Kaituo Li
U. Massachusetts, Amherst

Christoph Reichenbach
Goethe U. Frankfurt

Yannis Smaragdakis
U. Athens

Yanlei Diao
U. Massachusetts, Amherst

Christoph Csallner
U. Texas Arlington

# SEDGE: Symbolic Example Data Generation for Dataflow Programs

# Motivation

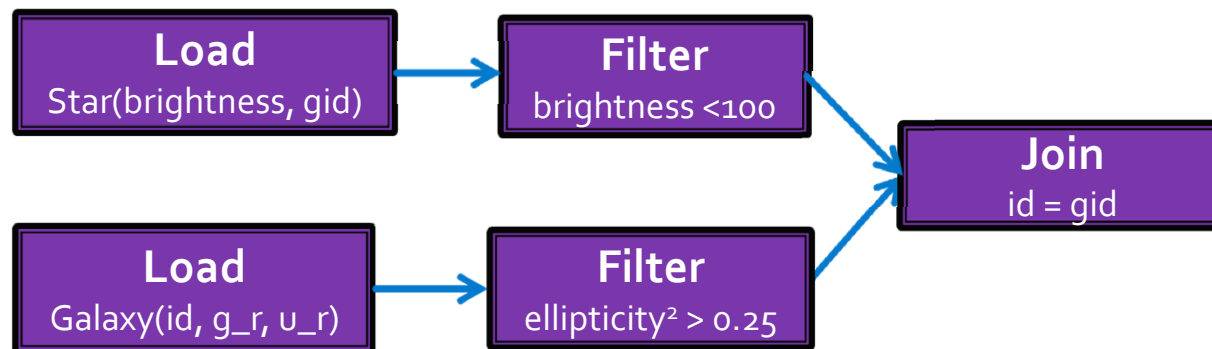# SEDGE: Symbolic Example Data GEnerator

- **Problem**: Generate **fewest** test cases to exercise **all** key behaviors of operators in a dataflow program
  - e.g. both passing and failing a filter

- **Our approach**: First dynamic-symbolic (aka "concolic") testing engine for a dataflow language

- **Results**:
  - Improved coverage and running time over industrial state-of-the-art
    - e.g. Pig Latin "illustrate", SIGMOD '09 best paper award

# Overview of Talk

- Background: Dataflow languages
- Metrics
- Our algorithm in action
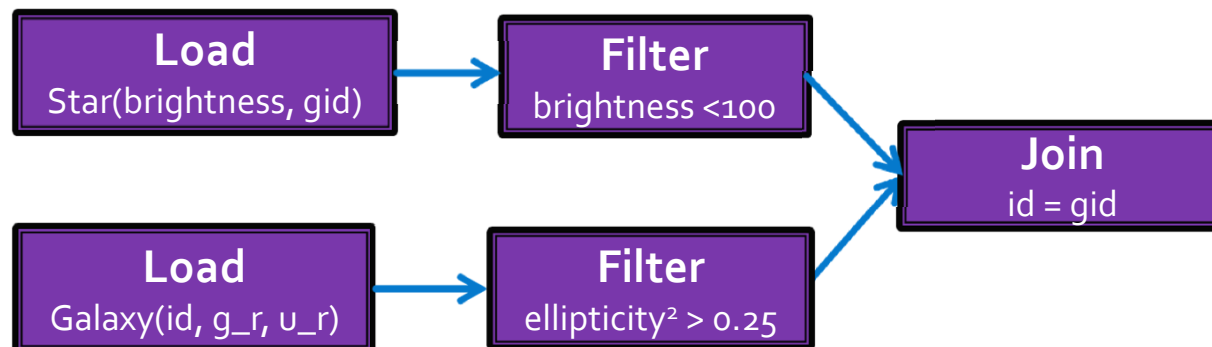- Comparison with state-of-the-art

# Example Dataflow Program

- Input: Galaxy profiles, Star profiles
- Find stars with a surface brightness less than 100 from galaxies with a squared ellipticity > 0.25
  - ellipticity$^2$ = g_r$^2$ + u_r$^2$



**Load**
Star(brightness, gid)

**Filter**
brightness <100

**Load**
Galaxy(id, g_r, u_r)

**Filter**
ellipticity$^2$ > 0.25

**Join**
id = gid

# In Pig Latin

```
A=LOAD 'Star' using PigStorage  AS (brightness:int, gid:int);
B=LOAD 'Galaxy' using PigStorage AS (id: int, g_r, u_r:double);
C=FILTER A BY brightness < 100;
D=FILTER B BY power(g_r, 2) + power(u_r, 2) > 0.25;
E=JOIN C ON gid, D ON id;
```
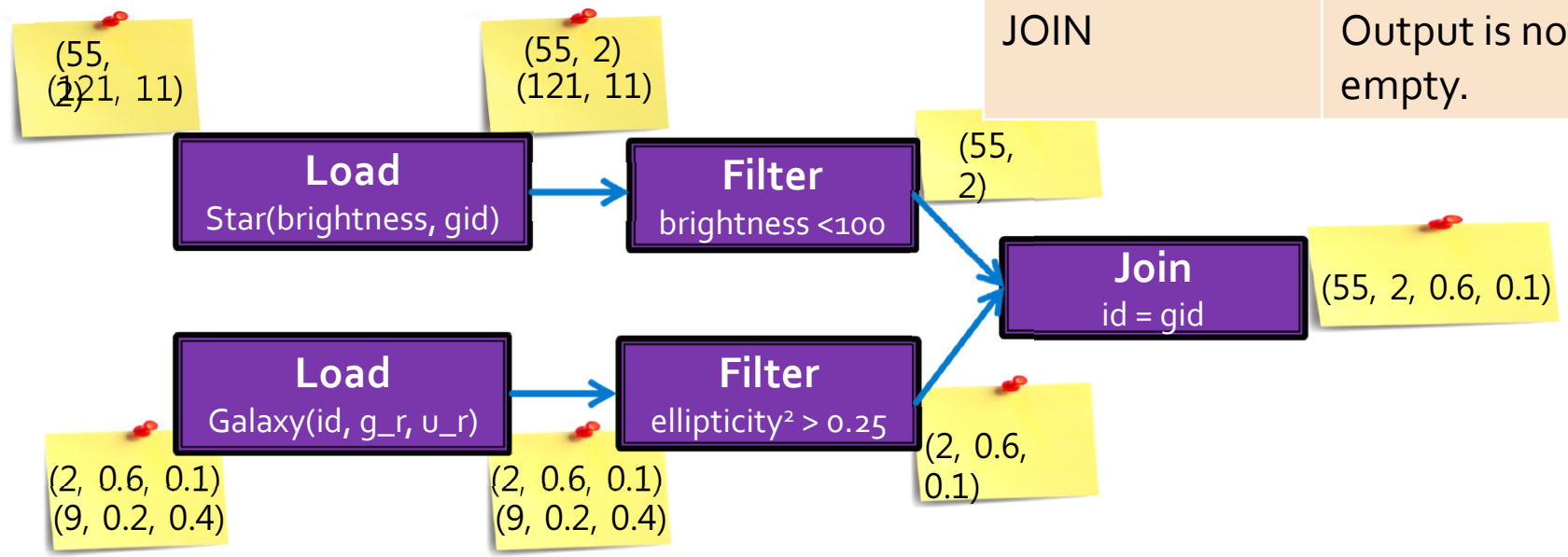
# Goal: Coverage with Fewest Tests

- *Completeness*
    - similar to "branch coverage" in imperative languages
- *Conciseness*
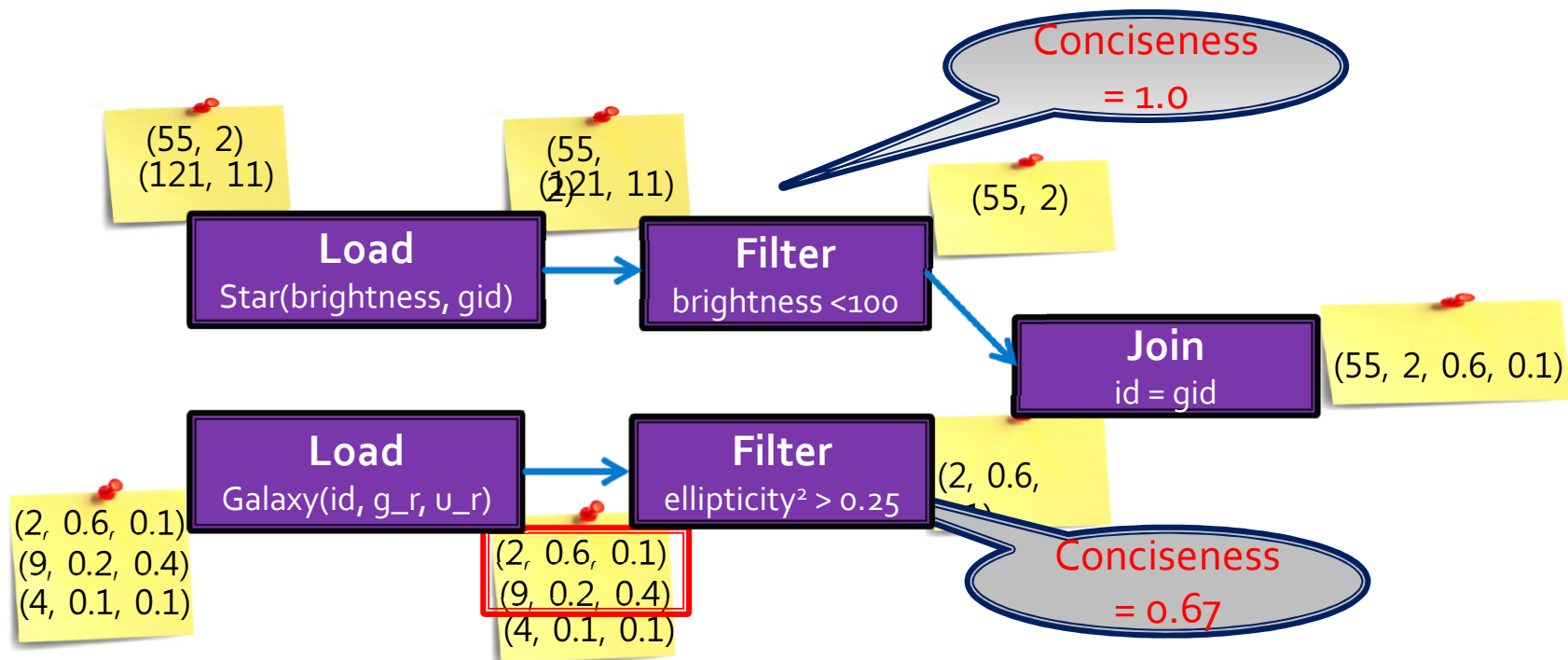    - with as few tests as possible

# But What Is Coverage?

| Operator | Coverage |
|----------|----------|
| LOAD | Input is not empty. |
| FILTER | **(a)** one record that passes the filter; **(b)** one that does not pass. |
| JOIN | Output is not empty. |

(55, 2)
(121, 11)

(55, 2)
(121, 11)

(55, 2)

**Load**
Star(brightness, gid)

**Filter**
brightness <100

**Join**
id = gid

(55, 2, 0.6, 0.1)

**Load**
Galaxy(id, g_r, υ_r)

**Filter**
ellipticity$^2$ > 0.25

(2, 0.6, 0.1)

(2, 0.6, 0.1)
(9, 0.2, 0.4)

(2, 0.6, 0.1)
(9, 0.2, 0.4)

# Our Algorithm

Dynamic execution

Pruning

Symbolic synthesis

# Our Algorithm

**Dynamic execution**

Pruning

Symbolic synthesis

Run sampled existing data

Cache concrete results of user-defined functions (UDF) across runs (if any)

# Dynamic Execution

| UDF | Parameters | Return |
|-----|-----------|--------|
| power | (0.6, 2) | 0.36 |
| power | (0.1, 2) | 0.01 |
| power | (0.2, 2) | 0.04 |
| power | (0.5, 2) | 0.25 |

(155, 2)
(121, 11)

**Load**
Star(brightness, gid)

**Filter**
brightness <100

**Join**
id = gid

**Load**
Galaxy(id, g_r, υ_r)

**Filter**
ellipticity² > 0.25

(4, 0.6, 0.1)
(9, 0.2, 0.5)

(4, 0.6, 0.1)
(9, 0.2, 0.5)
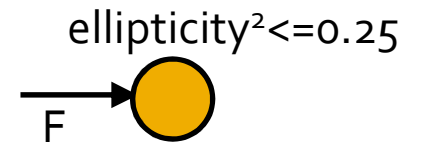
# Our Algorithm

Dynamic execution

**Pruning**

Eliminate redundancy

Symbolic synthesis
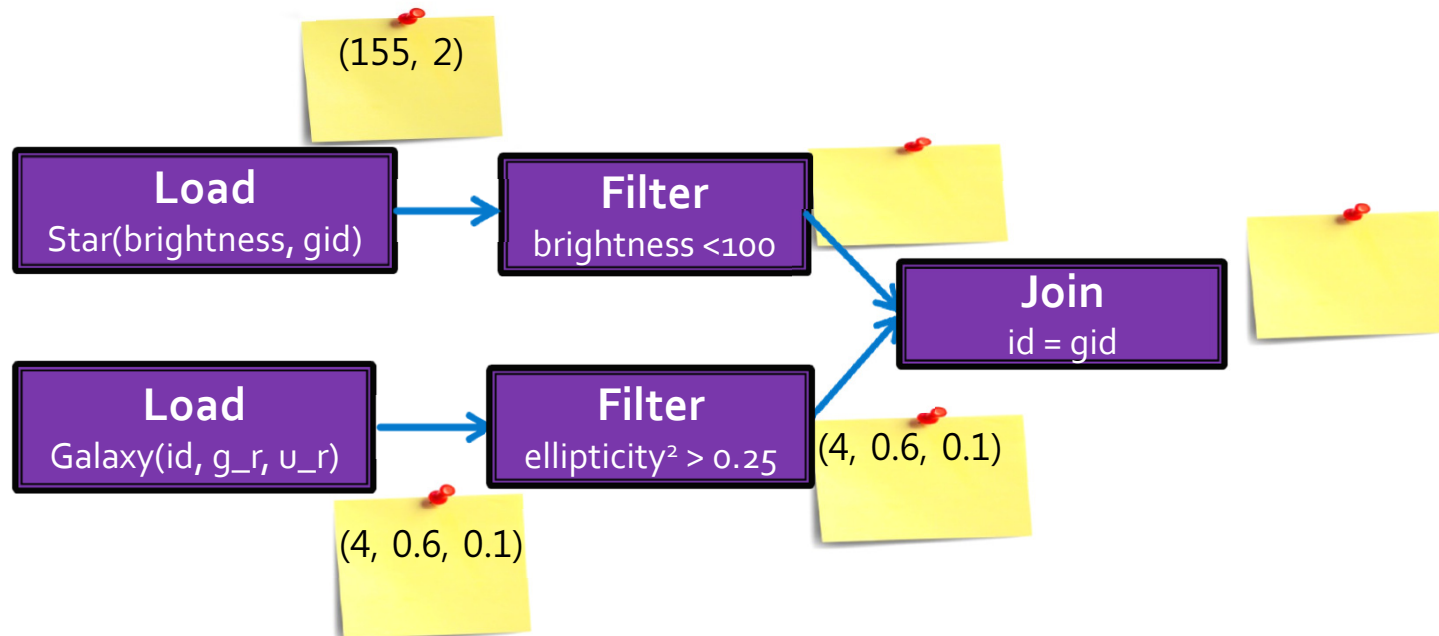
# Backward Symbolic Execution

gid==id && id==4

P

(155, 2)

**Load**
Star(brightness, gid)

**Filter**
brightness <100

**Join**
id = gid

**Load**
Galaxy(id, g_r, υ_r)

**Filter**
ellipticity² > 0.25

(4, 0.6, 0.1)

(4, 0.6, 0.1)

# Constraint Generation

P1: power(g_r, 2) + power(u_r, 2) <=0.25

$gid==id$ && $id==4$

P

$ellipticity^2<=0.25$

F

(155, 2)

**Load**
Star(brightness, gid)

**Filter**
brightness <100

**Join**
id = gid

**Load**
Galaxy(id, g_r, υ_r)

**Filter**
$ellipticity^2 > 0.25$
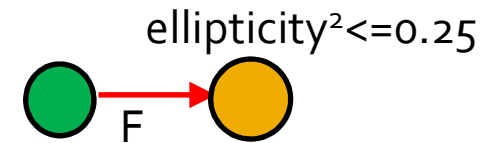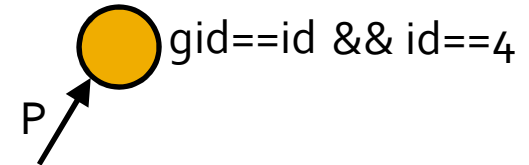
(4, 0.6, 0.1)

(4, 0.6, 0.1)
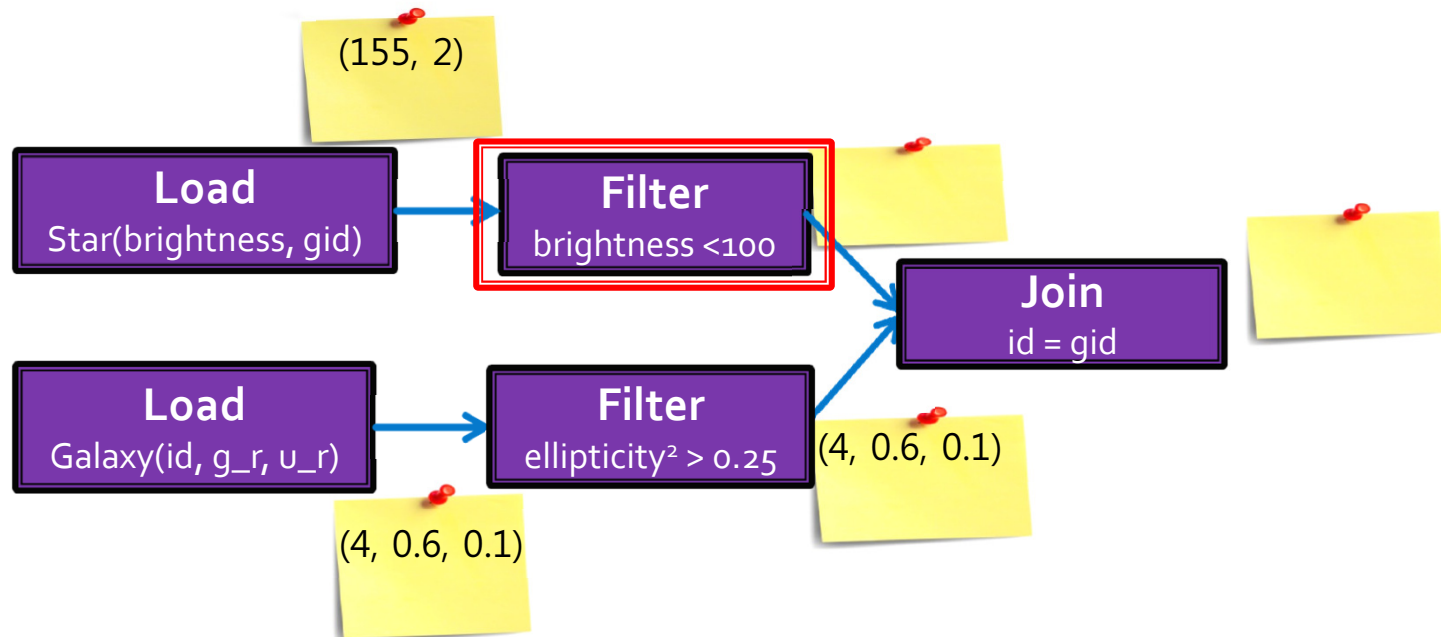
# Backward Symbolic Execution
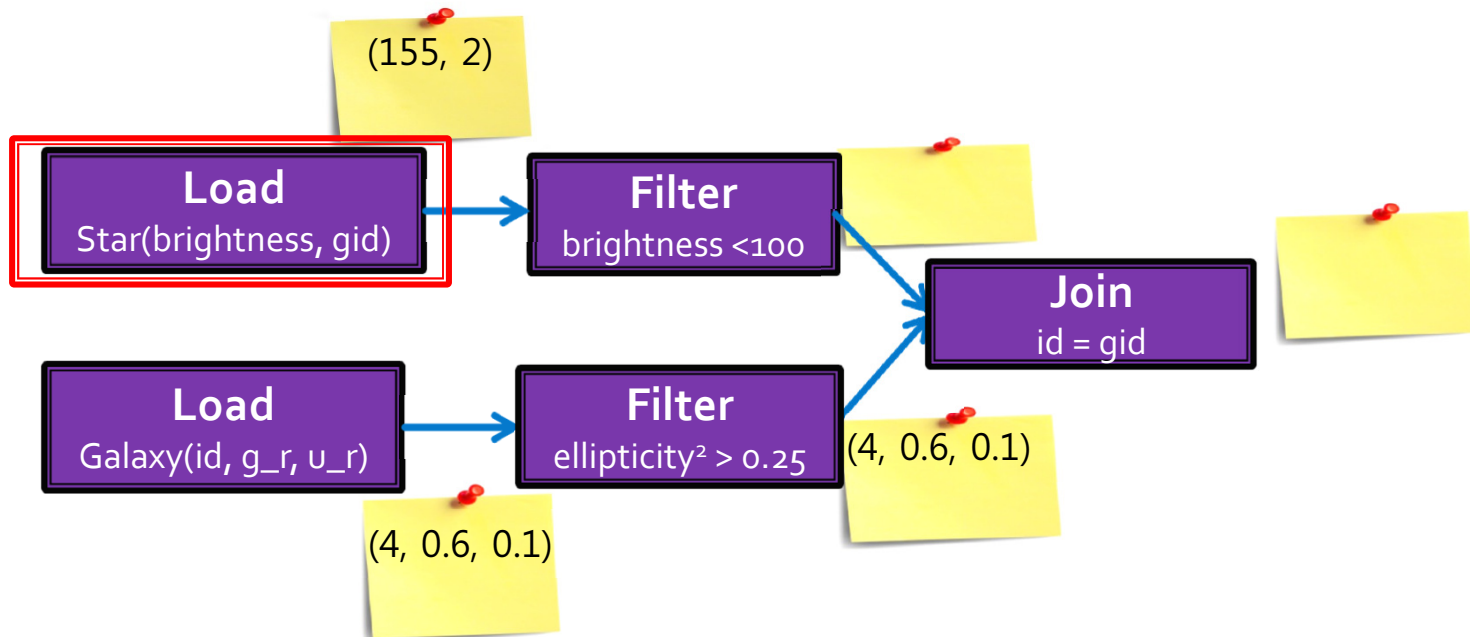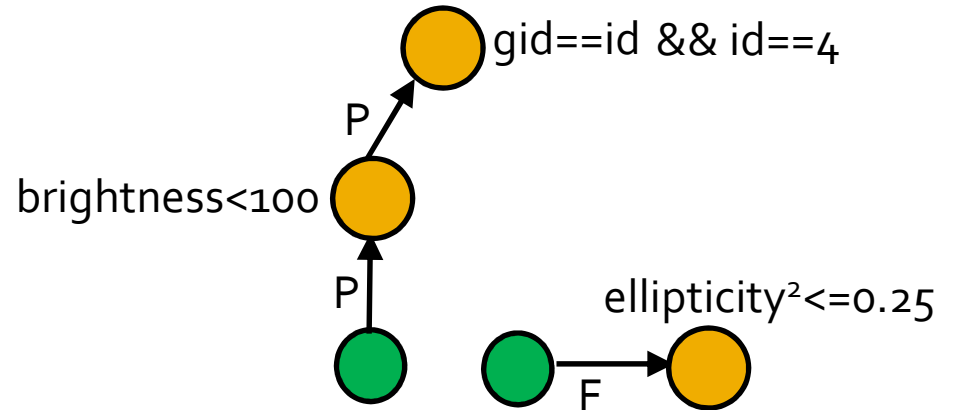


18

# Backward Symbolic Execution

# Constraint Generation



P2: gid==id && id==4
&& brightness<100

gid==id && id==4

brightness<100

P

P

ellipticity²<=0.25

F

(155, 2)

| Load | Filter | Join |
|------|--------|------|
| Star(brightness, gid) | brightness <100 | id = gid |
| Load | Filter | |
| Galaxy(id, g_r, υ_r) | ellipticity² > 0.25 | |

(4, 0.6, 0.1)

(4, 0.6, 0.1)

20

# Concretization

P1: power(g_r, 2) + power(u_r, 2) <=0.25

P2: gid==id && id==4 && brightness<100

(155, 2)

**Load**
Star(brightness, gid)

**Filter**
brightness <100

**Load**
Galaxy(id, g_r, u_r)

**Filter**
ellipticity$^2$ > 0.25

**Join**
id = gid

(4, 0.6, 0.1)

(4, 0.6, 0.1)

# Concretization

*Model the behavior of a user-defined function as its input-output behavior across all observed test cases.*

| UDF | Parameters | Return |
|-----|-----------|--------|
| power | (0.6, 2) | 0.36 |
| power | (0.1, 2) | 0.01 |
| power | (0.2, 2) | 0.04 |
| power | (0.5, 2) | 0.25 |

```
(define-fun power ((u Double) (v Int)) Double
  (ite (and (= u 0.6) (= v 2) 0.36
  (ite (and (= u 0.1) (= v 2) 0.01
  (ite (and (= u 0.2) (= v 2)  0.04
  (ite (and (= u 0.5) (= v 2) 0.25
   0)))
    (declare-const x Double)
    (declare-const y Int)
    (assert (not (= (power x y) 0)))
```
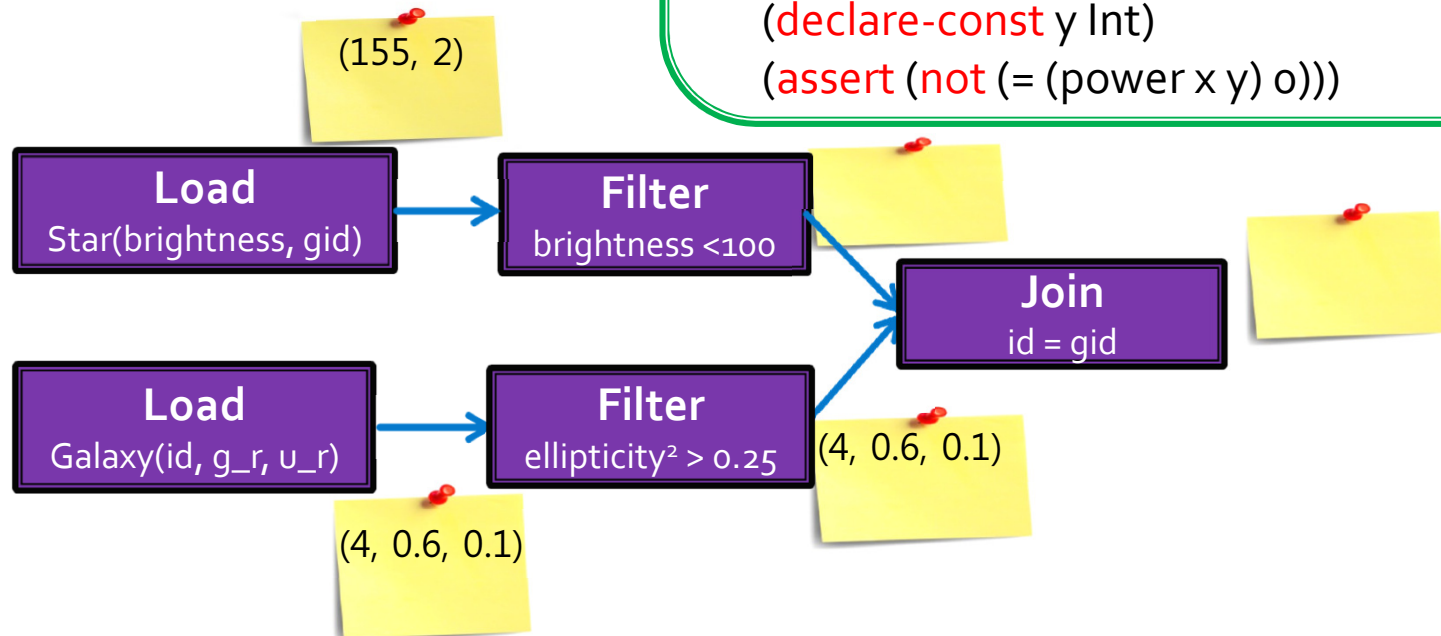
## Assumption

*UDFs return values only depend on argument value(s) in dataflow programs.*

# Constraint Solving

**Symbolic Constraints**

**Uninterpreted Functions**

**Constraint solvers**

Z3

**Coral**
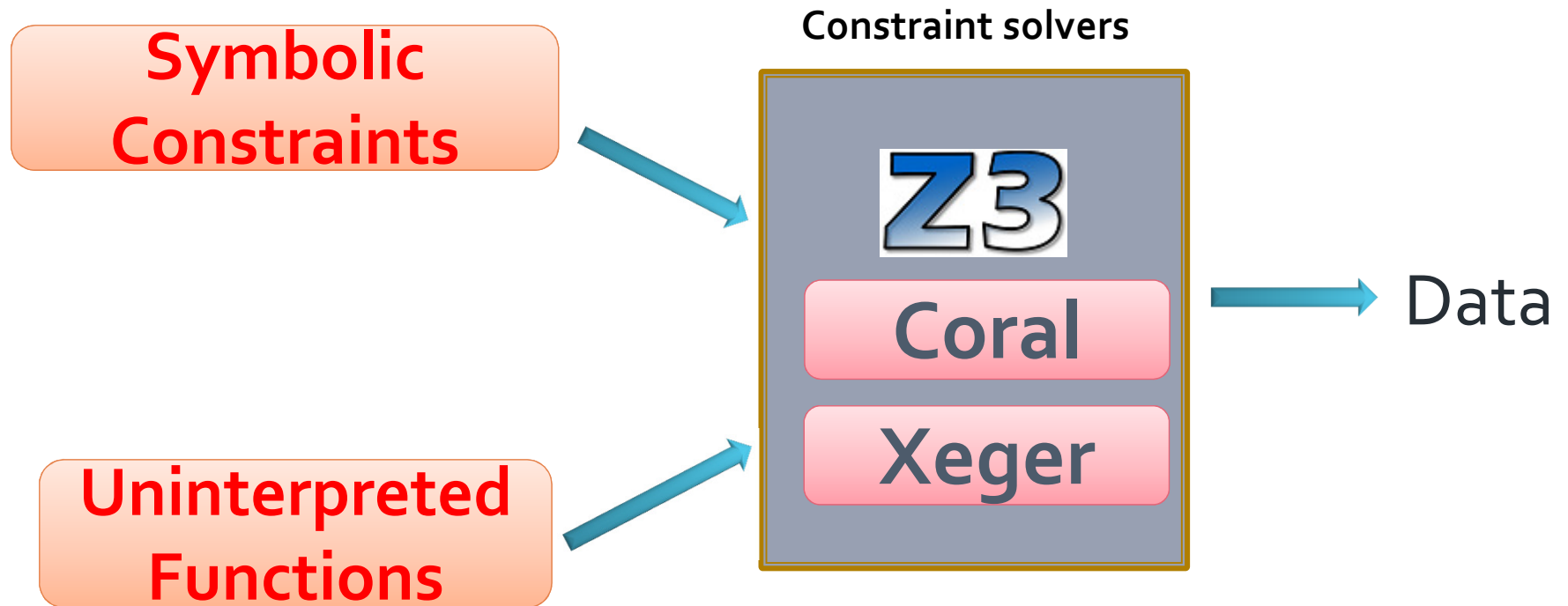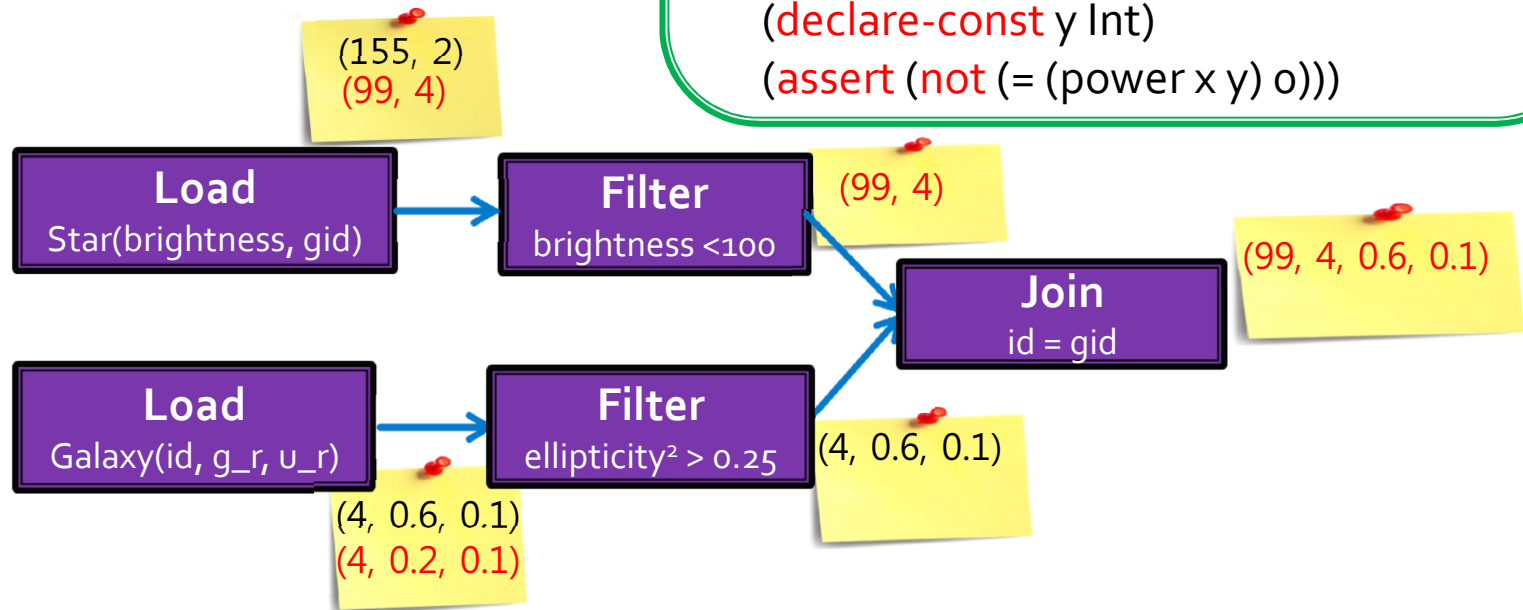
**Xeger**

Data

# Constraint Solving

P1: gid==id && id==4
&& brightness<100

P2: power(g_r, 2) + power(u_r, 2)
<=0.25

```
(define-fun power ((u Double) (v Int)) Double
  (ite (and (= u 0.6) (= v 2) 0.36
  (ite (and (= u 0.1) (= v 2) 0.01
  (ite (and (= u 0.2) (= v 2)  0.04
  (ite (and (= u 0.5) (= v 2) 0.25
   0)))
    (declare-const x Double)
    (declare-const y Int)
    (assert (not (= (power x y) 0)))
```

(155, 2)
(99, 4)

**Load**
Star(brightness, gid)

**Filter**
brightness <100

(99, 4)

**Join**
id = gid

(99, 4, 0.6, 0.1)

**Load**
Galaxy(id, g_r, u_r)

**Filter**
ellipticity$^2$ > 0.25

(4, 0.6, 0.1)

(4, 0.6, 0.1)
(4, 0.2, 0.1)

25

# Experiments

- Two benchmark suites
  - PigMix: 20 representative Pig programs from Pig community
  - SDSS: 11 Pig programs hand-translated from sample SQL queries
    - from the Sloan Digital Sky Survey (astronomy DB)

# SEDGE vs. Pig Latin Illustrate

## Pig Latin Illustrate

- ☹ Cannot handle UDFs
- ☹ Poor constraint solving ability
- ☹ Generate local constraint without looking ahead

## SEDGE

- ☺ Support UDFs
- ☺ Stronger constraint solving
- ☺ Generate inter-related constraints

Pig Latin Illustrate is the current state-of-the-art in example data generation for dataflow programs .
- Industrial tool ("illustrate" functionality in Pig, by Yahoo)
- SIGMOD '09 best paper

SEDGE can generate example data for dataflow programs better (completeness) and cheaper (running time)
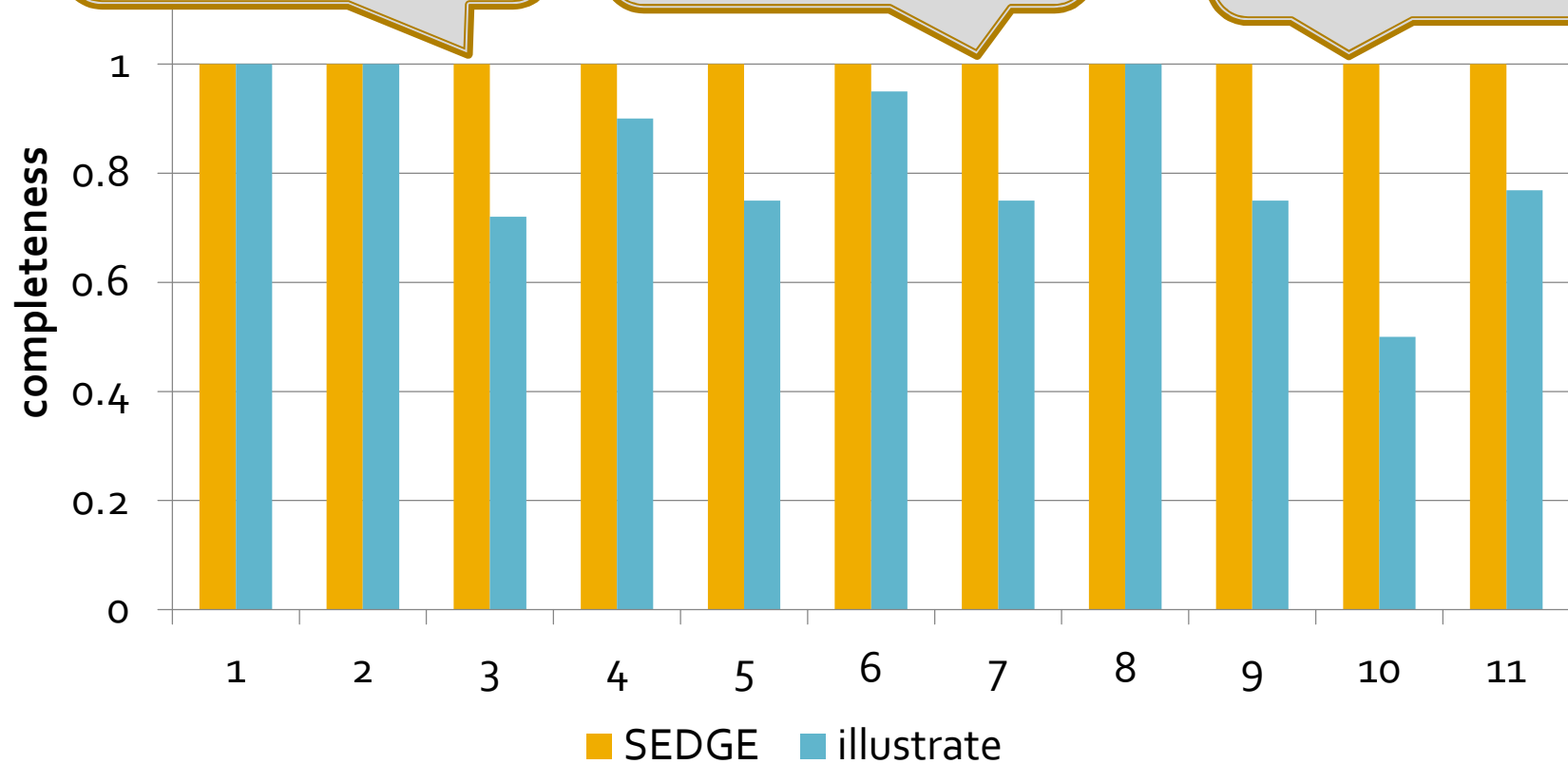
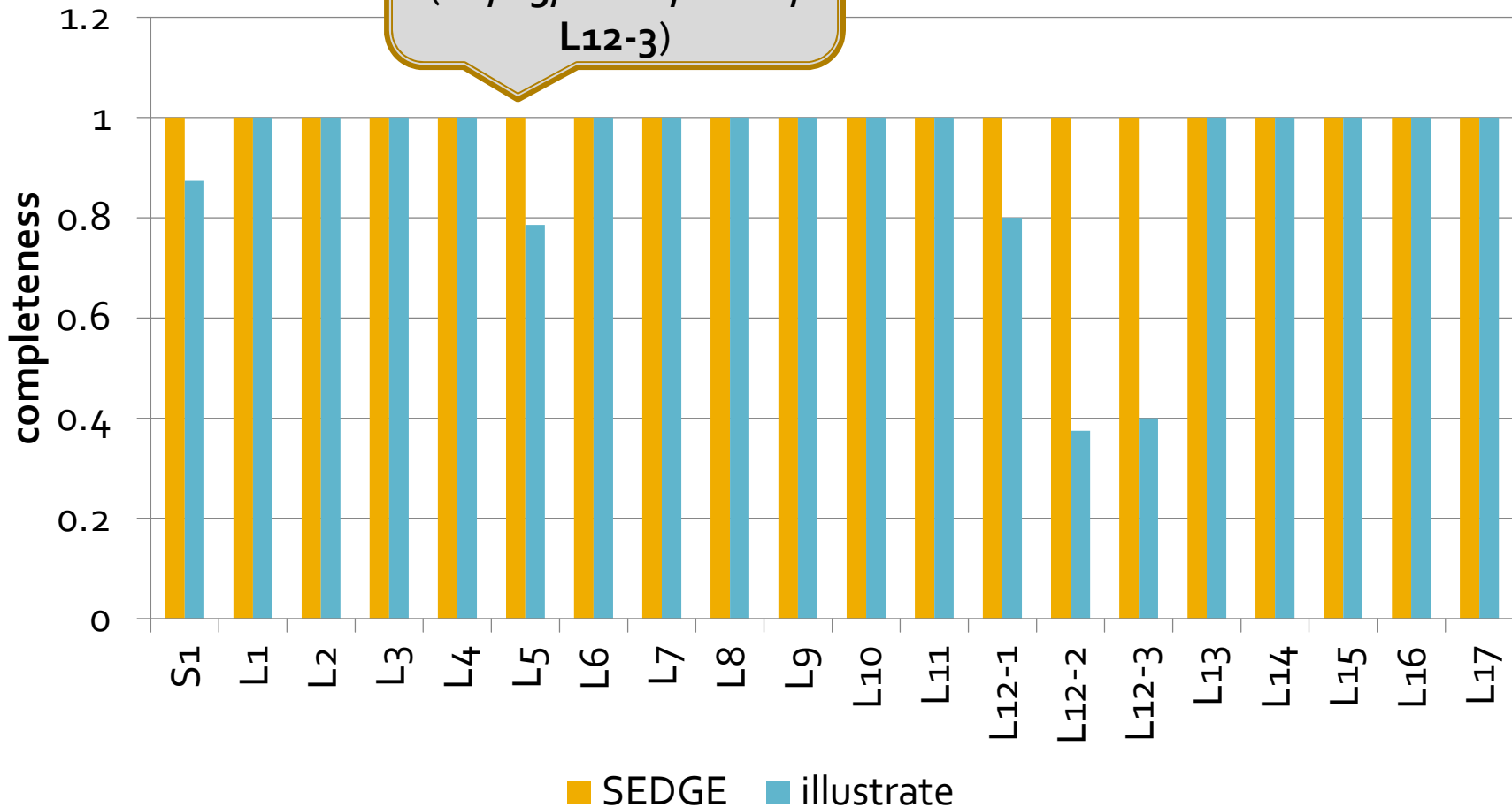Completeness Comparison Achieved for BigMix Benchmark

# Timing Comparison for SDSS Benchmark

# Timing Comparison for PigMix Benchmark

# Summary

- Created the first dynamic-symbolic (aka "concolic") testing engine for dataflow languages
- Suggested the use of concrete results across runs of a UDF to represent the UDF
- Proposed an approach that balances completeness, conciseness, and running time