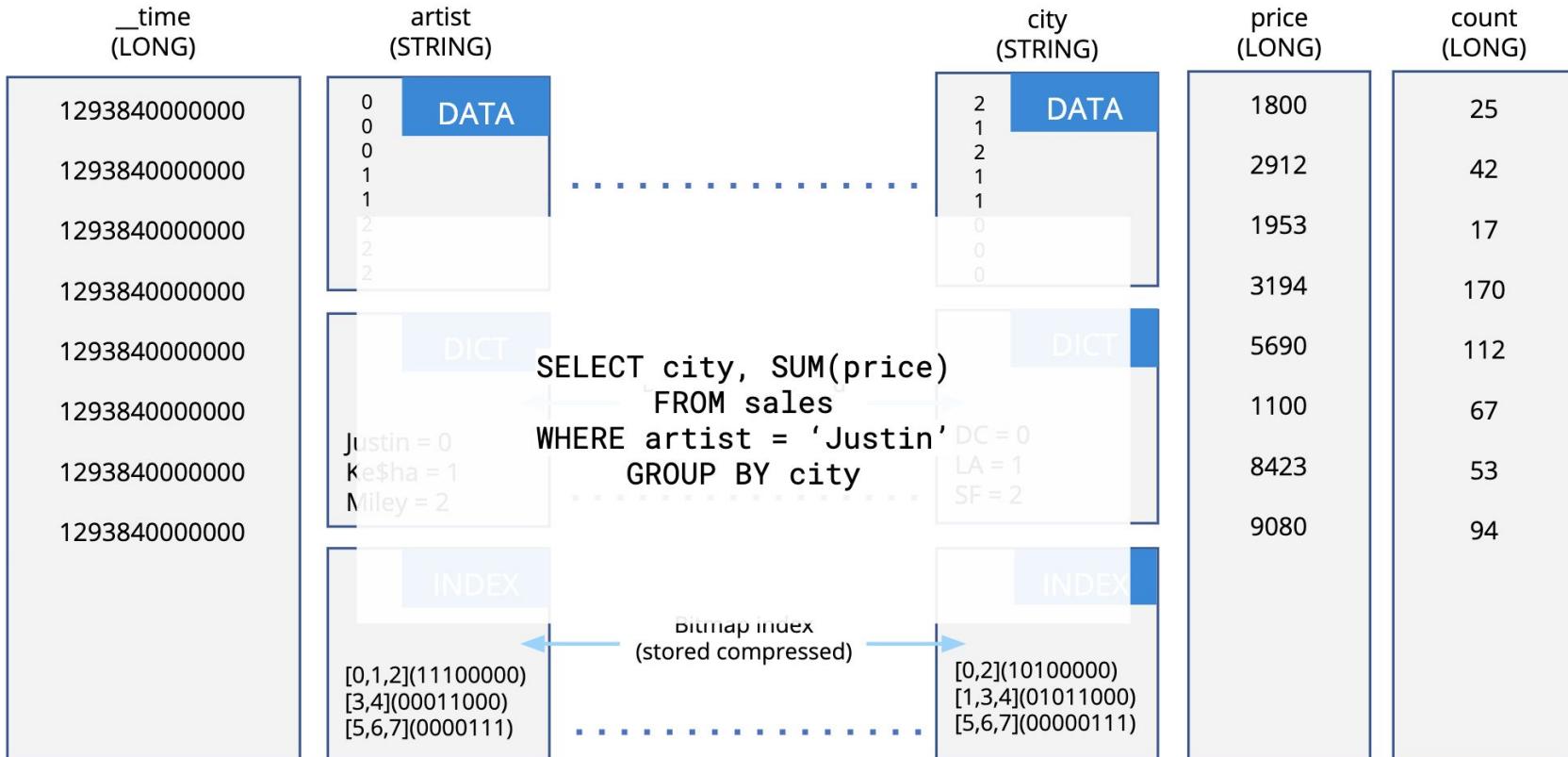


Integrated engine and data format



Query execution pipeline





```
time
SELECT city, SUM(price)
  FROM sales
 WHERE artist = 'Justin'
 GROUP BY city
```

1293840000000
1293840000000
1293840000000
1293840000000
1293840000000

artist
STRING)

0
1
1
2
2
2
2

DATA
DICT
Justin = 0
Ke\$ha = 1
Miley = 2

INDEX
[0,1,2](11100000)
[3,4](00011000)
[5,6,7](00001111)

city
STRING)

2
1
2
1
1
0
0
0

DATA
DICT
DC = 0
LA = 1
SF = 2

INDEX
[0,2](10100000)
[1,3,4](01011000)
[5,6,7](00000111)

price
LONG)

1800
2912
1953
3194
5690
1100
8423
9080

count
LONG)

25
42
17
170
112
67
53
94

Dictionary encoded
(sorted)

Bitmap index
(stored compressed)

```

SELECT city, SUM(price)
  FROM sales
 WHERE artist = 'Justin'
 GROUP BY city
  
```

1293840000000
1293840000000
1293840000000
1293840000000
1293840000000

time
("2011")
artist
(STRING)

DATA

DICT

INDEX

Justin = 0
Ke\$ha = 1
Miley = 2

[0,1,2](11100000)
[3,4](00011000)
[5,6,7](00000111)

city
(STRING)

DATA

DICT

INDEX

2
1
2
1
1
0
0
0

DC = 0
LA = 1
SF = 2

[0,2](10100000)
[1,3,4](01011000)
[5,6,7](00000111)

price
(LONG)

DATA

DICT

INDEX

1800
2912
1953
3194
5690
1100
8423
9080

count
(LONG)

DATA

DICT

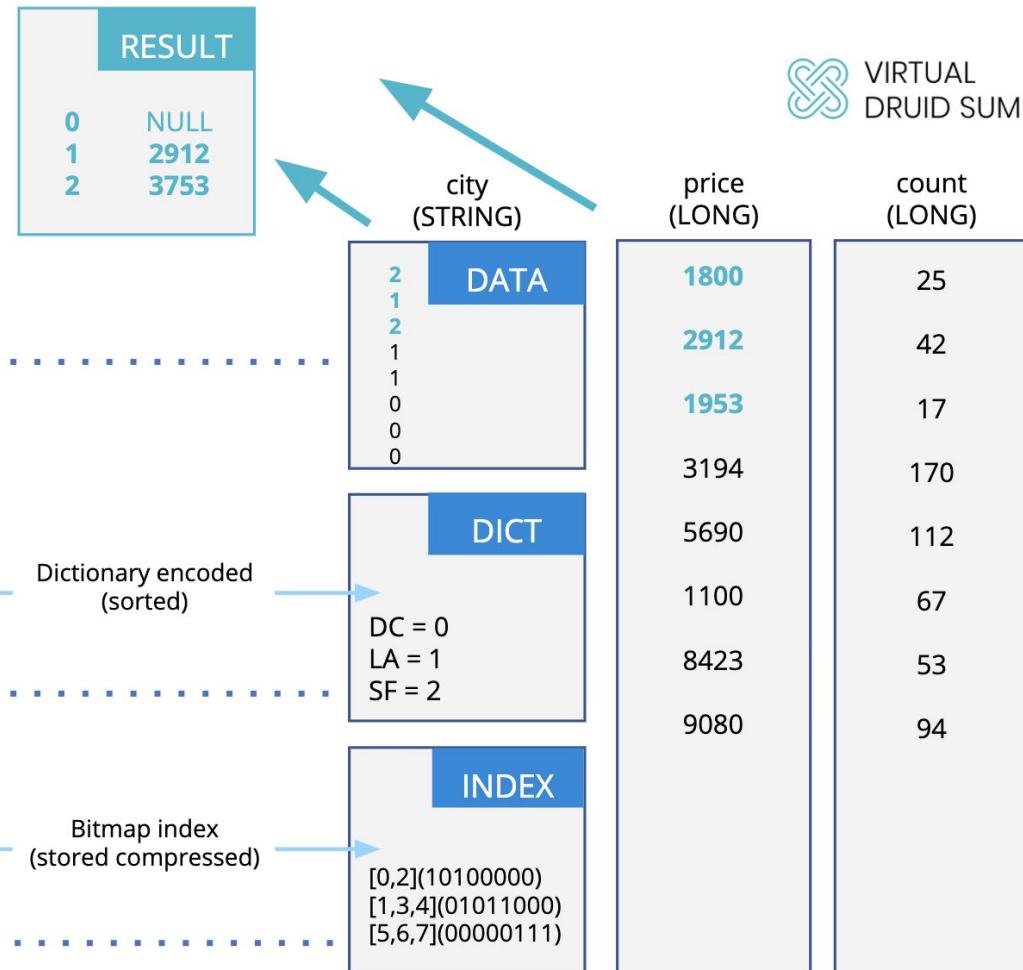
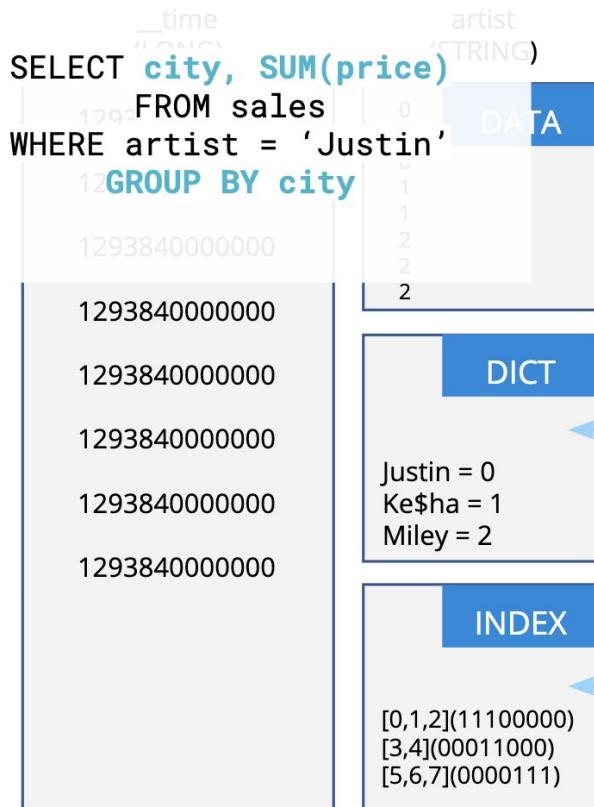
INDEX

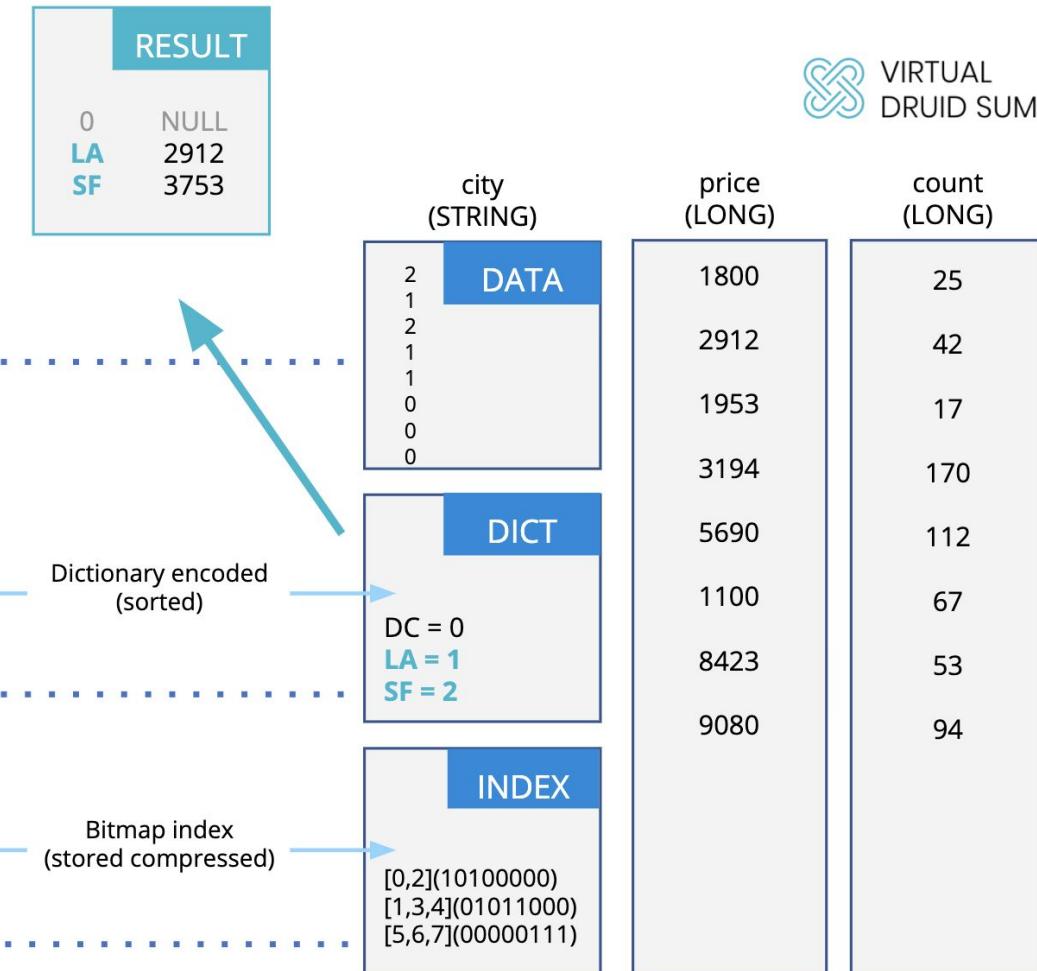
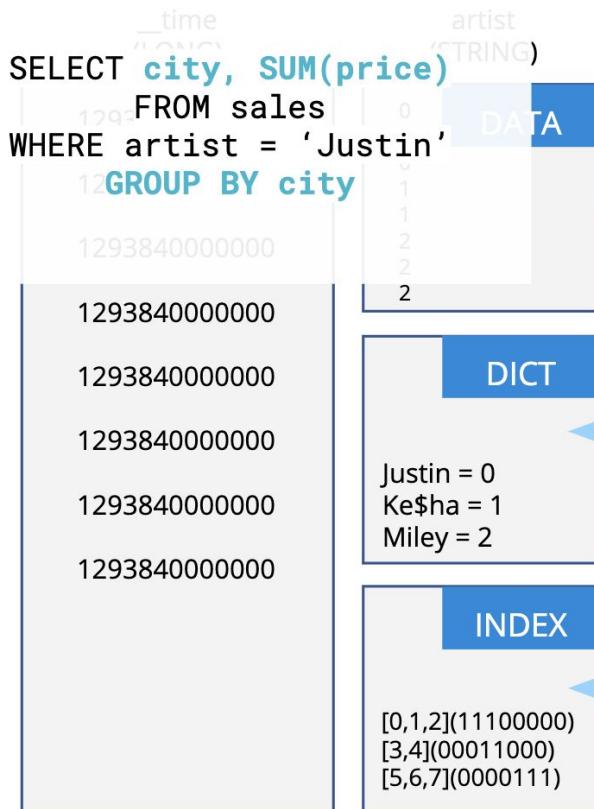
25
42
17
170
112
67
53
94

Dictionary encoded
(sorted)

Bitmap index
(stored compressed)







```

SELECT city, SUM(price)
  FROM sales
 WHERE artist = 'Justin'
 GROUP BY city
  
```

	time	artist	DATA
0	1293840000000		
1	1293840000000		
2	1293840000000		
3	1293840000000	Justin = 0	2
4	1293840000000	Ke\$ha = 1	1
5	1293840000000	Miley = 2	0
6	1293840000000		
7	1293840000000		
8	1293840000000		
9	1293840000000		

RESULT	
LA	2912
SF	3753

city (STRING)	price (LONG)	count (LONG)
2	1800	25
1	2912	42
2	1953	17
1	3194	170
0	5690	112
0	1100	67
0	8423	53
0	9080	94

Dictionary encoded
(sorted)

DC = 0
LA = 1
SF = 2

Bitmap index
(stored compressed)

[0,2](10100000)
[1,3,4](01011000)
[5,6,7](00000111)



Beyond the basics



VIRTUAL
DRUID SUMMIT

Getting my types right

Good segment sizes

Denormalization

Using approximations
when appropriate

Mastering data layout

Parachuting
computation behind
enemy lines



Getting my types right

Good segment sizes

Denormalization

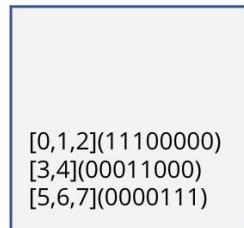
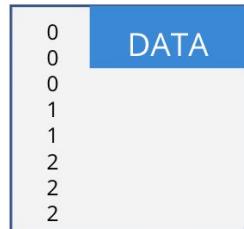
Using approximations
when appropriate

Mastering data layout

Parachuting
computation behind
enemy lines



Mastering data layout



Mastering data layout



Mastering data layout



Mastering data layout

```
“dimensionsSpec”: {  
    “dimensions”: [  
        “artist”,  
        “city”  
    ]  
}
```



Rows are sorted first by time,
then by dimensions in the order
provided in the dimensionsSpec

Placing “artist” first puts it
immediately after time in the
sort order

Mastering data layout

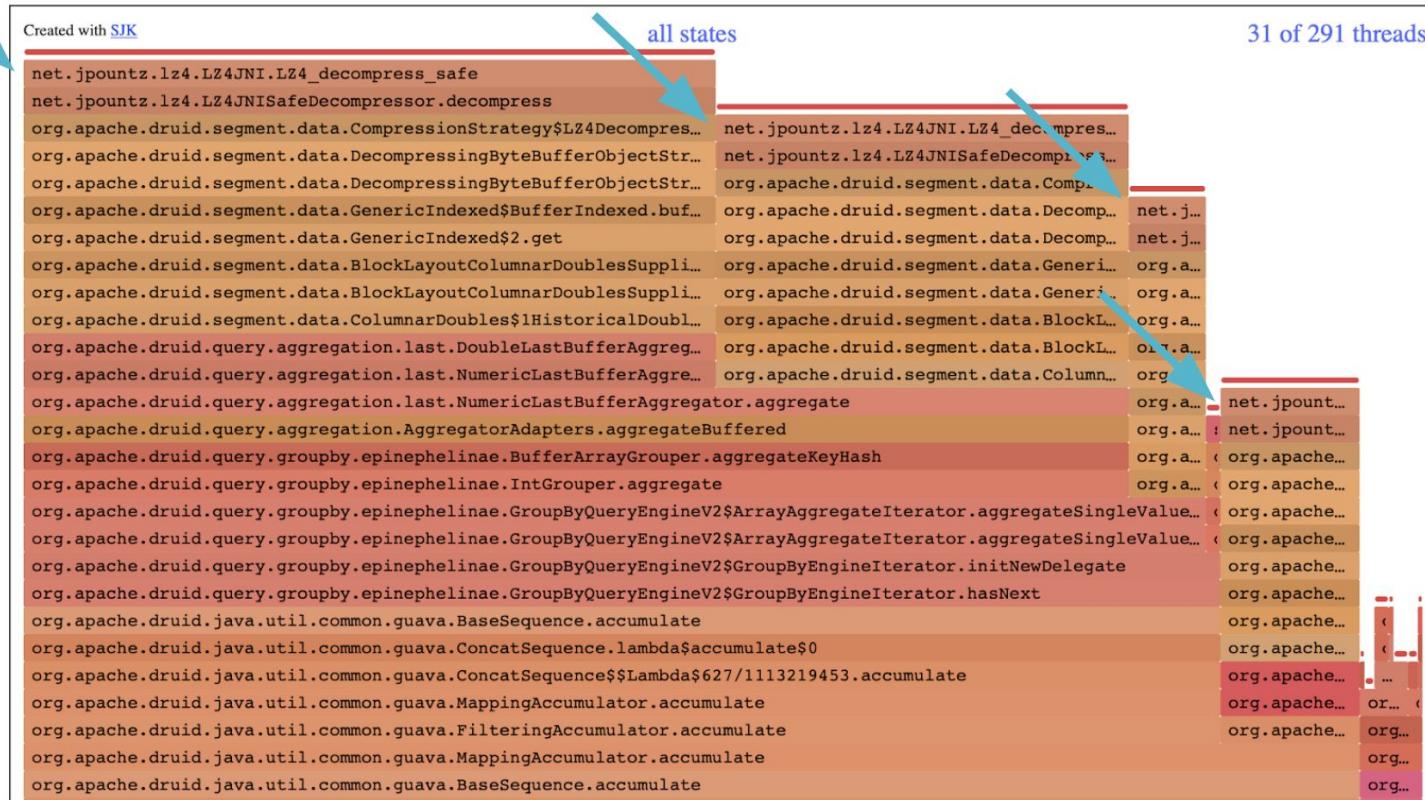
```
“dimensionsSpec”: {  
    “dimensions”: [  
        “artist”,  
        { “type” : “long”, “name”: “timestamp” },  
        “city”  
    ]  
}  
  
“timestampSpec”: {  
    “column”: “timestamp”,  
    “format”: “millis”  
}  
  
“granularitySpec”: {  
    “queryGranularity”: “day”,  
    “segmentGranularity”: “day”  
}
```



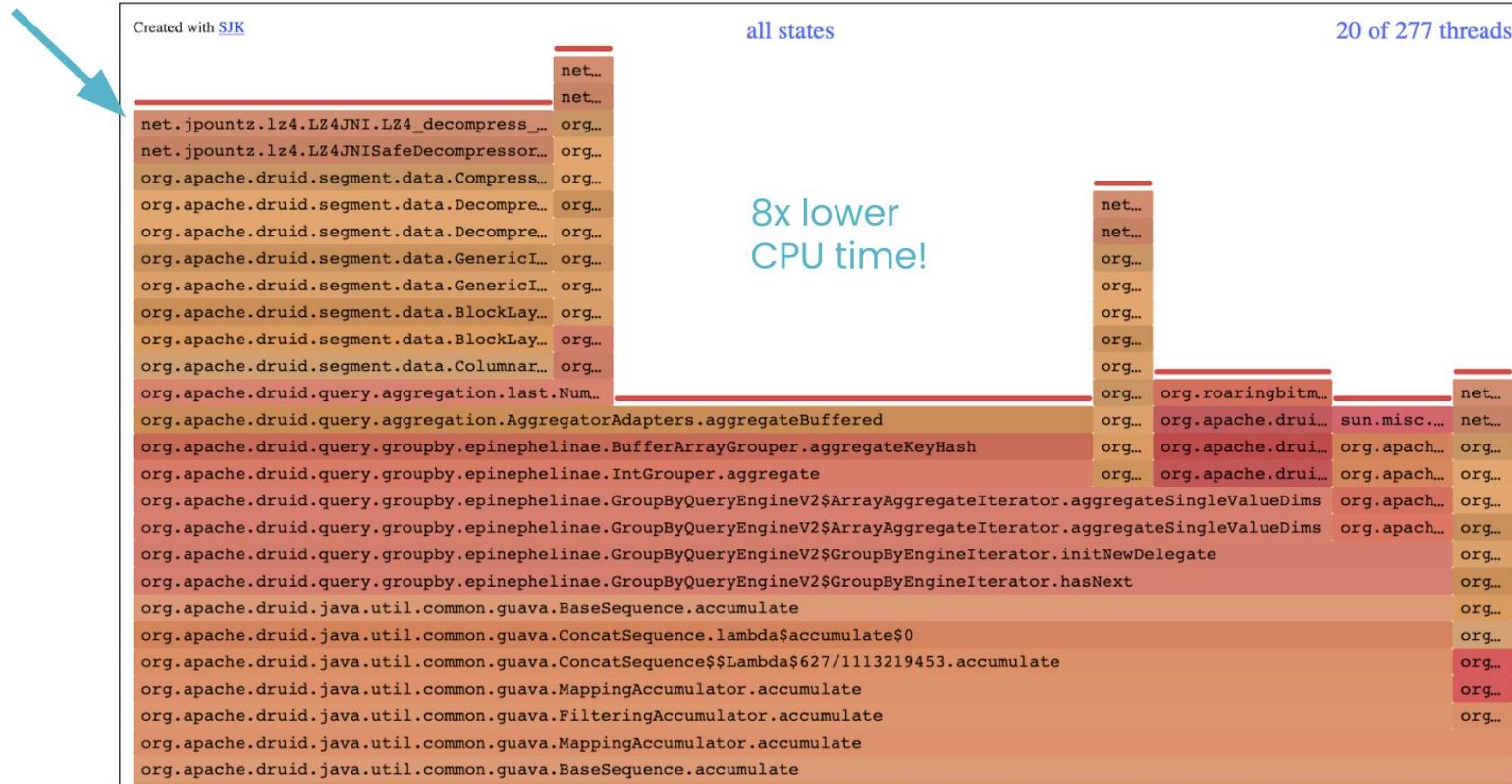
Can put “artist” before time in the sort order using this *secondary timestamp* trick!

(Example assumes “timestamp” column in millis format.)

Mastering data layout



Mastering data layout



Parachuting computations

Data size generally decreases in later stages of a group-by query

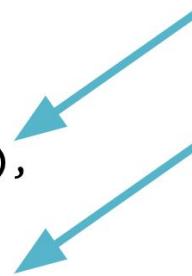
Move computations as late as possible



Parachuting computations

6 sec

```
SELECT
  LOWER(user),
  COUNT(*)
FROM traffic
GROUP BY 1
ORDER BY 2 DESC
```



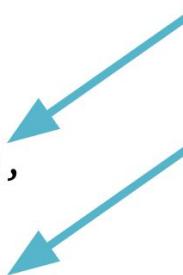
Grouping and selecting an expression of an underlying column

Parachuting computations

0.8 sec

(also 8x faster!)

```
SELECT
  LOWER(user),
  COUNT(*)
FROM traffic
GROUP BY user
ORDER BY 2 DESC
```



Group on column, select expression

Defers computation until after grouping

Not done automatically, because results will not be the same if the function maps two inputs to the same output

Stay in touch



Follow the Druid project on Twitter!

 @druidio

Join the community
(Mailing lists, Slack, meetups)
<https://druid.apache.org/community/>

Thank you!

Time for questions

@gianmerlino



VIRTUAL
DRUID SUMMIT