# Beyond Linearity

## Building reactive notebooks for data

Caitlin Colgrove, CTO @ Hex

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";
}
```
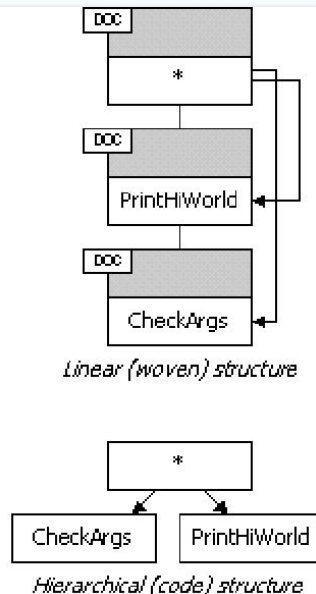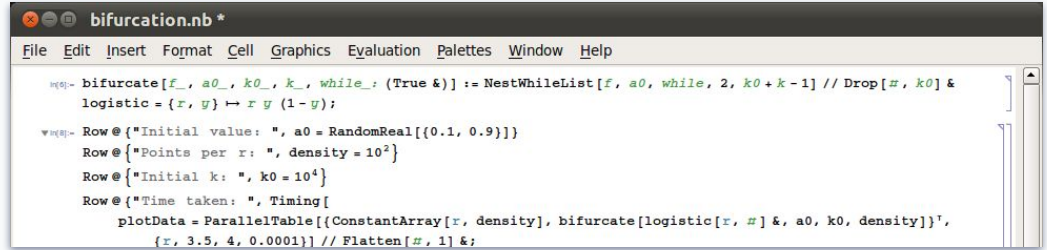
(a) Literate source.

(b) Linear and hierarchical views.

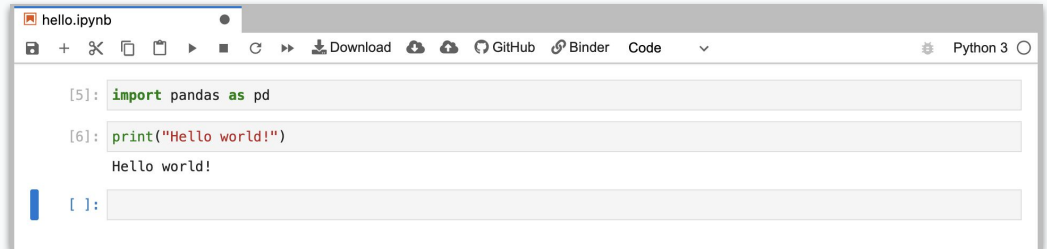*Linear (woven) structure*

*Hierarchical (code) structure*

# Fast forward to 2022

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

# 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.

2. **Reproducibility**

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

3. **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

DATA DEPARTMENT

HAVE YOU TRIED RESTARTING AND RUNNING FROM SCRATCH

# 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
```
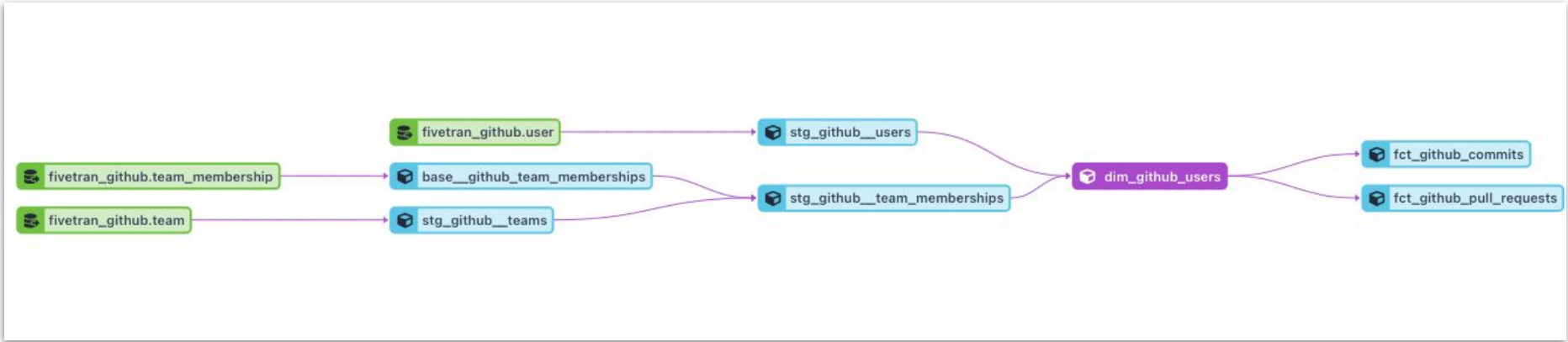
*Everyone's favorite reactive programming tool*

# 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

HEX | Flights Demo - Reactivity | </> Logic | [ ] | [ ] App

Caitlin Colgrove
Hex

Publish | Share

+ Add cell

Run mode | Auto ∨ | ▷ Run all ∨ | 🔲 Graph

○ 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
```

👁 View compiled

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 |

0:02

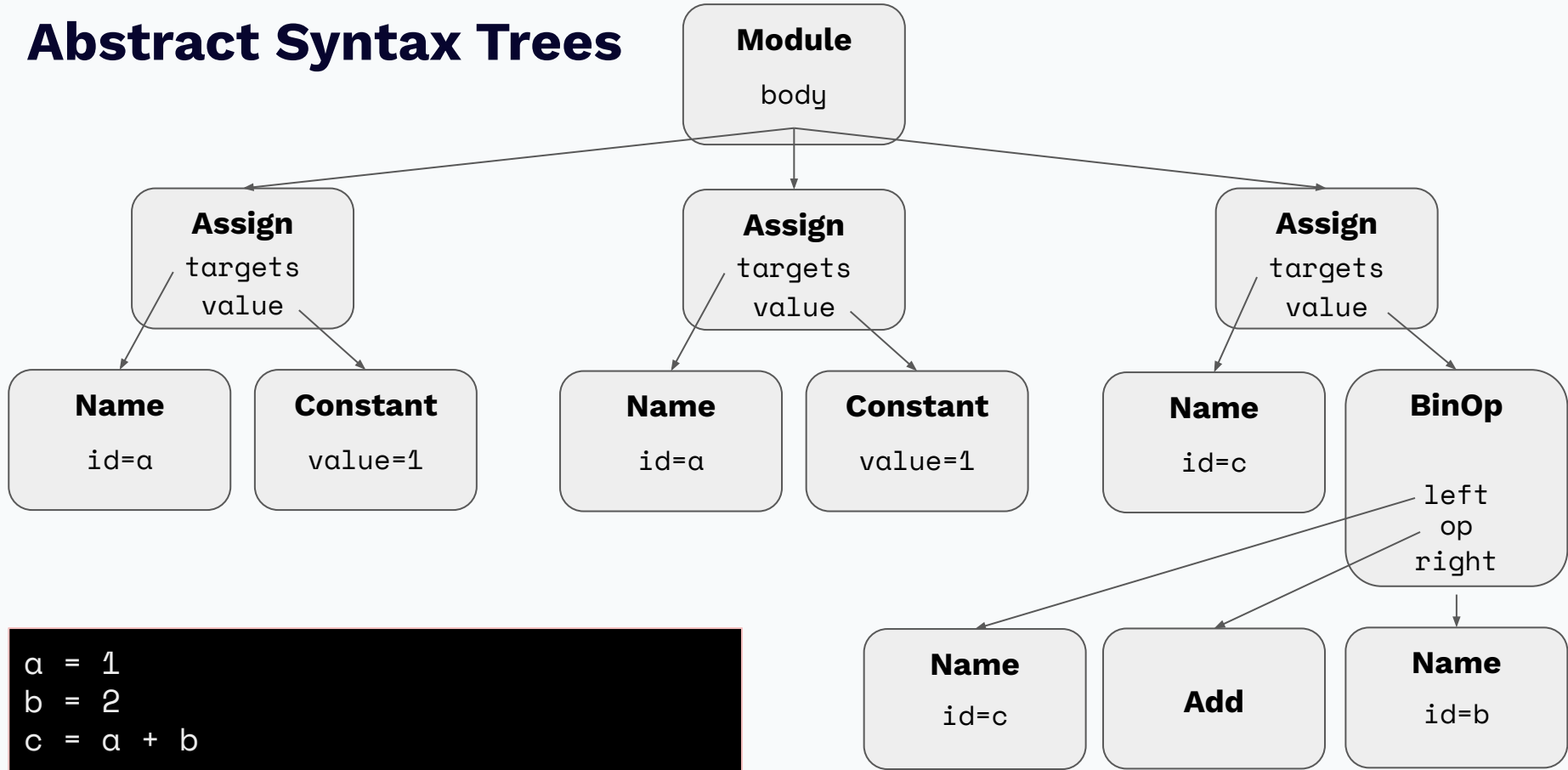# Under the hood: building the DAGs

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!

```
a = 1
b = a + 1
```

↓ ↑

```
b = 1
a = b + 1
```

The ordering is non-deterministic

```
a = 1
```

```
a = 2
```

↓ ↓

```
print(a)
```

# Solution: use notebook ordering

```
a = 1
b = a + 1
```

↓

```
a = 1
b = a + 1
```

```
a = 1
```

↓

```
a = 2
```

↓

```
print(a)
```

# 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

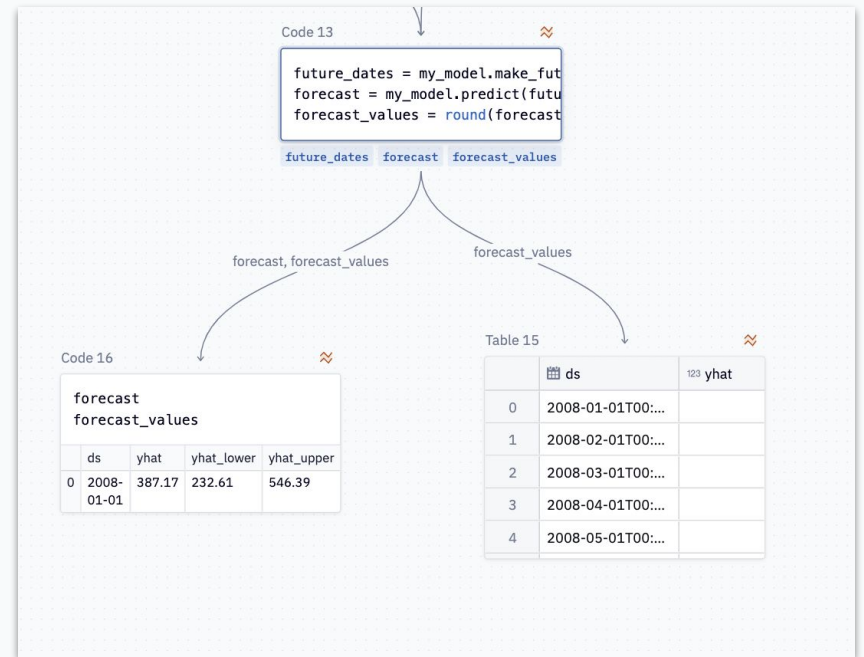    - Mark cell as "up to date" after successful run

# DAG usability cleanup

# Future exploration

# Future exploration

- Lambdas / better isolation

- Cell caching

- Performance & parallelism

Adam Storr
Design Lead

Melissa Carlson
Engineering Lead

Glen Takahashi
Chief Architect

# Interested?

Director, Platform Engineering
Backend Engineer
Cloud Engineer
Platform PM
Engineering Lead
... and many more

hex.tech/jobs

# Questions?