

# МНОГОЛИКАЯ PANDAS

**Павел Филонов**

# ИДЕЯ ДОКЛАДА



**Павел Филонов**  
Data Scientist  
[t.me/pavel\\_filonov](https://t.me/pavel_filonov)

- разнообразие DS проектов
- воспроизводимость
- быстрая проверка гипотез
- переиспользование кода

# THIS IS THE WAY

Small Data



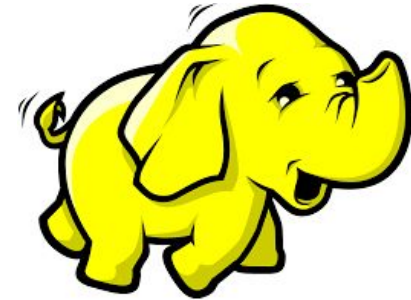
EDA

Enough Data



fit\_predict

Big Data



batch

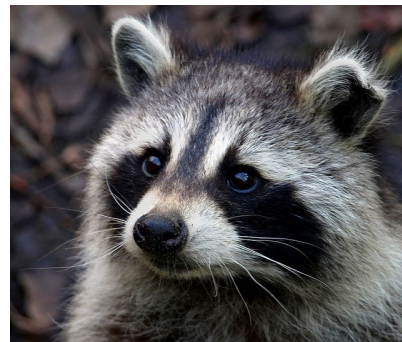
# НАШИ МИЛЫЕ АРІ



Pandas



Dask



Spark

# AMPLAB BIG DATA BENCHMARK

## Rankins

pageURL VARCHAR(300)  
pageRank INT  
avgDuration INT

## UserVisits

sourceIP VARCHAR(116)  
destURL VARCHAR(100)  
visitDate DATE  
adRevenue FLOAT  
userAgent VARCHAR(256)  
countryCode CHAR(3)  
languageCode CHAR(6)  
searchWord VARCHAR(32)  
duration INT

## Documents

*Unstructured HTML documents*

Size	Rankins (rows)	Rankings (bytes)	UserVisits (rows)	UserVisits (bytes)	Documents (bytes)
tiny	1200	77.6 KB	10000	1.7 MB	6.8MB
1node	18 Million	1.28 Gb	155 Million	25.4GB	29.0 GB
5node	90 Million	6.38 Gb	775 Million	126.8 GB	136.9 GB

# ROUND 1



Pandas



Tiny

# SIMPLE QUERIES

rankings

```
.loc[rankings.pageRank > 100]
```

uservisits

```
.groupby(uservisits.sourceIP.str.slice(0, 7))  
.adRevenue  
.sum()
```

Win

Win

# UDF

```
url_regex = re.compile("(?P<url>https?://[^\s]+)")
```

```
crawl
```

```
.pageSource
```

```
.apply(extract_url)
```

```
.value_counts()
```

Win

# JOIN

rankings

```
.merge(uservisits, left_on="pageURL", right_on="destinationURL", how="left")  
.query("visitDate > '1980-01-01' and visitDate < '1980-04-01'")  
.groupby("sourceIP")  
.agg({  
    "pageRank": "mean",  
    "adRevenue": "sum",  
})  
.sort_values("adRevenue", ascending=False)  
.head(10)
```

Flawless  
Victory

# ROUND 2



Pandas



1 Node

# SIMPLE QUERIES

rankings

```
.loc[rankings.pageRank > 100]
```

uservisits

```
.groupby(uservisits.sourceIP.str.slice(0, 7))  
.adRevenue  
.sum()
```

Win

Win

# UDF

```
url_regex = re.compile("(?P<url>https?://[^\s]+)")
```

```
crawl
```

```
.pageSource  
.apply(extract_url)  
.value_counts()
```

Too Long

# PANDARALLEL

```
pandarallel.initialize(progress_bar=True)
```

```
crawl
```

```
.pageSource
```

```
.parallel_apply(extract_url)
```

```
.value_counts()
```



# Win

# JOIN

rankings

```
.merge(uservisits, left_on="pageURL", right_on="destinationURL", how="left")  
.query("visitDate > '1980-01-01' and visitDate < '1981-01-01'")  
.groupby("sourceIP")  
.agg({  
    "pageRank": "mean",  
    "adRevenue": "sum",  
})  
.sort_values("adRevenue", ascending=False)  
.head(10)
```

Pandas  
loses

# РАЗБОР ПОЛЕТОВ

- идеально для маленьких объемов
- можно ускорить с помощью `pandarallel`
- нужно много ~~бамбука~~ RAM

# ROUND 3

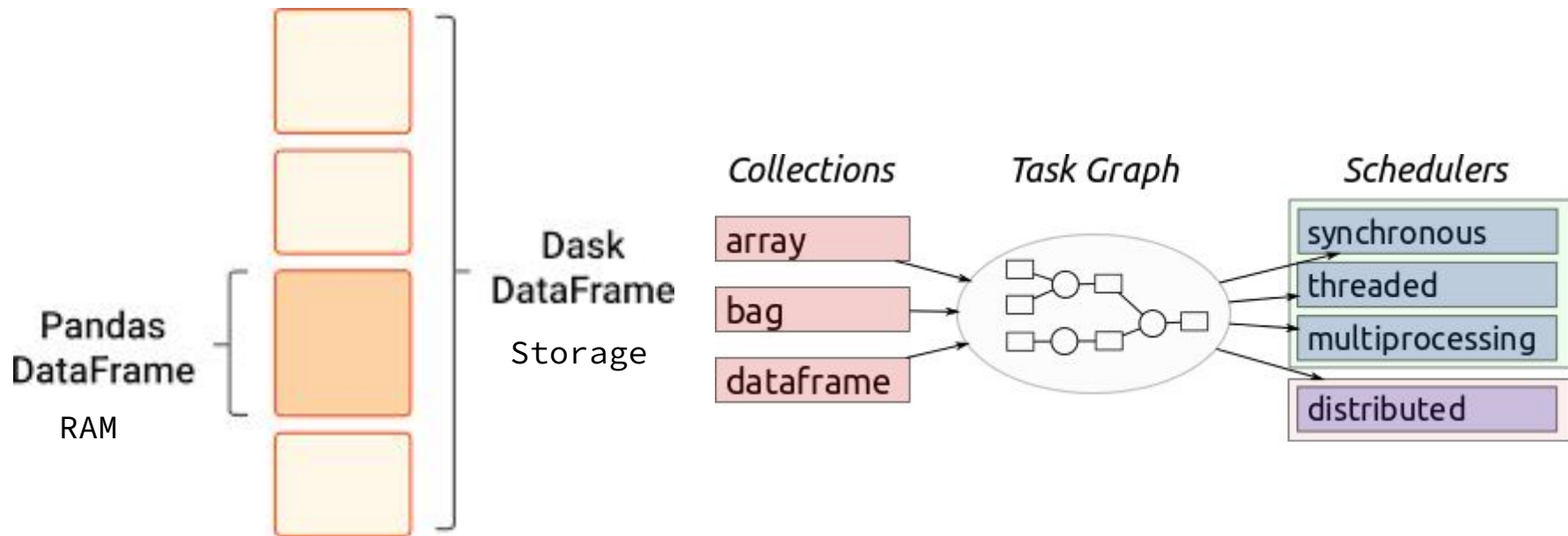


Dask



1 Node

# DASK



# SIMPLE QUERIES

rankings

```
.loc[rankings.pageRank > 100]
```

uservisits

```
.groupby(uservisits.sourceIP.str.slice(0, 7))  
.adRevenue  
.sum()
```

# UDF

```
url_regex = re.compile("(?P<url>https?://[^\s]+)")
```

```
crawl
```

```
    .pageSource
```

```
    .apply(extract_url)
```

```
    .value_counts()
```

# JOIN

rankings

```
.merge(uservisits, left_on="pageURL", right_on="destinationURL", how="left")  
.query("visitDate > '1980-01-01' and visitDate < '1990-01-01'")  
.groupby("sourceIP")  
.agg({  
    "pageRank": "mean",  
    "adRevenue": "sum",  
})  
.sort_values("adRevenue", ascending=False)  
.head(10)
```

Dask  
Wins

# ROUND 4



Dask



5Node

ХОРОШЕЕ ВРЕМЯ  
ДЛЯ ВОПРОСОВ

# ROUND 5



Spark



5Node

# SIMPLE QUERIES

rankings

```
.filter(rankings["pageRank"] > 100)
```

uservisits

```
.withColumn("sourceIPpart", substring("sourceIP", 1, 8))  
.groupBy("sourceIPpart")  
.sum("adRevenue")
```

# PANDAS API ON SPARK

```
import pyspark.pandas as ps
```

```
rankings
```

```
.loc[rankings.pageRank > 100]
```

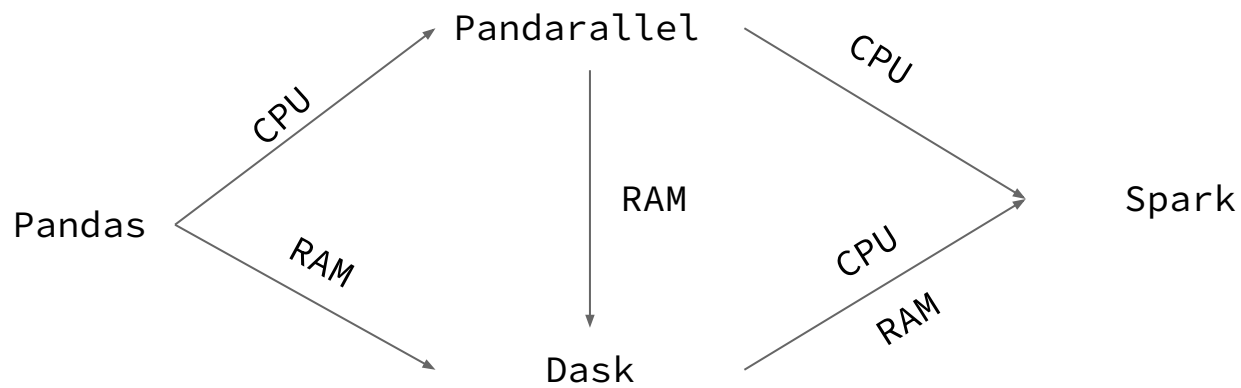
```
uservisits
```

```
.groupby(uservisits.sourceIP.str.slice(0, 7))
```

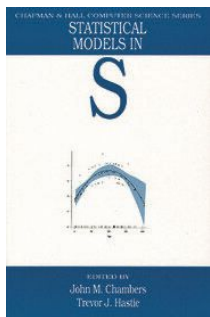
```
.adRevenue
```

```
.sum()
```

# КАК ЖЕ ВЫБРАТЬ?



# К ИСТОКАМ DATAFRAME API



1990



2000



2009



2015



2021  
(Pandas API)

# ИСТОЧНИКИ

- [AMPLabs Big Data Benchmark](#)
- [Pandarallel](#)
- [Dask](#)
- [Pandas API on Spark](#)
- [Preventing the Death of Dataframes](#)
- [Polars](#)
- [examples source codes](#)

A close-up photograph of a giant panda sitting against a tree trunk, holding a large piece of green bamboo in its mouth. The panda's black and white fur is clearly visible, and its teeth are engaged with the bamboo. The background is a soft-focus green forest.

**So long, and thanks  
for all the bamboo!**