

# Towards a Formal Type System for DataFrames

William W. Ma  
University of Chicago  
willamma@uchicago.edu

## ABSTRACT

While there has been a lot of work to improve the computational efficiency of data science, little work has been done to improve the data scientists ability to process data. The tools that they use, such as pandas, has evolved from single threaded to cloud scale. However, these tools still lack a consistency in their behavior. In this paper, we introduce the beginnings of a type system for dataframes. With this type system, we will be able to formally prove systems-level optimizations, but, more importantly, provide the data scientist with a consistent type system for them to do their analysis on.

## CCS CONCEPTS

• **Computer systems organization** → **Embedded systems**; *Redundancy*; Robotics; • **Networks** → Network reliability.

## KEYWORDS

datasets, neural networks, gaze detection, text tagging

### ACM Reference Format:

William W. Ma. 2018. Towards a Formal Type System for DataFrames. In *Woodstock '18: ACM Symposium on Neural Gaze Detection, June 03–05, 2018, Woodstock, NY*. ACM, New York, NY, USA, 3 pages. <https://doi.org/10.1145/1122445.1122456>

## 1 INTRODUCTION

With the rise of data science, dataframes rose in popularity as they were the go to data manipulation model for data analysis. In order to satisfy these data scientists, popular cloud-scale systems, such as Spark and dask, all provide dataframe interfaces. However, as dataframes have gone from single-threaded pandas dataframes to cloud-scale dataframes with SparkSQL, these dataframes all lack a consistent data model and type system. Although prior work has addressed the issue of a data model, there has been no attempt at a consistent data model [2]. While data models and type systems do not have catchy punchlines or provide 10x improvements in run time, they do provide clarity for the data scientist. Arguably, this clarity could be worth more than the 10x improvement that any system could provide.

Thus, in this paper, we introduce the beginnings of a consistent type system for dataframes. The type system on the data types is provided at both the dataframe level and individual cell level. Due

---

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

*Woodstock '18, June 03–05, 2018, Woodstock, NY*

© 2018 Association for Computing Machinery.

ACM ISBN 978-1-4503-XXXX-X/18/06...\$15.00

<https://doi.org/10.1145/1122445.1122456>

to the introductory nature of this work, some simplifications have been made but this does not detract from the main idea that we can have a consistent type system for dataframes.

We also included proofs of progress and preservation for both type systems.

The rest of this paper is structured as follows. We first introduce the logical models used and simplifications made in section 2. Then, we discuss the physical implementation in Coq 8.13 in section 3. Finally, we discuss future work in section 4.

## 2 LOGICAL MODEL

### 2.1 DataFrame Level Type System

For our purposes, we reduced a dataframe down to a tuple that consists of

- Row index, which consists of the position and label.
- Row data types (dtypes)
- Column index
- Column dtypes
- Matrix of values

Building on this data model, we reduce the possible functions on the dataframe to the following

- Transpose, which transposes the dataframe
- Mask, which gives the index positions to keep. This can be applied over the columns or the rows.
- Filter, which gives a function that returns a boolean when applied over the rows or columns. Then, only the true rows or columns are kept.
- ToLabels, which converts a given column or row index position to the corresponding axis' index.
- FromLabels, which adds a row or column that consists of the axis' index labels.
- Map, which applies a map over the rows or columns. We simplified this to also require the possible data types that each column or row could result into.
- Concat, which is the union, intersection, or difference of the columns or rows of two dataframes.
- InferDTypes, which infers the data types from the values. This can be applied to the columns or rows.

For this introductory work, we limit the data types that a dataframe can have to be only strings and natural numbers. Thus, the data values at the dataframe level are

```
DFValue :=
| DFV_String
| DFV_Nat
```

The terms for the type system are

```

Axis :=
| Columns
| Rows

Index := ℕ * DFValue

ConcatType :=
| Union
| Intersection
| Difference

DataFrame :=
| DataFrame
| Transpose DataFrame
| Mask (list ℕ)
    Axis DataFrame
| Filter (list DFValue → ℬ)
    Axis DataFrame
| ToLabels ℕ Axis DataFrame
| FromLabels Axis DataFrame
| Map (list DFValue → list DFValue)
    (list DFValue) DataFrame
| Concat ConcatType Axis
    DataFrame1 DataFrame2
| InferDTypes Axis DataFrame

```

The only value in the dataframe level type system is the dataframe constructor itself. Because every function can take in a value or step, we only list the step relation for transpose

$$\text{TransValue} \frac{\text{value } df}{\text{Transpose } df \rightarrow \text{DataFrame } (df.\text{transpose})}$$

$$\text{TransStep} \frac{df \rightarrow df'}{\text{Transpose } df \rightarrow \text{Transpose } df'}$$

We also include three optimizations, which are listed below:

$$\text{TransTwice} \frac{\text{value } df}{\text{Transpose}(\text{Transpose } df) \rightarrow df}$$

$$\text{FilterNone} \frac{\text{value } df}{\text{Filter } (\text{fun row} \Rightarrow \text{true}) \text{ axis } df \rightarrow df}$$

$$\text{ToFromLabels} \frac{\text{value } df}{\text{ToLabels } \emptyset \text{ a } (\text{FromLabels } a \text{ } df) \rightarrow df}$$

For the type definitions, we follow the following format for all of the terms.

$$\frac{\vdash df \in \text{row}_{dtypes}, \text{col}_{dtypes} \quad df' = \text{Transpose } df}{\vdash df' \in \text{col}_{dtypes}, \text{row}_{dtypes}}$$

## 2.2 Data Level Type System

At the data level, we adjust the data model to use a matrix of values that is infinitely large. Thus, we have to introduce a “empty” cell value. Because we switch over to the data level, we have to define the casting rules between data types. We allow all data types to be cast into the “empty” and string data types. In addition to the above

list of functions on dataframes, we also introduce a cast function for the data level type system.

Additionally, we change the semantics of some of the functions to work independently of the other data values in the row and column. For mask, filter, and map, we assume that the functions passed in can work on the individual data values independent of the other values in the row or column. It is worth noting that at the data level, Transpose, ToLabels, FromLabels, Concat, and InferDTypes are all given a new value to replace the old value with. Thus, for simplicity, we replace these functions with an Update function.

Thus, the terms, values, and types are

```

dterm :=
| DEmpty dterm
| DString dterm
| DNat dterm
| CastDType dterm dterm dterm
| Mask ℬ dterm dterm
| Filter (dterm → ℬ) dterm dterm
| Map (dterm → dterm) dterm dterm
| Update dterm dterm dterm

dvalue :=
| dv_empty DEmpty
| dv_nat DNat
| dv_string DString

dtype :=
| Empty
| String
| Nat

```

The step relations are as follows:

$$\text{CastEmpty} \frac{}{\text{CastDType } d \text{ DEmpty} \rightarrow d}$$

$$\text{CastToString} \frac{}{\text{CastDType } DString \text{ } d \rightarrow DString}$$

$$\text{CastNat} \frac{}{\text{CastDType } DNat \text{ } DNat \rightarrow DNat}$$

$$\text{MaskFalse} \frac{}{\text{Mask } \text{false } d \rightarrow DEmpty}$$

$$\text{MaskTrue} \frac{}{\text{Mask } \text{true } d \rightarrow d}$$

$$\text{FilterTrue} \frac{\text{CastDType } d_{in} \text{ } d \rightarrow d'}{\text{Filter } d_{in} \text{ true } d' \rightarrow d'}$$

$$\text{FilterFalse} \frac{\text{CastDType } d_{ind} \text{ } d' \rightarrow d'}{\text{Filter } \text{true } d_{in} \text{ false } d' \rightarrow DEmpty}$$

$$\text{Map} \frac{\text{CastDType } d_{in} \text{ } d \rightarrow d'}{\text{Map } d_{in} \text{ } d_{out} \text{ } d' \rightarrow d_{out}}$$

$$\text{Update} \frac{}{\text{Update } d_{new} \text{ } d_{old} \rightarrow d_{new}}$$

The type definitions are listed below.

$$\begin{array}{c}
\frac{}{\text{T\_Empty} \vdash \text{DEmpty} \in \text{Empty}} \\
\frac{}{\text{T\_String} \vdash \text{DString} \in \text{String}} \\
\frac{}{\text{T\_Nat} \vdash \text{DNat} \in \text{Nat}} \\
\frac{}{\text{T\_CastEmpty} \vdash \text{CastDType } d \text{ DEmpty} \in D} \\
\frac{}{\text{T\_CastToString} \vdash \text{CastDType } \text{DString } d \in \text{String}} \\
\frac{}{\text{T\_CastNat} \vdash \text{CastDType } \text{DNat } \text{DNat} \in \text{Nat}} \\
\frac{}{\text{T\_MaskFalse} \vdash \text{Mask } \text{false } d \in \text{Empty}} \\
\frac{}{\text{T\_MaskTrue} \vdash \text{Mask } \text{true } d \in D} \\
\frac{\vdash d \in D \quad \vdash d_{in} \in D_{in} \quad \vdash d' \in D_{in} \quad \text{CastDType } d_{in} d \rightarrow d'}{\text{T\_FilterTrue} \vdash \text{Filter } d_{in} \text{ true } d' \in D_{in}} \\
\frac{\vdash d \in D \quad \vdash d_{in} \in D_{in} \quad \vdash d' \in D_{in} \quad \text{CastDType } d_{in} d \rightarrow d'}{\text{T\_FilterFalse} \vdash \text{Filter } d_{in} \text{ false } d' \in \text{Empty}} \\
\frac{\vdash d \in D \quad \vdash d_{in} \in D_{in} \quad \vdash d_{out} \in D_{out} \quad \text{CastDType } d_{in} d \rightarrow d'}{\text{T\_Map} \vdash \text{Map } d_{in} d_{out} d' \text{ false } d' \in D_{out}} \\
\frac{\vdash d_{new} \in D}{\text{T\_Update} \vdash \text{Update } d_{new} d_{old} \in D}
\end{array}$$

### 3 PHYSICAL IMPLEMENTATION

We implemented the type systems and dataframe operations within Coq version 8.13<sup>1</sup>. Proofs of progress and preservation are within their respective `TypeSystem_*.v` files. The implementations of the dataframe operations are within `Dataframe.v`. We do include an `Optimizations.v` file but it is empty and is there mainly for thinking about the organization of the project.

### 4 FUTURE WORK

For future work, we intend to extend this type system to fully represent the behavior of pandas. Given the foundation in this paper, adding additional functions in the dataframe level type system will only require implementing them within Coq. For the data level, additional operators will look a lot like the existing step relations and type definitions. For instance joins will most likely look like a filter and groupbys will end up looking a lot like maps.

This work done for this paper will lay the ground work for formal verification of dataframe optimizations. These optimizations include the fact that the transpose of an operation done on a transpose is essentially the operation with a different axis argument, assuming that the operation has an axis argument. Some of the helper theorems have been defined in `Dataframe.v` to demonstrate the type of theorems that will be needed to prove the optimizations that we use.

This project is also on github [1].

## REFERENCES

- [1] William W. Ma. 2021. dataframe-proofs. [https://github.com/williamma12/dataframe\\_proofs](https://github.com/williamma12/dataframe_proofs).
- [2] Devin Petersohn, William W. Ma, Doris Jung Lin Lee, Stephen Macke, Doris Xin, Xiangxi Mo, Joseph E. Gonzalez, Joseph M. Hellerstein, Anthony D. Joseph, and Aditya G. Parameswaran. 2020. Towards Scalable Dataframe Systems. *CoRR* abs/2001.00888 (2020). arXiv:2001.00888 <http://arxiv.org/abs/2001.00888>

<sup>1</sup>I couldn't get auto to use a custom database so I added hints to core. I will fix this once I have time.