



Scaling up your pandas workflows with Modin

Devin Petersohn

Co-founder and CTO, Ponder



PONDER

About me

- Active Duty U.S. Marine
 - Korean Crypto-Linguist
 - 3d Radio BN (2008-2012)
- BS - University of Missouri (2016)
- MS - UC Berkeley (2018)
- PhD - UC Berkeley (2021)
 - NSF Graduate Research Fellow
 - Chancellor's Fellow
 - Modin started as my PhD project
- Currently: Cofounder and CTO of Ponder





A dataframe built from first
principles



Data Science Organization problems

What problems do data science teams face?

Data Science has a scalability problem



Data Science has a scalability problem



Data Science Teams* have a scalability problem



Data Science Teams* have a scalability problem



MORE DATA SCIENTISTS != MORE INSIGHTS

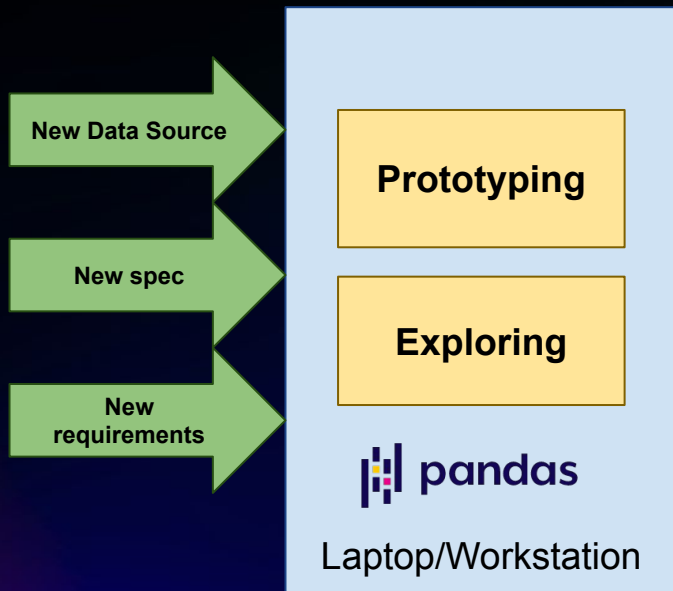
Data Science Teams* have a scalability problem



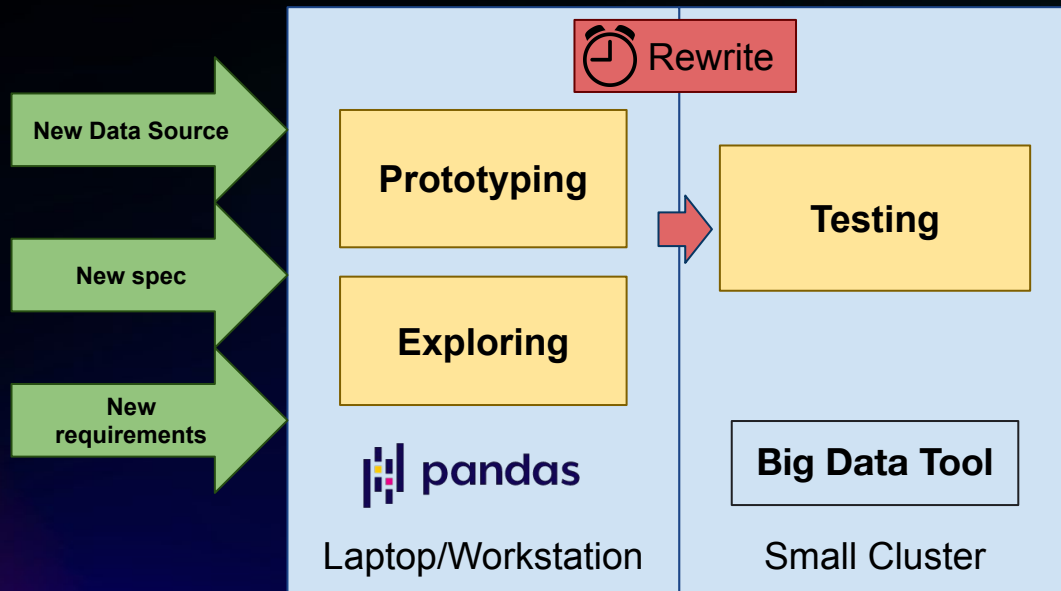
MORE DATA SCIENTISTS != MORE INSIGHTS

MORE DATA SCIENTISTS != MORE PRODUCTION MODELS/JOBS

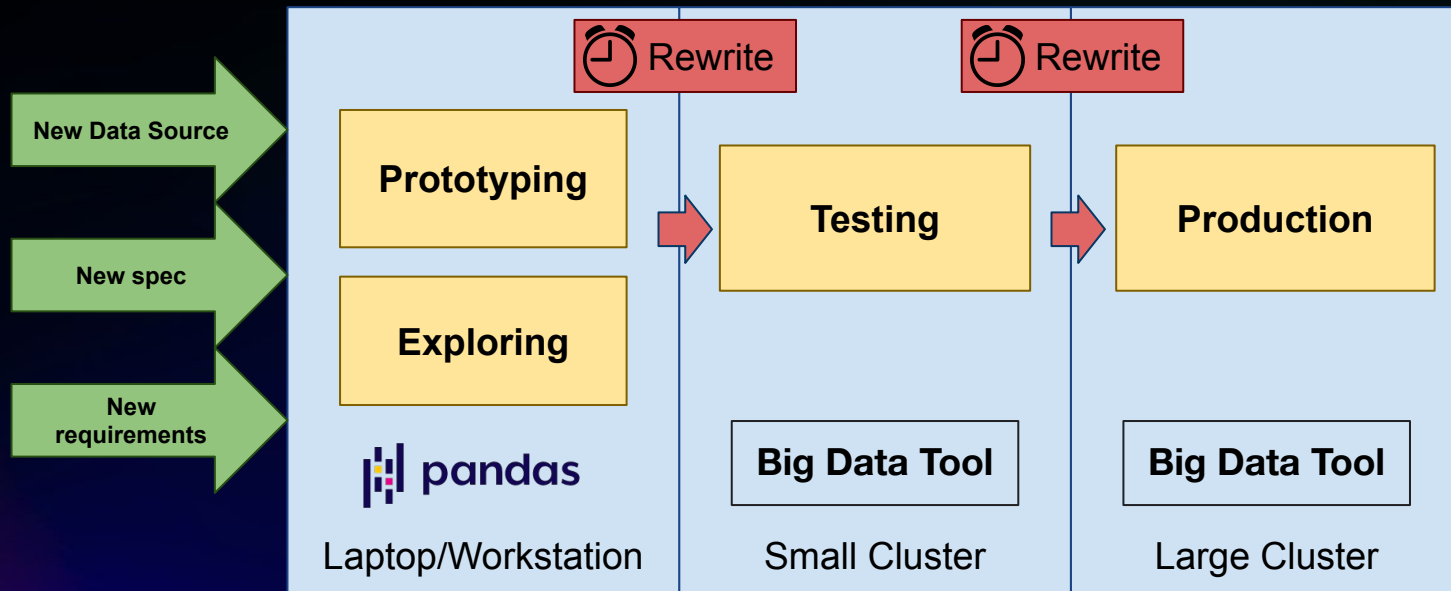
Many organizations look like this



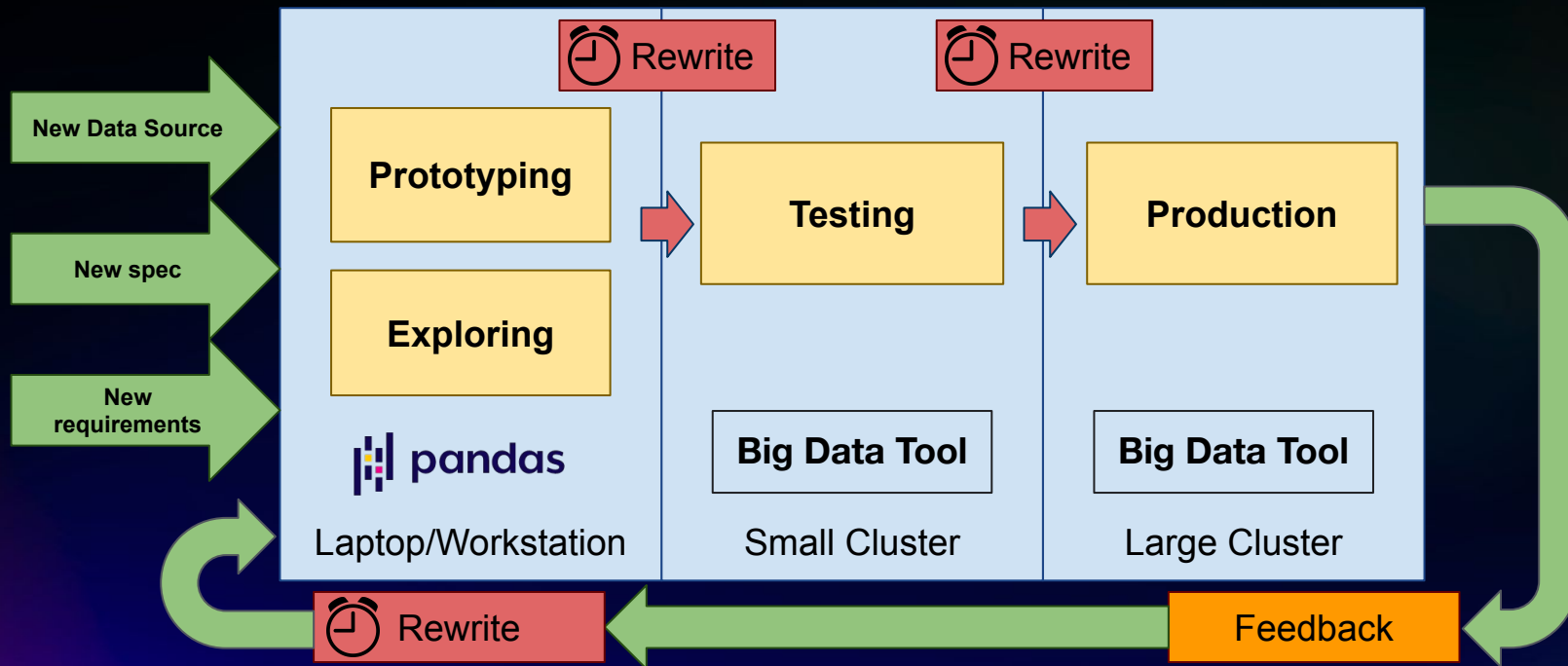
Many organizations look like this



Many organizations look like this



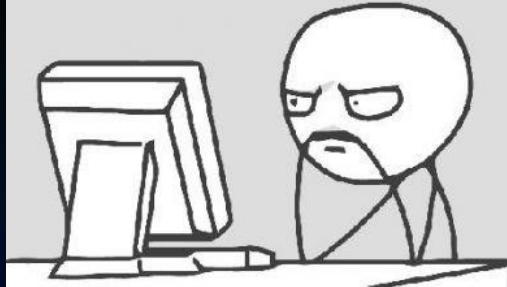
Many organizations look like this



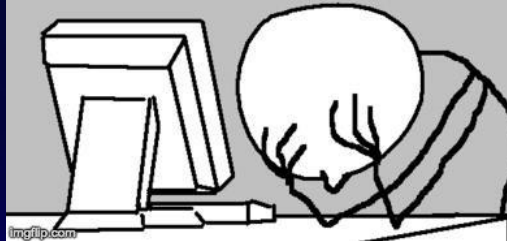


PONDER

**FINALLY FINISHED MY PANDAS
NOTEBOOK ON THIS SAMPLE DATA**



**NOW TO REWRITE
IT FOR THE CLOUD**



imgflip.com



Data Science scalability is human scalability

There is no service that can spin up more **Data Scientists**,
so we must treat them like the finite resource they are

Data Science scalability is a human scalability

But why?

There is no way you can spin up more **Data Scientists**,
so we must treat them like the finite resource they are

Ease of use

expressiveness, flexibility, agility

pandas

?

Performance

scalability, robustness, efficiency



Ease of use

expressiveness, flexibility, agility

?

big data
frameworks

Performance

scalability, robustness, efficiency

Shifting the focus from the machine to the user



Data Scientists shouldn't have to work for their tools

Tools should work for data scientists

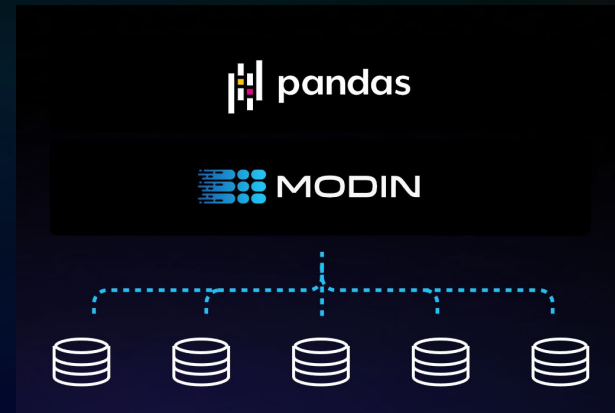
Ponder's work:

Transparently scale existing tools

Abstract away all of the components of the system that data scientists don't care about, only expose details they do care about.

Ponder's work:

Transparently scale existing tools



Abstract away all of the components of the system that data scientists don't care about, only expose details they do care about.



Ease of use

expressiveness, flexibility, agility

pandas



PONDER

(powered by Modin)

**big data
frameworks**

Performance

scalability, robustness, efficiency



Let's start from first principles!

Performance

scalability, robustness, efficiency

But the Pandas API doesn't scale!



The API is simply an expression of what to do

Solving dataframes from first principles

- Dataframe Data Model
- Dataframe Algebra
- Parallelism / Decomposition Rules
- Type System
- Operator Semantics
- Implementation

First Steps: Formalize the Dataframe



Towards Scalable Dataframe Systems

Devin Petersohn, William Ma, Doris Lee, Stephen Macke, Doris Xin, Xiangxi Mo
Joseph E. Gonzalez, Anthony D. Joseph, Joseph M. Hellerstein, Aditya Parameswaran
UC Berkeley

{devin.petersohn, williamma, dorislee, smacke, dorx, xmo, jgonzal, adj, hellerstein, adityagp}@berkeley.edu

ABSTRACT

Dataframes are a popular and convenient abstraction to represent, structure, clean, and analyze data during exploratory data analysis. Despite the success of dataframe libraries in R and Python (pandas), dataframes face performance issues even on moderately large datasets. In this vision paper, we take the first steps towards formally defining dataframes, characterizing their properties, and outlining a research agenda towards making dataframes more interactive at scale. We draw on tools and techniques from the database community, and describe ways they may be adapted to serve dataframe systems, as well as the new challenges therein. We also describe our current progress toward a scalable dataframe system, MODIN, which is already up to 30× faster than pandas in preliminary case studies, while enabling unmodified pandas code to run as-is. In its first 18 months, MODIN is already used by over 60 downstream projects, has over 250 forks, and 3,900 stars on GitHub, indicating the pressing need for pursuing this agenda.

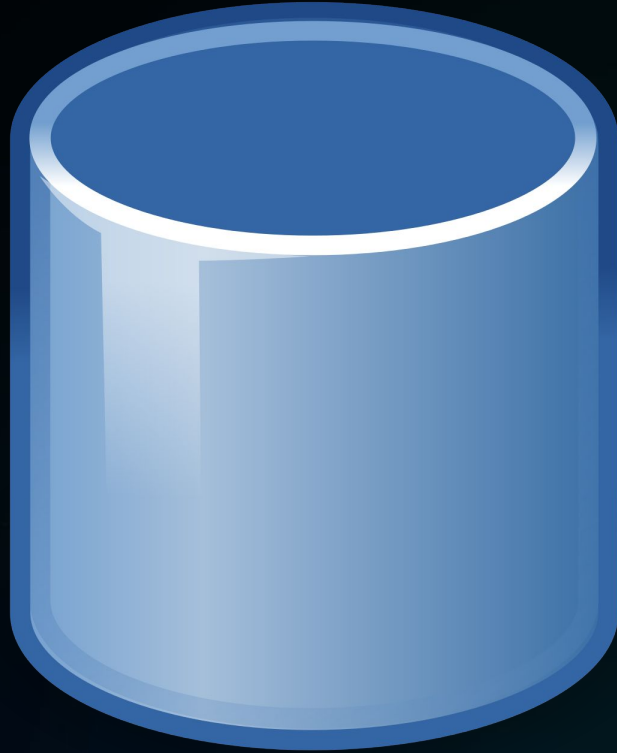
1. INTRODUCTION

Characteristics such as these have helped dataframes become incredibly popular for EDA; for instance, the dataframe abstraction provided by pandas within Python (pandas.pydata.org), has, as of 2019, been downloaded over 200 million times, served as a dependency for over 160,000 repositories in GitHub, and starred on GitHub more 22,000 times. Python's own popularity has been attributed to the success of pandas for data exploration and data science [7,9]. Due to its ubiquity, we focus on pandas for concreteness.

Pandas has been developed from the ground up via open-source contributions; pandas DataFrame API operators and their implementations have been provided by dozens of contributors to satisfy immediate or ad-hoc needs, spanning capabilities that mimic relational algebra, linear algebra, and spreadsheet formula computation. To date, the pandas DataFrame API has ballooned to over 200 distinct operators [14]. R, which is both more mature and more carefully curated than pandas, has only 70 operators—but this is still far more than, say, relational and linear algebra combined [15].

While this rich API is sometimes cited as a reason for pandas' attractiveness, the set of operators has significant redundancies, often with different performance implications. These redundancies place a considerable burden on users to select the optimal one of

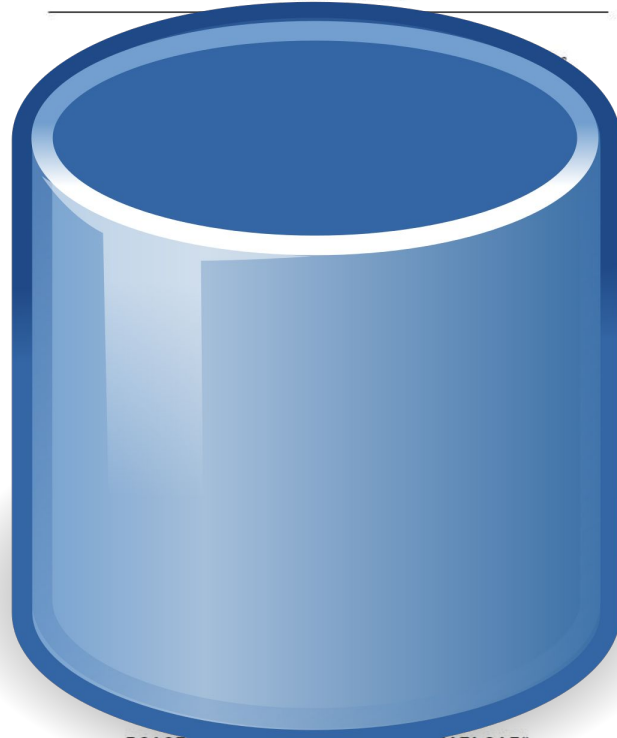
888v2 [cs.DB] 6 Jan 2020



THE FRENCH LAUNDRY

TASTING OF VEGETABLES

28 MAY 2021



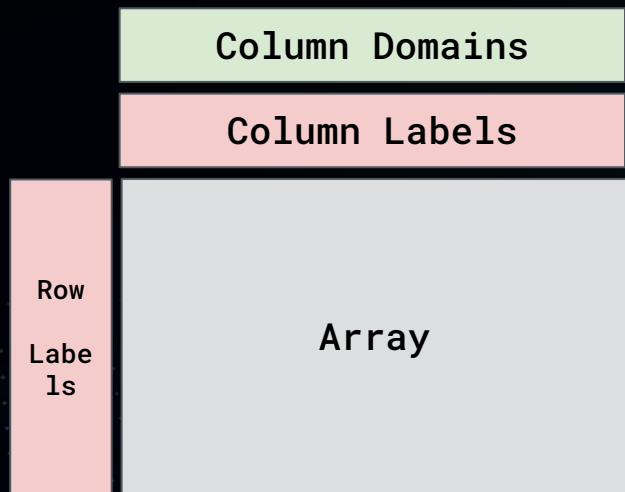
ROASTED VEGETABLES "MEATLOAF"
*La Ratte Potato Purée, Crispy Cipolini Onions, Smoked Tomato Glaze,
Watercress Leaves and Whole Grain Mustard Gravy*

"GOUGÈRE"

Andante Dairy "Etude" and Preserved Australian Black Winter Truffle "Fondue"



Dataframes: A New Data Model and API



`df.groupby(..)`

`df.drop(..)`

`pd.merge(..)`

`df.describe(..)`

`pd.concat(..)`

`df.pivot(..)`

`df.explode(..)`

Dataframes: Mixed type-array, w/ row and column labels

600+ functions to clean, reshape, explore, and summarize data spanning rel., linear, & spreadsheet algebra





Convenience

Entire query at once

Incremental + inspection

Flexible

Strict schema

Mixed types, R/C and
data/metadata equiv.

Versatility

SFW or bust

600+ functions

Everybody loves pandas!



data school

Launch a data science career!

December 12, 2018 - PYTHON

What's the future of the pandas library?

pandas is a powerful, open source Python library for **data analysis, manipulation, and visualization**. I've been teaching data scientists to use pandas since 2014, and in the years since, it has grown in popularity to an estimated **5 to 10 million users** and become a "must-use" tool in the Python data science toolkit.

SIGN IN The Register

{* DEVOPS *}

Python explosion blamed on pandas

Data science fad just won't die

Thomas Claburn in San Francisco

Thu 14 Sep 2017 // 20:02 UTC

34

Not content to bait developers by declaring that Python is the **fastest-growing major programming language**, coding community site Stack Overflow has revealed the reason for its metastasis.



DISCOVER LATEST OBSESSIONS

QUARTZ

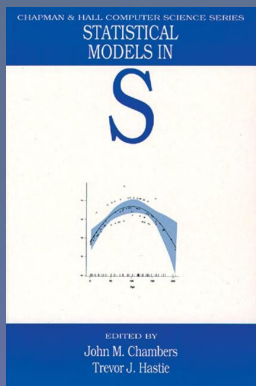
FEATURED EMAILS BECOME A MEMBER

PANDAS

Meet the man behind the most important tool in data science

Inherited Data Model

Small group development



1990

~70 operators



1995
2000 (stable)

Large community development

>200 operators



2008



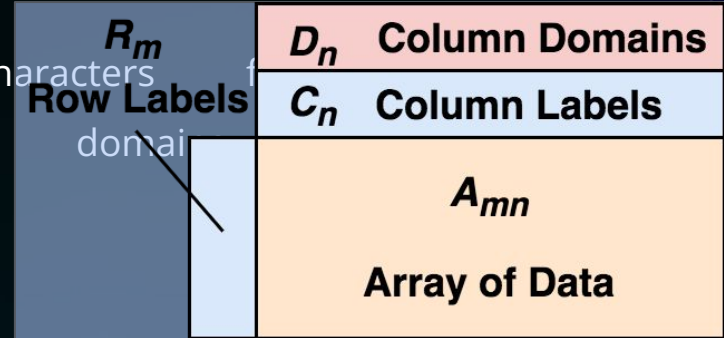
PONDER

Dataframe formal definition (VLDB 2020)

Let Σ^* be the finite set of characters

Let Dom be a finite set of domain

Let each $dom_i \in Dom$ have a mapping $p_i: \Sigma^* \rightarrow dom_i$.

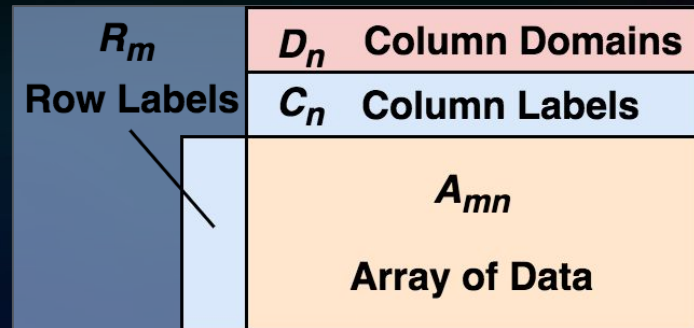


A **dataframe** is a tuple (A_{mn}, R_m, C_n, D_n) , where A_{mn} is an arrangement of entries in columns and rows from the domain Σ^* , R_m is a vector of row labels from Σ^* , C_n is a vector of column labels from Σ^* , and D_n is a vector of n domains from some finite set of domains Dom , one per column, each of which can also be left unspecified. We call D_n the schema of the dataframe. If any of the n entries within D_n is left unspecified, then that domain can be induced by applying a schema induction function $S(\cdot)$ to the corresponding column of A_{mn} . The schema induction function $S: \Sigma^* \rightarrow Dom$, assigns an arrangement of m strings to a domain in Dom .



Dataframe data model

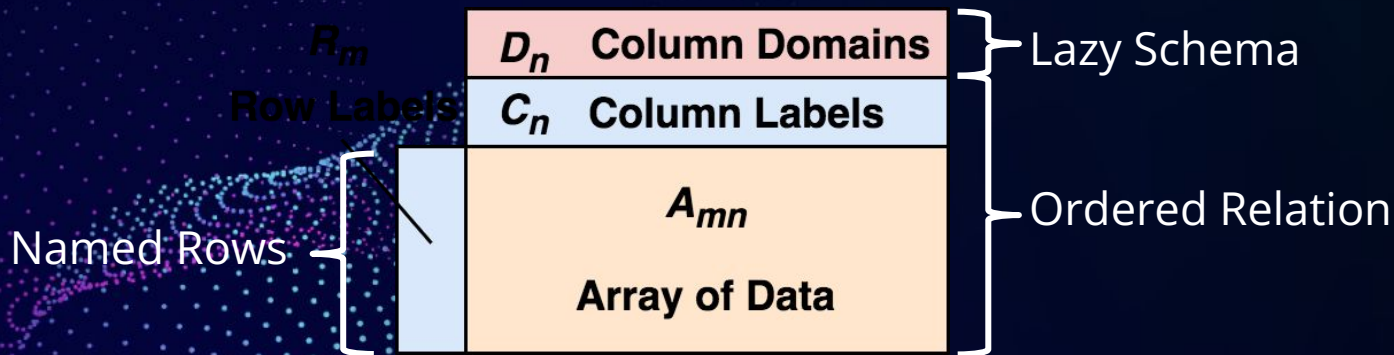
- Ordered, but not necessarily sorted
 - Rows **and** columns
- No predefined schema necessary
 - Types can be induced at runtime
- Typed Row/column labels
 - Labels can become data
- Indexing by label or by row/column numeric index
 - “Named notation” or “Positional notation”



Dataframes from two perspectives

From a **relational algebra** perspective, dataframes contain:

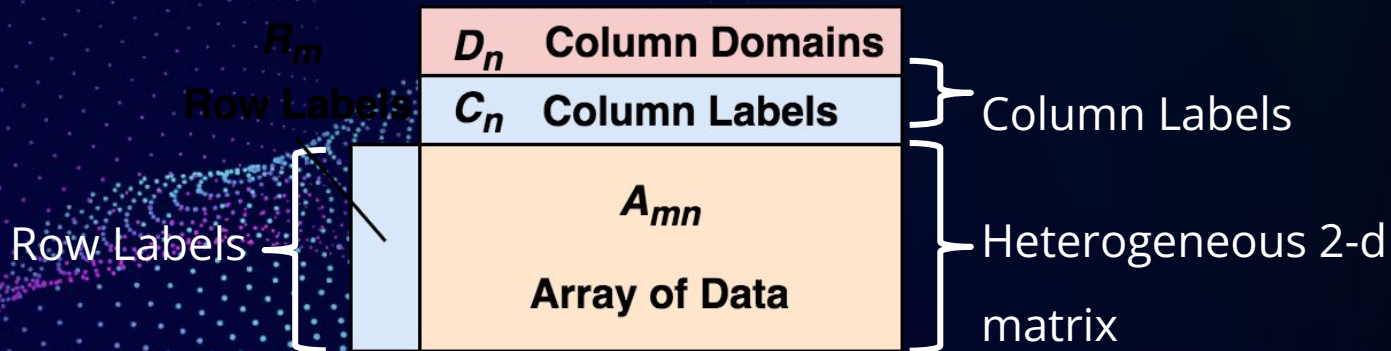
- An ordered table
- Named rows of arbitrary type
- A lazily-induced schema
- Column names of arbitrary type
- Column and row symmetry
- Support for linear algebra operators (e.g. matrix multiply)



Dataframes from two perspectives

From a **linear algebra** perspective:

- Heterogeneous matrix-like data structure
- Both numeric and non-numeric types
- Explicit row and column labels
- Indexing by label in addition to position
- Support for relational algebra operators (e.g. join)



- Dataframe Data Model ✓
- Dataframe Algebra
- Parallelism / Decomposition Rules
- Type System
- Operator Semantics
- Implementation

Next: What can a dataframe do?

First Principles next steps:
define an algebra

What can Pandas Do?



pd.DataFrame	280+ methods
pd.Series	280+ methods
Convenience methods (e.g. concat)	40+ APIs

Total: 600+ operators in pandas



PONDER

Example
Pandas
APIs

T
abs
add
all
corr
dropna
get
groupby
head
join
median
pivot
reindex
query
rolling
tail
where

Operator
SELECTION
PROJECTION
UNION
DIFFERENCE
CROSS-PRODUCT / JOIN
DROP DUPLICATES
GROUPBY
SORT
RENAME
WINDOW
TRANSPOSE
MAP
TOLABELS
FROMLABELS

Example
R
APIs

t()
melt()
merge()
head()
slice()
distinct()
arrange()
rename()
summary()
group_by()
aggregate()
union()
with()
subset()
match()
sample_n()
row_labels()

Proof by exhaustion that all pandas APIs are covered

Operator Pattern	Modin Syntax	Exhaustive list of pandas public API	Count
Applying a user-defined function uniformly element-wise, column-wise, or row-wise	<code>map(df, axis, f)</code>	<code>abs, astype, clip, eval, isna, isnull, notna, notnull, fillna, isin, apply, applymap, transform, fillna, rank, round, unary version of [add, subtract, etc.], string manipulations (str), datetime manipulation (dt), replace</code>	42
Binary function between two dataframes	<code>map(join(df1, df2), f)</code>	<code>add, subtract, multiply, divide, eq, floordiv, where, update, corr, corrwith, combine, combine_first, cov, dot, ge, gt, lt, pow, radd, rdiv, rfloordiv, rmod, rmul, rpow, rsub, rtruediv, sub, truediv</code>	28
Indexing: Queries on the row labels	<code>mask(df, row_indices=query)</code>	<code>as_freq, asof, at, at_time, between_time, drop_level, drop, first, first_valid_index, iloc, last, last_valid_index, loc, mask, pop, get, head, tail, iat, lookup, resample, sample, select, take, truncate</code>	26
Treating metadata as data, metadata manipulation, and querying	<code>to_labels(map(from_labels(df), f), label)</code>	<code>add_prefix, add_suffix, swap_level, melt, reorder_levels, shift, slice_shift, tshift, tz_convert, tz_localize</code>	10
Reshaping, transposing, pivot	<code>[map/groupby](transpose(df), f)</code>	<code>pivot, pivot_table, stack, unstack, transpose, T</code>	6
One-hot (dummy) encoding	<code>transpose(to_labels(transpose(map(df, f, axis=1))))</code>	<code>get_dummies</code>	1
User-defined aggregation of values per-column	<code>reduction(df, f)</code>	<code>all, any, count, agg, aggregate, compound, describe, duplicated, equals, idxmax, idxmin, kurtosis, mad, max, mean, median, min, mode, memory_usage, prod/product, nunique, quantile, sem, skew, std, sum, var</code>	28
Aligning and joining two dataframes on row or column labels	<code>join(from_labels(df1), from_labels(df2), on="index", axis)</code>	<code>align, concat, join, reindex, reindex_axis, reindex_like</code>	6
Window user-defined functions (window size < length of dataframe)	<code>window(df, axis, f, size)</code>	<code>bfill, cumsum, cumprod, cummax, cummin, diff, ewm, expanding, ffill, fillna, interpolate, nlargest, nsmallest, pct_change, rolling</code>	15
Conditional filter	<code>filter(df, f)</code>	<code>dropna, drop_duplicates, filter, query</code>	4
Queries on the schema of the dataframe	<code>filter_by_types(df, types)</code>	<code>select_dtypes</code>	1
Type Inference and induction	<code>infer_types(df, columns)</code>	<code>infer_objects, convert_objects</code>	2
Column/row insertion and assignment, appending columns/rows	<code>concat(df1, [df2])</code>	<code>append, assign, concat, insert</code>	4
Groupby with a user-defined aggregation or function	<code>groupby(df, by, f)</code>	<code>groupby</code>	
Join on an attribute	<code>join(df1, df2, condition, how)</code>	<code>merge, merge_asof</code>	2
Sorting on labels or column values	<code>sort_by(df)</code>	<code>sort_index, sort_values</code>	2
Expand the number of rows or columns	<code>explode(df, f)</code>	<code>explode</code>	1

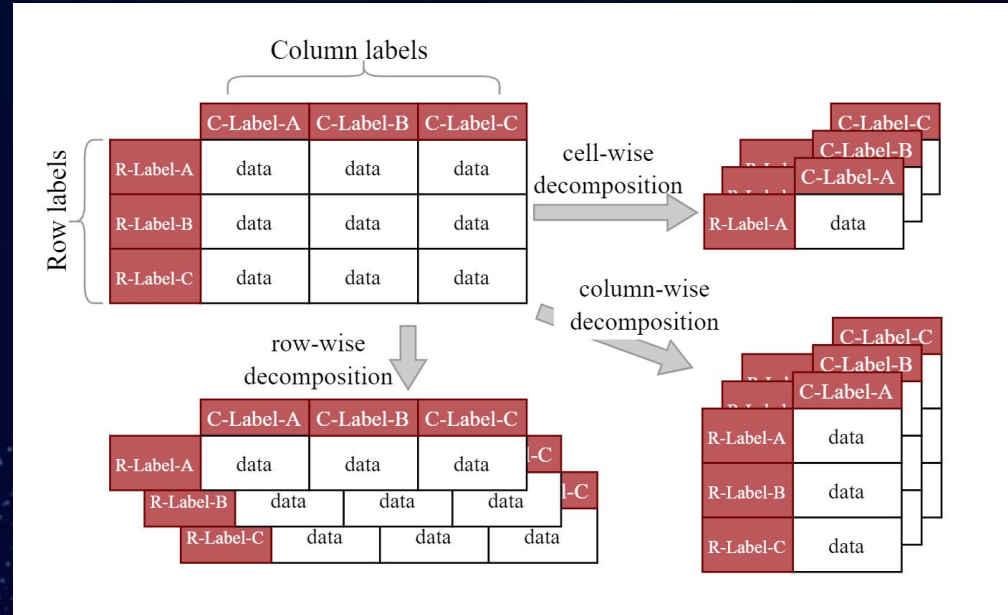
- Dataframe Data Model ✓
- Dataframe Algebra ✓
- Parallelism / Decomposition Rules
- Type System
- Operator Semantics
- Implementation

Decomposition Rules -> Formalize parallelism

Cell wise: An operator can be applied to a “unit dataframe” independently

Row-wise: An operator can be applied to each row independently

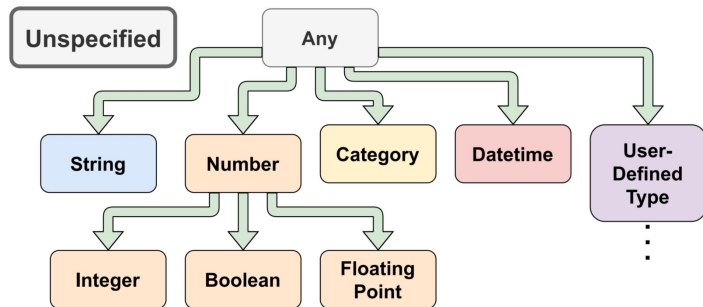
Column-wise: An operator can be applied to each column independently



- Dataframe Data Model ✓
- Dataframe Algebra ✓
- Parallelism / Decomposition Rules ✓
- Type System
- Operator Semantics
- Implementation

Formalization of other components of the dataframe

Type System



Order Semantics

Operator	Input Order	Position	Output Order & Position
mask	N	Y*	Parameter-Dependent
filter_by_types	N	N	Inherited Updated
map	N	Y◊	Inherited from Inputs
filter	N	Y◊	Inherited Updated
explode	N	Y◊	Inherited Updated
reduce	N	Y◊	Inherited
window	Y	N	Inherited
groupby	N	N	Data-dependent
infer_types	N	N	Inherited
join	N	N	Inherited*
concat	N	N	Inherited*
transpose	N	N	Inherited
to_labels	N	N	Inherited
from_labels	N	Y	Inherited
sort	N	N	Data-dependent
rename	N	N	Inherited

- Dataframe Data Model ✓
- Dataframe Algebra ✓
- Parallelism / Decomposition Rules ✓
- Type System ✓
- Operator Semantics ✓
- Implementation

So now we know:

- What a dataframe is (formally)
- What operators a dataframe supports
- How these operators map back to pandas
- How to handle dataframe types
- How to decouple logical and physical order
- How to maximally parallelize each operator



A dataframe built from first
principles

Modin

Accelerate your pandas workloads by changing one line of code



To use Modin:

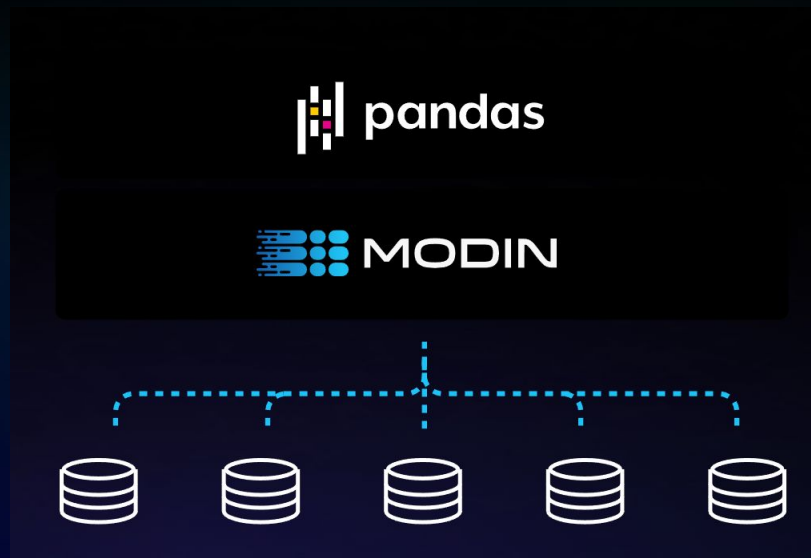
```
import pandas as pd
```

To install Modin:

```
pip install modin
```


What can we do with the formalism?

- Data model
 - Expose the “feel” of pandas without the baggage
- Dataframe Algebra
 - Smaller surface to implement
 - Map operators to other systems



Pandas on everything with Ponder



PostgreSQL



Your DB Here!

Databases



RAY



Compute Engine

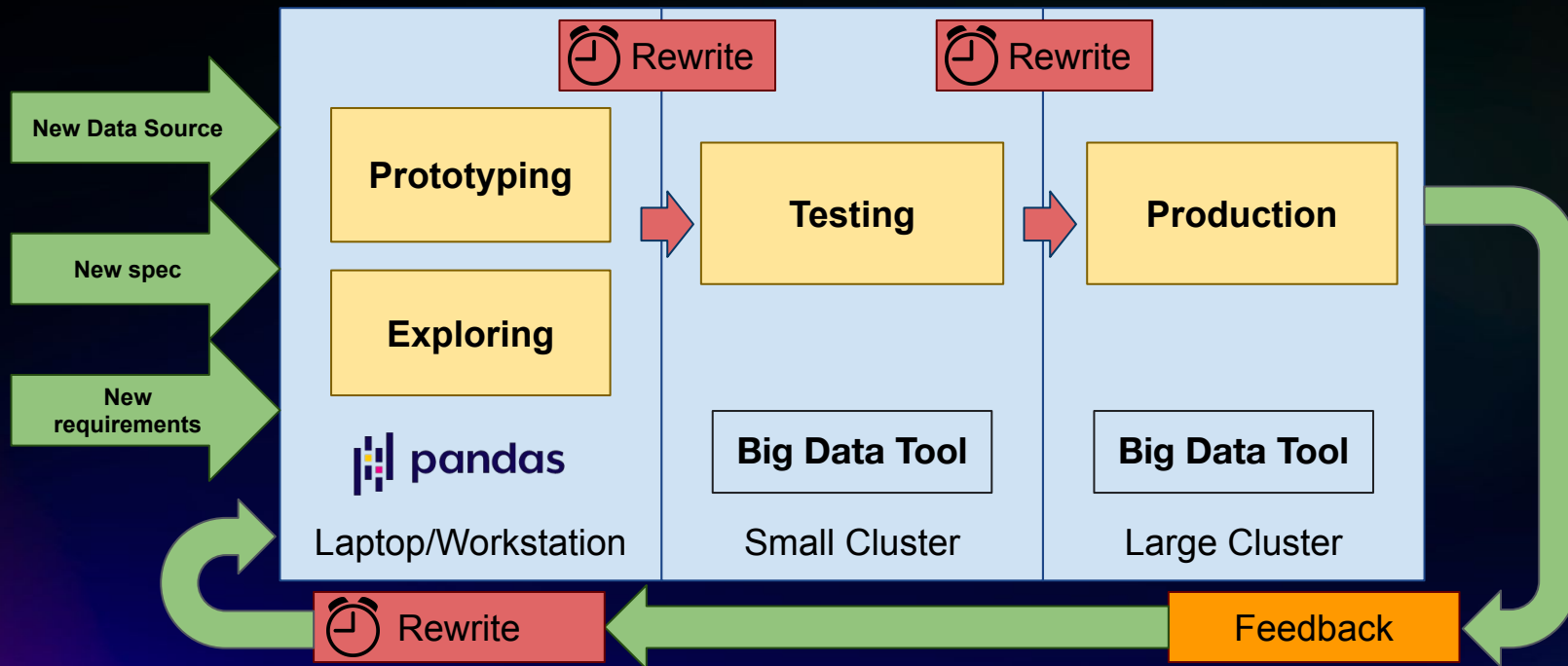


Azure

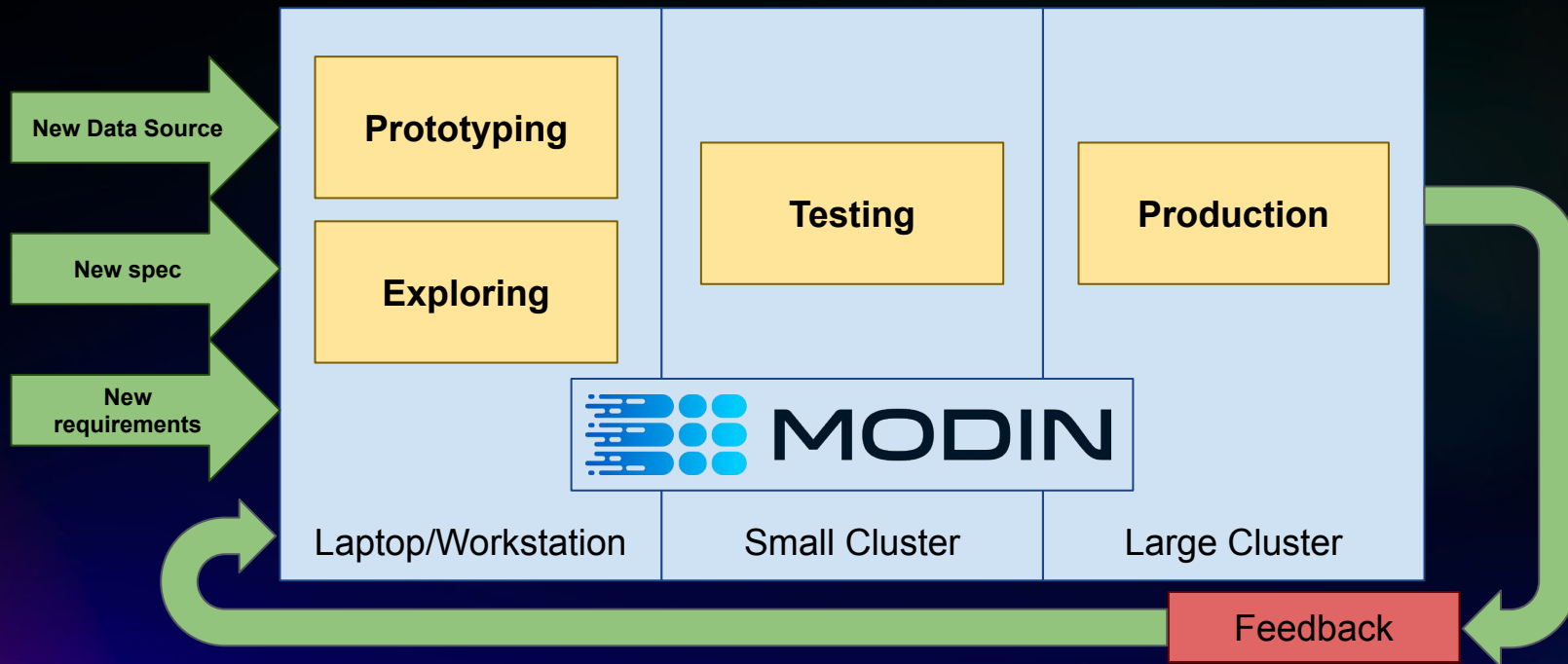


Infrastructure

Many organizations look like this



How Modin is being used in practice:





Modin Open Source Impact

Used by 10% Fortune 100 companies & more!



2.6M+

Downloads to date



7k

GitHub stars



90+

Open Source
Contributors

E-Commerce Case Study



Current State

Pricing + Recommendation
Pipelines in Production

150M Events
Streamed Daily

50k Records
Pandas Limit



With Ponder



Modin

1000X
More Data with Modin

< 20 Sec.
Faster with Modin

Finance Case Study



Current State

Regulatory Reporting
Pipelines Migration

10k+ Lines
Of Code to change est.

2k+ Hours
Human time to rewrite
est.



With Ponder



Modin

100X
Reduction in code
reduction w/ Modin

100X
Reduction in human
time with Modin

Bring your database!



PostgreSQL



amazon
REDSHIFT



Google BigQuery



Your DB Here!

Databases



RAY



Compute Engine



Infrastructure

In summary...

- Formal dataframe data model & algebra
 - Dataframes are newly defined structures with a lot of open problems!
 - Dataframe algebra can express all of pandas
- Reference implementation: Modin
 - High impact -> problems are real and pressing need
- Deep technical problems still exist!
 - It is an exciting time to be working on dataframes

Thank you!

Devin Petersohn
devin@ponder.io

Backup

Current state of affairs

Current state of affairs

- Coverage of pandas API
 - pandas.DataFrame - 83% (93% based on usage)
 - pandas.Series - 77% (86% based on usage)
 - pandas.read_* - 42% (>90% based on usage)

-

API Vision - 5 years

- More complete pandas API - 95%
- All common interactive data processing modalities
 - Spreadsheet API
 - SQL API
 - New Modin API
- Hooks for SQL systems to implement parts of the API
 - A partial pandas API for relational systems
- Preliminary numpy API - more implementation proof that dataframes can act as matrices
 - Plug in to sklearn

Engineering Vision - 5 years

- Query planning/optimization
 - Optimize for user's time
 - Some research involved here (more later)
 - Extending Calcite?
- GPU integration
- MPI, other compute engines
- Serverless



Case study: Pivot

PONDER

transpose +
groupby + map

Narrow Table (SALES)

Year	Month	Sales
2001	Jan	100
2001	Feb	110
2001	Mar	120
2002	Jan	150
2002	Feb	200
2002	Mar	250
2003	Jan	300
2003	Feb	310

Wide Table of MONTHS

Month	2001	2002	2003
Jan	100	150	300
Feb	110	200	310
Mar	120	250	NULL

Pivot →

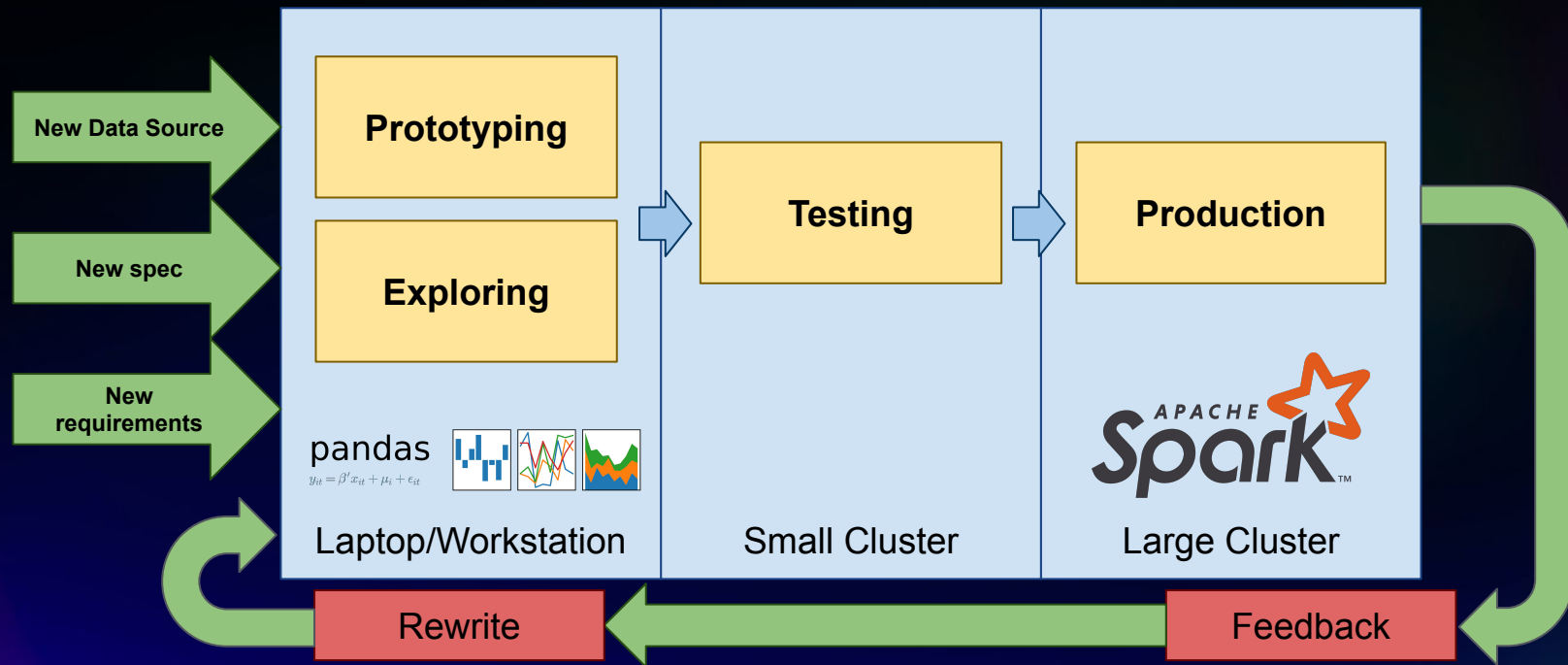
← Unpivot

Year	Jan	Feb	Mar
2001	100	110	120
2002	150	200	250
2003	300	310	NULL

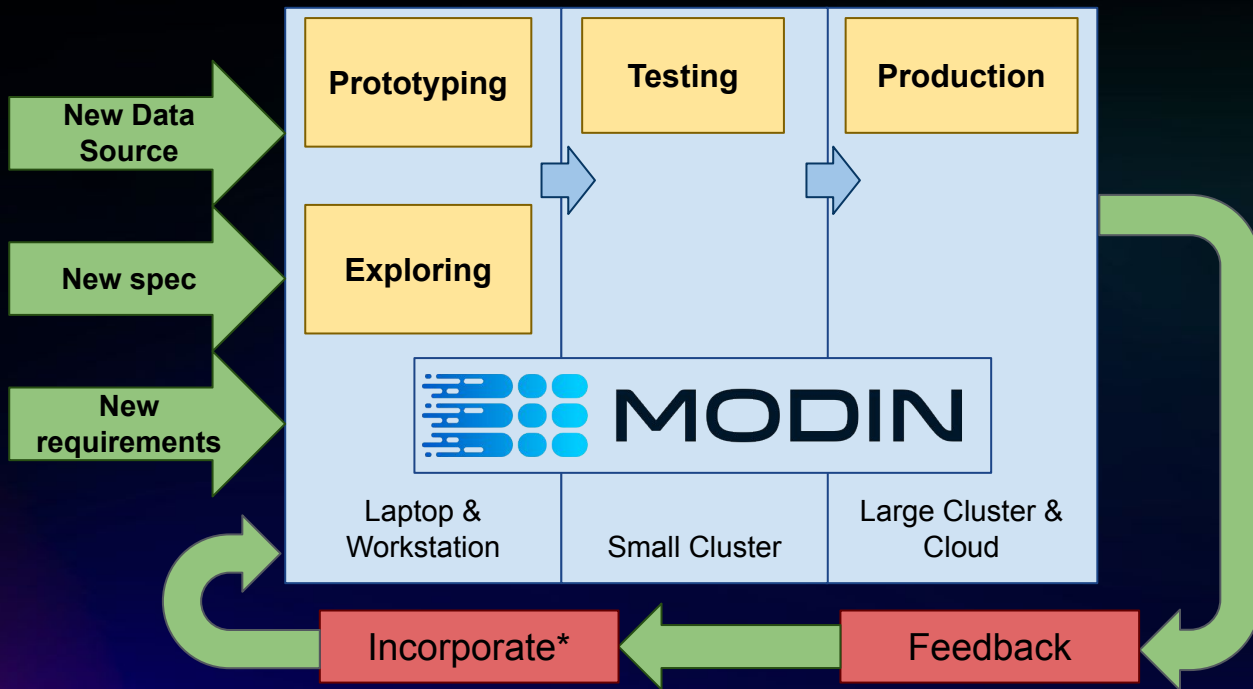
Wide Table of YEARS

Transpose + map

Current Machine Learning Lifecycle



One API, all scales (think SQL)



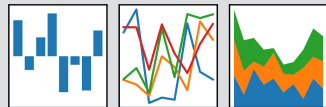
Data Science Landscape: Today

Tools efficient for O(1MB)

Usable but not
scalable

pandas

$$y_{it} = \beta' x_{it} + \mu_i + \epsilon_{it}$$



Tools efficient for O(100s GB+)

Scalable but not
usable



Data Science Landscape: Today



Tools efficient for O(1MB)

- Follow typical programming styles
- Tools are widely used and understood - in production
- The majority of college graduates will already know these tools
- No scalability



Tools efficient for O(100s GB+)

- Difficult to debug
- Requires specialized knowledge
 - Must understand partitioning
 - Lazy evaluation
- Designing systems for large data
 - New APIs that do the same thing



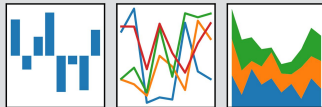
Data Science Landscape: Today

Tools efficient for O(1MB)

Usable and
scalable

pandas

$$y_{it} = \beta' x_{it} + \mu_i + \epsilon_{it}$$

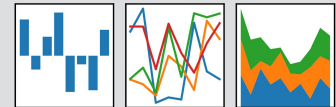


Tools efficient for O(100s GB+)

Scalable and
usable

pandas

$$y_{it} = \beta' x_{it} + \mu_i + \epsilon_{it}$$



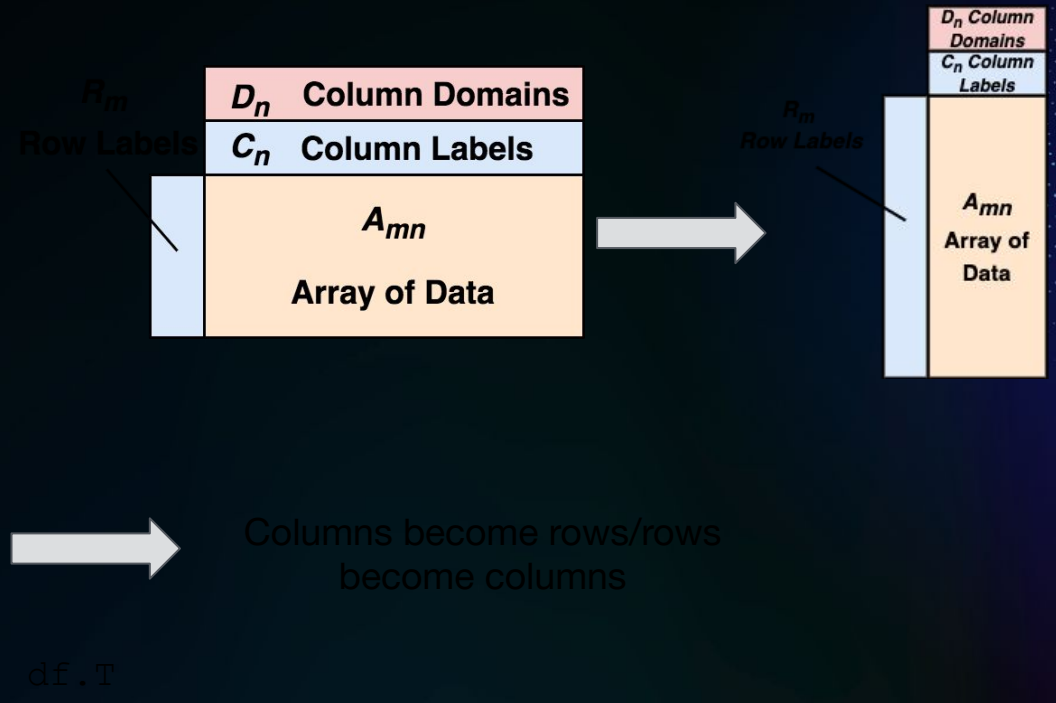
pandas++



Dataframe Algebra

PONDER

Operator
SELECTION
PROJECTION
UNION
DIFFERENCE
CROSS-PRODUCT / JOIN
DROP DUPLICATES
GROUPBY
SORT
RENAME
WINDOW
TRANSPOSE
MAP
TOLABELS
FROMLABELS

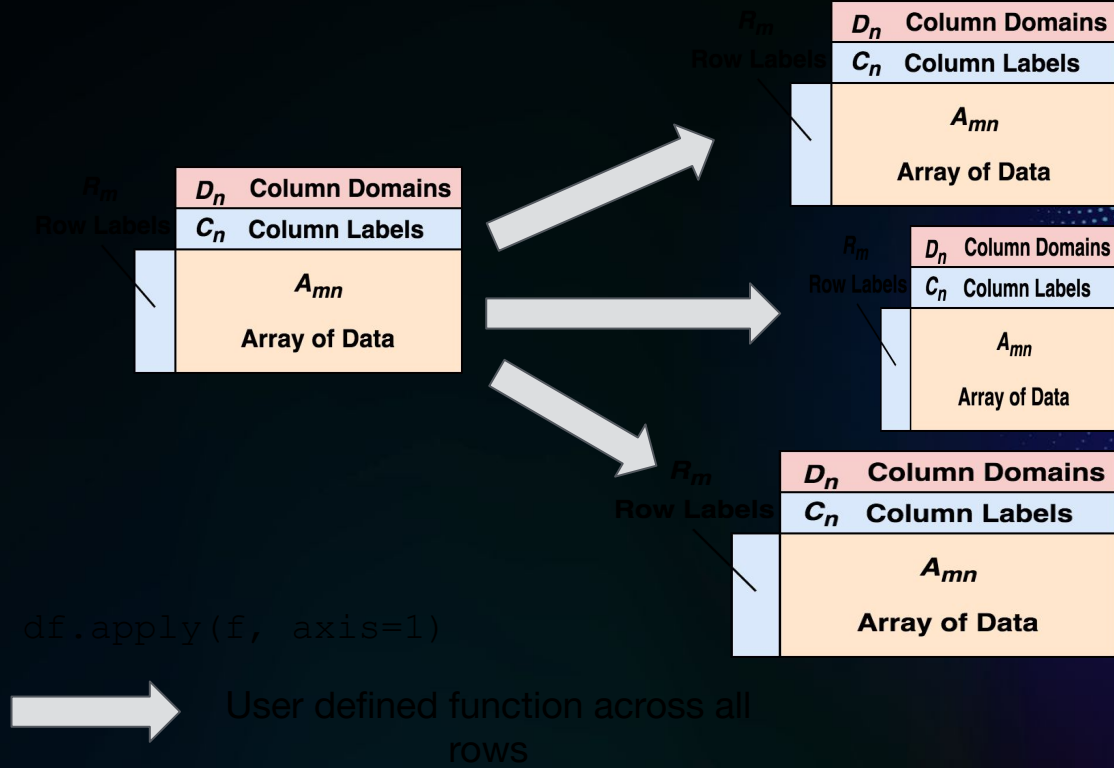




Dataframe Algebra

PONDER

Operator
SELECTION
PROJECTION
UNION
DIFFERENCE
CROSS-PRODUCT / JOIN
DROP DUPLICATES
GROUPBY
SORT
RENAME
WINDOW
TRANSPOSE
MAP
TOLABELS
FROMLABELS

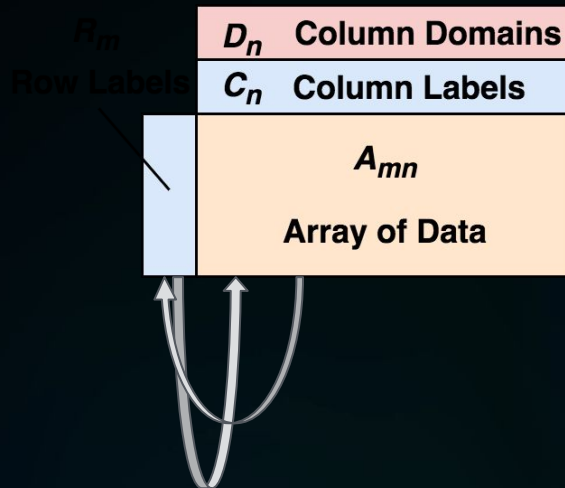




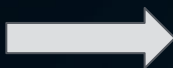
Dataframe Algebra

PONDER

Operator
SELECTION
PROJECTION
UNION
DIFFERENCE
CROSS-PRODUCT / JOIN
DROP DUPLICATES
GROUPBY
SORT
RENAME
WINDOW
TRANSPOSE
MAP
TOLABELS
FROMLABELS



```
df.set_index(df["col"])  
df.reset_index(drop=False)
```



Elevating data into metadata, or
moving metadata into the data

Working with tabular data

A (very) brief history

Relational Databases

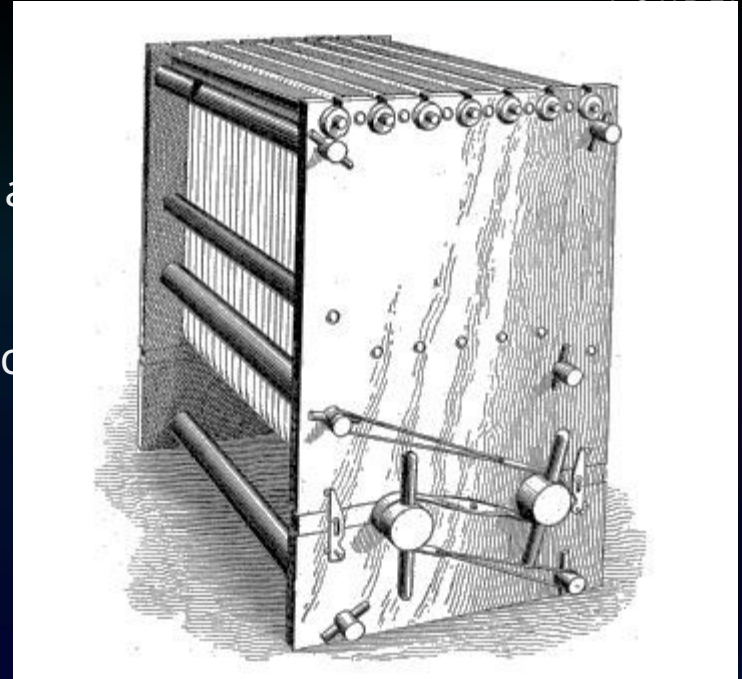
- Invented in the 1970's by Edgar F. Codd
 - Defined a data model for **structured data**
 - Schema must be known before input into DBMS
 - Decoupling of the physical representation from logical
- Popularized in the 1980's
 - Data model has stood the test of time, still in use!
 - SQL

Census - the original “Big Data”



PONDER

- 1880's - Seaton Device, manual intervention and 10 years
- 1890's - Herman Hollerith's machine, punch cards
- 1940's - First use of electronic computers



Dataframes



- Emerged from a real-world need at the time
 - No way of handling unstructured or semi-structured data
 - Matrices and Tables did not fit their need, something new needed
- Not formalized!
- Dataframes have an origin in S

Focus: pandas



pandas

pandas success



- 23 million installs/month
- Over 300 million total installs
- Used by over 209k projects in GitHub
- 24.6k GitHub stars

pd.DataFrame	280+ methods
pd.Series	280+ methods
Convenience methods (e.g. concat)	40+ APIs



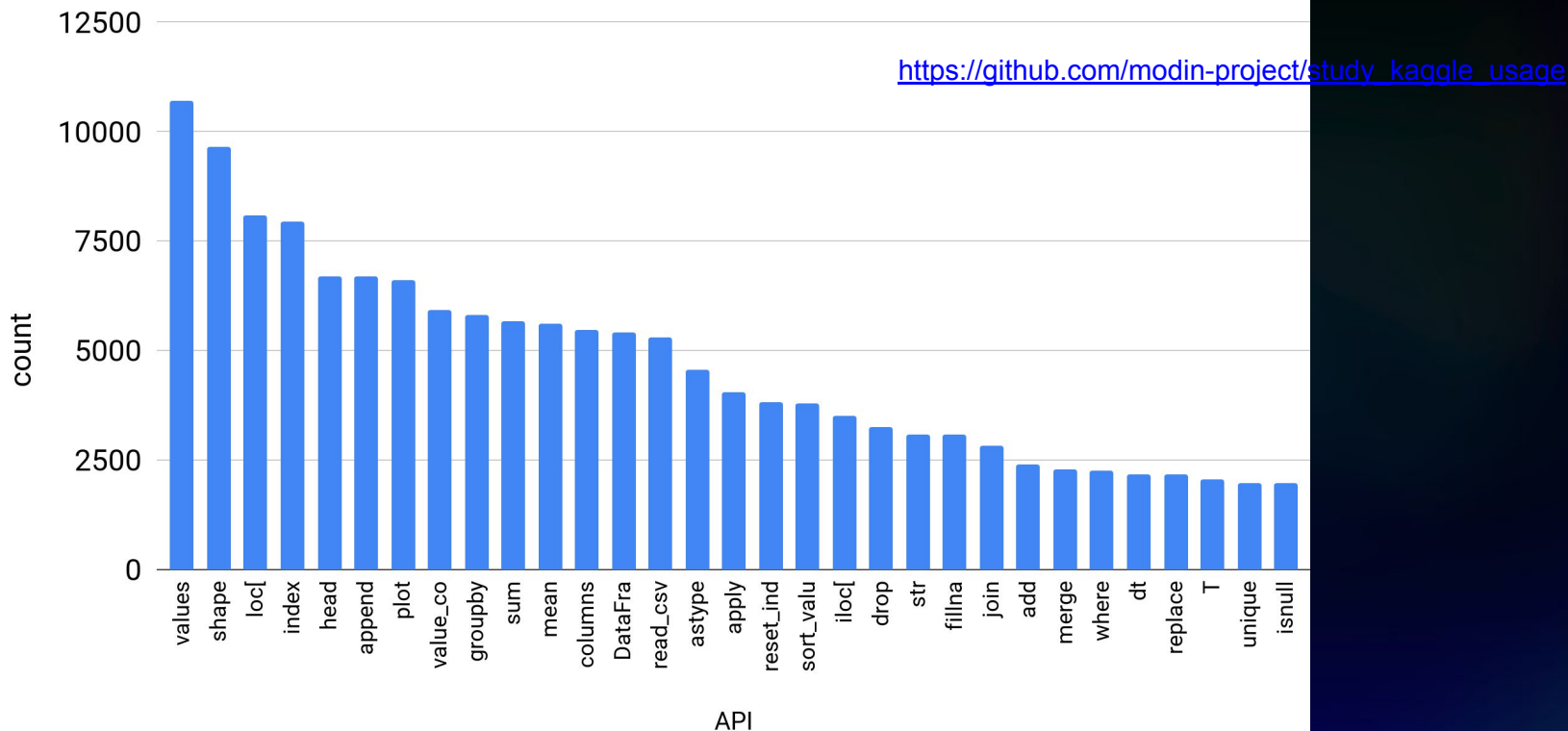
pandas API is huge and expressive

PONDER

pd.DataFrame	280+ methods
pd.Series	280+ methods
Convenience methods (e.g. concat)	40+ APIs

What do people use within pandas API?

Top 30 most used pandas APIs by count -- Kaggle notebooks dataset





Definition: Dataframes

PONDER

A_{mn}

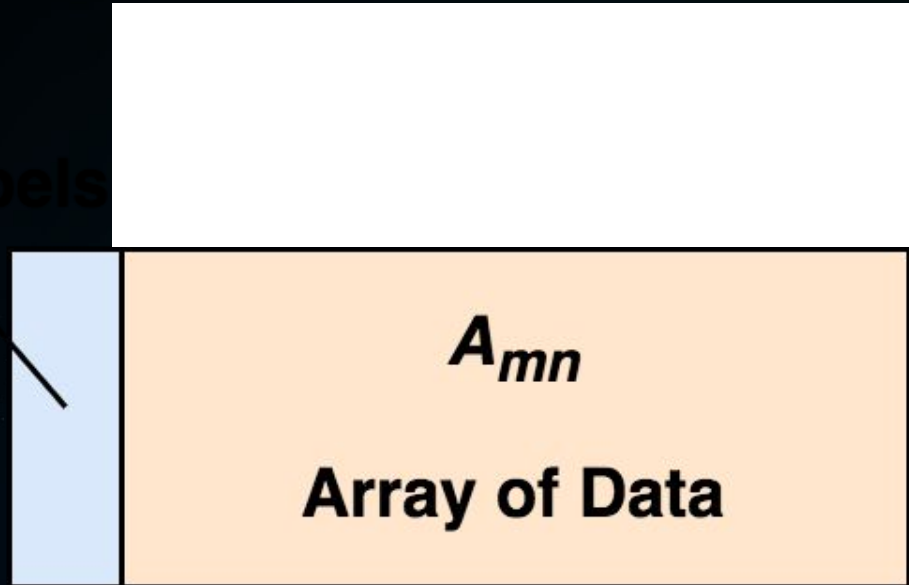
Array of Data



Definition: Dataframes

PONDER

A_m
Row Labels



A_{mn}
Array of Data



Definition: Dataframes

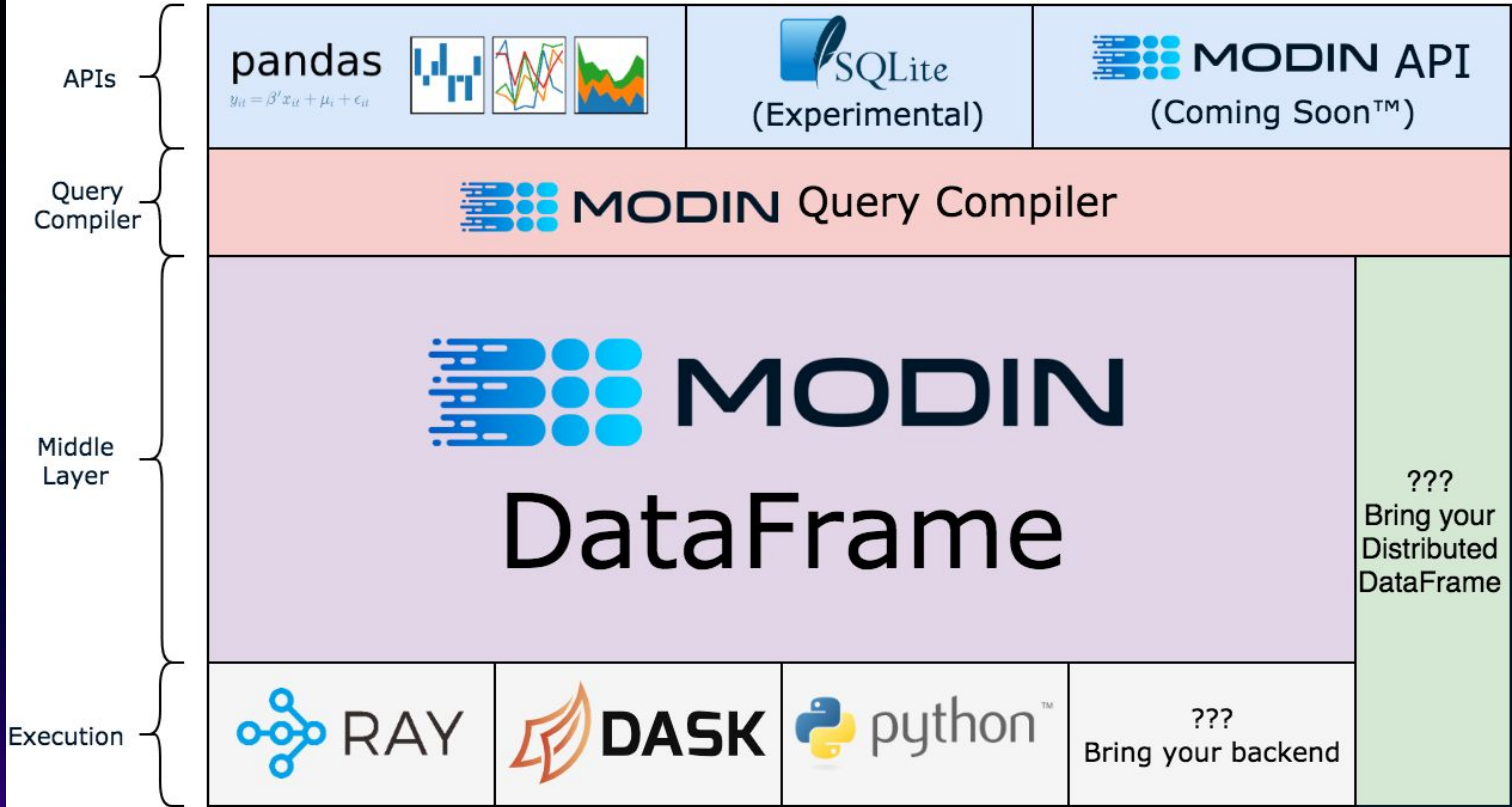
PONDER

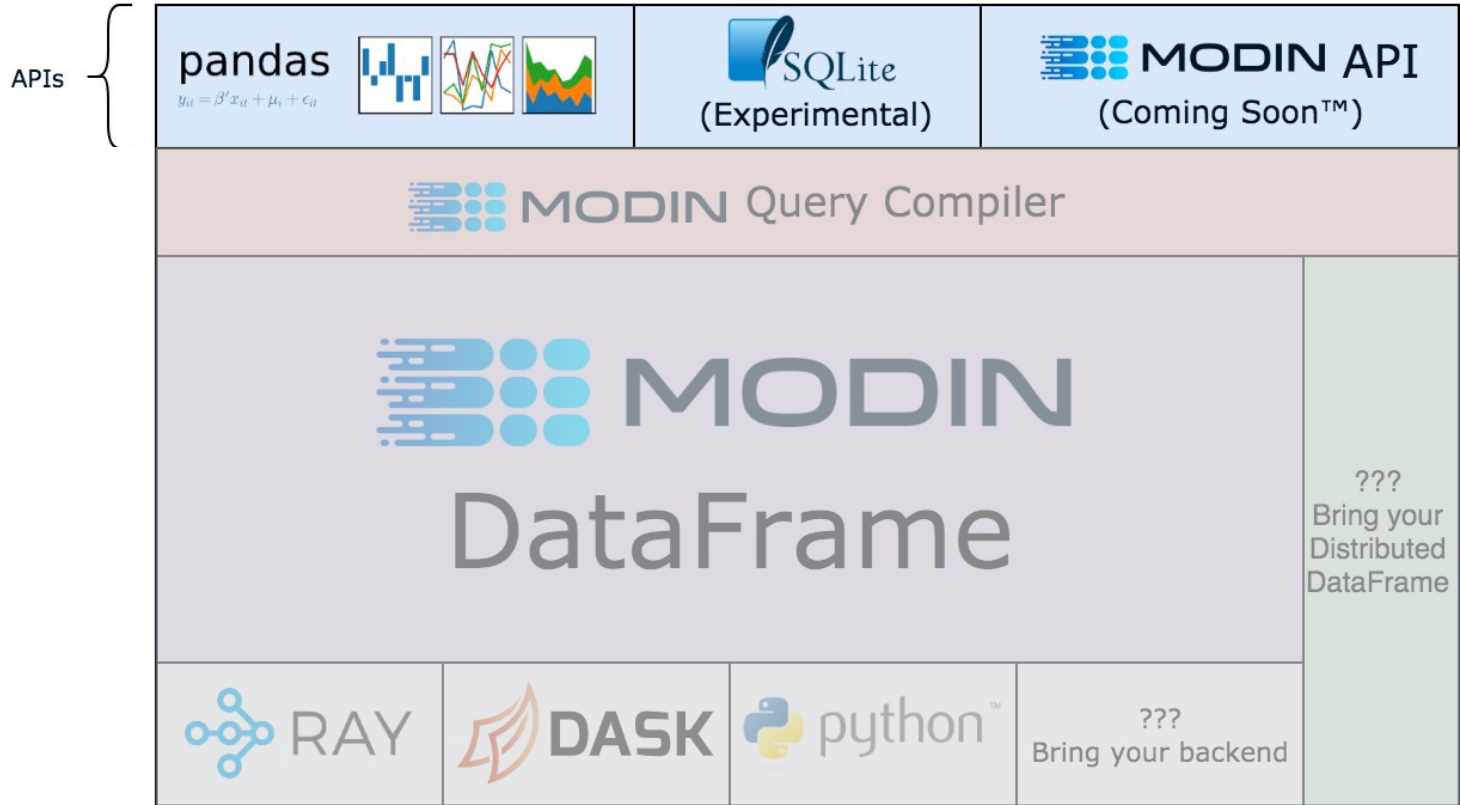
R_m
Row Labels

C_n Column Labels

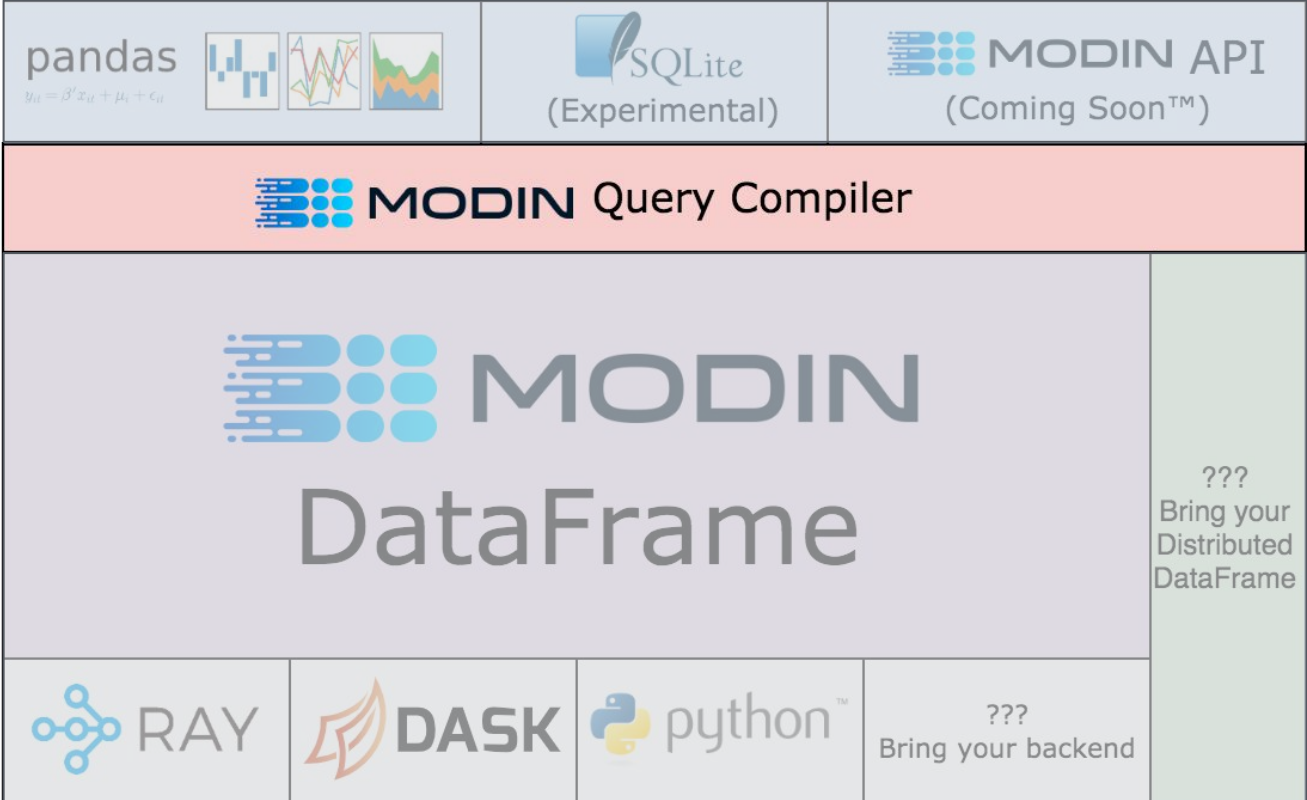
A_{mn}

Array of Data

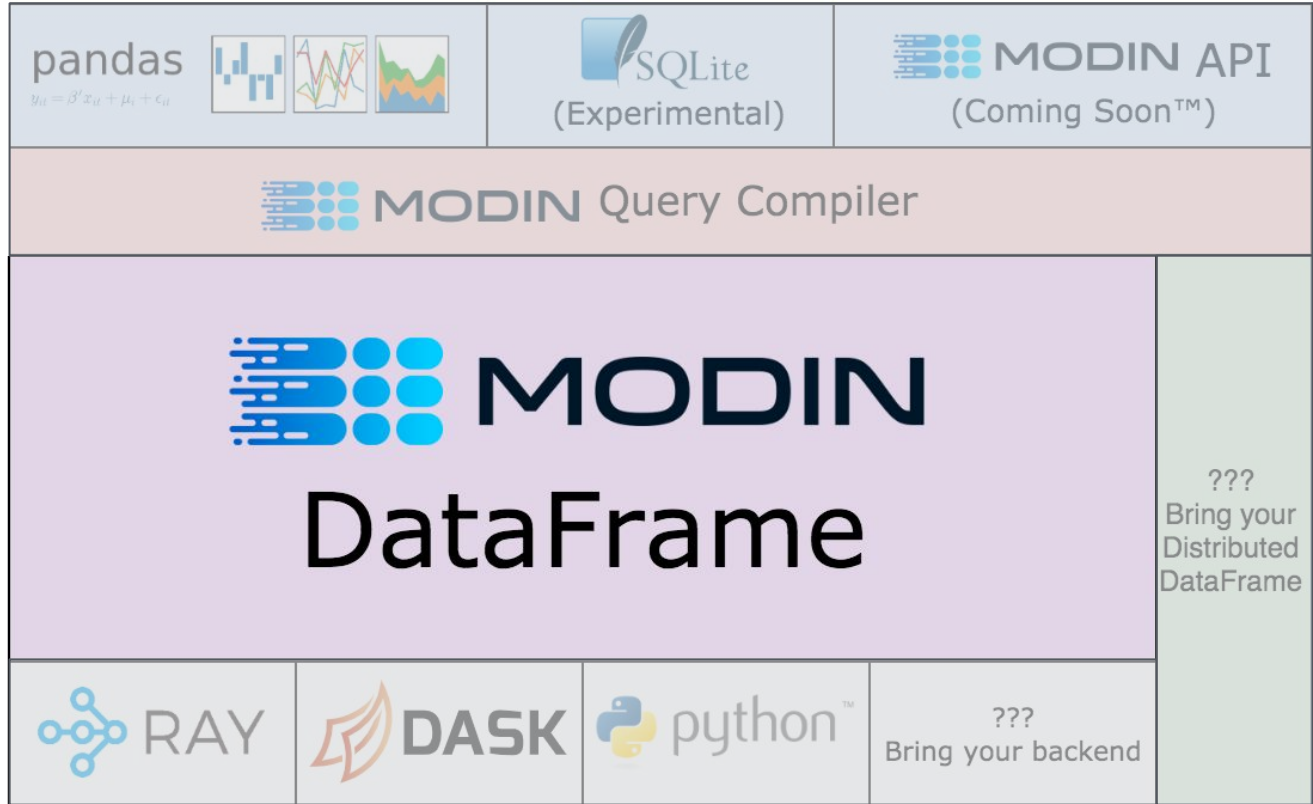


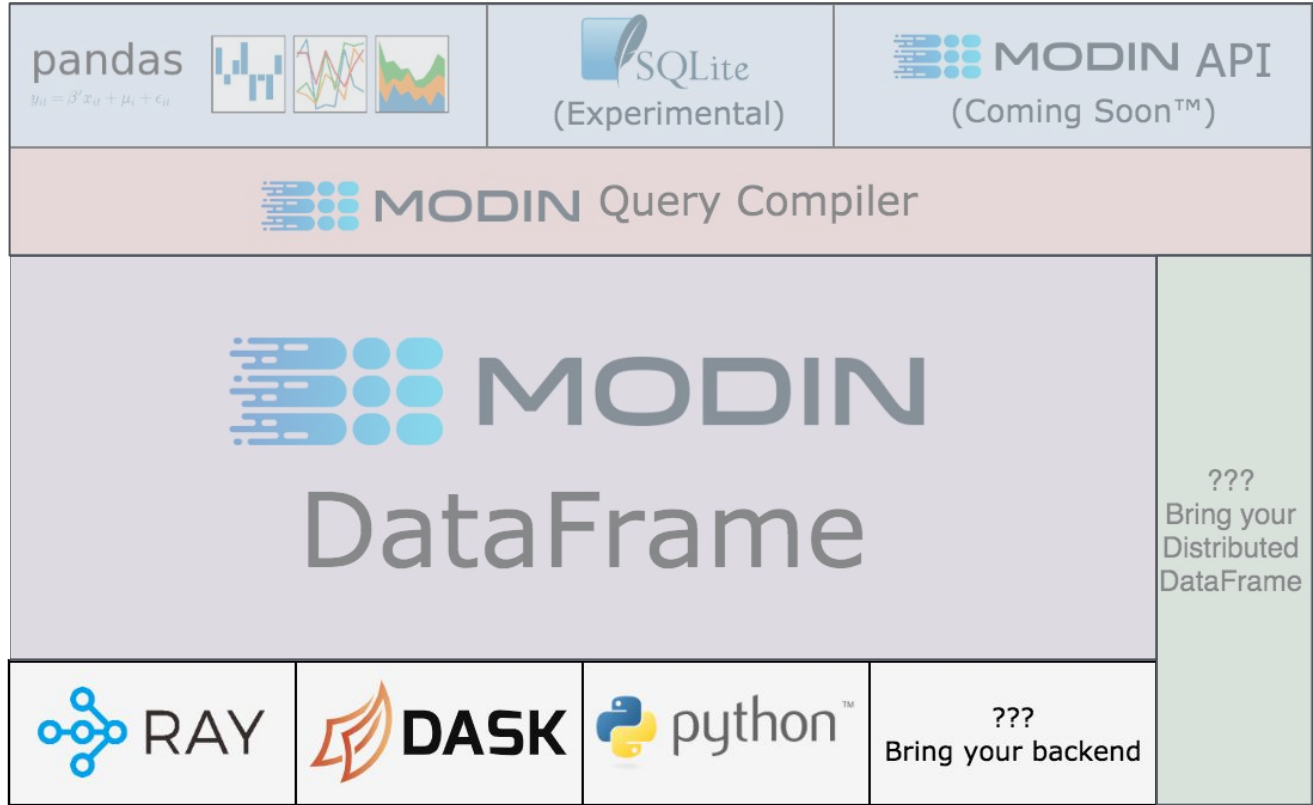


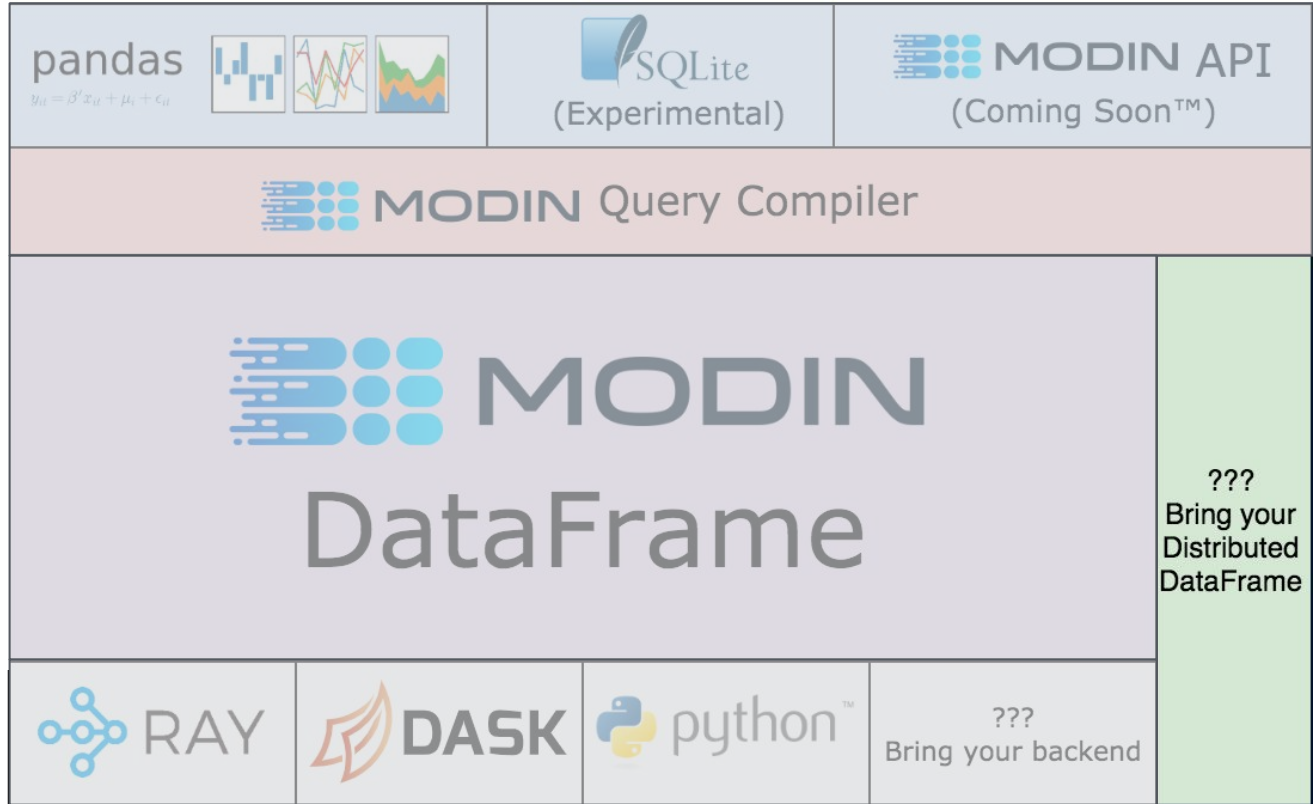
Query Compiler

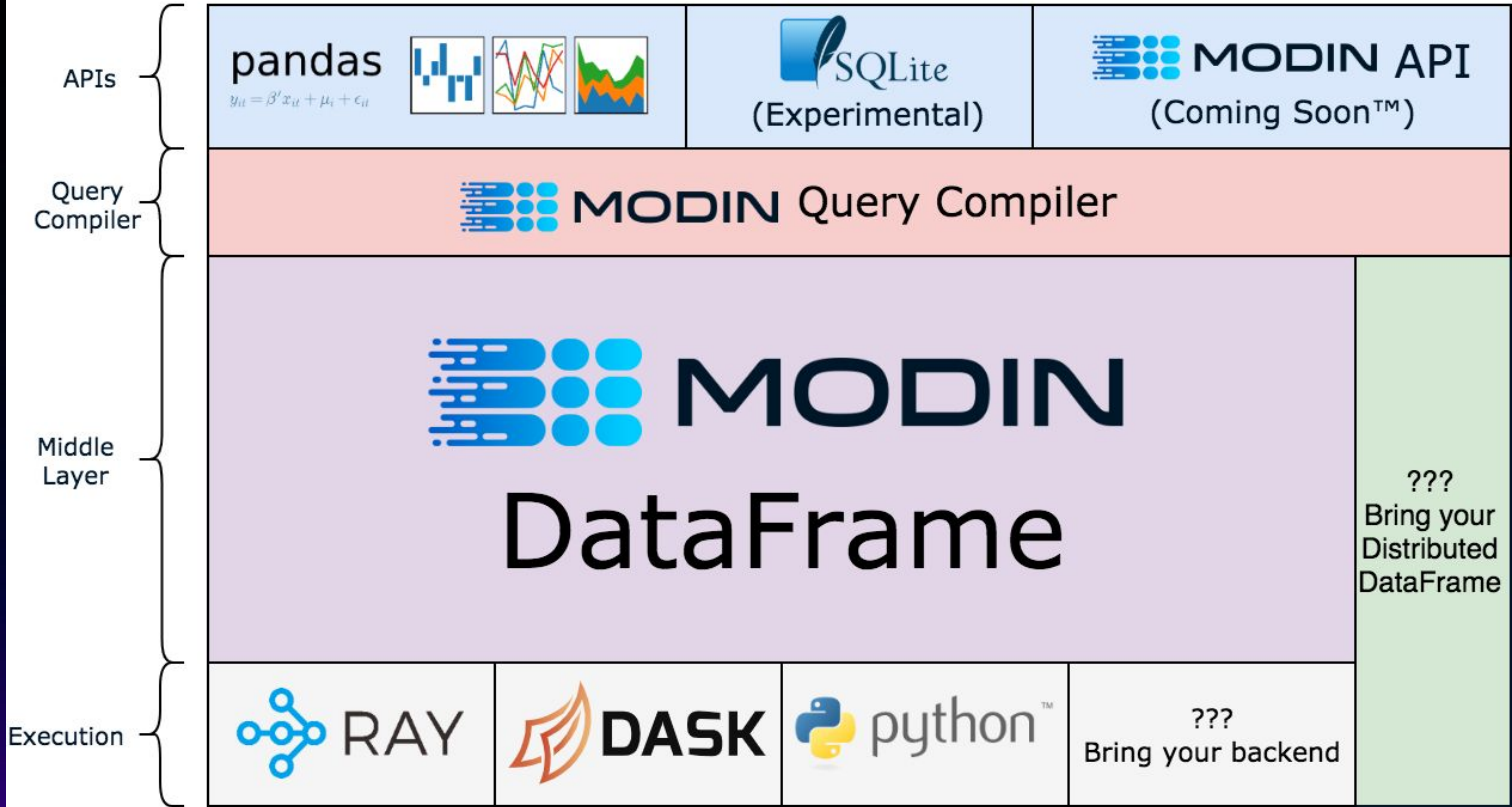


Middle Layer











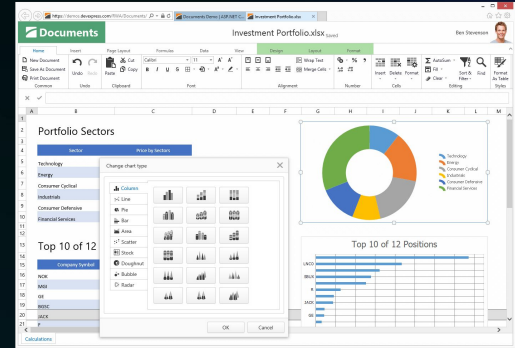
PONDER

Dataframes emerged from a need to hybridize

Relational Table

Name	FName	City	Age	Salary
Smith	John	3	35	\$280
Doe	Jane	1	28	\$325
Brown	Scott	3	41	\$265
Howard	Shemp	4	48	\$359
Taylor	Tom	2	22	\$250

Spreadsheet





Dataframes emerged from a need to hybridize

PONDER

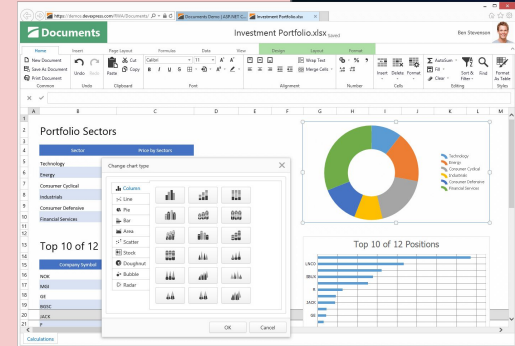
Matrix

$$\begin{pmatrix} 0 & 3 & 1 & 0 & 2 & 3 & 8 & 1 & 1 & 3 \\ 1 & 1 & 0 & 0 & 7 & 1 & 2 & 2 & 3 & 3 \\ 1 & 2 & 2 & 0 & 0 & 6 & 7 & 1 & 2 & 2 \\ 1 & 2 & 3 & 10 & 0 & 4 & 6 & 1 & 0 & 5 \\ 3 & 2 & 2 & 1 & 4 & 3 & 2 & 1 & 6 & 0 \\ 7 & 4 & 4 & 5 & 3 & 9 & 6 & 1 & 6 & 1 \\ 7 & 1 & 1 & 5 & 2 & 8 & 9 & 1 & 3 & 6 \\ 5 & 0 & 1 & 6 & 2 & 0 & 0 & 0 & 1 & 5 \\ 1 & 6 & 3 & 3 & 4 & 6 & 2 & 0 & 1 & 1 \\ 1 & 2 & 2 & 4 & 1 & 1 & 3 & 0 & 8 & 2 \end{pmatrix}$$

Relational Table

Name	FName	City	Age	Salary
Smith	John	3	35	\$280
Doe	Jane	1	28	\$325
Brown	Scott	3	41	\$265
Howard	Shemp	4	48	\$359
Taylor	Tom	2	22	\$250

Spreadsheet



Dataframes



Dataframe Algebra

PONDER

Operator
SELECTION
PROJECTION
UNION
DIFFERENCE
CROSS-PRODUCT / JOIN
DROP DUPLICATES
GROUPBY
SORT
RENAME
WINDOW
TRANSPOSE
MAP
TOLABELS
FROMLABELS



“Position notation”
or
“Named notation”

```
df.iloc[row_pos, col_pos]  
df.loc[row_lab, col_lab]
```



Dataframe Algebra

PONDER

Operator
SELECTION
PROJECTION
UNION
DIFFERENCE
CROSS-PRODUCT / JOIN
DROP DUPLICATES
GROUPBY
SORT
RENAME
WINDOW
TRANSPOSE
MAP
TOLABELS
FROMLABELS



Ordered union

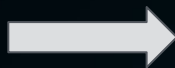
`df.append(df2)`



Dataframe Algebra

PONDER

Operator
SELECTION
PROJECTION
UNION
DIFFERENCE
CROSS-PRODUCT / JOIN
DROP DUPLICATES
GROUPBY
SORT
RENAME
WINDOW
TRANSPOSE
MAP
TOLABELS
FROMLABELS



Ordered Joins bring fundamentally new challenges

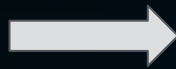
```
df.merge(df2, how="inner")  
df2.merge(df, how="inner")
```



Dataframe Algebra

PONDER

Operator
SELECTION
PROJECTION
UNION
DIFFERENCE
CROSS-PRODUCT / JOIN
DROP DUPLICATES
GROUPBY
SORT
RENAME
WINDOW
TRANSPOSE
MAP
TOLABELS
FROMLABELS



This is a generalized definition of groupby

```
df.count()  
df.groupby(df.columns).count()
```



Dataframe Algebra

PONDER

Operator
SELECTION
PROJECTION
UNION
DIFFERENCE
CROSS-PRODUCT / JOIN
DROP DUPLICATES
GROUPBY
SORT
RENAME
WINDOW
TRANSPOSE
MAP
TOLABELS
FROMLABELS



Conceptually a rolling function,
can output same table shape or
smaller (groupby)

```
df.cumsum()  
df.rolling
```

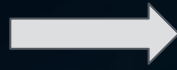


Dataframe Algebra

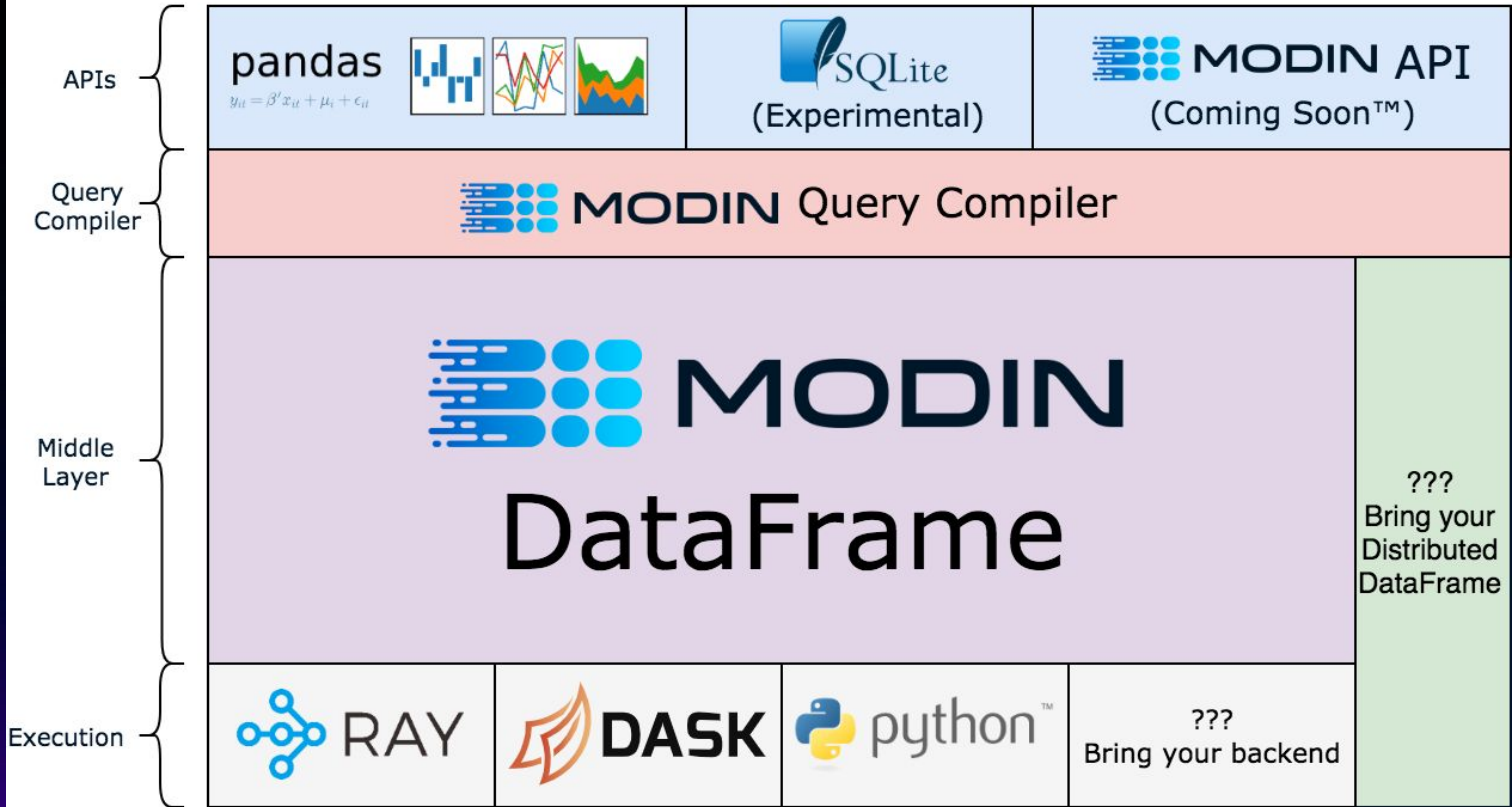
PONDER

Operator
SELECTION
PROJECTION
UNION
DIFFERENCE
CROSS-PRODUCT / JOIN
DROP DUPLICATES
GROUPBY
SORT
RENAME
WINDOW
TRANSPOSE
MAP
TOLABELS
FROMLABELS

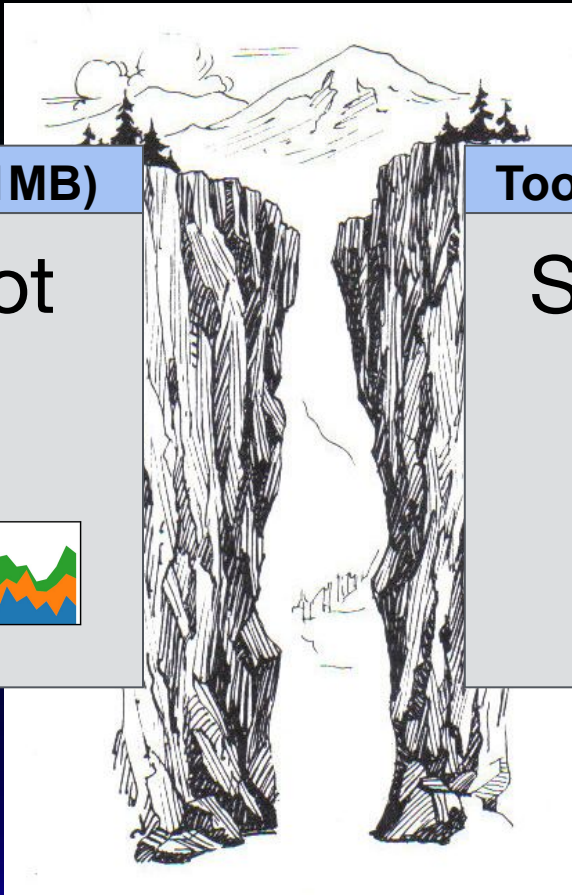
```
df.reset_index(drop=False)
```



Insert the row labels into the data
and reset the row labels to the
positional notation



Data Science Landscape: Today

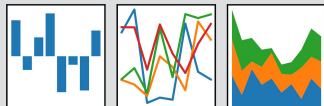


Tools efficient for O(1MB)

Usable but not
scalable

pandas

$$y_{it} = \beta' x_{it} + \mu_i + \epsilon_{it}$$



Tools efficient for O(100s GB+)

Scalable but not
usable



Data Science Landscape: Today

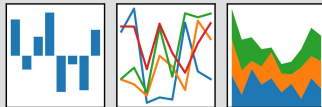
Tools efficient for O(1MB)

Usable and
scalable



pandas

$$y_{it} = \beta' x_{it} + \mu_i + \epsilon_{it}$$

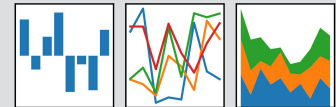


Tools efficient for O(100s GB+)

Scalable and
usable

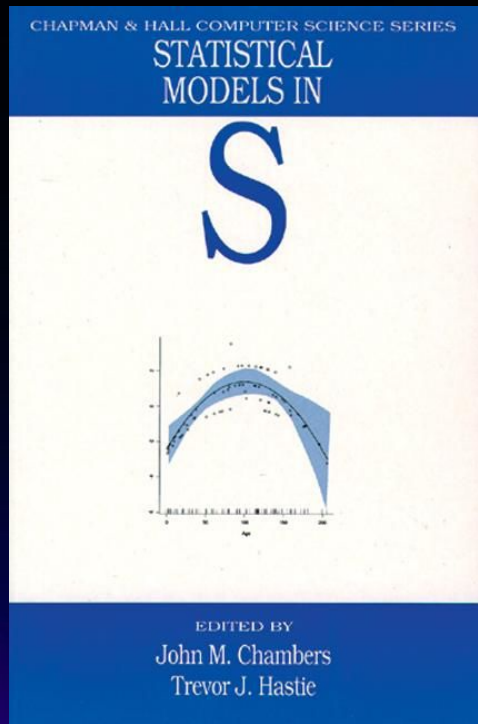
pandas

$$y_{it} = \beta' x_{it} + \mu_i + \epsilon_{it}$$



Dataframe origin

A (not so) long time ago, at Bell Labs

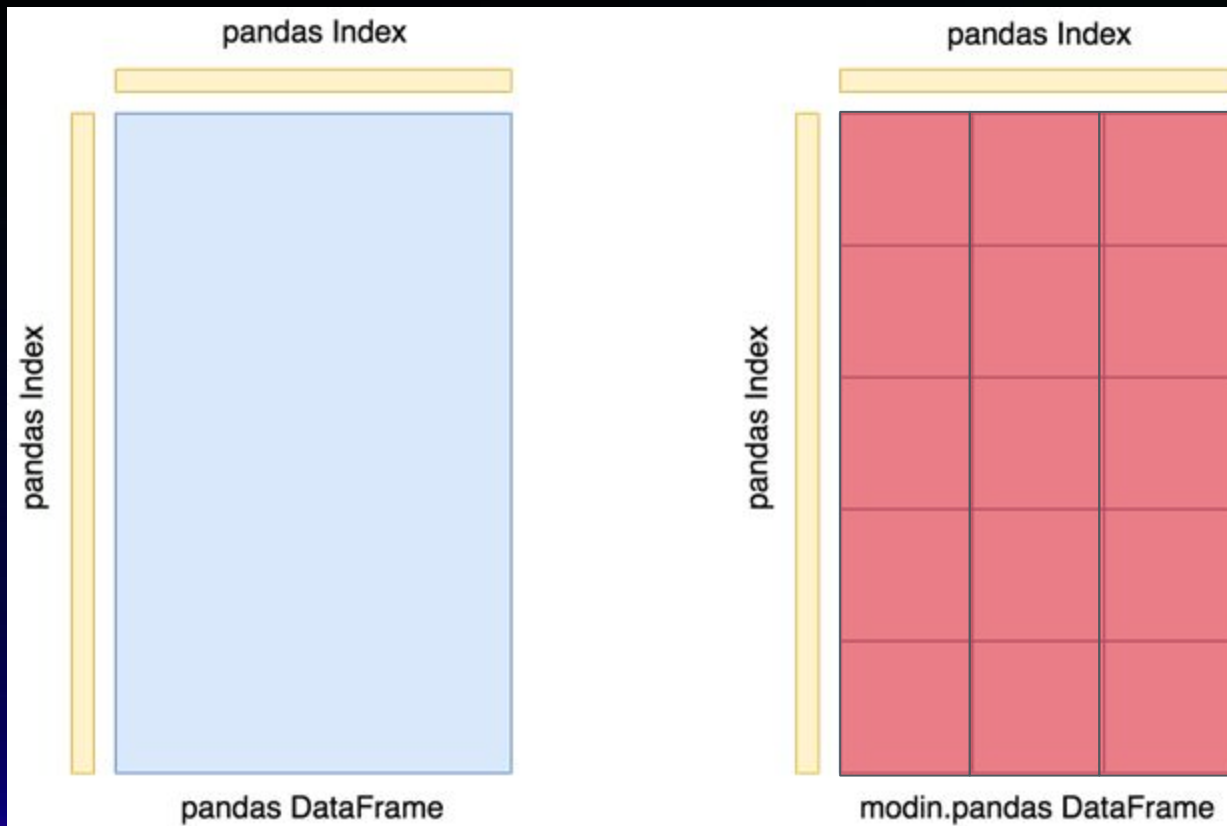


Chapter 3: Data for Models

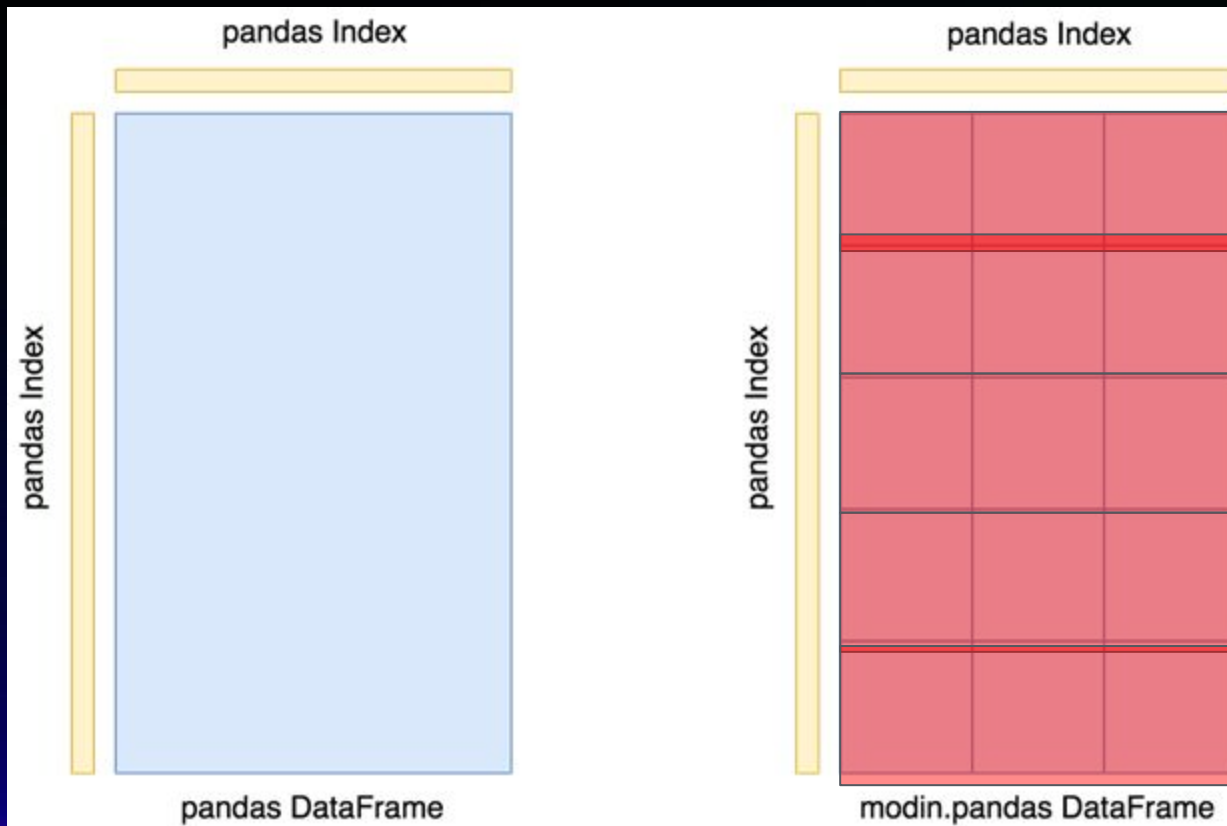
“Dataframes are more general than matrices in the sense that matrices in S assume all elements to be of the same mode—all numeric, all logical, all character string, etc.”

“... data frames support matrix-like computation, with variables as columns and observations as rows, and, in addition, they allow computations in which the variables act as separate objects, referred to by name.”

Modin partitioning - logical column partitioning



Modin partitioning - logical row partitioning



Modin architecture

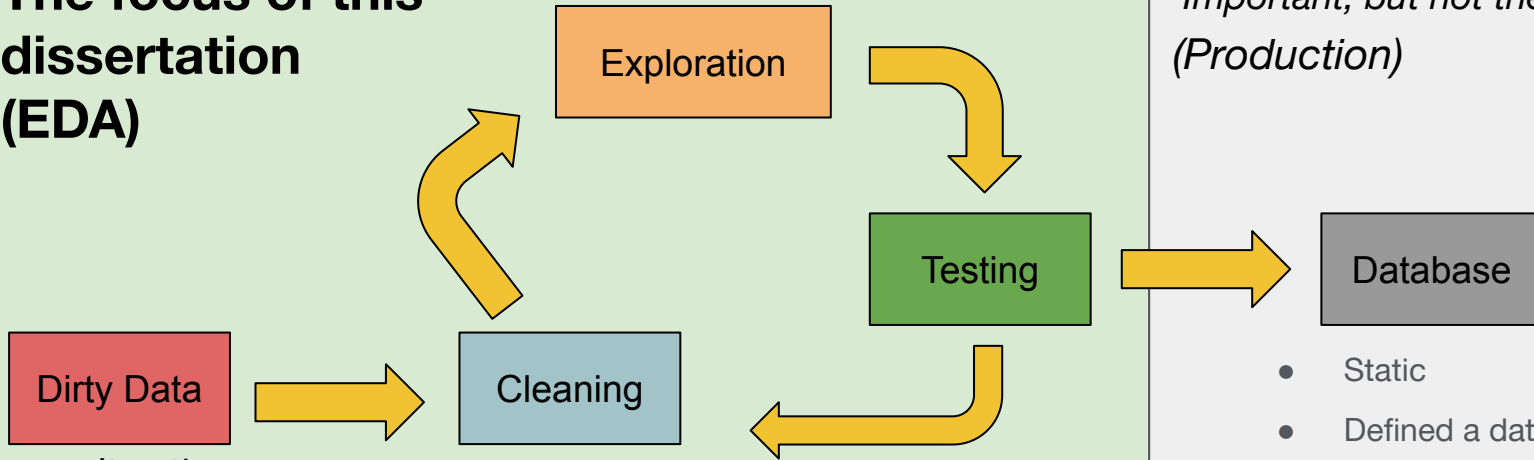
- Highly flexible
- Layered architecture
 - New optimizations can be implemented as they are developed
- Support for dataframe algebra
- Partitioning approach lends itself to allowing us to use optimizations from multiple domains

Data Science LifeCycle



PONDER

The focus of this dissertation (EDA)



- Iterative
- **Unstructured** data model
- Schema is lazily induced
- No decoupled physical representation from logical

Important, but not the focus (Production)

- Static
- Defined a data model for **structured data**
- Schema must be known before input into DBMS
- Decoupling of the physical representation from logical

Data Science LifeCycle



PONDER

The focus of this dissertation (EDA)

Dirty Data

- Iterative
- **Unstructured** data model
- Schema is lazily induced
- No decoupled physical representation from logical

Dataframes

Exploration

Testing

Important, but not the focus (Production)

Databases

- Prepared a data model for **structured data**
- Schema must be known before input into DBMS
- Decoupling of the physical representation from logical