

# Deconstructing Feather

---

Bill Lattner (@wlattner)  
PyData Chicago, 2016



**CIVIS**<sup>™</sup>  
ANALYTICS

Building a Data-Driven World<sup>™</sup>

# Why?

---

- Exchange tabular data between Python, R, and others
- Fast read/write
- Represent categorical features
- *It's about the metadata<sup>1</sup>*

1. <http://wesmckinney.com/blog/feather-its-the-metadata/>



CSV Files

# com·plex·i·ty

/kəm'pleksədē/

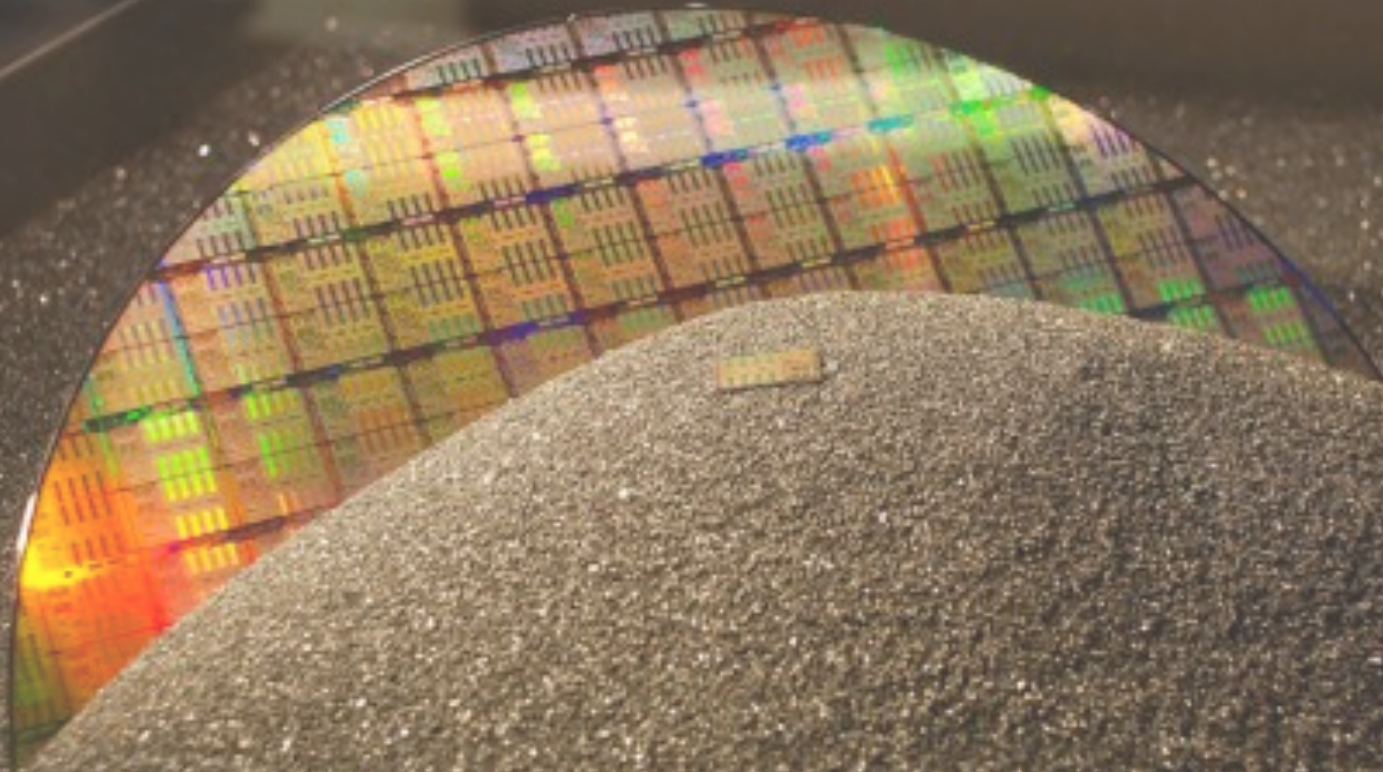
*noun*

```
1 pandas.read_csv(filepath_or_buffer, sep=', ', delimiter=None, header='infer',
2                 names=None, index_col=None, usecols=None, squeeze=False, prefix=None,
3                 mangle_dupe_cols=True, dtype=None, engine=None, converters=None,
4                 true_values=None, false_values=None, skipinitialspace=False,
5                 skiprows=None, skipfooter=None, nrows=None, na_values=None,
6                 keep_default_na=True, na_filter=True, verbose=False,
7                 skip_blank_lines=True, parse_dates=False, infer_datetime_format=False,
8                 keep_date_col=False, date_parser=None, dayfirst=False, iterator=False,
9                 chunksize=None, compression='infer', thousands=None, decimal='.',
10                lineterminator=None, quotechar='"', quoting=0, escapechar=None,
11                comment=None, encoding=None, dialect=None, tupleize_cols=False,
12                error_bad_lines=True, warn_bad_lines=True, skip_footer=0,
13                doublequote=True, delim_whitespace=False, as_recarray=False,
14                compact_ints=False, use_unsigned=False, low_memory=True,
15                buffer_lines=None, memory_map=False, float_precision=None)
```



Computers!

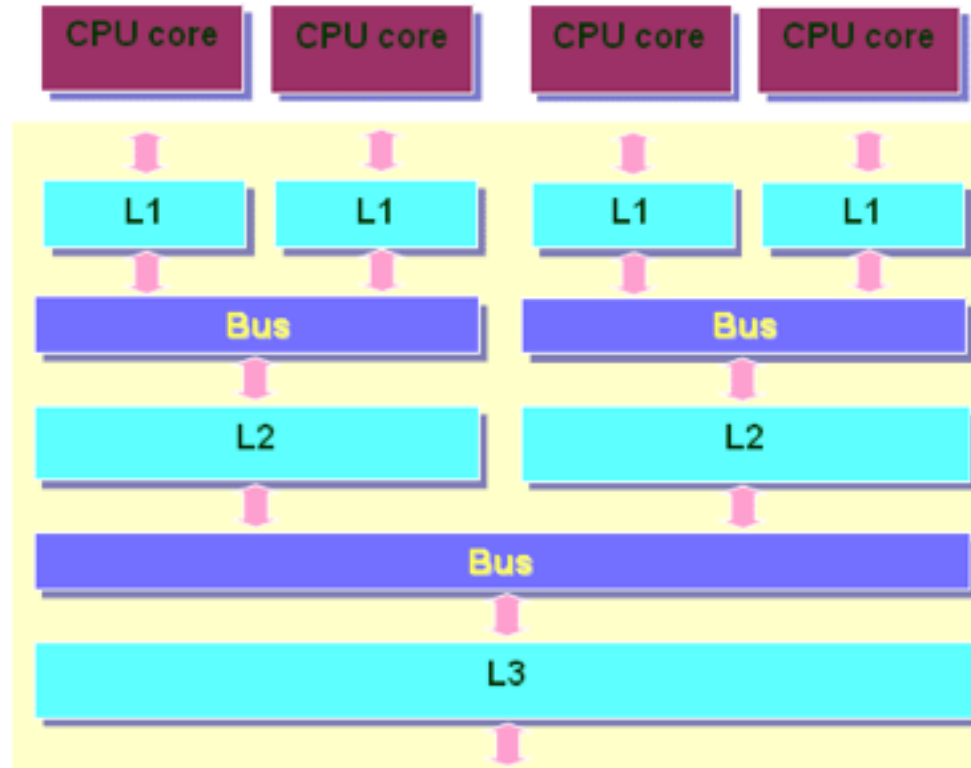
It all starts with sand...



$O(1)$  all the memory access



# How they actually work



[https://software.intel.com/sites/default/files/m/d/4/1/d/8/196578\\_196578.gif](https://software.intel.com/sites/default/files/m/d/4/1/d/8/196578_196578.gif)



## How they actually work

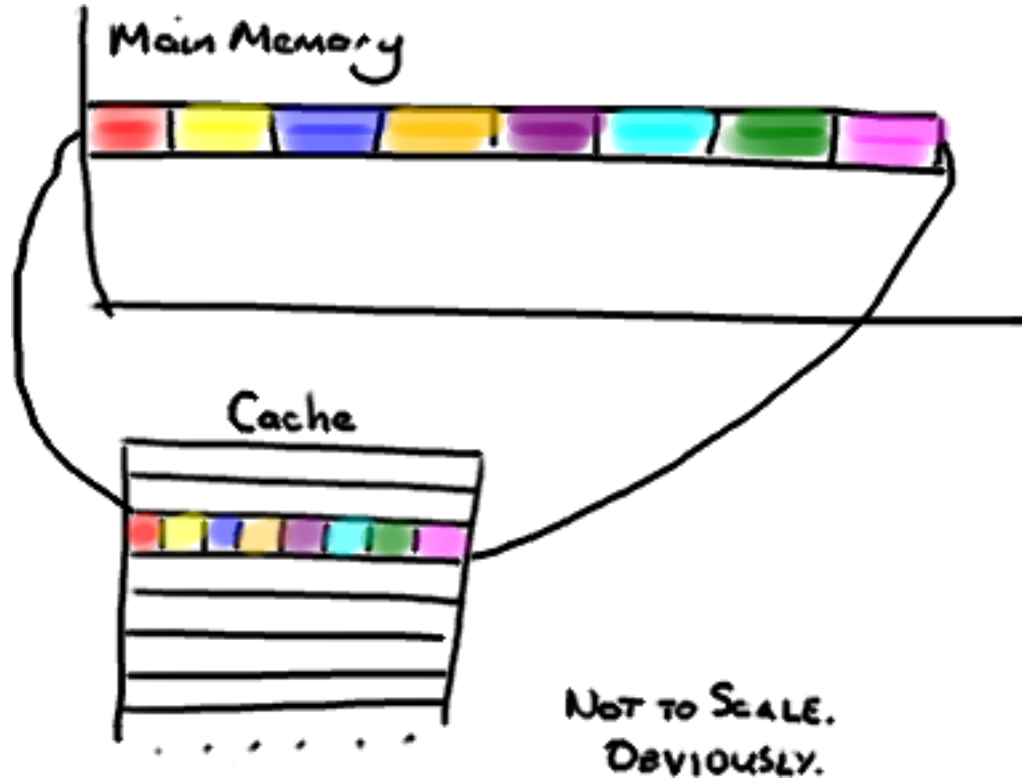
---

|                                     |                |
|-------------------------------------|----------------|
| L1 cache reference                  | 0.5 ns         |
| Branch mispredict                   | 5 ns           |
| L2 cache reference                  | 7 ns           |
| Mutex lock/unlock                   | 100 ns         |
| Main memory reference               | 100 ns         |
| Compress 1K bytes with Zippy        | 10,000 ns      |
| Send 2K bytes over 1 Gbps network   | 20,000 ns      |
| Read 1 MB sequentially from memory  | 250,000 ns     |
| Round trip within same datacenter   | 500,000 ns     |
| Disk seek                           | 10,000,000 ns  |
| Read 1 MB sequentially from network | 10,000,000 ns  |
| Read 1 MB sequentially from disk    | 30,000,000 ns  |
| Send packet CA->Netherlands->CA     | 150,000,000 ns |

<http://static.googleusercontent.com/media/research.google.com/en//people/jeff/stanford-295-talk.pdf>

# How they actually work

---



[http://mechanitis.blogspot.com/2011/07/dissecting-disruptor-why-its-so-fast\\_22.html](http://mechanitis.blogspot.com/2011/07/dissecting-disruptor-why-its-so-fast_22.html)

## What this means

---

Data layout needs to be tailored for expected read/write operations.

Memory access cost (latency) depends on location ***and*** predictability.

Sequential access FTW!!!



Feather (<https://github.com/wesm/feather>)

## The idea

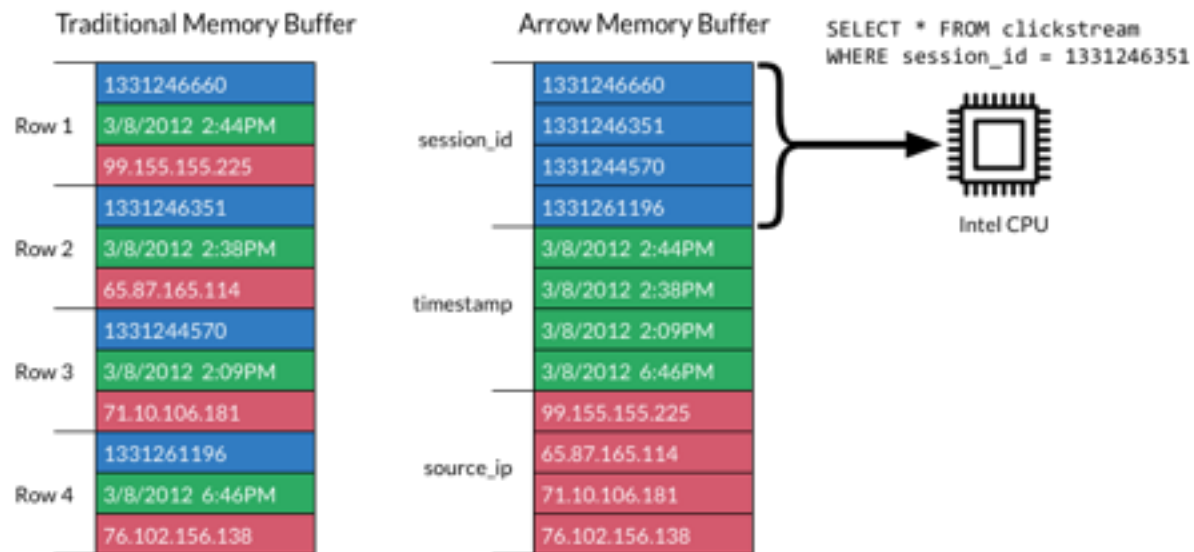
---

On disk representation of tabular data should be similar to the in memory representation.

Columnar layout is a good fit for analytic workflows.

# The idea

|       | session_id | timestamp       | source_ip      |
|-------|------------|-----------------|----------------|
| Row 1 | 1331246660 | 3/8/2012 2:44PM | 99.155.155.225 |
| Row 2 | 1331246351 | 3/8/2012 2:38PM | 65.87.165.114  |
| Row 3 | 1331244570 | 3/8/2012 2:09PM | 71.10.106.181  |
| Row 4 | 1331261196 | 3/8/2012 6:46PM | 76.102.156.138 |



<https://arrow.apache.org/>

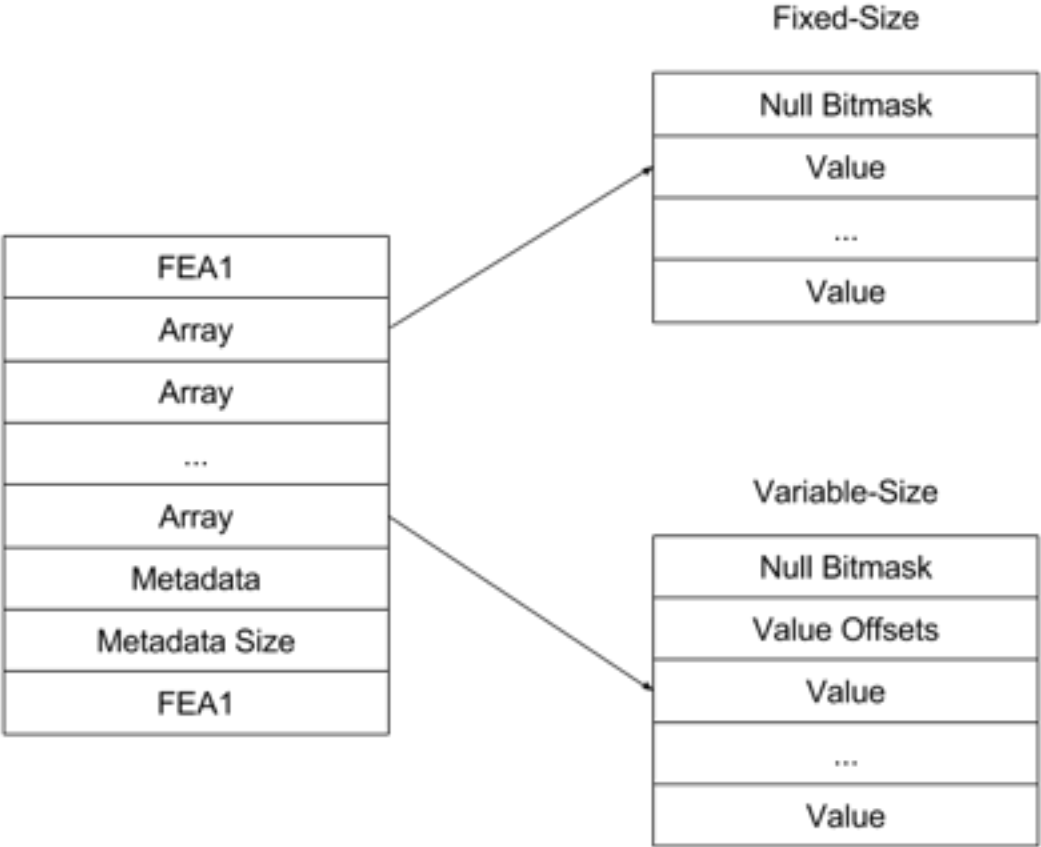
# sim·plic·i·ty

/sim'plisədē/

*noun*

```
1 feather.read_dataframe(path, columns=None)
```

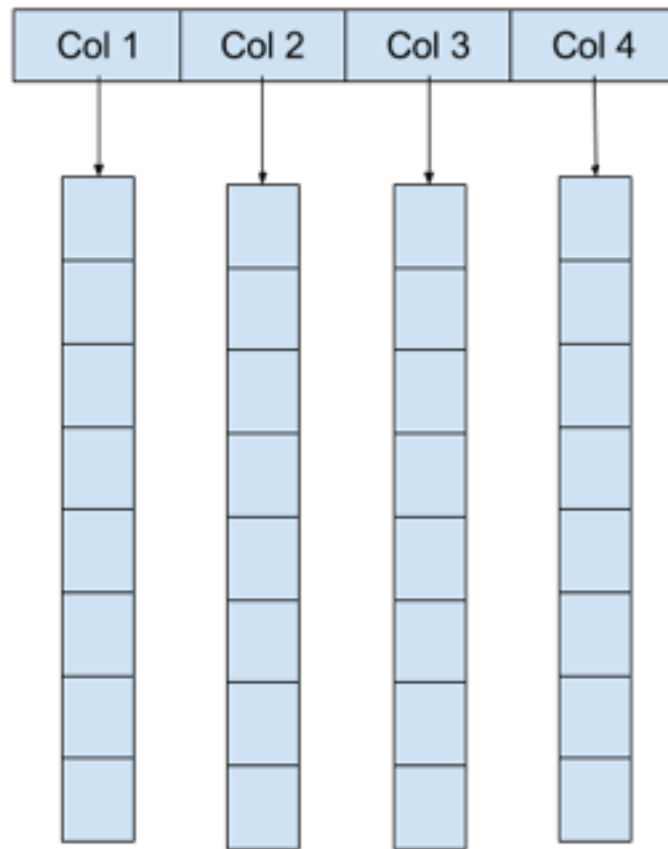
# The details





# Compare to a dataframe in R

---





**DANGER**

Live Code



The future

- In-place operations
- Share operational code between languages
- Zero parsing or copying to Pandas memory representation, mmap the feather file

## De facto interchange format

---

- input to tools like Scikit-Learn or StatsModels
- output from like PostgreSQL



Thanks