

Beyond Linearity

Building reactive notebooks for data

Caitlin Colgrove, CTO @ Hex



Poll: how do code notebooks make you *feel*?

- A. I use notebooks for everything! Analysis, text editing, email... all notebooks!
- B. They're useful sometimes but they have their drawbacks.
- C. I will literally quit my job if they make me use a notebook.
- D. You mean, like... to write in?

Historical background: literate programming

In 1984, Donald Knuth introduced the concept of "literate programming", a way of developing that mixes code, explanation, and outputs together in a way that's meant to be more interpretable by humans.

```
@ Here is a Perl program that simply
prints out |Hello, world!| the number of
times specified in the first argument.

<<*>>=
#!/usr/bin/perl
  <<CheckArgs>>
  <<PrintHiWorld>>

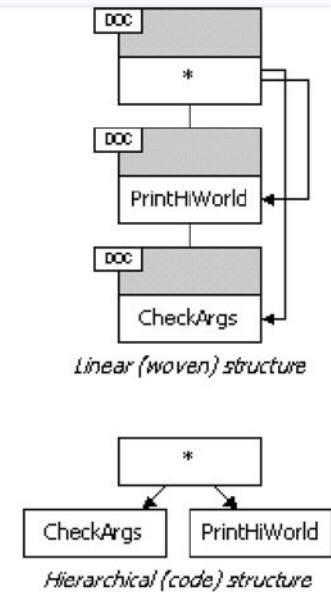
@ Printing involves a simple loop. Line
breaks are added for clarity.

<<PrintHiWorld>>=
for ($i = 0; $i < $ARGV[0]; $i++) {
    print "Hello, world!\n";
}

@ We \emph{must} make sure, however,
that an argument was specified.

<<CheckArgs>>=
if (@ARGV != 1) {
    die "No argument specified";
}
```

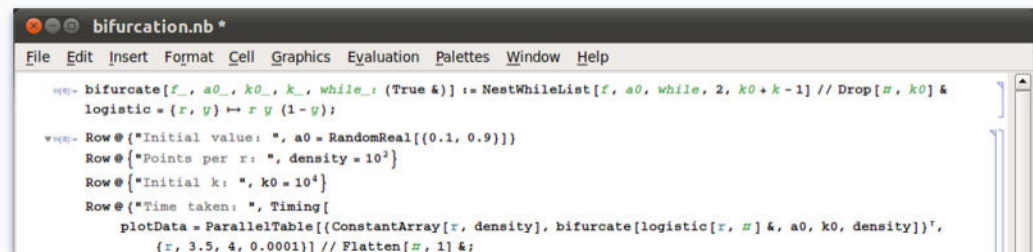
(a) Literate source.



(b) Linear and hierarchical views.

Fast forward to 2022

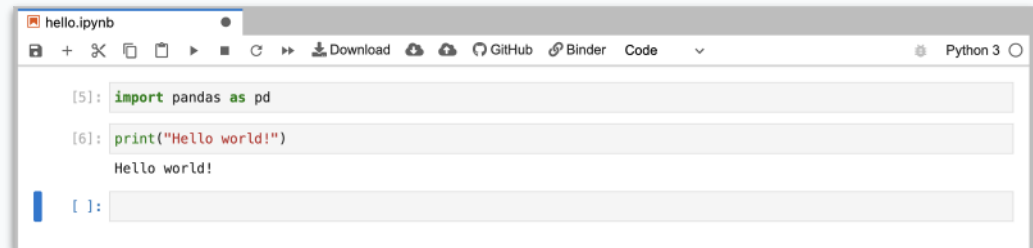
Notebooks are the most widely-used example of literate programming in practice.



```
bifurcation.nb *
File Edit Insert Format Cell Graphics Evaluation Palettes Window Help

bifurcate[f_, a0_, k0_, k_, while_: (True &)] := NestWhileList[f, a0, while, 2, k0 + k - 1] // Drop[#, k0] &
logistic = {x, y} -> x y (1 - y);

Row@{"Initial value: ", a0 = RandomReal[{0.1, 0.9}]}
Row@{"Points per x: ", density = 10^3}
Row@{"Initial k: ", k0 = 10^4}
Row@{"Time taken: ", Timing[
  plotData = ParallelTable[{ConstantArray[x, density], bifurcate[logistic[x, #] &, a0, k0, density]}',
    {x, 3.5, 4, 0.0001}] // Flatten[#, 1] &;
```



```
hello.ipynb
Download GitHub Binder Code Python 3

[5]: import pandas as pd

[6]: print("Hello world!")
Hello world!

[ ]:
```



```
Published Feb 17

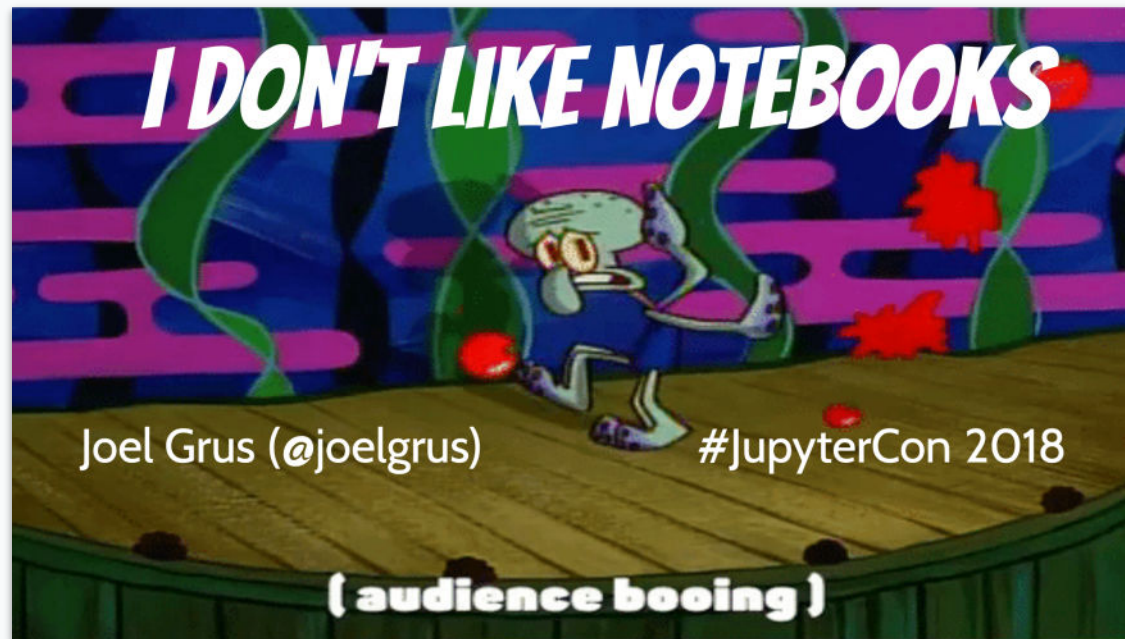
JS lists and arrays
md # JS lists and arrays

mylist = > Array(5) ["tekst", "b", "car", 4, 12.5]
```

Why notebooks?

- Mix code and outputs together
- Great for iterating on smaller chunks of code; well-suited to exploration
- Linear, narrative layout that is great for storytelling

But notebooks have... issues



The State Problem

```
a = 1
```

```
a = 2
```

```
print(a)
```

What does this print?

imperative programming

a programming paradigm that uses statements that change a program's state.

Notebook state causes 3 major problems

1. Interpretability

It's hard to reason about what's happening in a notebook, especially someone else's.

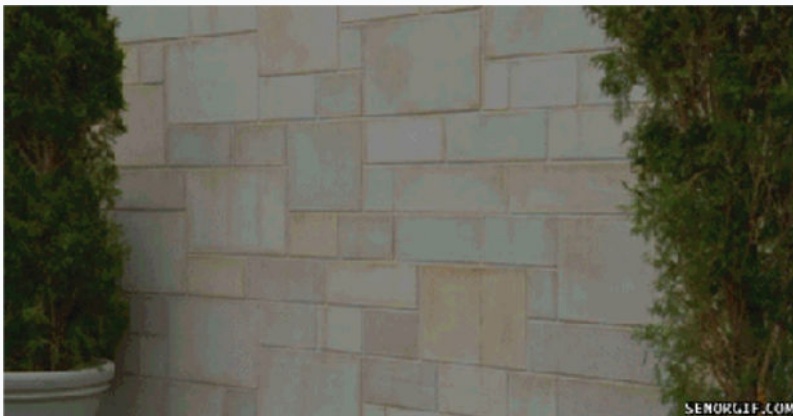
1. Reproducibility

Out of order cells make it hard to reproduce work without frequent restart-and-run-alls.

1. Performance

Re-runs are wasteful and time-consuming... especially in Hex :(

Another barrier to entry



This is exactly the kind of thing that scares people off from analytics and data science, and gives code a bad name.

The state of state



imgflip.com



imgflip.com

Re-thinking state

reactive programming

a programming paradigm oriented around data flows and the propagation of change.

In practice, this means that reactive objects maintain references to their dependencies and update automatically when their dependencies change.

Why reactive programming?

- State consistency
- Performance
- Nice abstractions for async and concurrent data flows

Imperative

```
>> a = 4
>> b = 10
>> c = a + b
>> c
14
>> a = 25
>> c
14
```

Reactive

```
>> a = 4
>> b = 10
>> c = a + b
>> c
14
>> a = 25
>> c
35
```


FILE HOME INSERT PAGE LAYOUT **FORMULAS** DATA REVIEW VIEW KUTOOLS ENTERPRISE

fx Insert Function ∑ AutoSum Used ★ Recently Used Financial Logical ? Text A Date & Time 🔍 Lookup & Reference θ Math & Trig ⋮ More Functions

Name Manager Define Name Use in Formula Create from Selection **Trace Precedents** Trace Dependents Remove Arrows

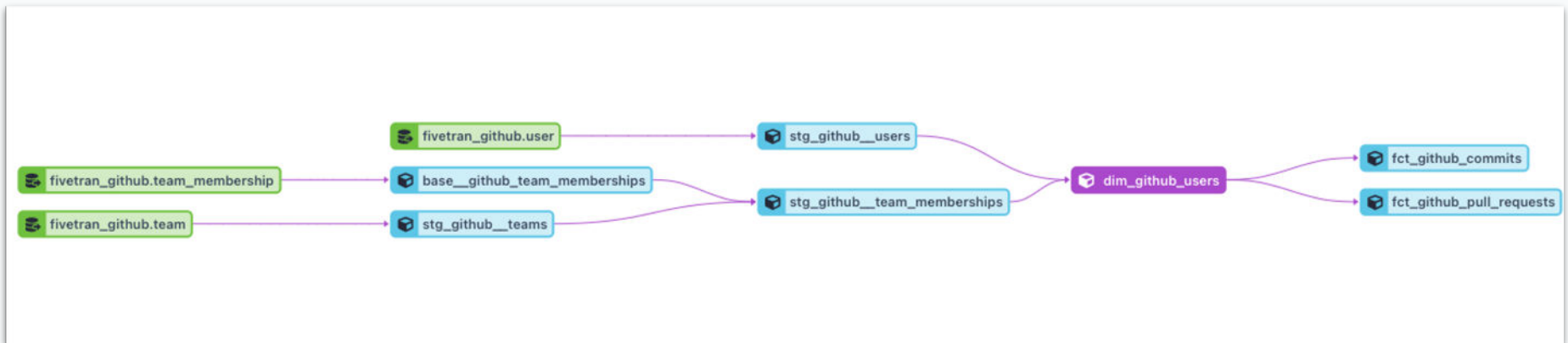
Function Library Defined Names Form

E5 : × ✓ fx =1:1048576

	A	B	C	D	E	F	G	H	I	J
1										
2		23								
3	40	0.88	95.04							
4	45									
5	50	2.4	108		0					
6										
7	100	230	456							
8		0.9								
9		291.111111								
10	111.111111		659.04							

Everyone's favorite reactive programming tool

DAGs!

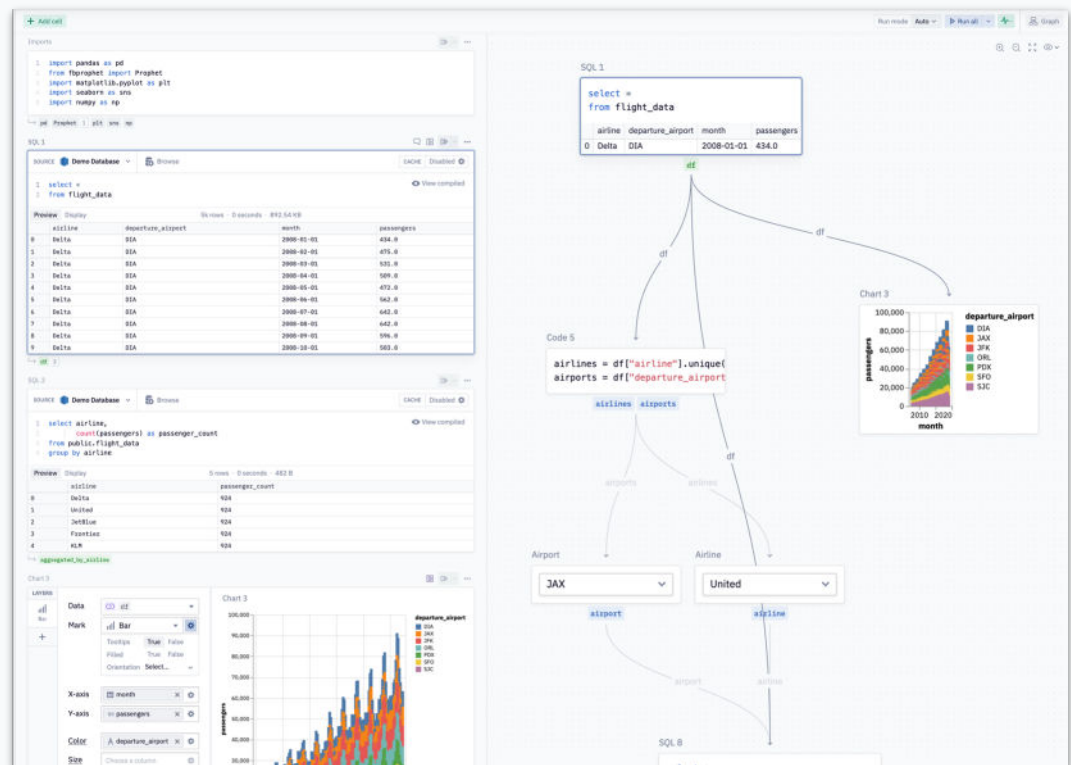


a DAG in dbt

Bringing reactivity and DAGs to notebooks

We introduced a **fully-reactive, DAG-based execution model** in Hex 2.0, which solves for all 3 problems we discussed earlier:

- Interpretability
- Reproducibility
- Performance



Demo

Production Demo Internal

Flights Demo - Reactivity

This forecast takes in historic flight volumes, and generates a prediction going forward some number of months into the future.

Imports

```
1 import pandas as pd
2 from fbprophet import Prophet
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5 import numpy as np
```

pd Prophet 1 plt sns np

SQL 1

SOURCE Demo Database Browse CACHE Disabled

```
1 select *
2 from flight_data
```

Preview Display 5k rows · 0 seconds · 892.54 KB

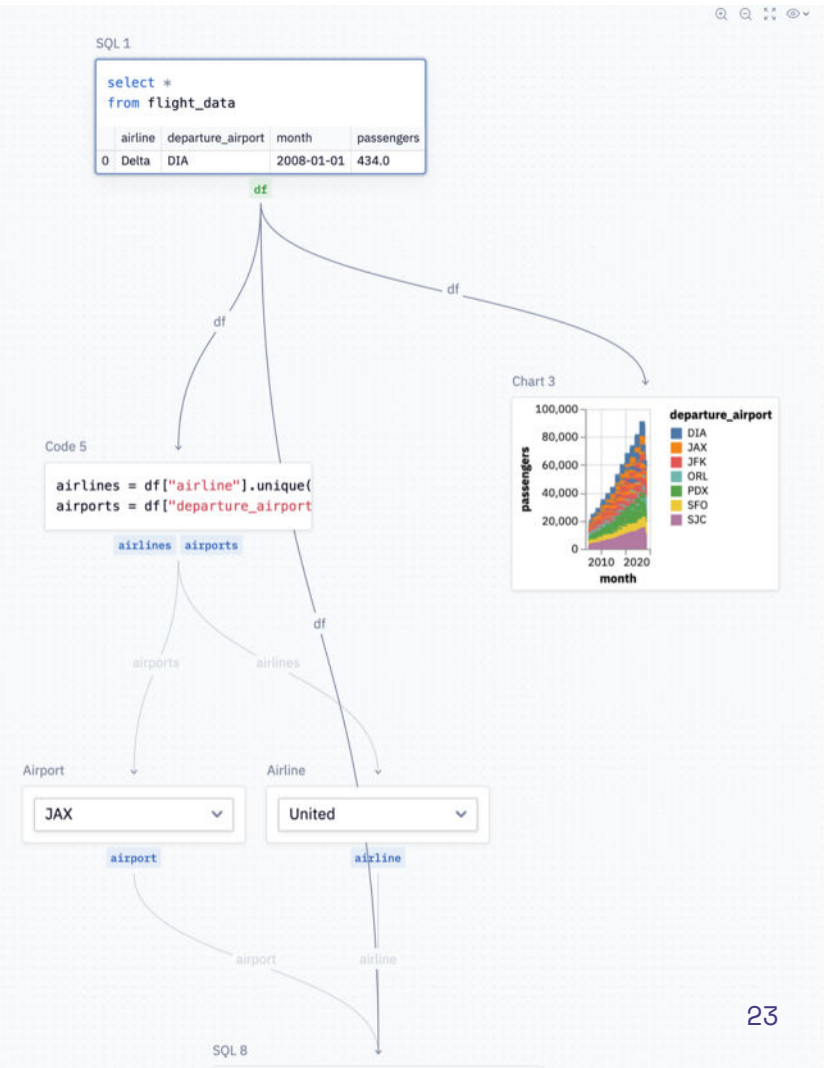
	airline	departure_airport	month	passengers
0	Delta	DIA	2008-01-01	434.0
1	Delta	DIA	2008-02-01	475.0
2	Delta	DIA	2008-03-01	531.0
3	Delta	DIA	2008-04-01	509.0

Under the hood: building the DAGs

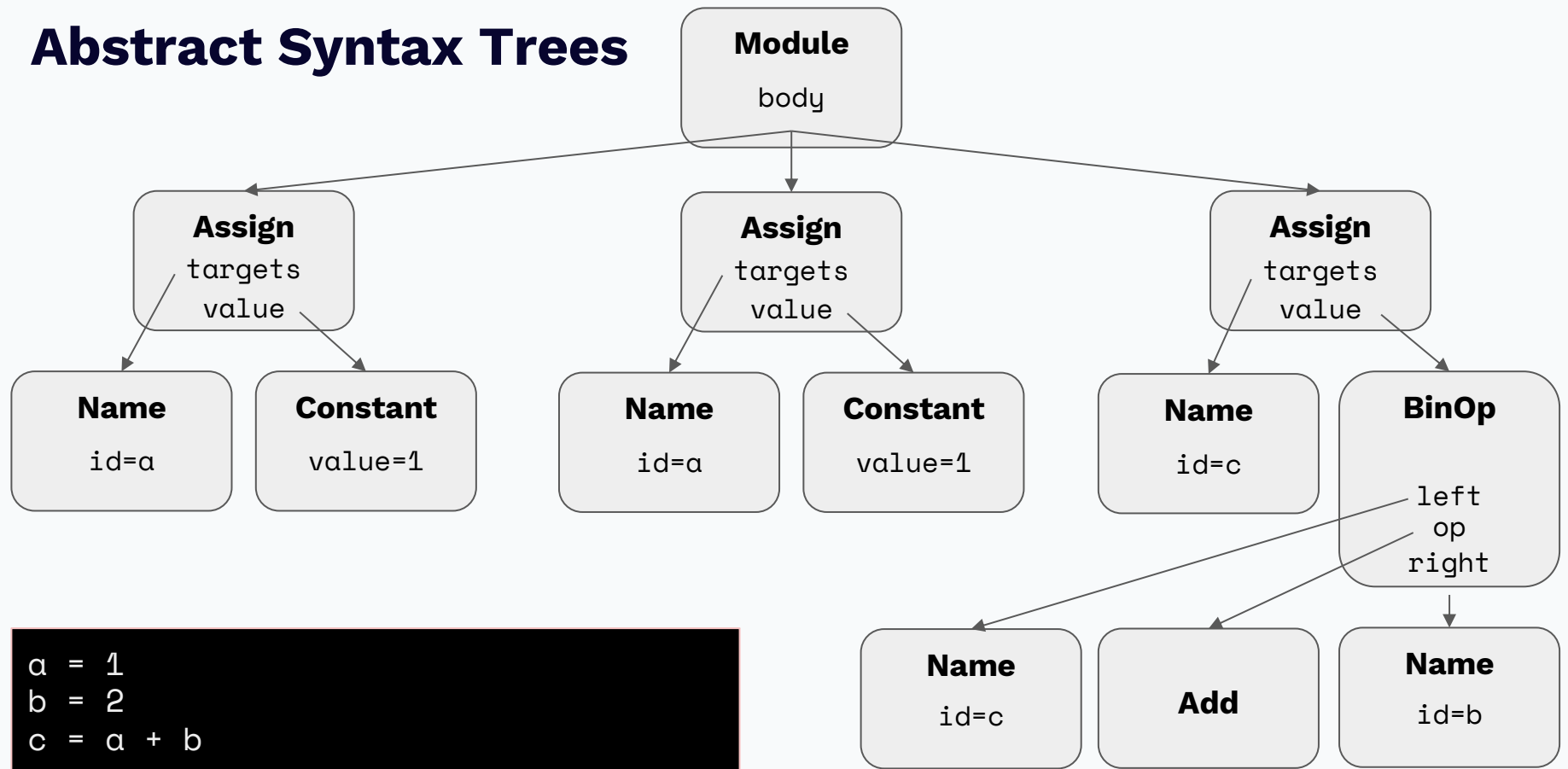
Graphs have **Nodes** and **Edges**:

- Nodes = Cells
- In edges: Variable references
- Out edges: Variable assignments

How do we determine relationships?



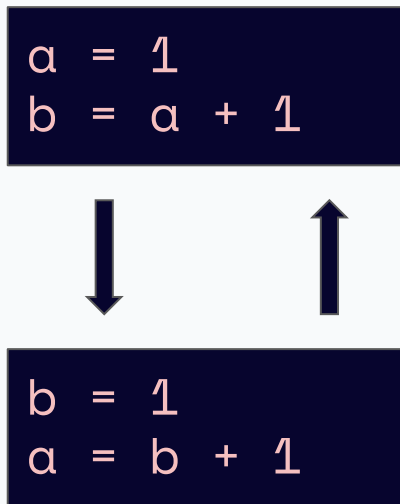
Abstract Syntax Trees



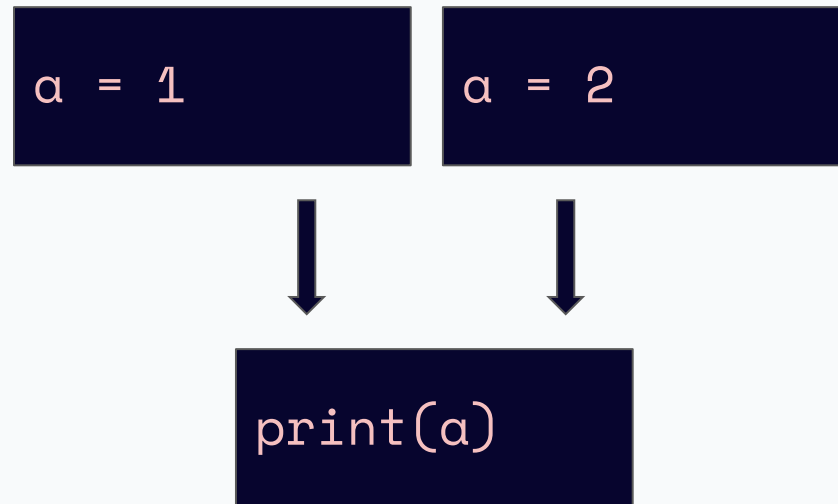
```
a = 1
b = 2
c = a + b
```

Issues with this approach

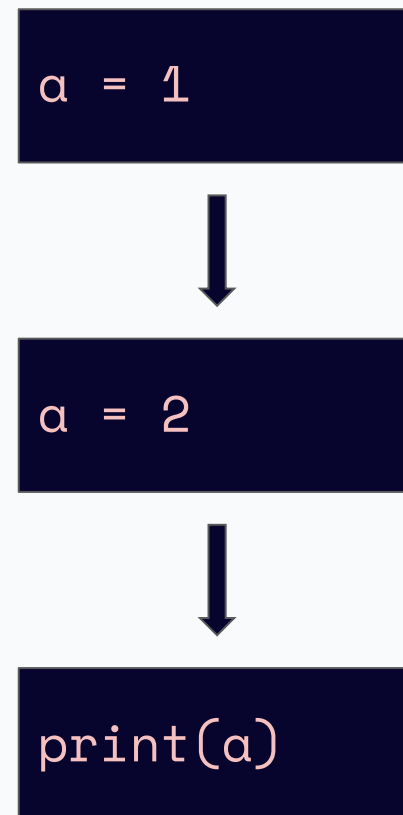
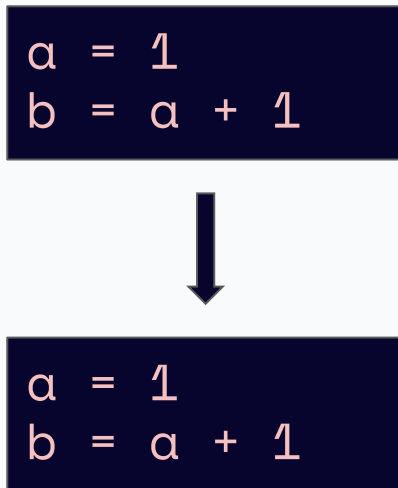
It's not actually a DAG!



The ordering is non-deterministic



Solution: use notebook ordering



**Pulling it all together:
bringing DAGs into Hex
notebooks**

Determining “staleness”

In order to know which cells to recompute, we track a condition called *staleness*.

A cell is *stale* if:

- It hasn't been run yet this kernel session
- An upstream cell has been **edited** and it hasn't been re-run
- An upstream cell has been **run** and it hasn't been re-run
- An upstream cell has **become stale**

Implementing Reactivity with iPython

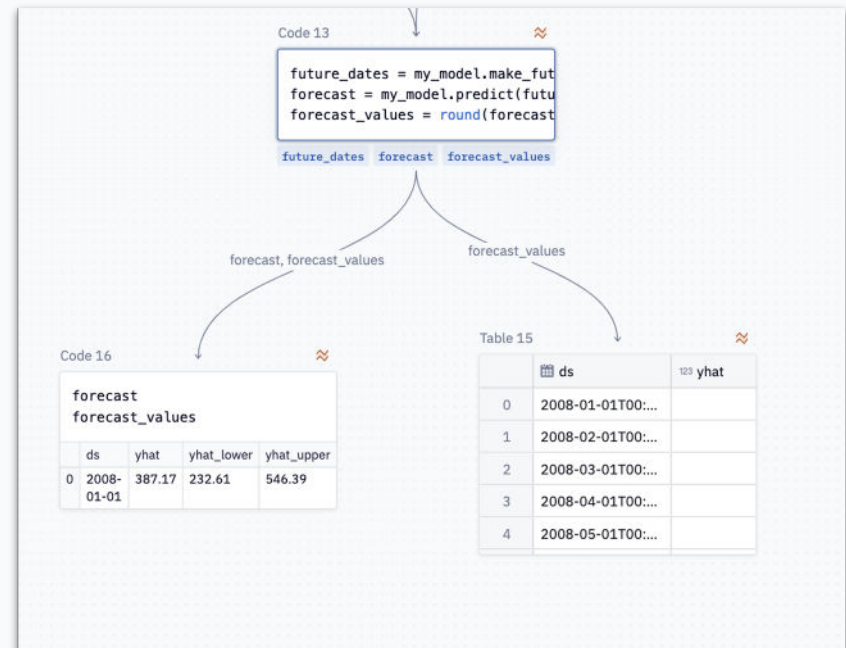
On each edit:

- Run each cell through an AST parser to compute inputs and outputs
- Re-compute the cell DAG
- Traverse graph upstream **and** downstream to determine list of cells needed to be run
 - Upstream, filter out cells that are already “up to date”
 - Downstream, mark as “stale”
- Queue all remaining stale cells in notebook order into the kernel
 - Mark cell as “up to date” after successful run

DAG usability cleanup

```
Code 0  
1 import pandas as pd  
2 from fbprophet import Prophet  
3 import matplotlib.pyplot as plt  
4 import seaborn as sns  
5 sns.set()
```

```
Markdown 1  
# Flight Traffic Forecast  
  
Flight Traffic Forecast
```



Future exploration

Future exploration

- Lambdas / better isolation
- Cell caching
- Performance & parallelism



Adam Storr
Design Lead



Melissa Carlson
Engineering Lead



Glen Takahashi
Chief Architect

Questions?