

Co-dependence Aware Fuzzing for Dataflow-based Big Data Analytics

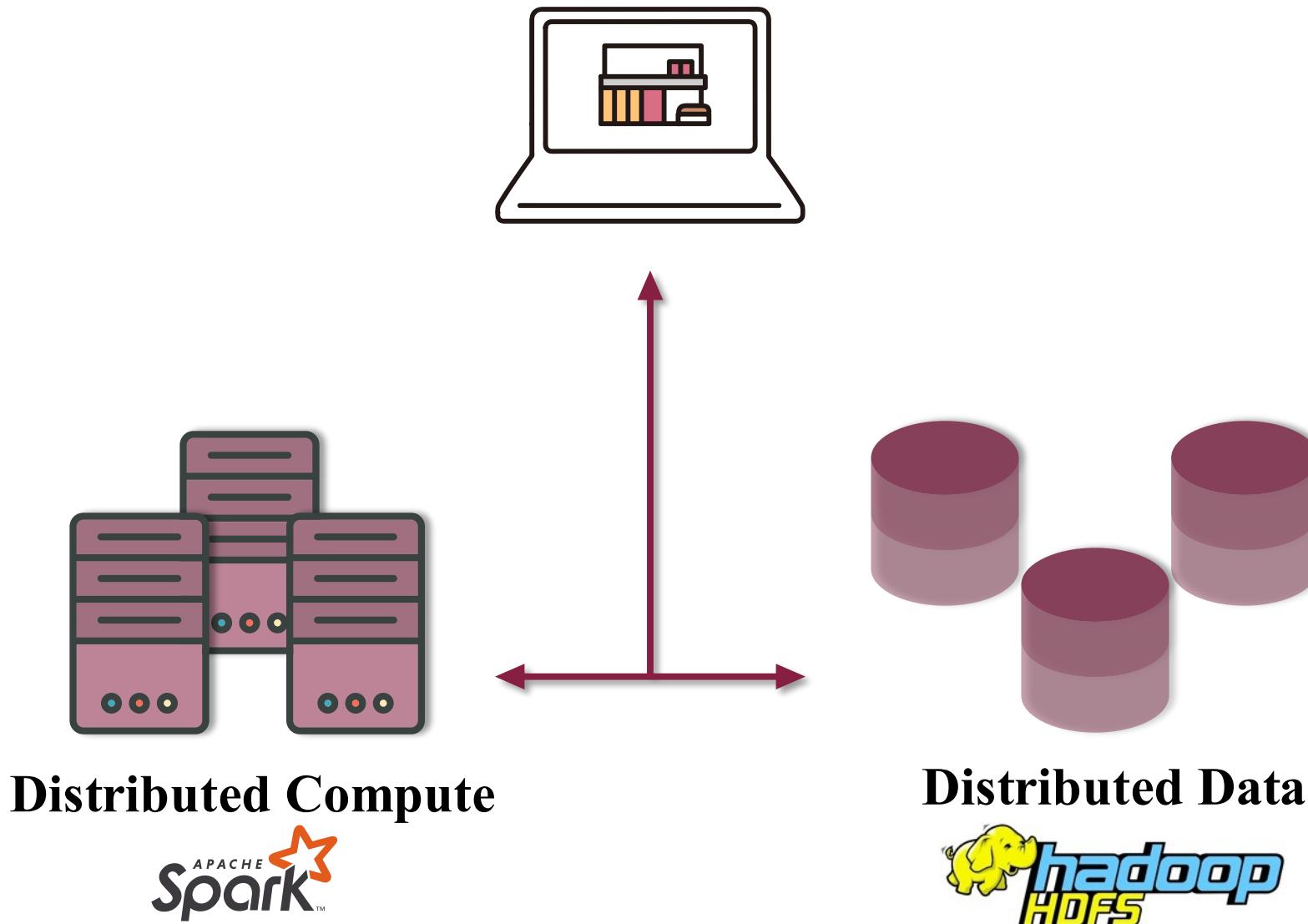
Ahmad Humayun, Miryung Kim, Muhammad Ali Gulzar



Outline

- Problem: Test input generation for Big Data (DISC) applications
 - Programs have complex interactions with data that create inter-data dependencies
 - A general solution for data dependent programs
-
- Motivating Example
 - Technical Challenges
 - Approach Overview
 - Evaluation

Data Intensive Scalable Computing (DISC)



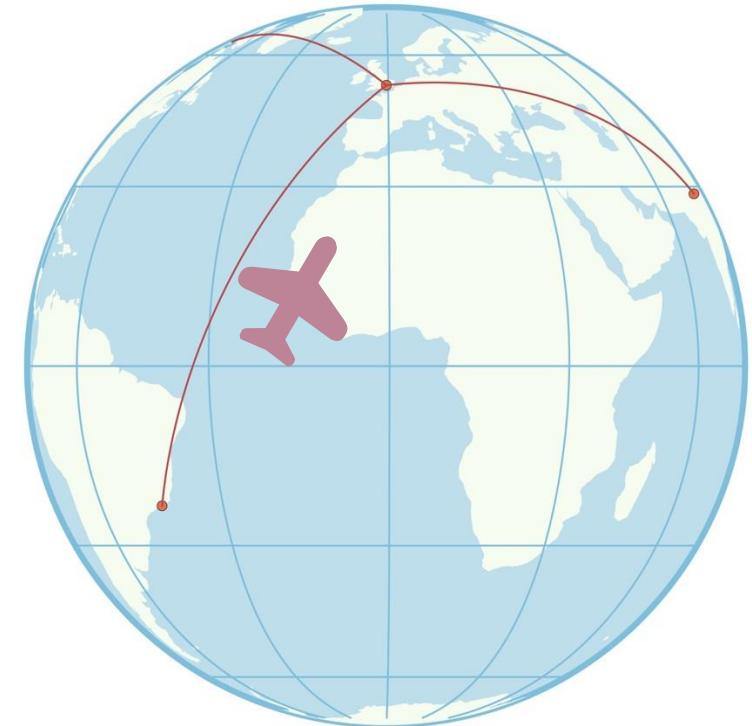
DISC Application Example

- Compute the distance travelled by a flight, between two airports



32556, 08-22 07:25, 08-22 09:50, LAX, SFO
22019, 09-03 15:00, 09-03 02:50, ROA, SFO
31522, 09-03 08:55, 09-03 10:20, IAD, SFO

flights.csv



DISC Application Example

- Compute the distance travelled by a flight, between two airports



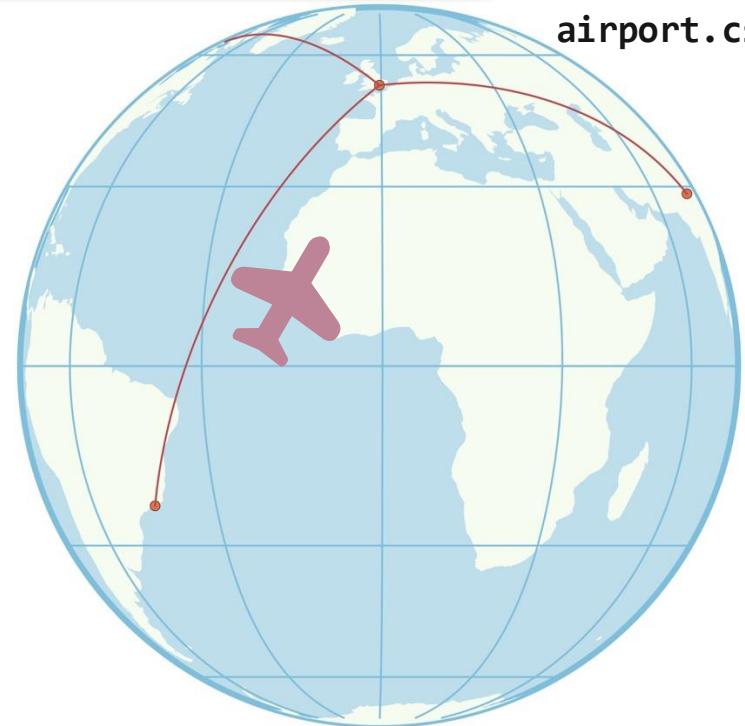
flights.csv

```
32556, 08-22 07:25, 08-22 09:50, LAX, SFO  
22019, 09-03 15:00, 09-03 02:50, ROA, SFO  
31522, 09-03 08:55, 09-03 10:20, IAD, SFO
```

```
LAX, 33.94° N, 118.41° W, Los Angeles  
SFO, 37.61° N, 122.38° W, San Francisco  
IAD, 38.95° N, 77.45° W, Dulles
```



airport.csv



DISC Application Example

- Step 1: Join departure airport codes with airport codes in airports.csv



flights.csv

32556, 08-22 07:25, 08-22 09:50, LAX, SFO
22019, 09-03 15:00, 09-03 02:50, ROA, SFO
31522, 09-03 08:55, 09-03 10:20, IAD, SFO

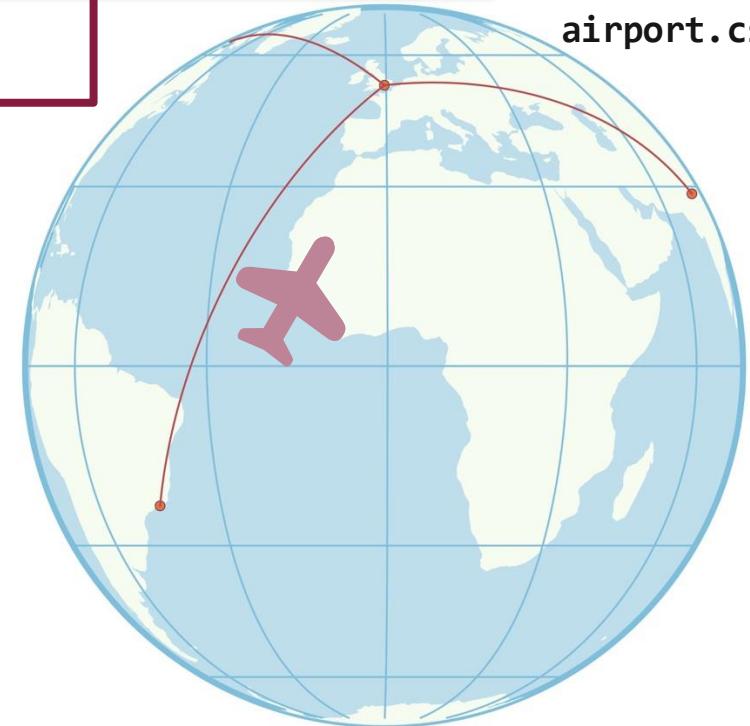
LAX, 33.94° N, 118.41° W, Los Angeles
SFO, 37.61° N, 122.38° W, San Francisco
IAD, 38.95° N, 77.45° W, Dulles



airport.csv

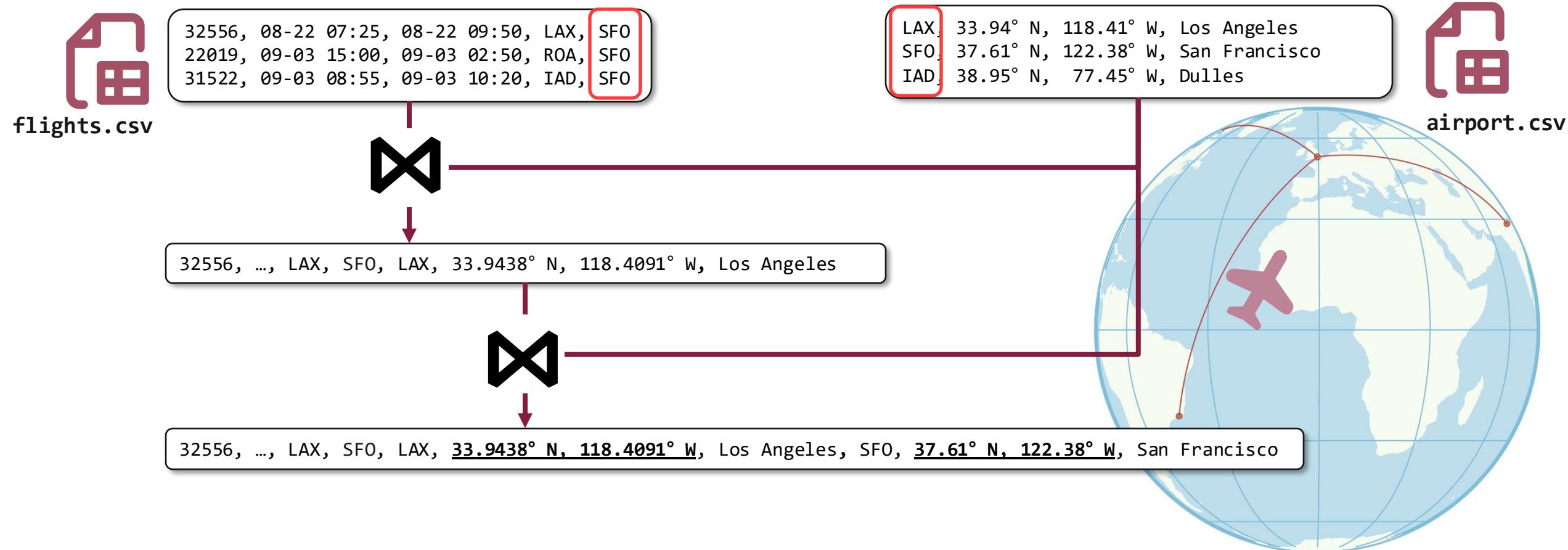


32556, ..., LAX, SFO, LAX, 33.9438° N, 118.4091° W, Los Angeles



DISC Application Example

- Step 2: Join arrival airport codes with airport codes in airports.csv



DISC Application Example

- Step 3: Use a map function to apply a formula using the longitude and latitude values



32556, 08-22 07:25, 08-22 09:50, LAX, SFO
22019, 09-03 15:00, 09-03 02:50, ROA, SFO
31522, 09-03 08:55, 09-03 10:20, IAD, SFO



LAX, 33.94° N, 118.41° W, Los Angeles
SFO, 37.61° N, 122.38° W, San Francisco
IAD, 38.95° N, 77.45° W, Dulles

flights.csv

airport.csv



32556, ..., LAX, SFO, LAX, 33.9438° N, 118.4091° W, Los Angeles



32556, ..., LAX, SFO, LAX, 33.9438° N, 118.4091° W, Los Angeles, SFO, 37.61° N, 122.38° W, San Francisco



32556, LAX, SFO, 543185.787



DISC Application Example

- Step 3: Use a map function to apply a formula using the longitude and latitude values



32556, 08-22 07:25, 08-22 09:50, LAX, SFO
22019, 09-03 15:00, 09-03 02:50, ROA, SFO
31522, 09-03 08:55, 09-03 10:20, IAD, SFO

flights.csv

LAX, 33.94° N, 118.41° W, Los Angeles
SFO, 37.61° N, 122.38° W, San Francisco
IAD, 38.95° N, 77.45° W, Dulles



airport.csv

DISC bugs are very **costly**.

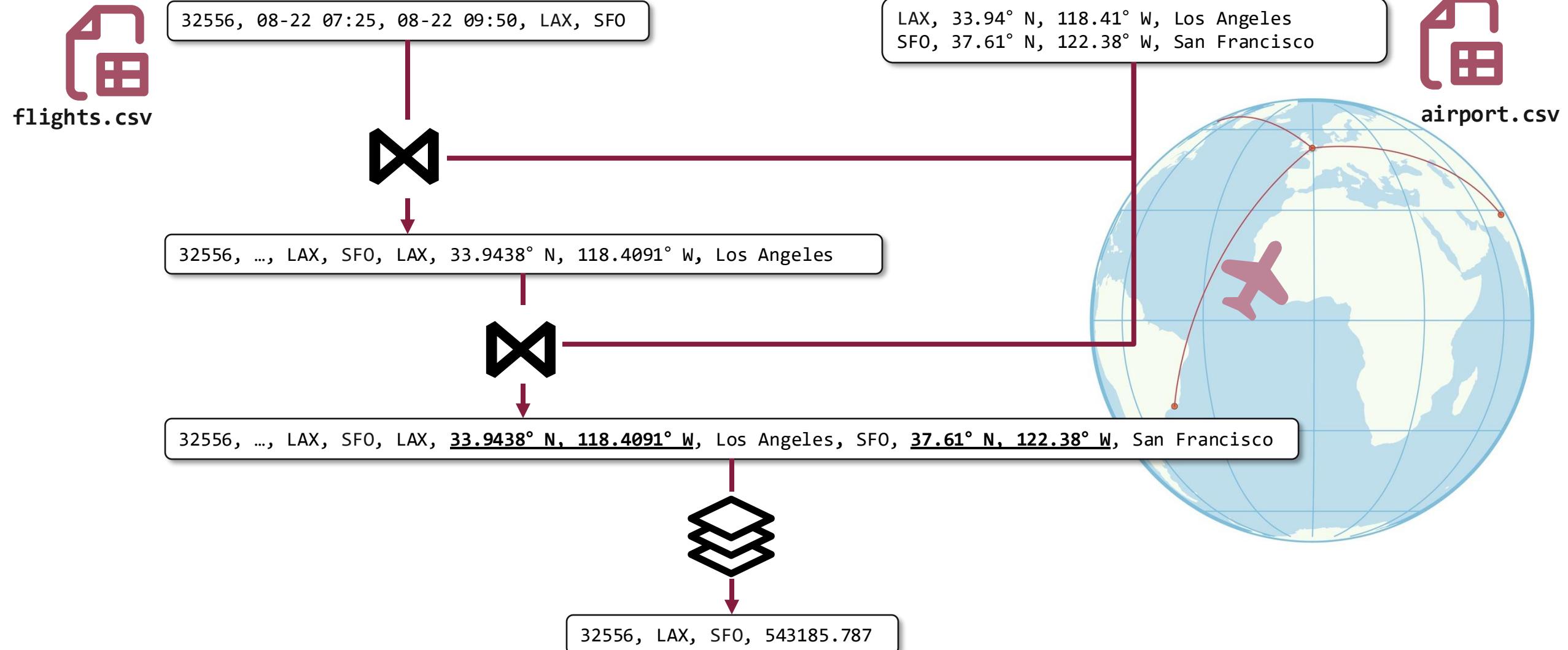
32556, ..., LAX, SFO, LAX, 33.9438° N, 118.4091° W, Los Angeles, SFO, 37.61° N, 122.38° W, San Francisco



32556, LAX, SFO, 543185.787



Challenges – Dataflow



Challenges – Dataflow

- Normal Execution:



32556, 08-22 07:25, 08-22 09:50, LAX, SFO

Column 3

Column 0

LAX, 33.94° N, 118.41° W, Los Angeles
SFO, 37.61° N, 122.38° W, San Francisco



32556, ..., LAX, SFO, LAX, 33.9438° N, 118.4091° W, Los Angeles



32556, ..., LAX, SFO, LAX, 33.9438° N, 118.4091° W, Los Angeles, SFO, 37.61° N, 122.38° W, San Francisco

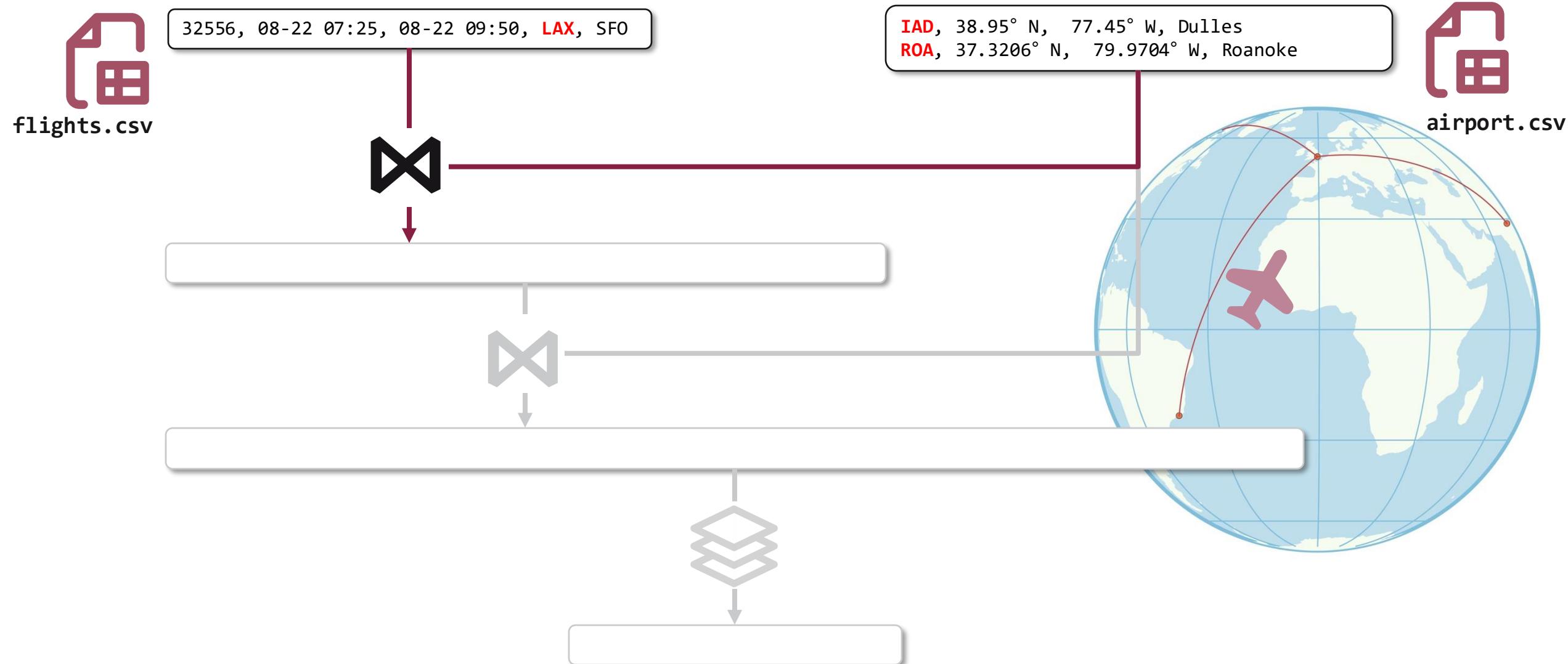


32556, LAX, SFO, 543185.787



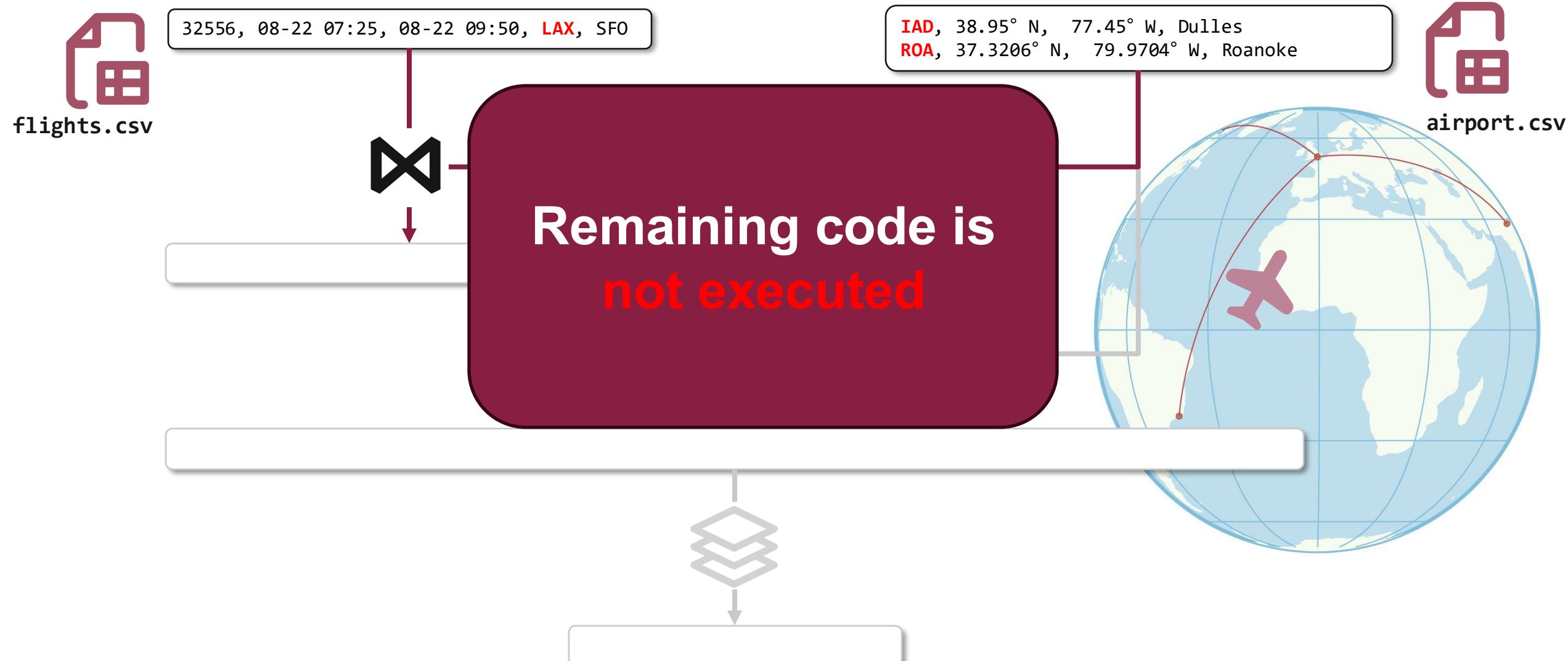
Challenges – Dataflow

- Dataflow terminates prematurely



Challenges – Dataflow

- Dataflow terminates prematurely



Challenges



flights.csv

Column 3



32556, 08-22 07:25, 08-22 09:50, LAX, SFO

Column 0

LAX, 33.94° N, 118.41° W, Los Angeles
SFO, 37.61° N, 122.38° W, San Francisco



airport.csv



32556, ..., LAX, SFO, LAX, 33.9438° N, 118.4091° W, Los Angeles



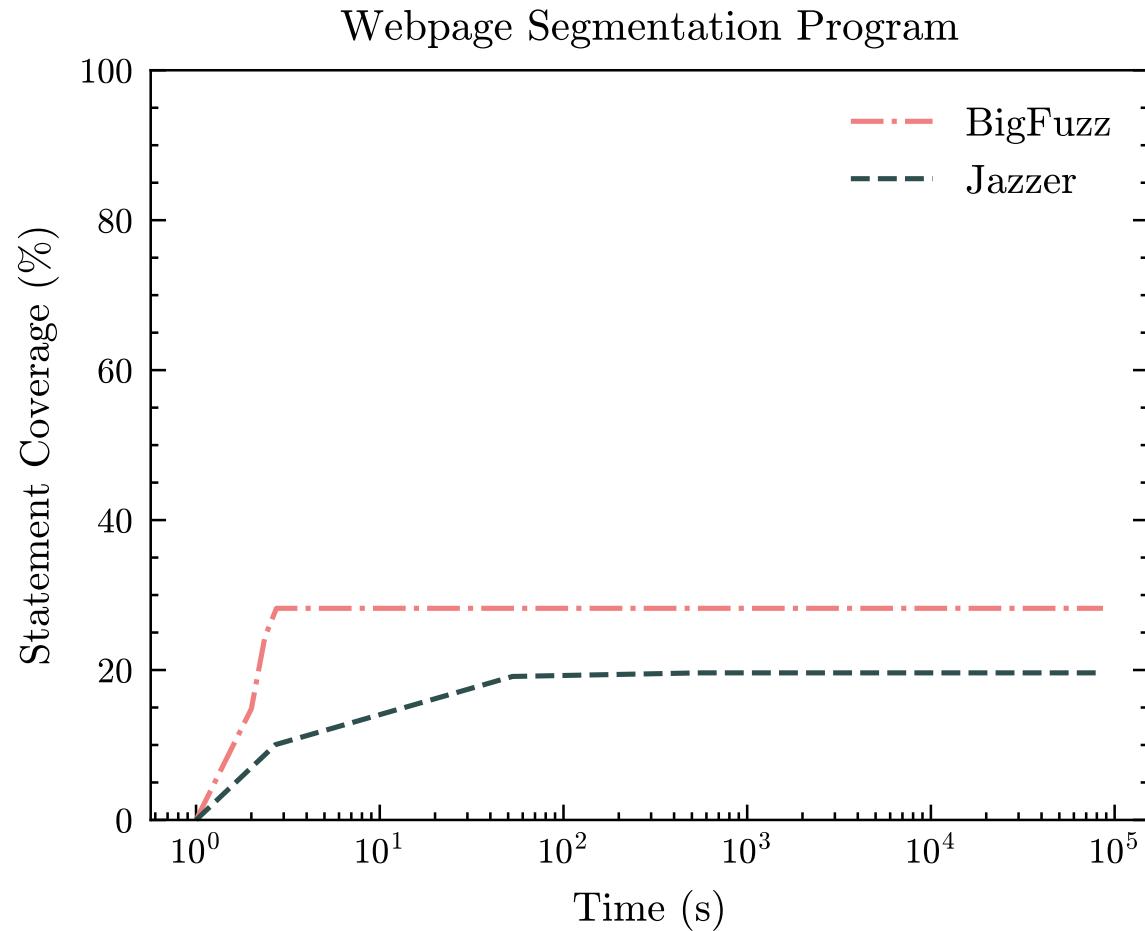
32556, ..., LAX, SFO, LAX, 33.9438° N, 118.4091° W, Los Angeles, SFO, 37.61° N, 122.38° W, San Francisco



32556, LAX, SFO, 543185.787



Current State of the Art



Challenges



flights.csv

Column 3



32556, 08-22 07:25, 08-22 09:50, LAX, SFO

Column 0



LAX, 33.94° N, 118.41° W, Los Angeles
SFO, 37.61° N, 122.38° W, San Francisco



airport.csv

Challenge 1

Strict dataflow constraints



32556, ..., LAX, SFO, LAX, 33.9438° N, 118.4091° W, Los Angeles



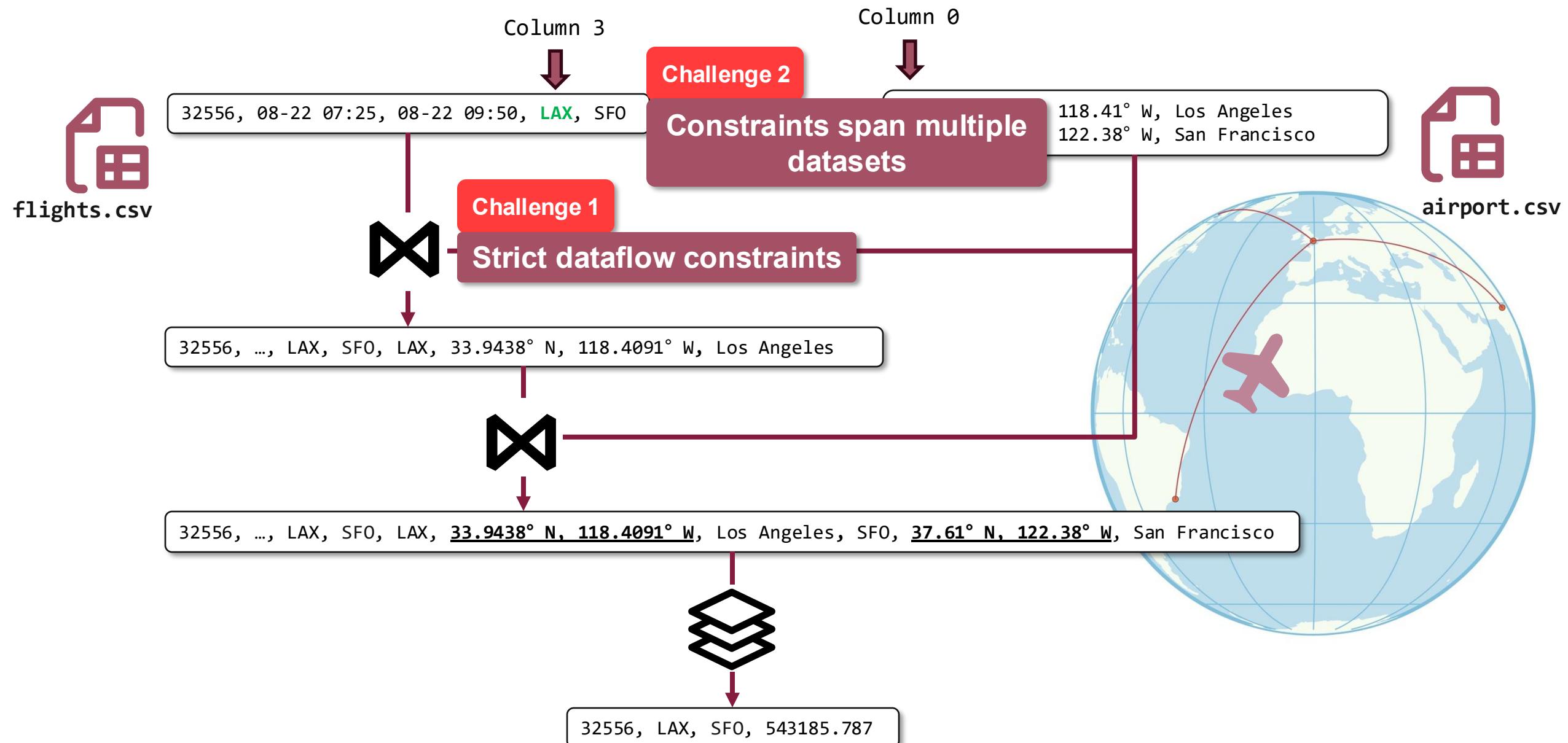
32556, ..., LAX, SFO, LAX, 33.9438° N, 118.4091° W, Los Angeles, SFO, 37.61° N, 122.38° W, San Francisco



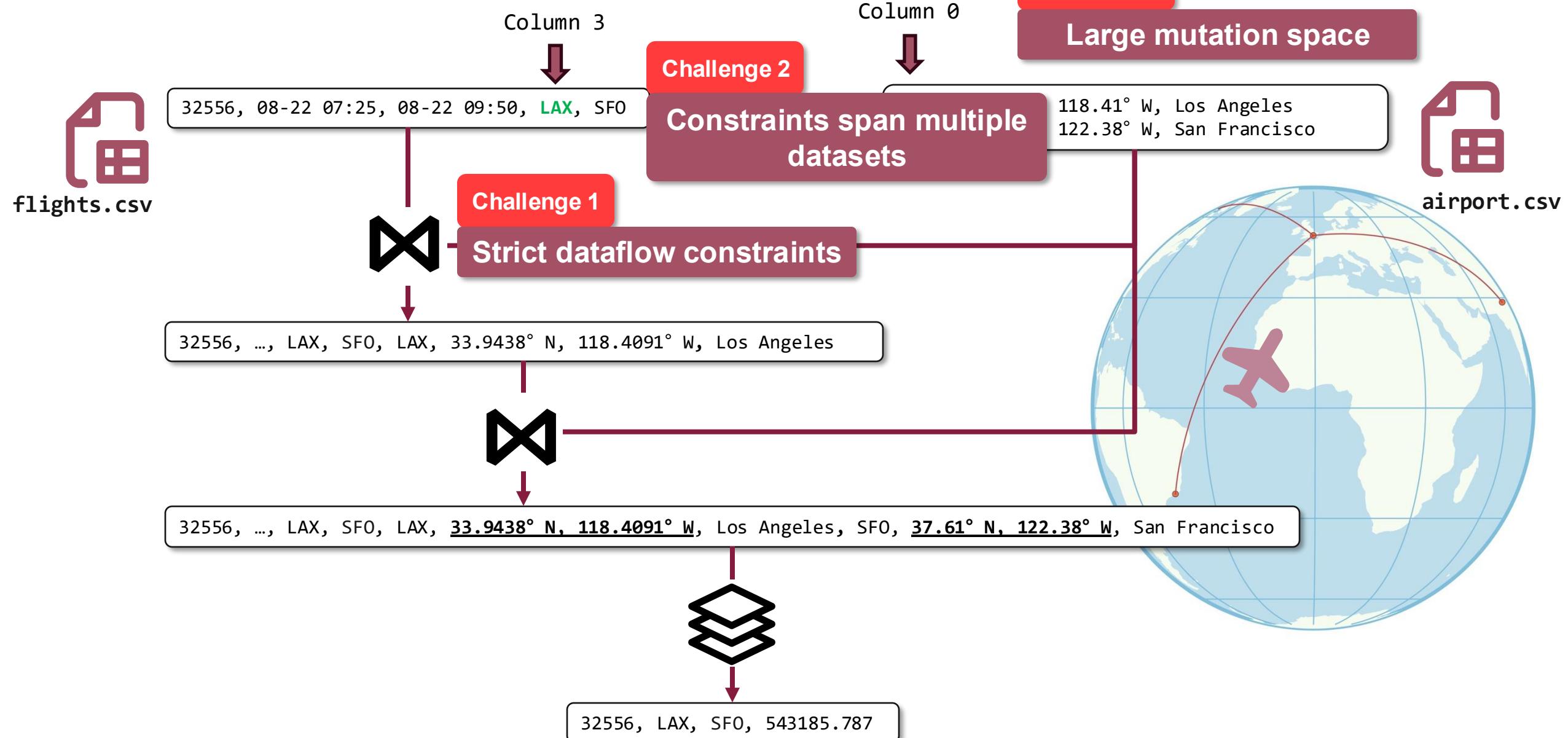
32556, LAX, SFO, 543185.787



Challenges



Challenges



What is Co-Dependence? Intuition

- Equality constraint



flights.csv

32556, 08-22 07:25, 08-22 09:50, **LAX**, SFO

Column 3



Column 0

IAD, 38.95° N, 77.45° W, Dulles
ROA, 37.3206° N, 79.9704° W, Roanoke

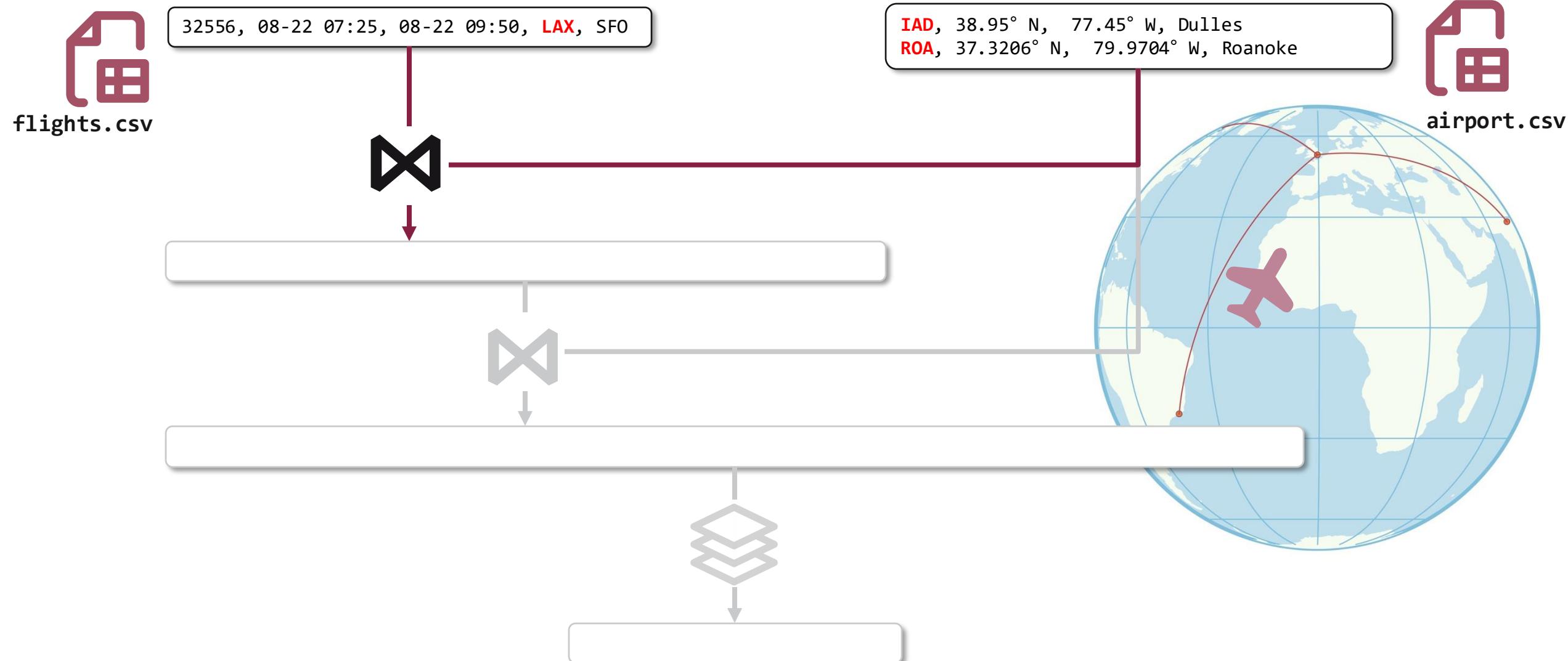


airport.csv



What is Co-Dependence?

- For optimal performance, fuzzers should respect this constraint



What is Co-Dependence?

- w.r.t mutations



flights.csv

32556, 08-22 07:25, 08-22 09:50, LAX, SFO

Column 3

Column 0

IAD, 38.95° N, 77.45° W, Dulles
ROA, 37.3206° N, 79.9704° W, Roanoke



airport.csv



The two columns are
Co-dependent



What is Co-Dependence?

Column 3

Column 0

Can we **generalize** this
idea?



flights.csv

32556, 0

Bulles
N, Roanoke



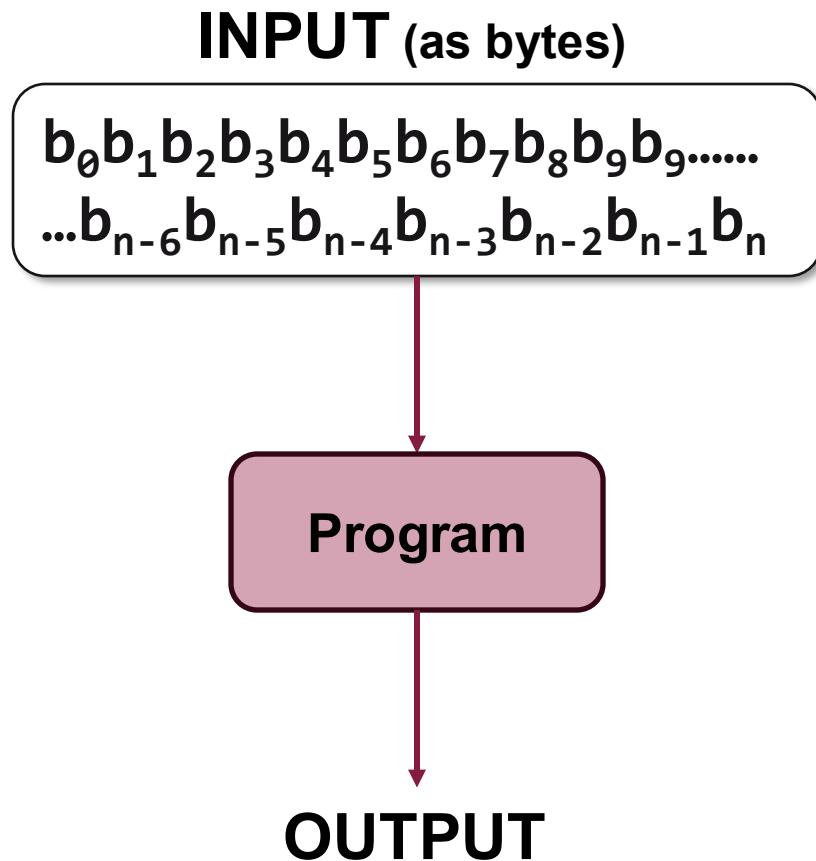
airport.csv

The two columns are
Co-dependent



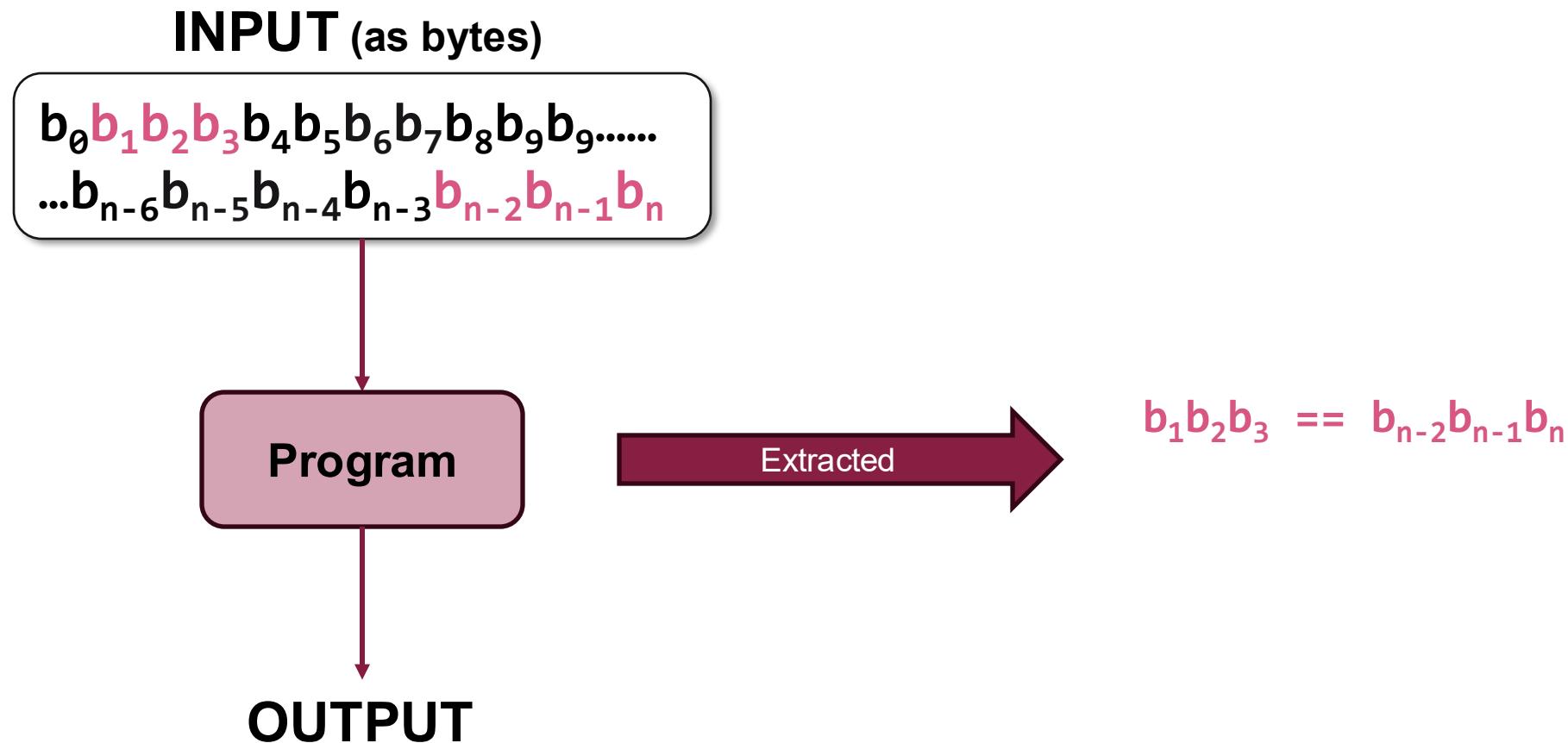
Generalizing Co-Dependence

- An input is a blob of bytes



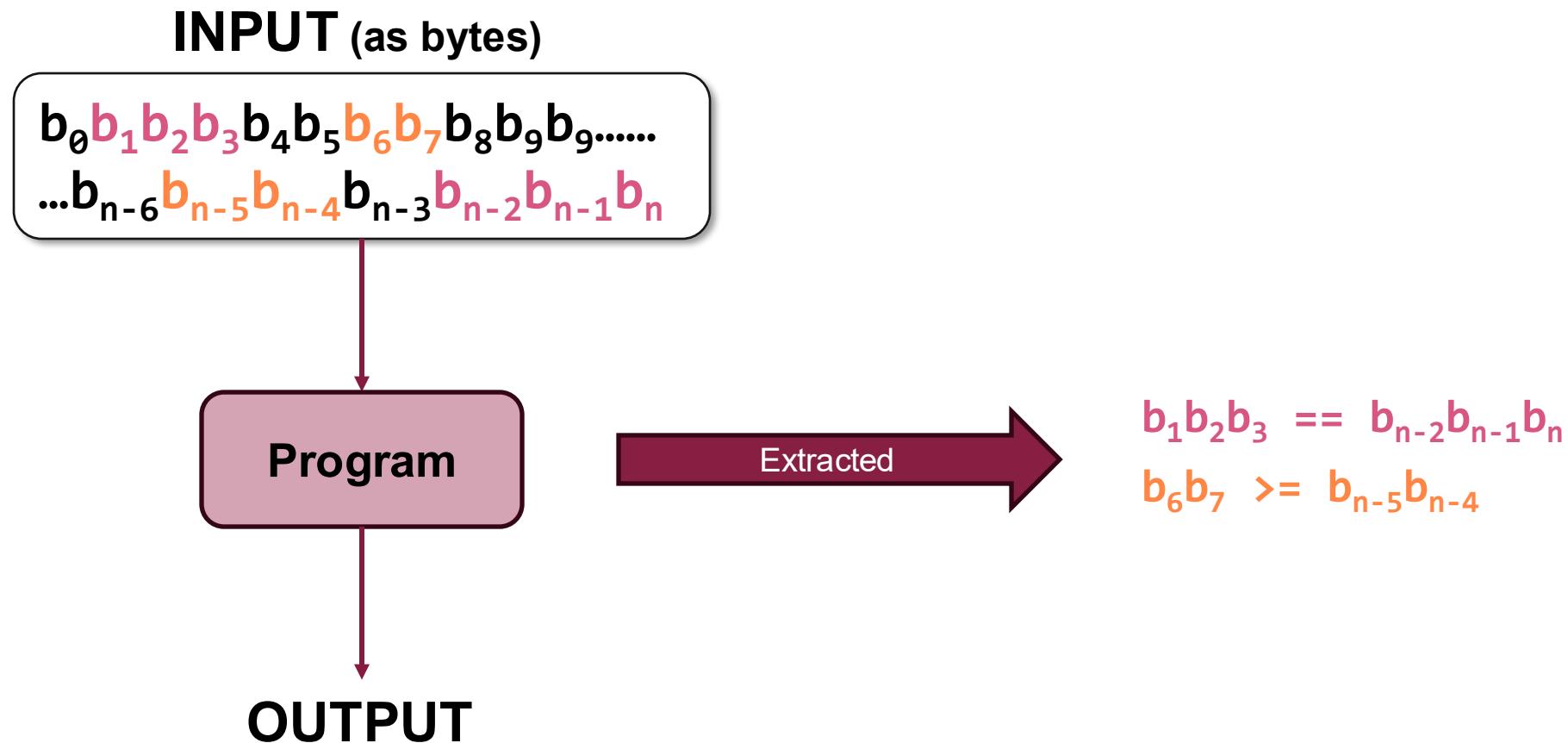
Generalizing Co-Dependence

- Two byte-regions of the input that are related to each other by some operation



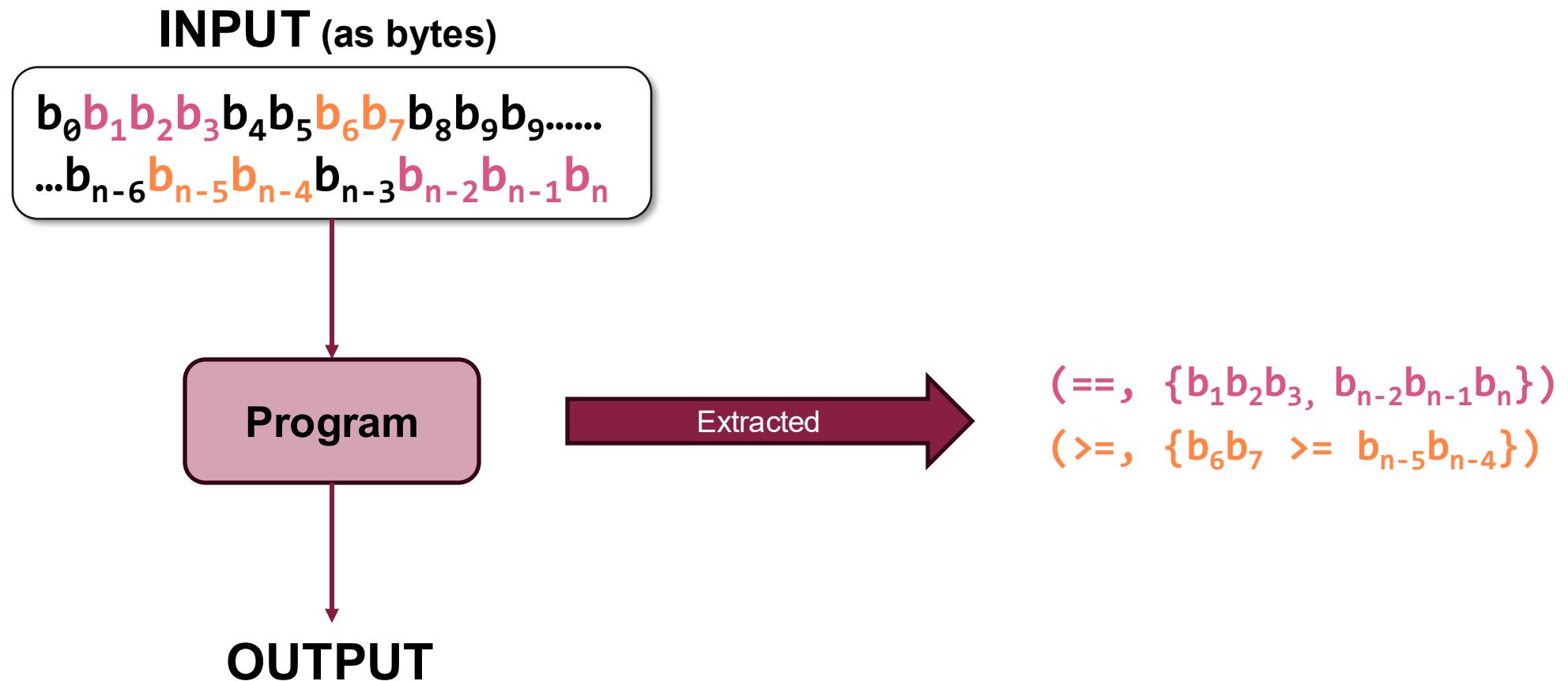
Generalizing Co-Dependence

- Two byte-regions of the input that are related to each other by some operation



Generalizing Co-Dependence

- Two byte-regions of the input that are related to each other by some operation



Generalizing Co-Dependence

other (i.e., input and output of an operator). Thus, we define a DISC application's DFG, G , as

$$G = \langle O \cup N, E \rangle$$

where $O = \{o_1, o_2, \dots, o_n\}$ is a set of operators, $N = \{n_1, n_2, \dots, n_m\}$ is a set of datanodes. $E \subseteq (O \times N) \cup (N \times O)$ is a set of directed edges connecting operators with datanodes. An atomic unit of this DFG has three nodes and two edges i.e., an operator with an incoming edge from an input datanode and an outgoing edge to an output datanode. Furthermore, a datanode n holds data in the form of a *byte sequence* $b_1 b_2 \dots b_k$. Let $D(n)$ be the set of all possible subsequences of the byte sequence in a datanode n , i.e., $\{b_1, b_2, \dots, b_1 b_2, \dots, b_1 \dots b_k\}$. Input datasets of DISC applications are defined as S , a set of initial datanodes which are external inputs to the DFG. We combine regions in input datasets in S' , a union of $D(n)$ across all input datanodes.

$$S' = \bigcup_{n \in S} D(n)$$

Finally, we characterize co-dependency among input regions as a set of tuples, C

$$C = \{(o, R) \mid R \subseteq S', \forall o \in O\}$$

The first element, o , is an operator in the dataflow graph; and the second element, R , is a subset of the regions in the input datasets that are co-dependent due to operator o . Let $I(o)$ be the incoming data to an operator o . Since co-dependence can only occur between regions of the original datasets, we must extract R from $I(o)$, which can be any arbitrary byte sequence in the incoming datanode of operator o . To extract such information, we define *monitors* that are concretely explained in Section 3.1.

$$M_o : I(o) \rightarrow \mathcal{P}(S')$$

where $\mathcal{P}(S')$ is the powerset of S' . A monitor, M_o , is an operator-specific function that takes $I(o)$ as input and outputs a set of byte

input

bytes)

$b_8 b_9$
 \vdots
 $b_{n-2} b_n$



T

Table 1: Summary of how each class of operators produces co-dependent regions in the input dataset. For simplicity, we use $\text{row}[1].\text{col}[3]$ as a human-readable representation of input byte region $b_i \dots b_j$, where $0 < i < j < \text{size}(\text{dataset})$

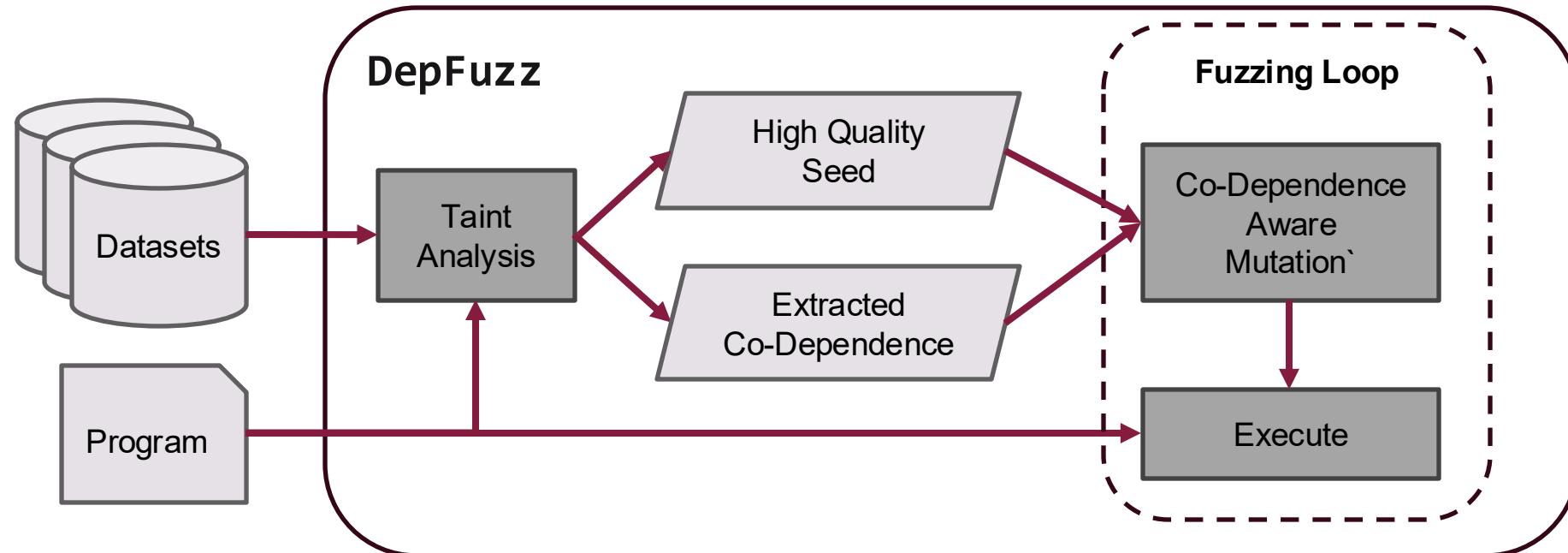
Operator Class	Sample Operators	Example of Identified Constraint	Mutation strategy
Fusions	<code>data1.join(data2)</code> <code>data1.intersection(data2)</code> <code>data1.cogroup(data2)</code>	$M_{join}(\{\text{data1}, \text{data2}\})$ Possible output of M_{join}: <code>{data1.row[1].col[3], data2.row[23].col[0]}</code> <code>{data1.row[31].col[3], data2.row[52].col[0]}</code> Co-dependence tuples: <code>(==, {data1.row[1].col[3], data2.row[23].col[0]})</code> <code>(==, {data1.row[31].col[3], data2.row[52].col[0]})</code>	Any mutation applied to <code>data1.row[1].col[3]</code> must also be applied to <code>data2.row[23].col[0]</code> . If no rows with matching keys exists, select a row from <code>data1</code> and copy <code>data1.col[3]</code> to <code>data2.col[0]</code> .
Aggregations	<code>data.aggregateByKey(udf)</code> <code>data.reduceByKey(udf)</code> <code>data.groupByKey()</code> <code>data.countByKey()</code>	$M_{reduceByKey}(\{\text{data}\})$ Possible output of $M_{reduceByKey}$: <code>{data.row[2].col[2],</code> <code>data.row[43].col[2],</code> <code>data.row[63].col[2]}</code> Co-dependence tuple: <code>(==, {data.row[2].col[2],</code> <code>data.row[43].col[2],</code> <code>data.row[63].col[2]})</code>	Any mutation applied to <code>data[row=2,col=2]</code> must also be applied to <code>data[row=43,col=2]</code> . Duplicate rows and apply same mutation to key columns of duplicates.
Filters	<code>data.filter(col0 > col5)</code>	$M_{filter}(\{\text{data}\})$ Possible output of M_{filter}: <code>{data.row[23].col[0], data2.row[23].col[5]}</code> <code>{data.row[31].col[0], data2.row[31].col[5]}</code> Co-dependence tuples: <code>(>, {data1.row[23].col[0], data2.row[23].col[5]})</code> <code>(>, {data1.row[31].col[3], data2.row[31].col[5]})</code>	Any mutation applied to <code>data.col[0]</code> and <code>data.col[5]</code> must ensure that there is a true and false row for the predicate.
UDF Operators	<code>if(a.contains(b))</code> <code>if(a != b)</code> <code>if(a > b)</code>	$M_{contains}(\{a, b\})$ Possible output of $M_{contains}$: <code>{data.row[1].col[0], data2.row[1].col[2]}</code> Co-dependence tuples: <code>(contains, {data1.row[1].col[0], data2.row[1].col[2]})</code>	Any mutation applied to <code>data.col[0]</code> and <code>data.col[2]</code> must ensure that string <code>a</code> contains <code>b</code> for some mutations. It must also ensure it occasionally creates inputs that violate this.

Generalizing Co-Dependence

- Two byte-regions of the input that are related to each other by some operation

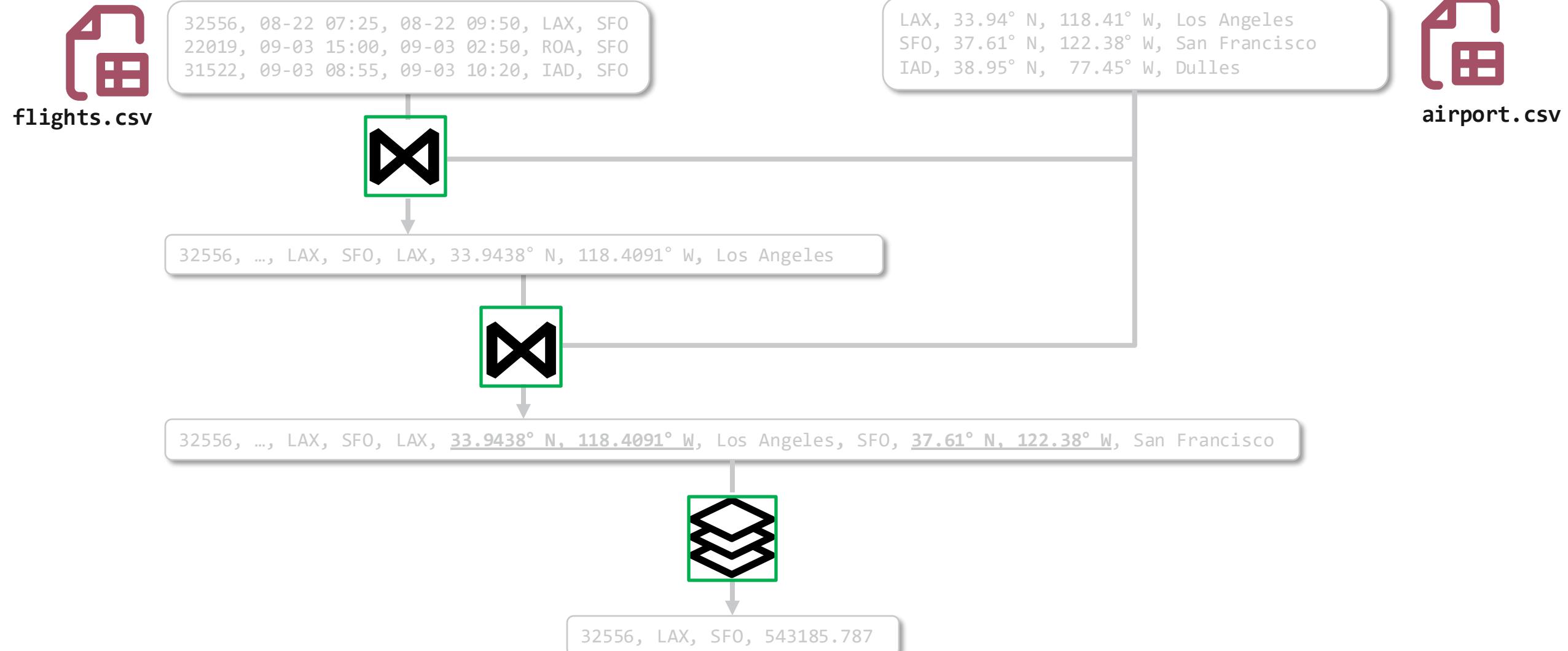


Approach Overview



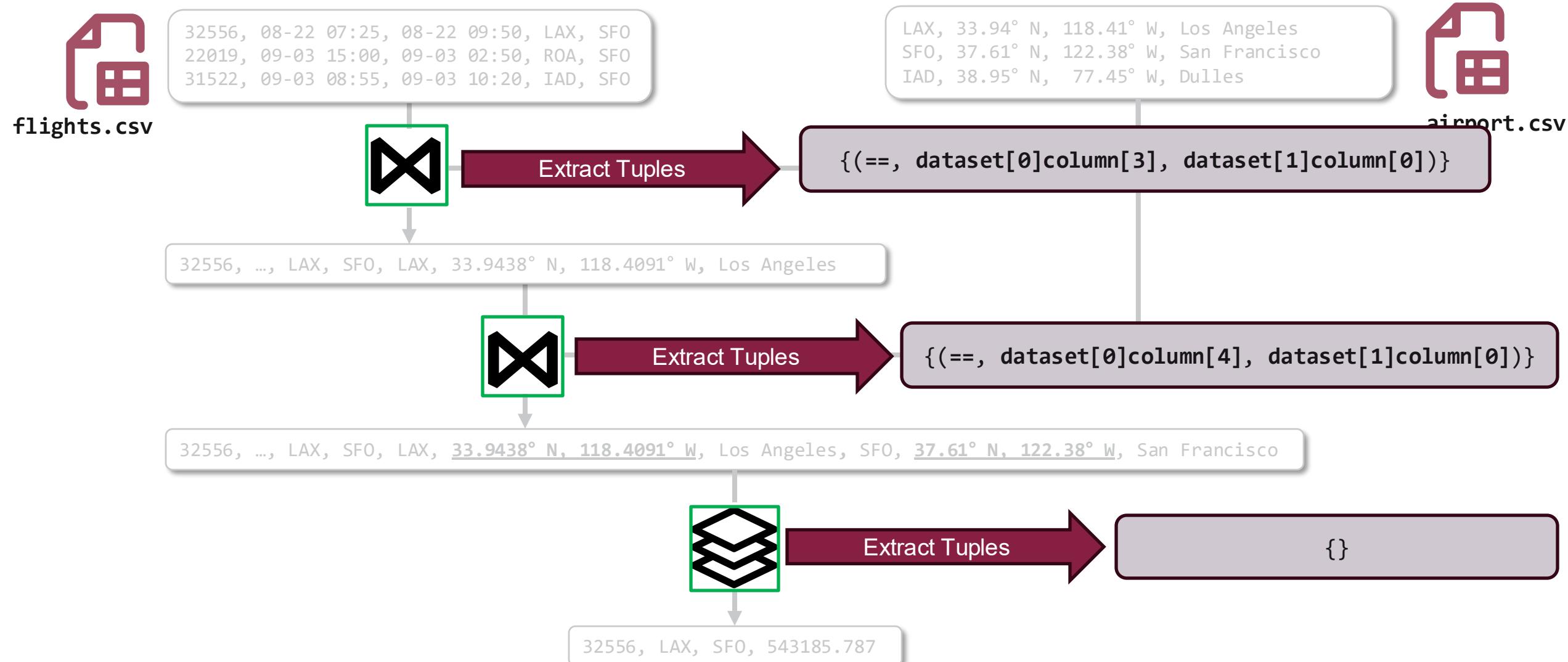
Capturing Co-Dependence

The **monitor** abstraction



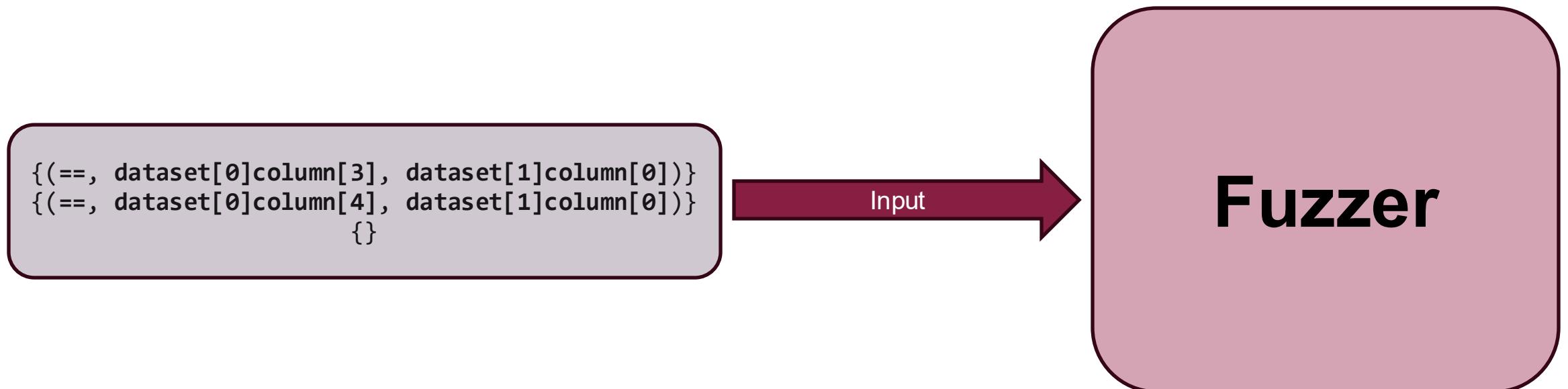
Capturing Co-Dependence

The **monitor** abstraction



Fuzzing

The fuzzer can use this representation perform better mutations



Representation

- This representation abstracts the code from the fuzzer

```
(Operator1, {Region A, Region B})  
(Operator2, {Region B, Region C})
```

Representation

- Transitive co-dependencies can be merged

```
(==, {Region A, Region B})  
(==, {Region B, Region C})
```

Representation

- Transitive co-dependencies can be merged

```
(==, {Region A, Region B})  
(==, {Region B, Region C})
```

```
{(==, {Region A, Region B, Region C})}
```

Representation

- Transitive co-dependencies can be merged

```
(==, {Region A, Region B})  
(==, {Region B, Region C})
```

```
{(==, {Region A, Region B, Region C})}
```

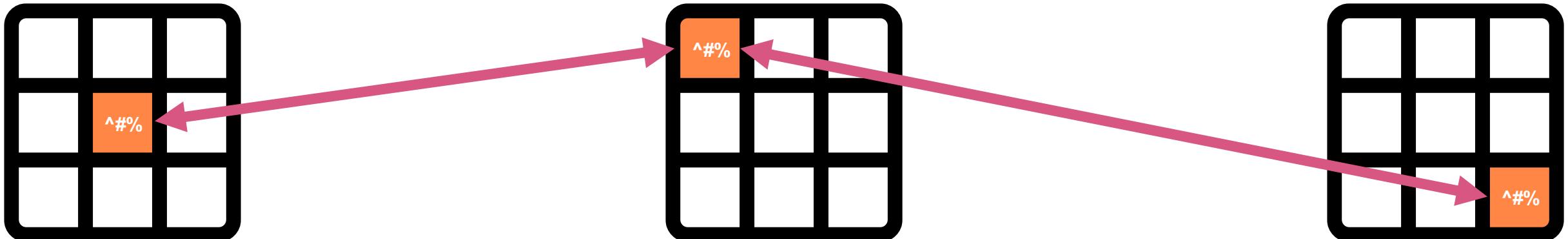


Representation

- Transitive co-dependencies can be merged

```
(==, {Region A, Region B})  
(==, {Region B, Region C})
```

```
{(==, {Region A, Region B, Region C})}
```

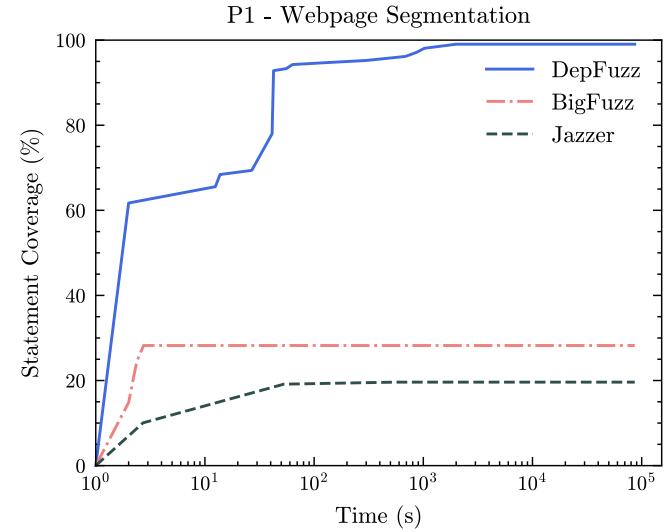


DepFuzz Evaluation

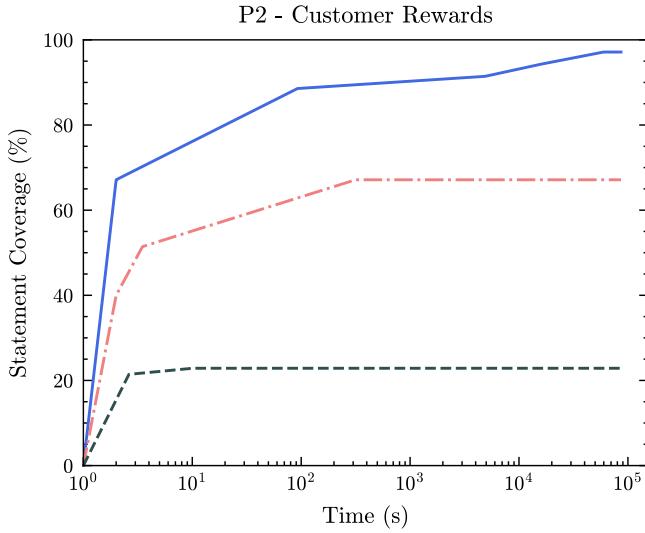
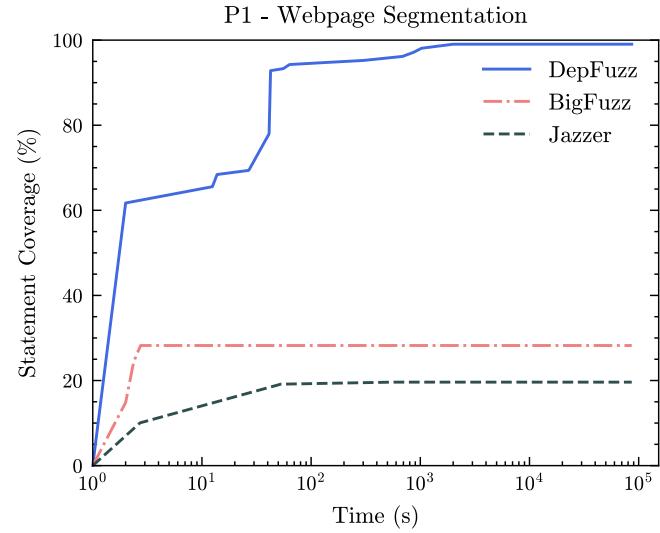
- DepFuzz is available on github: <https://github.com/SEED-VT/DepFuzz>
- Dockerized and easily runnable
- Baselines: BigFuzz and Jazzer
- Benchmark of 17 programs



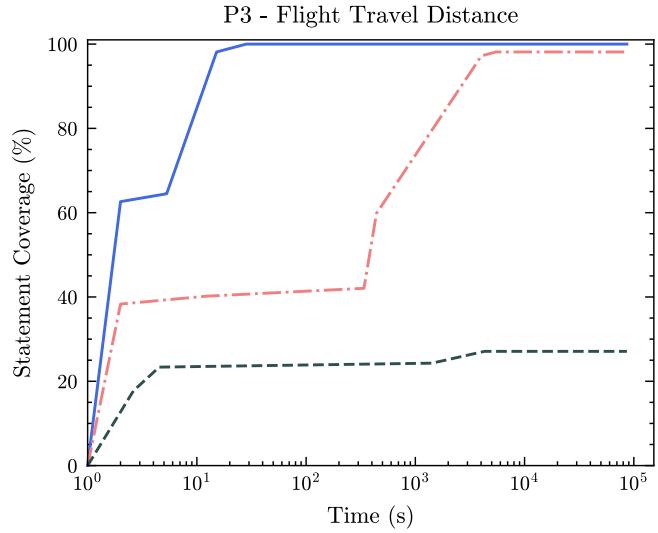
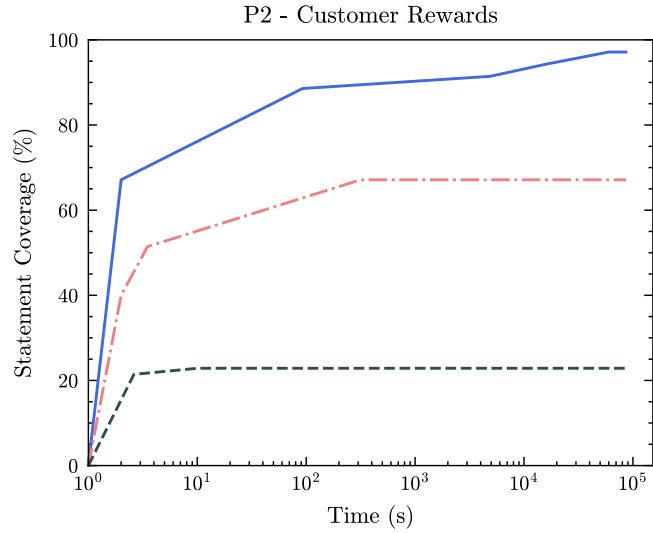
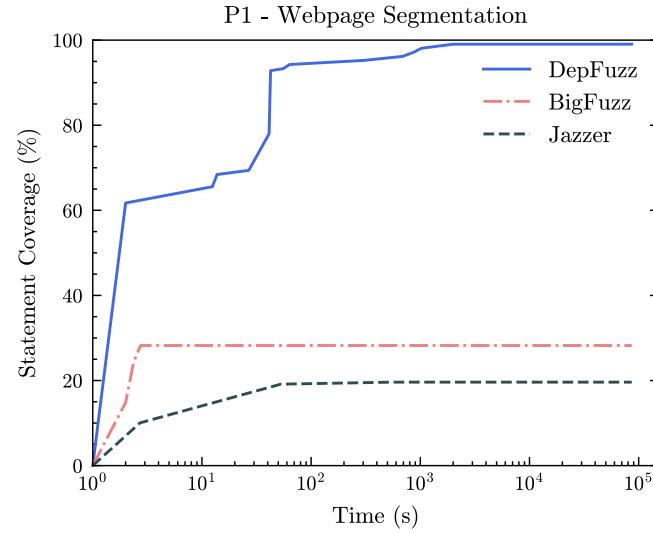
Results – Coverage



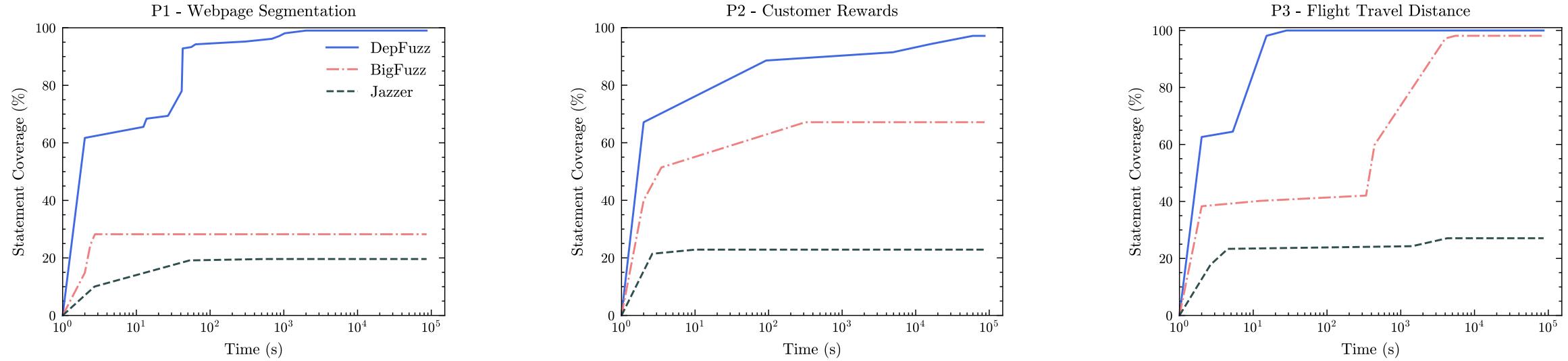
Results – Coverage



Results – Coverage

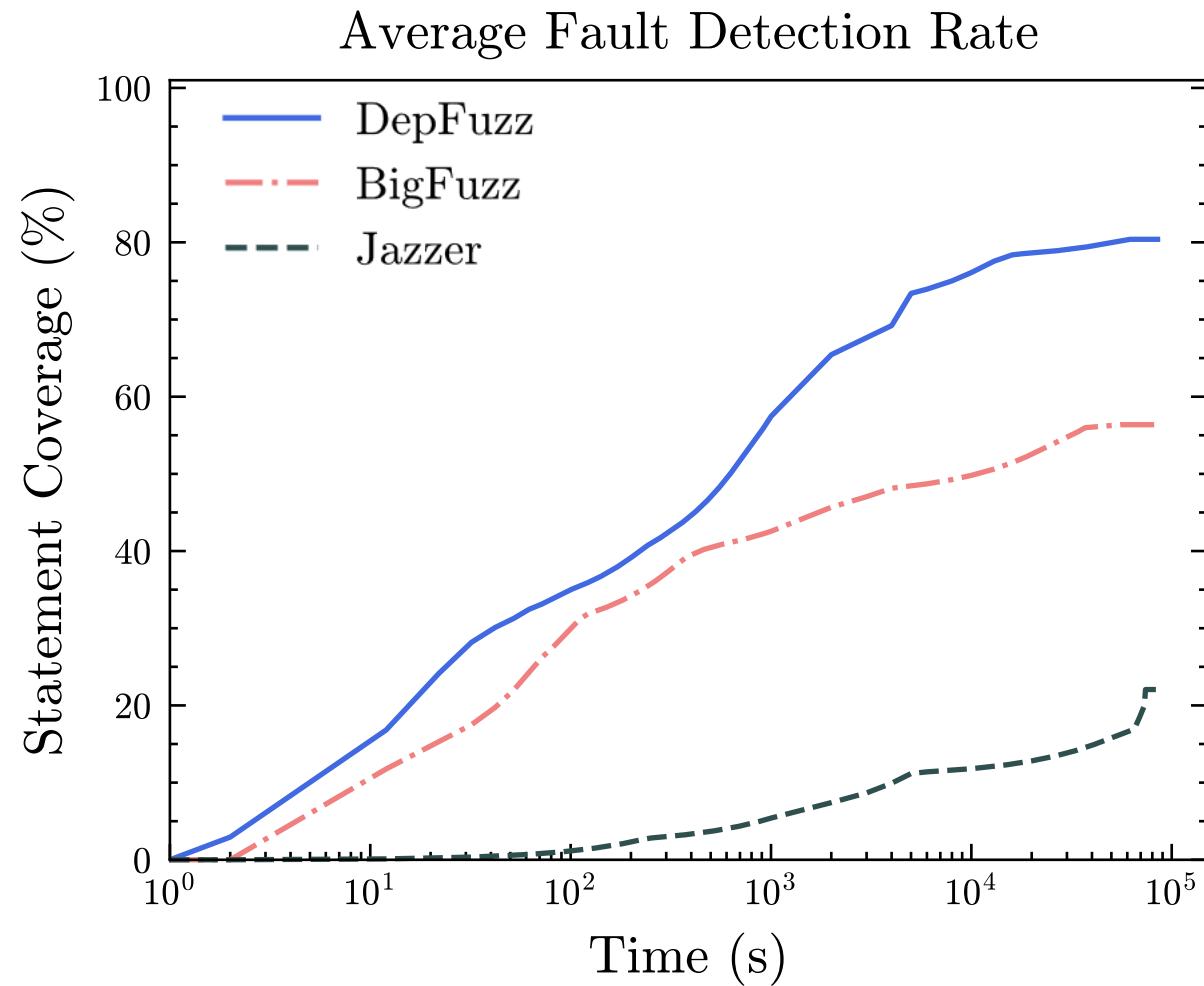


Results – Coverage

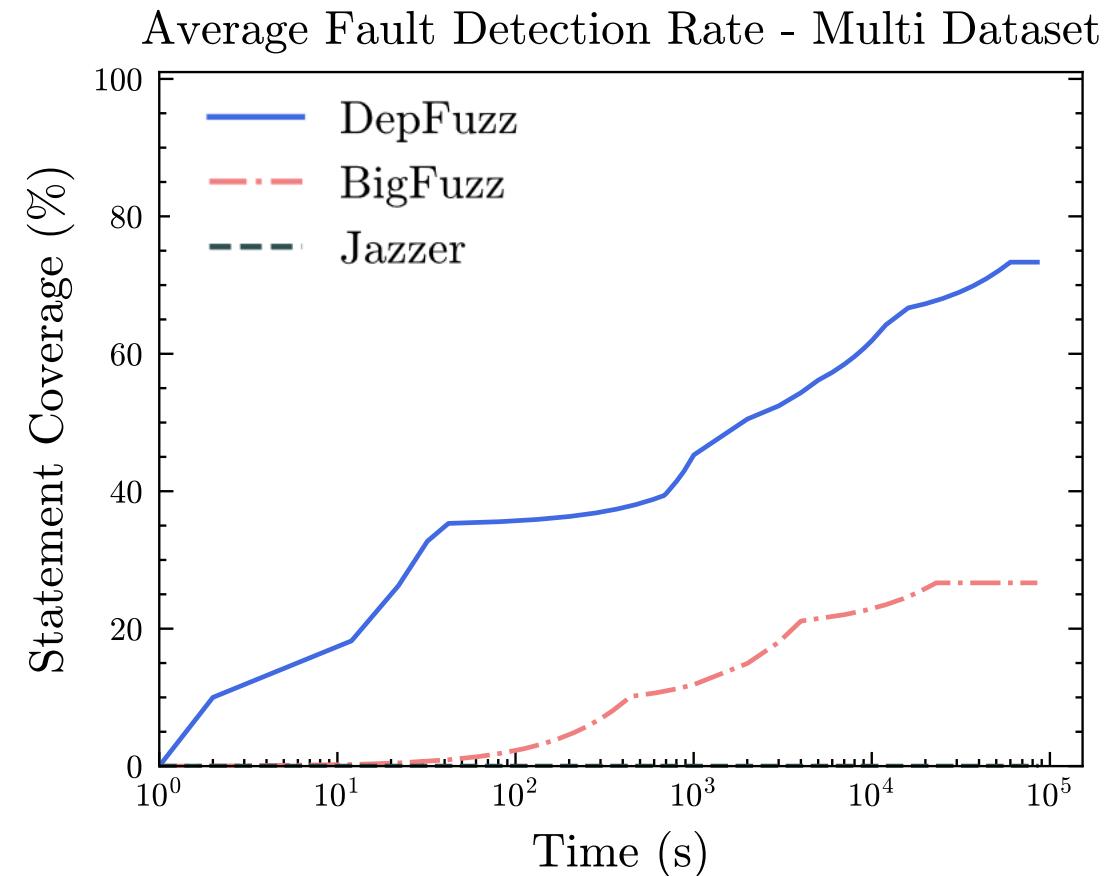
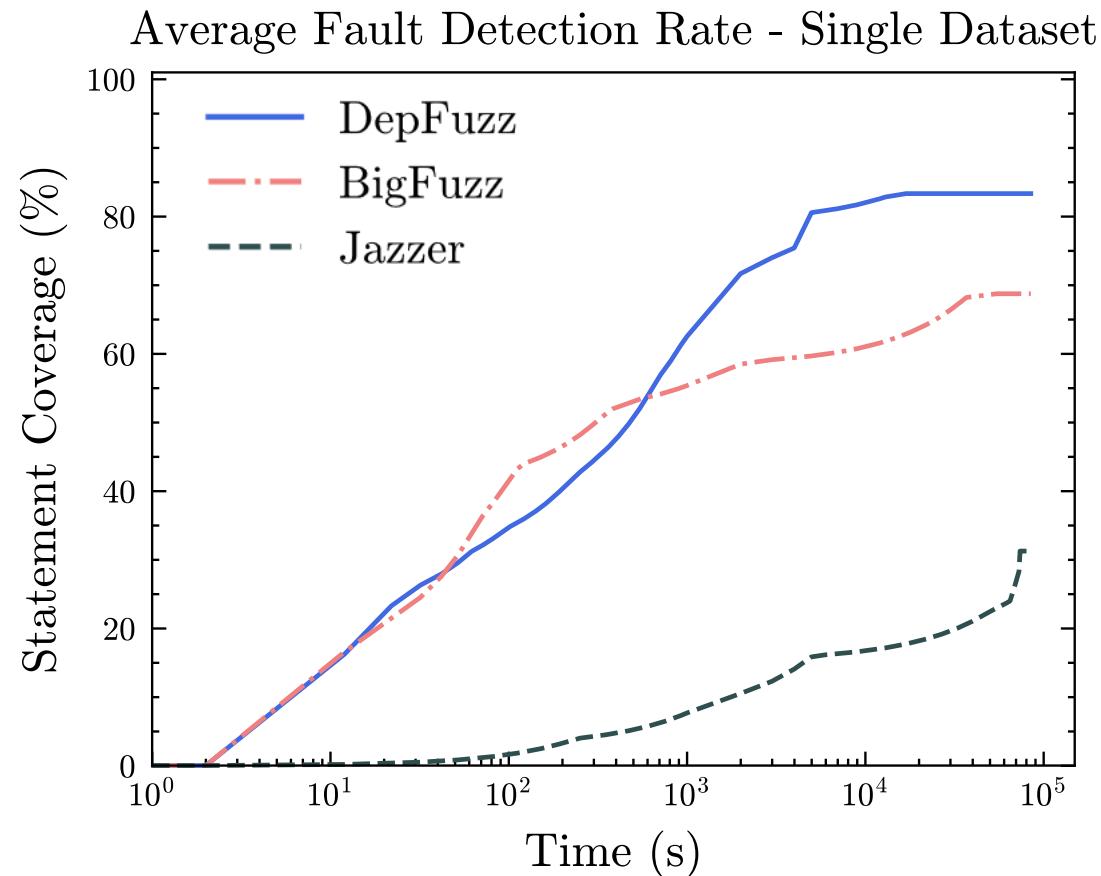


- On average, DepFuzz achieves 21% higher coverage than baselines.

Results – Fault Detection

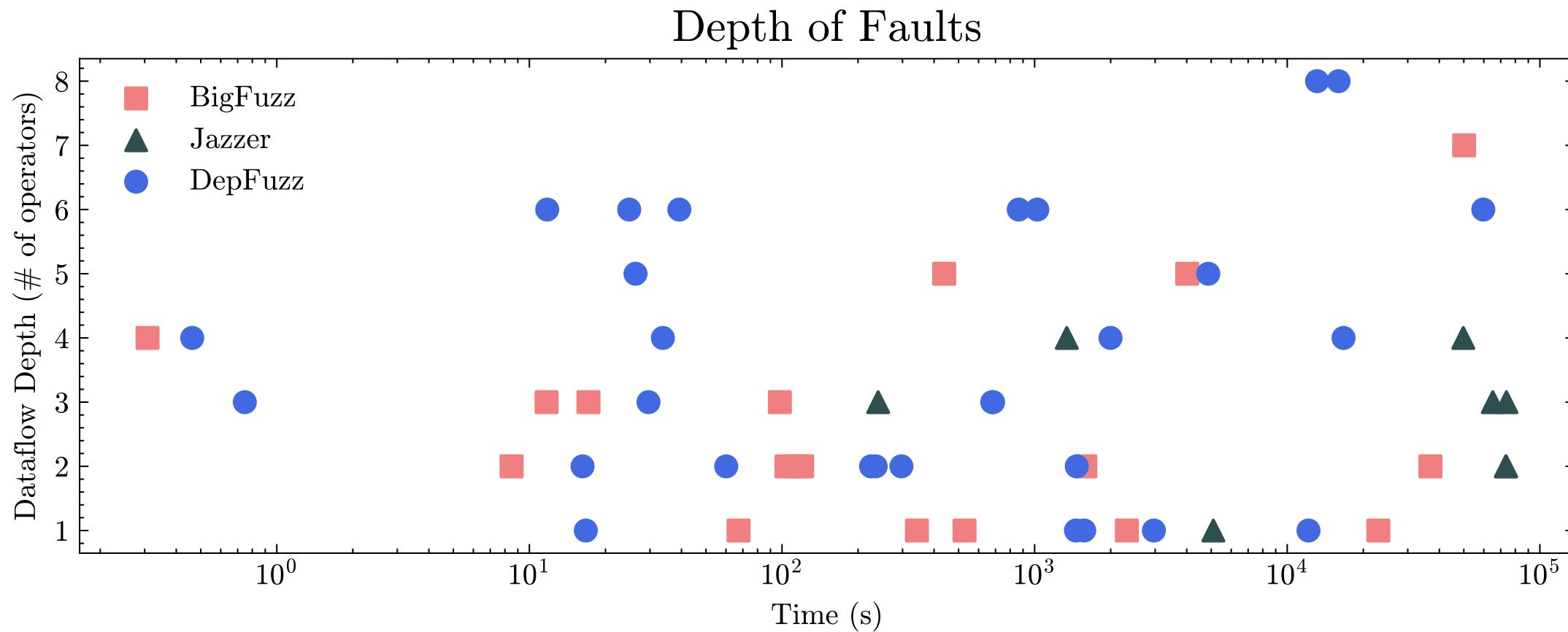


Results – Fault Detection



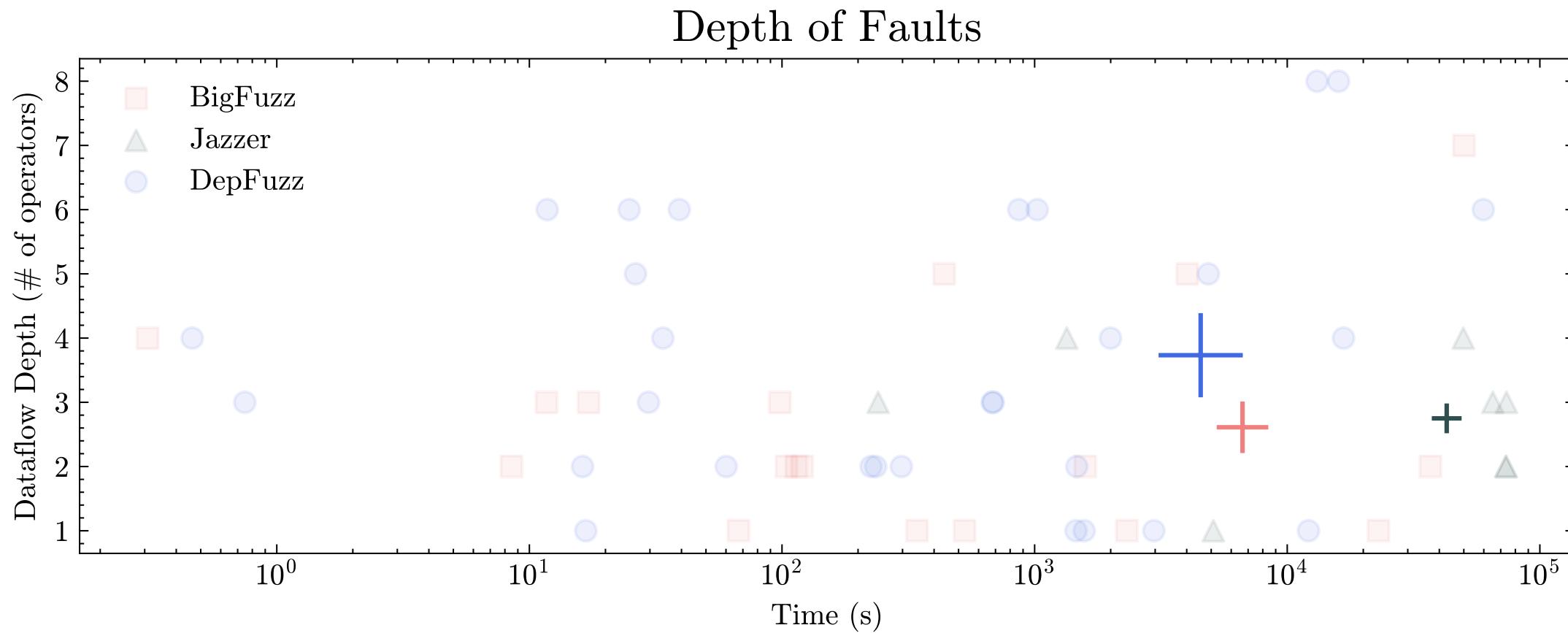
- On average, DepFuzz finds 2.6X more faults than baselines.

Results – Detected Fault Depth



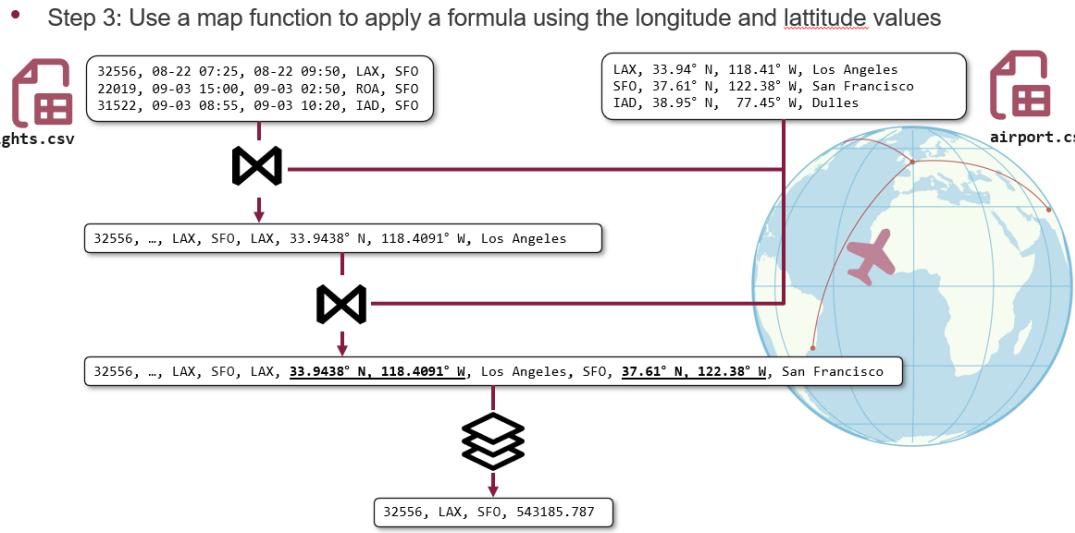
- On average DepFuzz detects faults that are deeper

Results – Detected Fault Depth



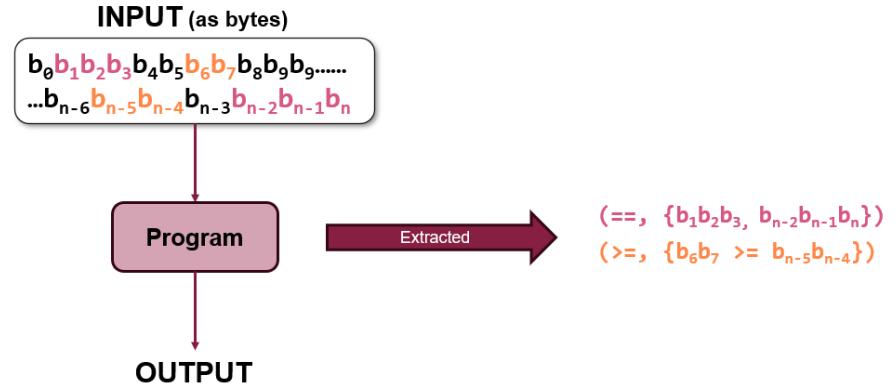
- On average DepFuzz detects faults that are deeper

DISC Application Example



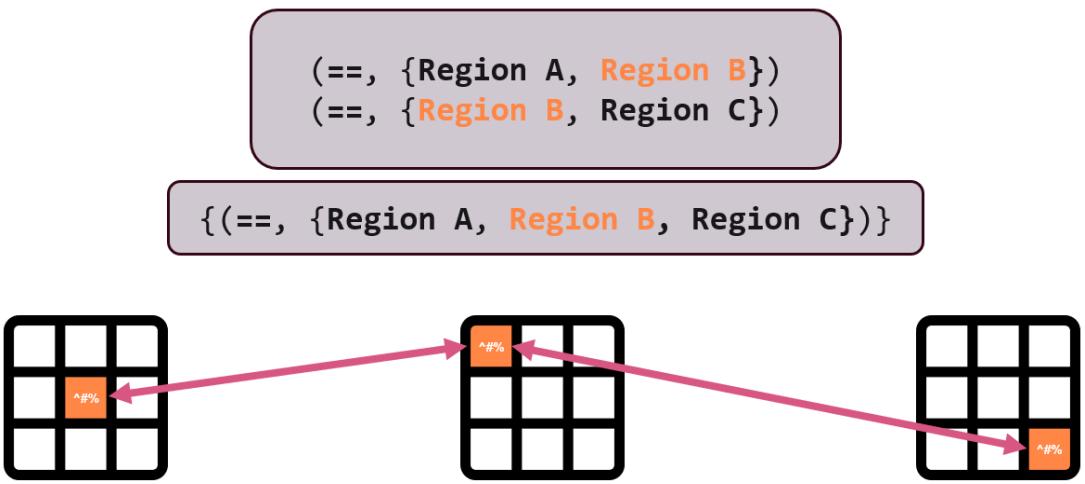
Generalizing Co-Dependence

- Two byte-regions of the input that are related to each other by some operation

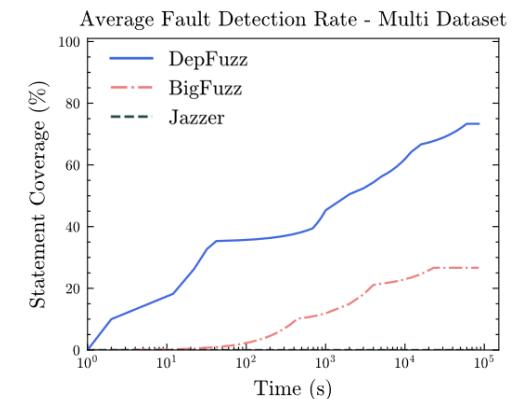
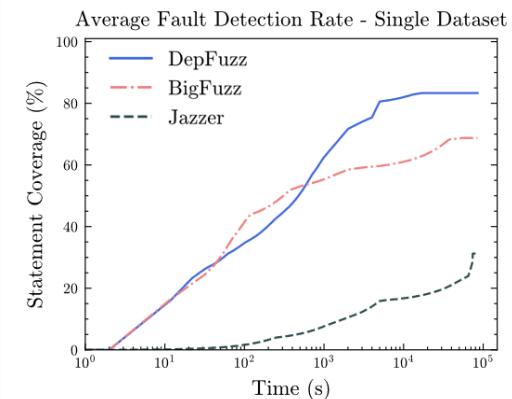


Representation

- Transitive co-dependencies can be merged



Results – Fault Detection



- On average, DepFuzz finds 2.6X more faults than baselines.

Co-dependence Aware Fuzzing for Dataflow-based Big Data Analytics

Ahmad Humayun, Miryung Kim, Muhammad Ali Gulzar

