

# Eugene Wu

Assistant Professor

[ewu@cs.columbia.edu](mailto:ewu@cs.columbia.edu)

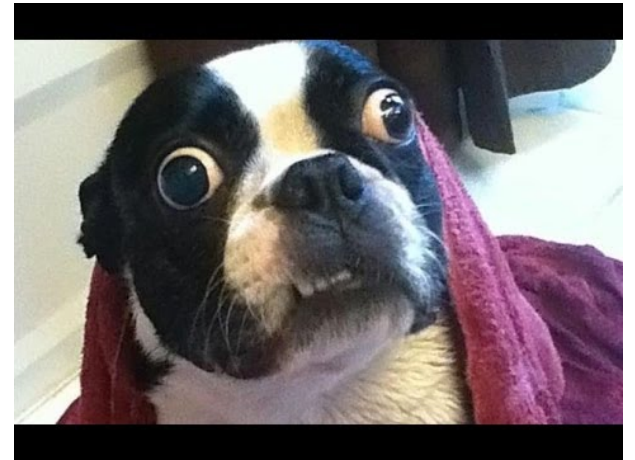




Art



Legal Auditing



Showdog Pedigree  
(stud book)

# Provenance

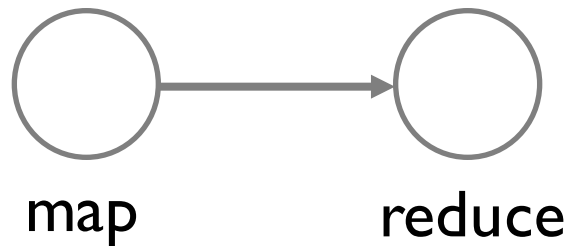
What is it?

Projects@Columbia

Open Problems

*What*  
*Provenance Information*  
*is in a*  
*Map-Reduce job?*

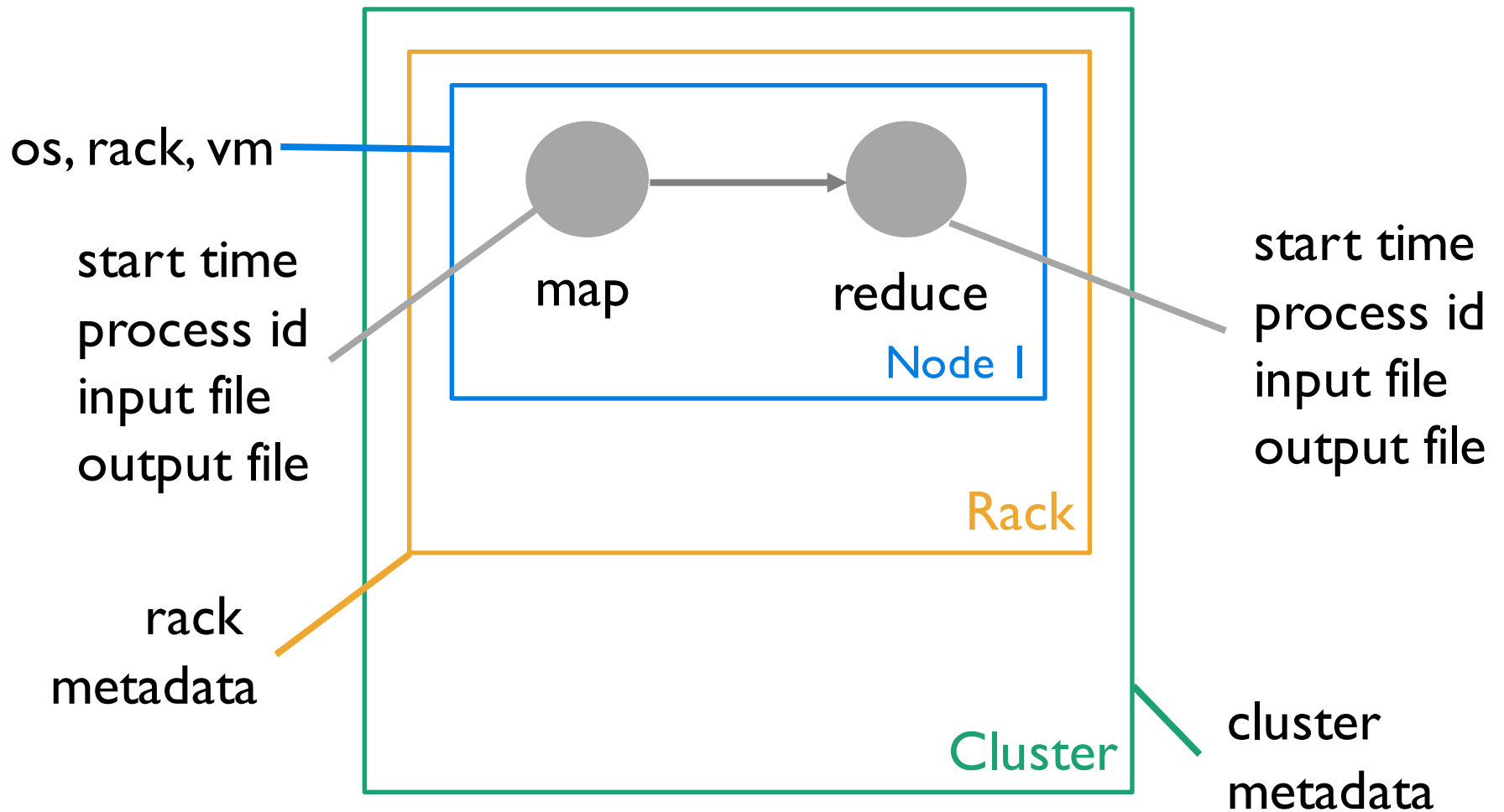
# Static Information



```
UPDATE CUSTOMER SET C_BALANCE =  
C_BALANCE + ? WHERE C_ID = ? AND C_D_ID  
= ? AND C_W_ID = ? UPDATE CUSTOMER SET  
C_BALANCE = ?, C_YTD_PAYMENT = ?,  
C_PAYMENT_CNT = ?, C_ORDER = ? WHERE  
C_W_ID = ? AND C_D_ID = ? AND C_ID = ?  
UPDATE WAREHOUSE SET W_YTD = W_YTD + ?  
WHERE W_ID = ?  
UPDATE = ? WHERE D_ID = ? AND D_W_ID =  
?  
INSERT INTO HISTORY VALUES (?, ?, ?, ?,  
?, ?, ?, ?) D_ID, OL_W_ID, OL_NUMBER,  
OL_I_ID, OL_SUPPLY_W_ID, OL_DELIVERY_D,  
OL_QUANTITY, OL_AMOUNT, OL_DIST_INFO)  
VALUES (?, ?, ?, ?, ?, ?, ?, ?)  
UPDATE ORDER LINE SET OL_DELIVERY_D = ?  
WHERE OL_O_ID = ? AND OL_I_ID = ? AND  
OL_W_ID = ?
```

wordcount.java  
file version

# Dynamic Information



# Answerable Queries

## About Metadata

What processes failed?

What machines were they running on?

What were inputs to these processes?

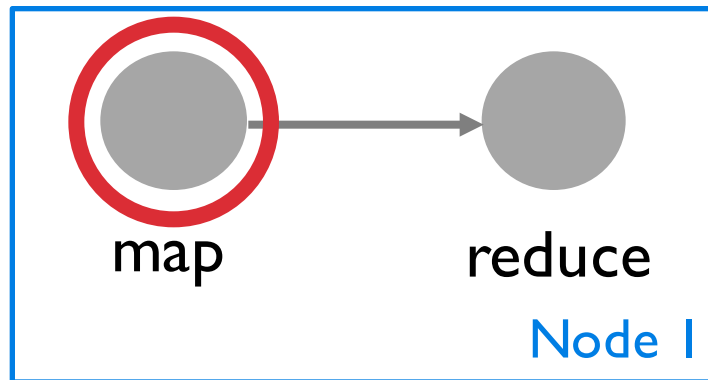
What was their source code?

## About Data

This record is wrong, results are contaminated?

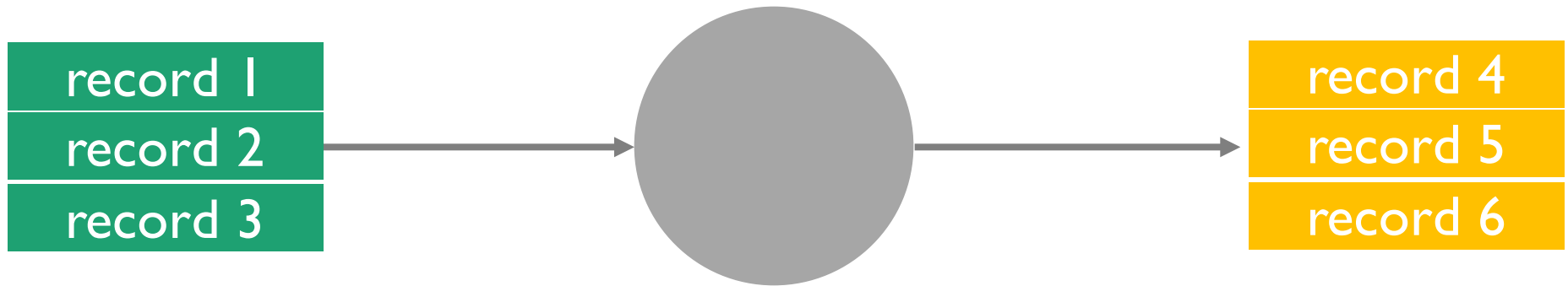
This result val is wrong, what input records to check?

# Lineage

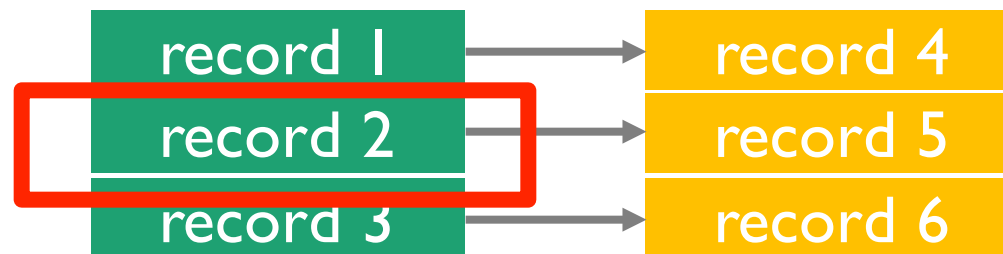




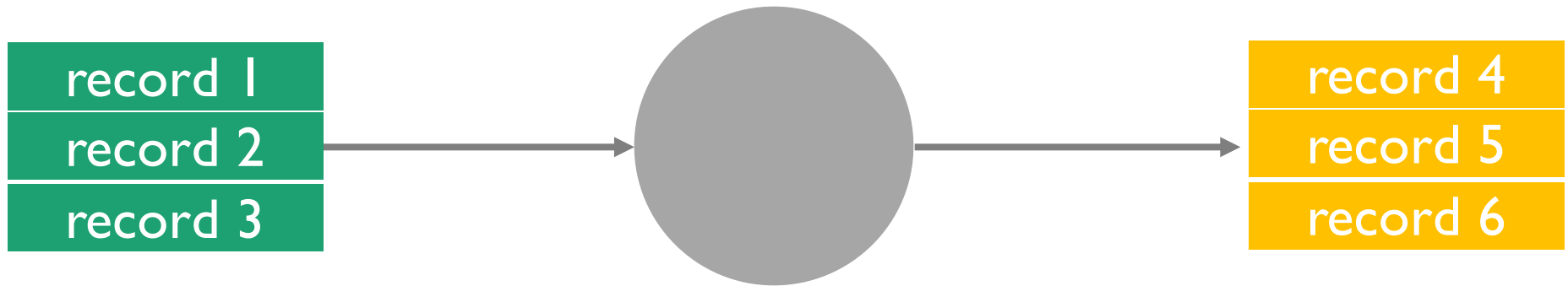
# Lineage



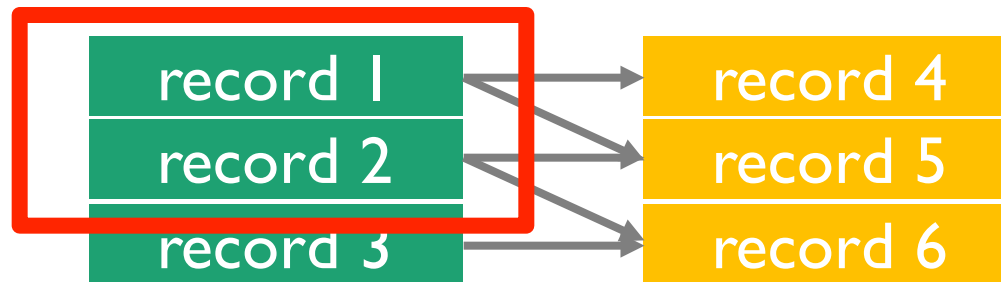
Which input records generated **record 5** ?



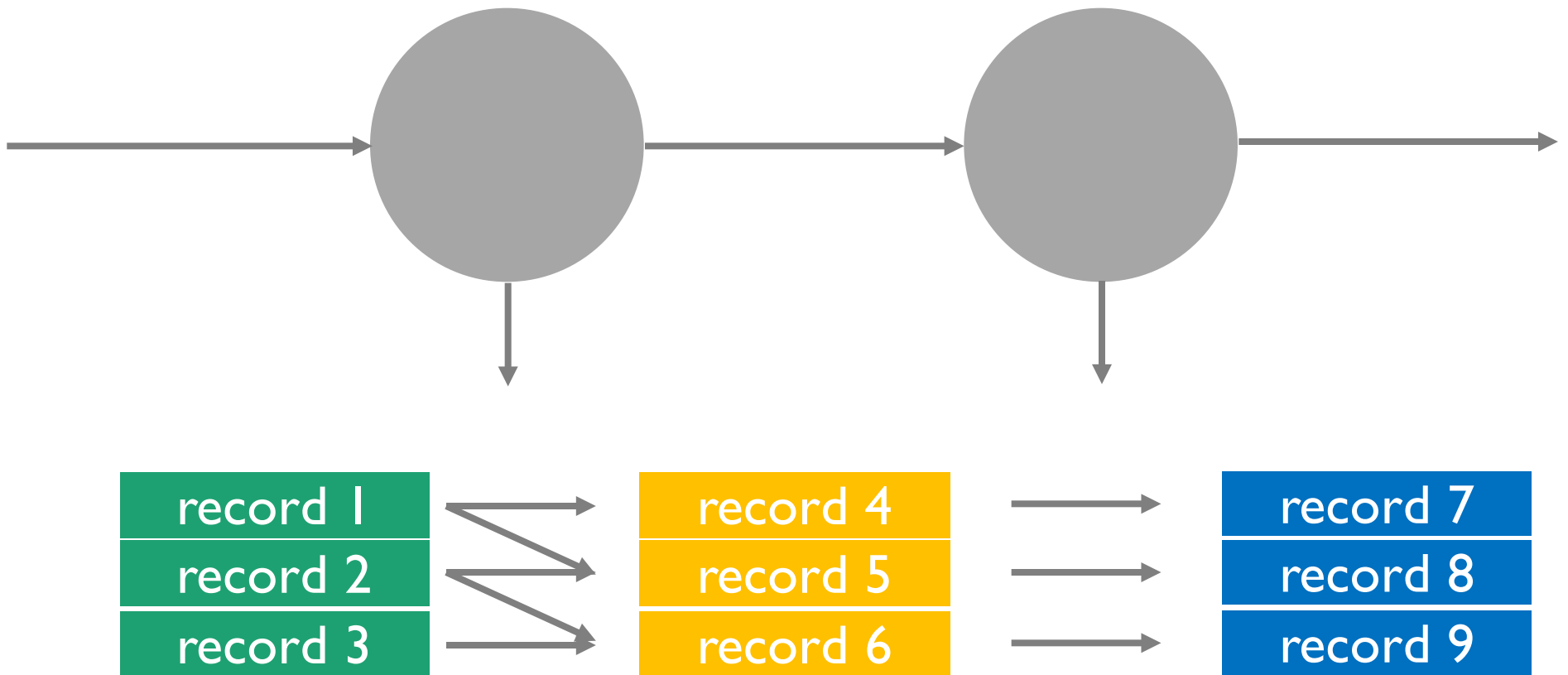
# Lineage



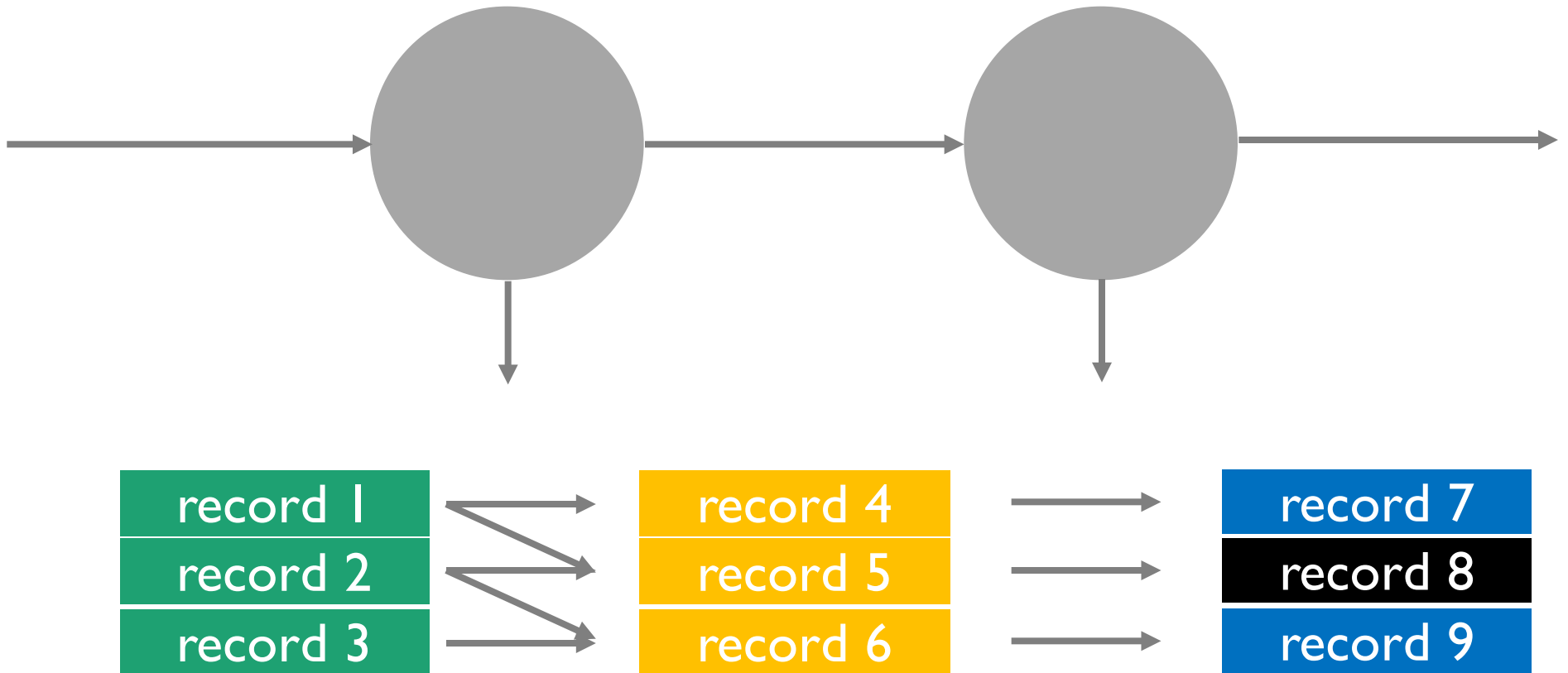
Which input records generated **record 5** ?



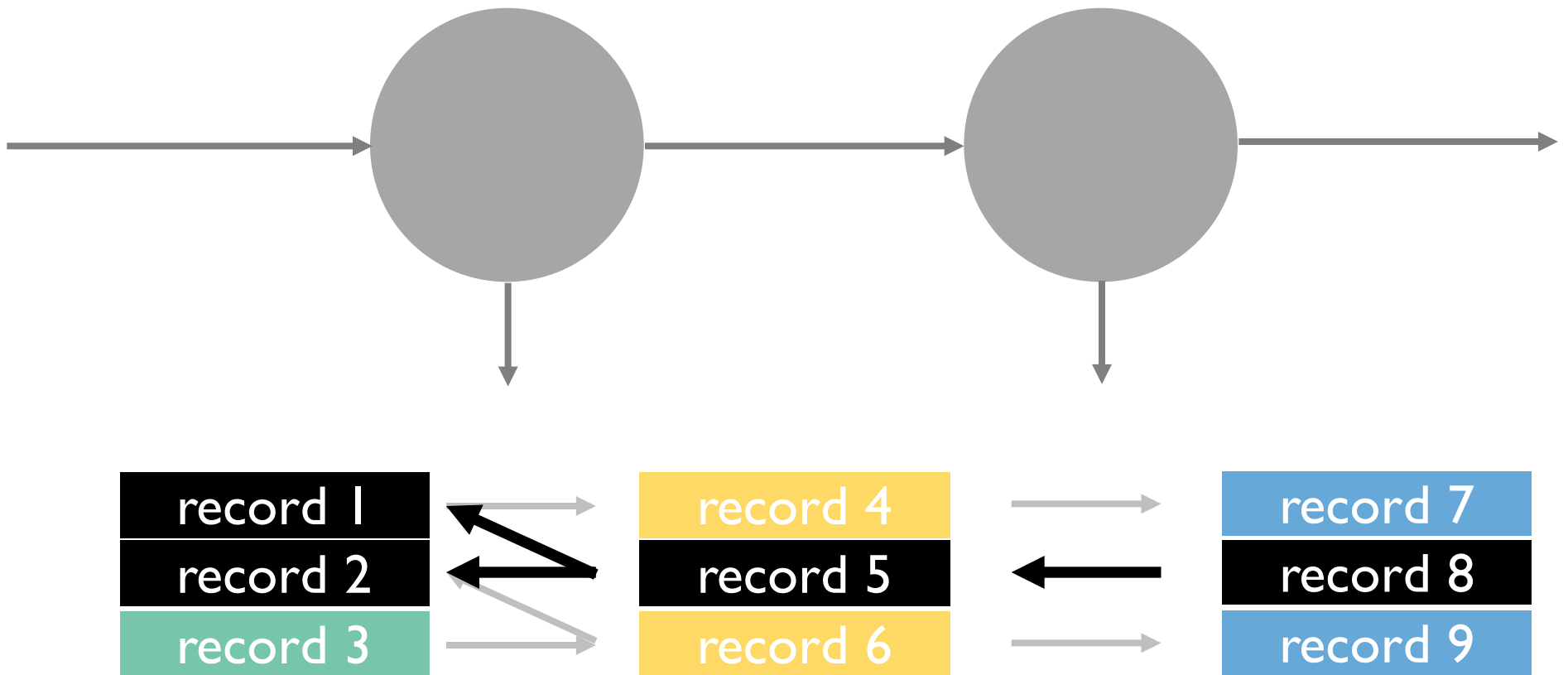
# Lineage



# Lineage



# Lineage



```
sum = 0
for cell in input:
    if cell.type == 'valid':
        sum += cell.value
return sum
```

What cells were read to compute the result?

Return all cell

What values were added to compute the result?

Return all valid cells

Compute  
framework can

```
UPDATE CUSTOMER SET  
C_BALANCE = C_BALANCE + ?  
WHERE C_ID = ? AND C_D_ID  
= ? AND C_W_ID = ? UPDATE  
CUSTOMER SET C_BALANCE =  
?, C_YTD_PAYMENT = ?,  
C_PAYMENT_CNT = ?, C_ORDER  
= ? WHERE C_W_ID = ? AND  
C_D_ID = ? AND
```

# Provenance refers to all of this

Depends on  
program and  
lineage queries  
(hard)



# Provenance

What is it?

**Projects@Columbia**

**Open Problems**



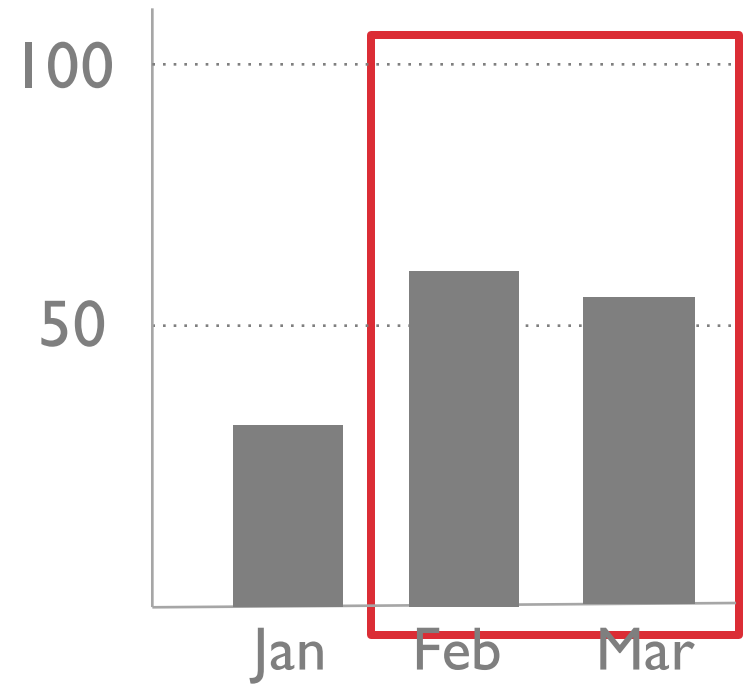
# Scorpion: Explain Data



analysis



SUM(sales)



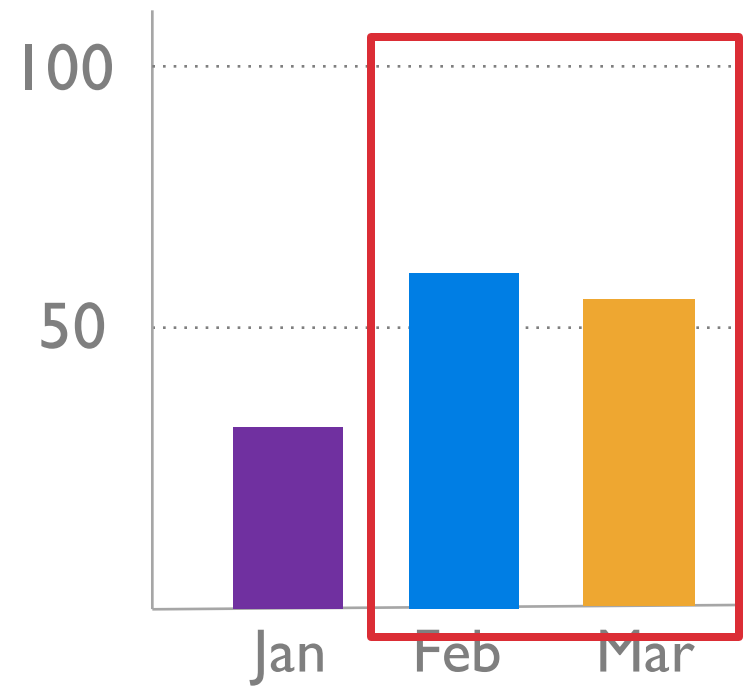
# Scorpion: Explain Data



analysis



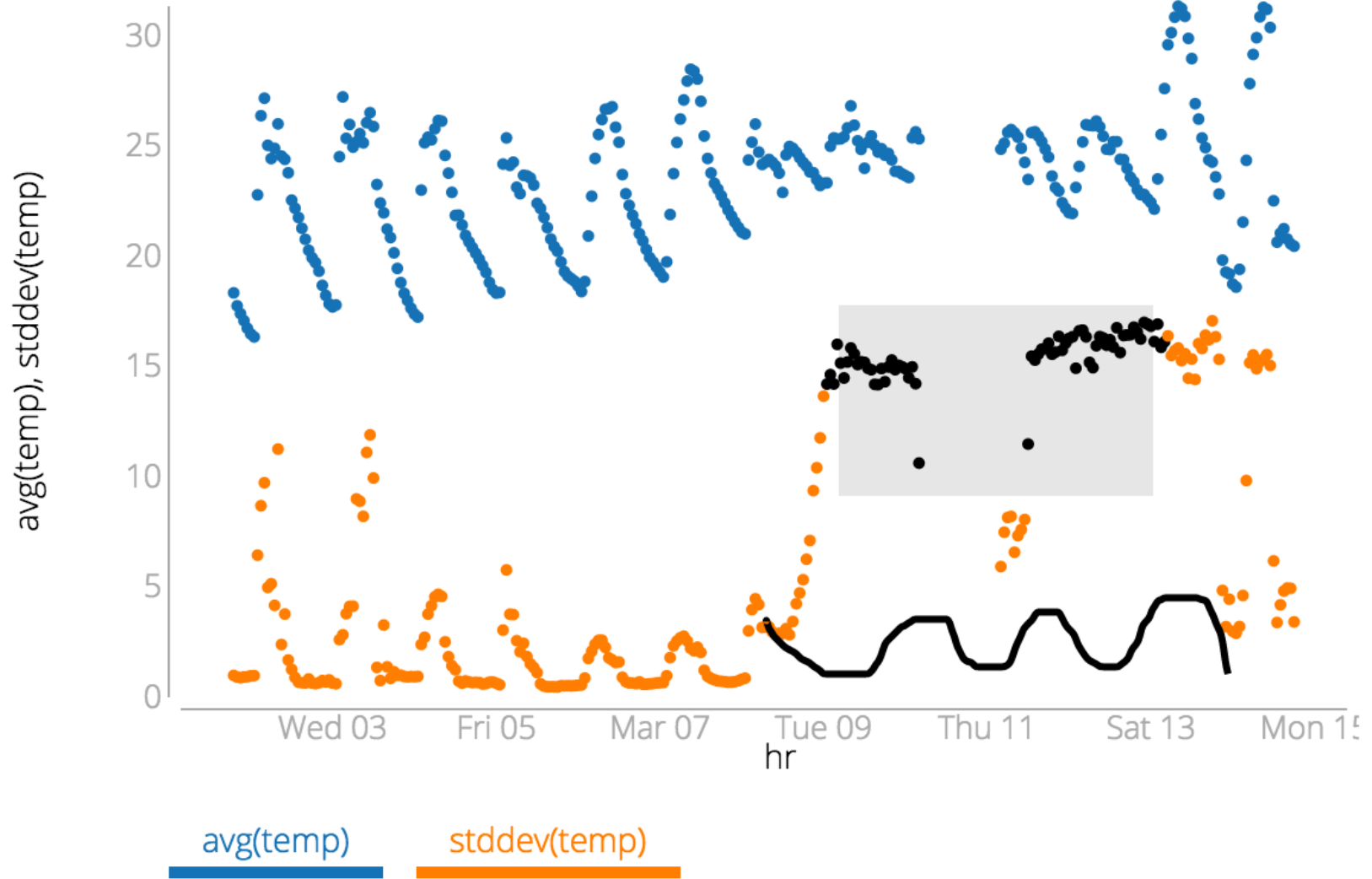
SUM(sales)



Visualization

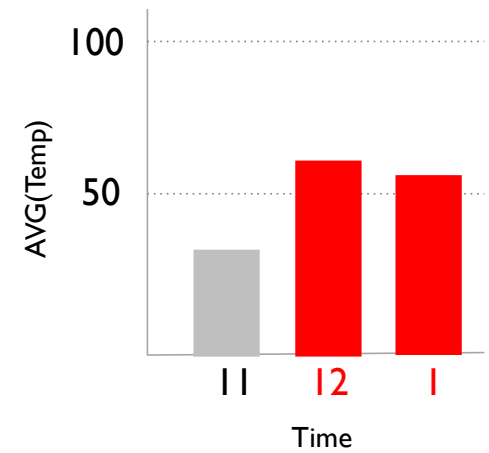
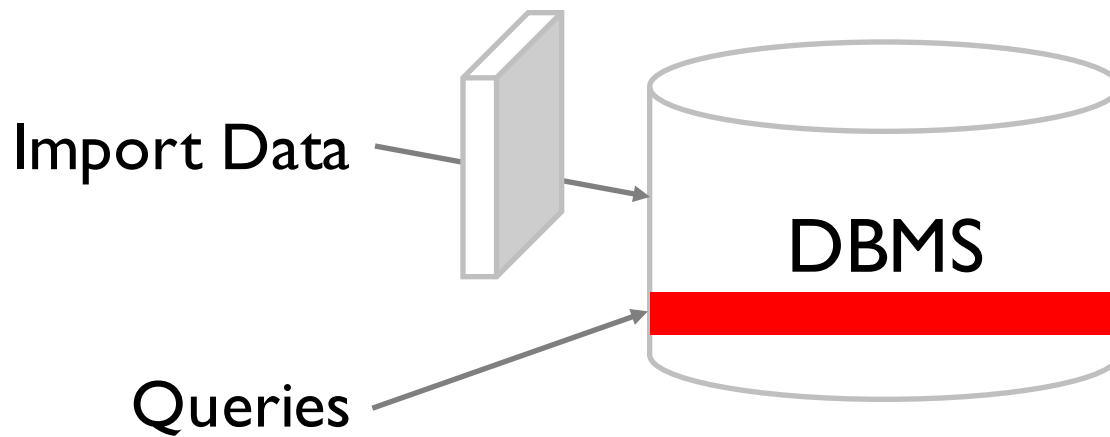
Show Rows

Show Query Form

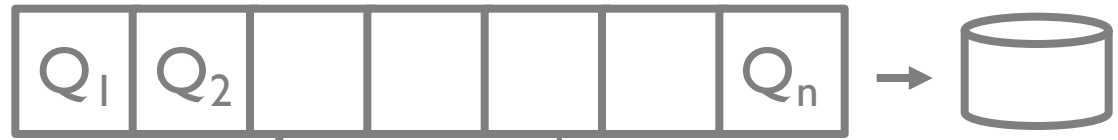


“Sensors near the window”

# QueryFracker: Explain Queries



True Query Log  
( $Qlog^*$ )



Corruption

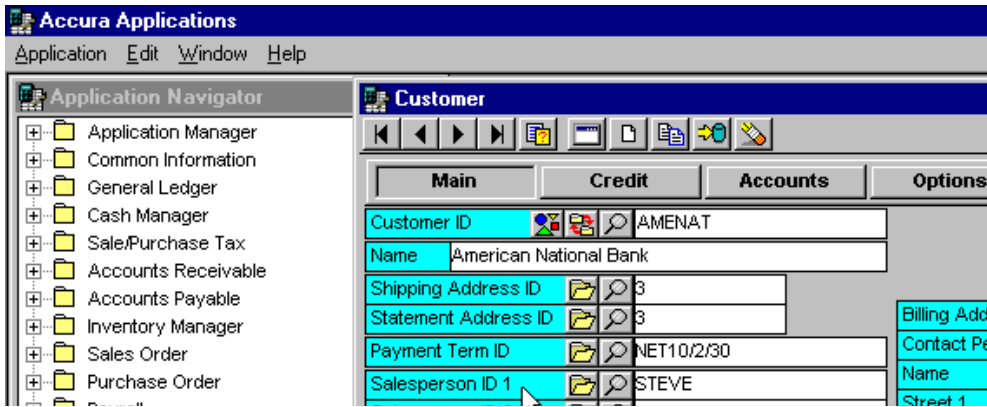
QueryFracker

Corrupt Query  
Log ( $Qlog$ )

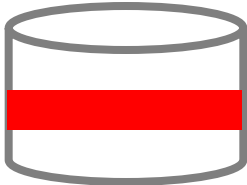


Goal

Given ( $Qlog$ , complaints)  $\rightarrow$  recover  $Qlog^*$



Application Server



Corrupt Query  
Log (Qlog)



provenance of the  
complaint records

Filtered Query  
Log

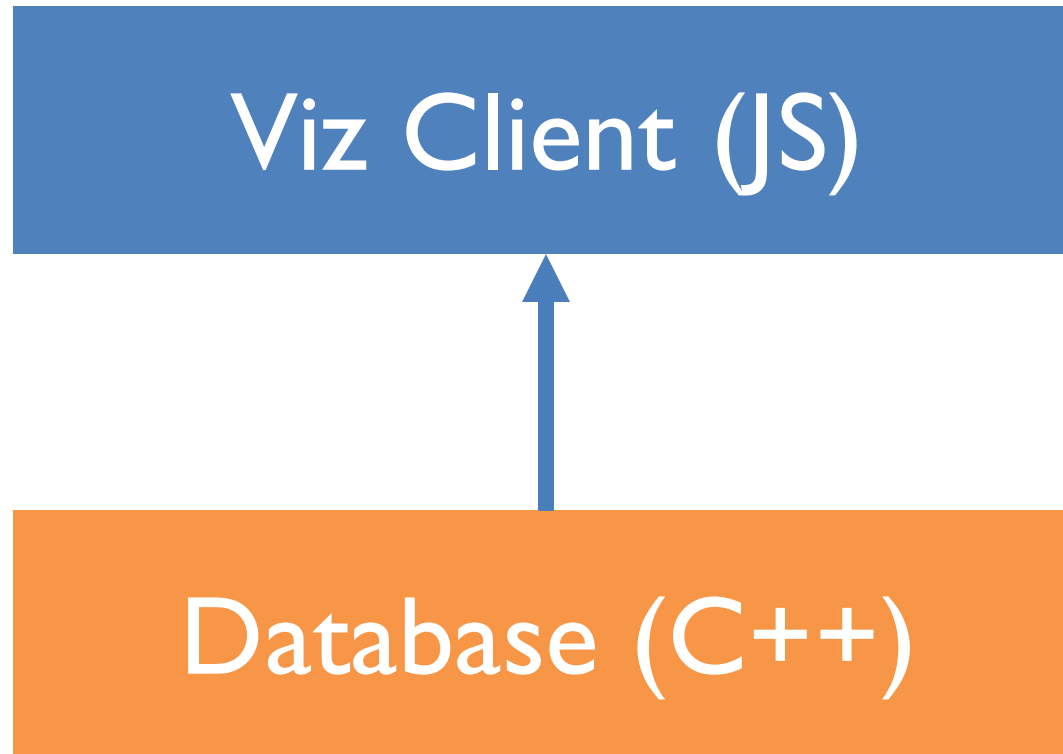


find repairs using  
MILP

Repair  
Recommendation



# Data Visualization Management System

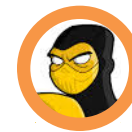




# Data Visualization Management System



Analysis Tools



Provenance

Performance-aware Design

DEclarative Visual Interaction Language (DEVIL)

Design-aware Optimization

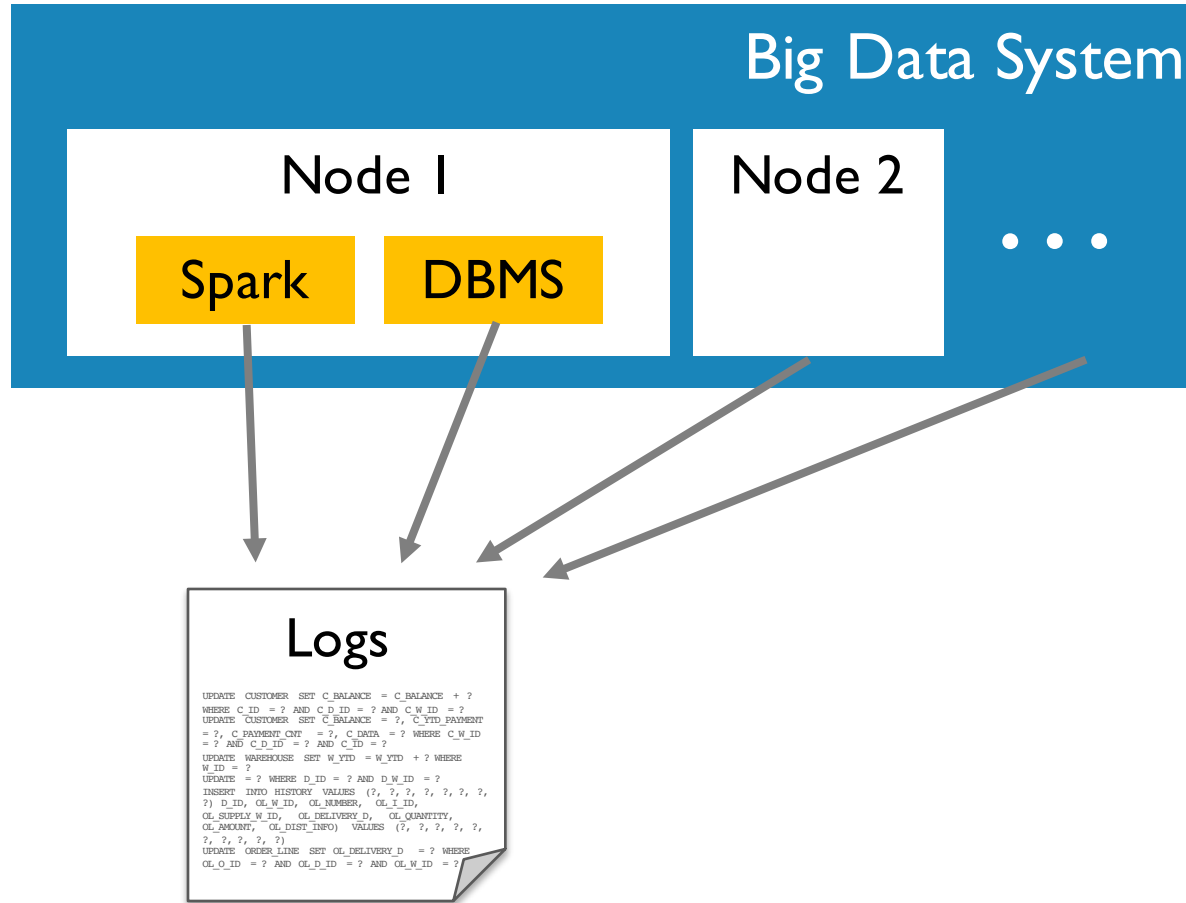
# Provenance

What is it?

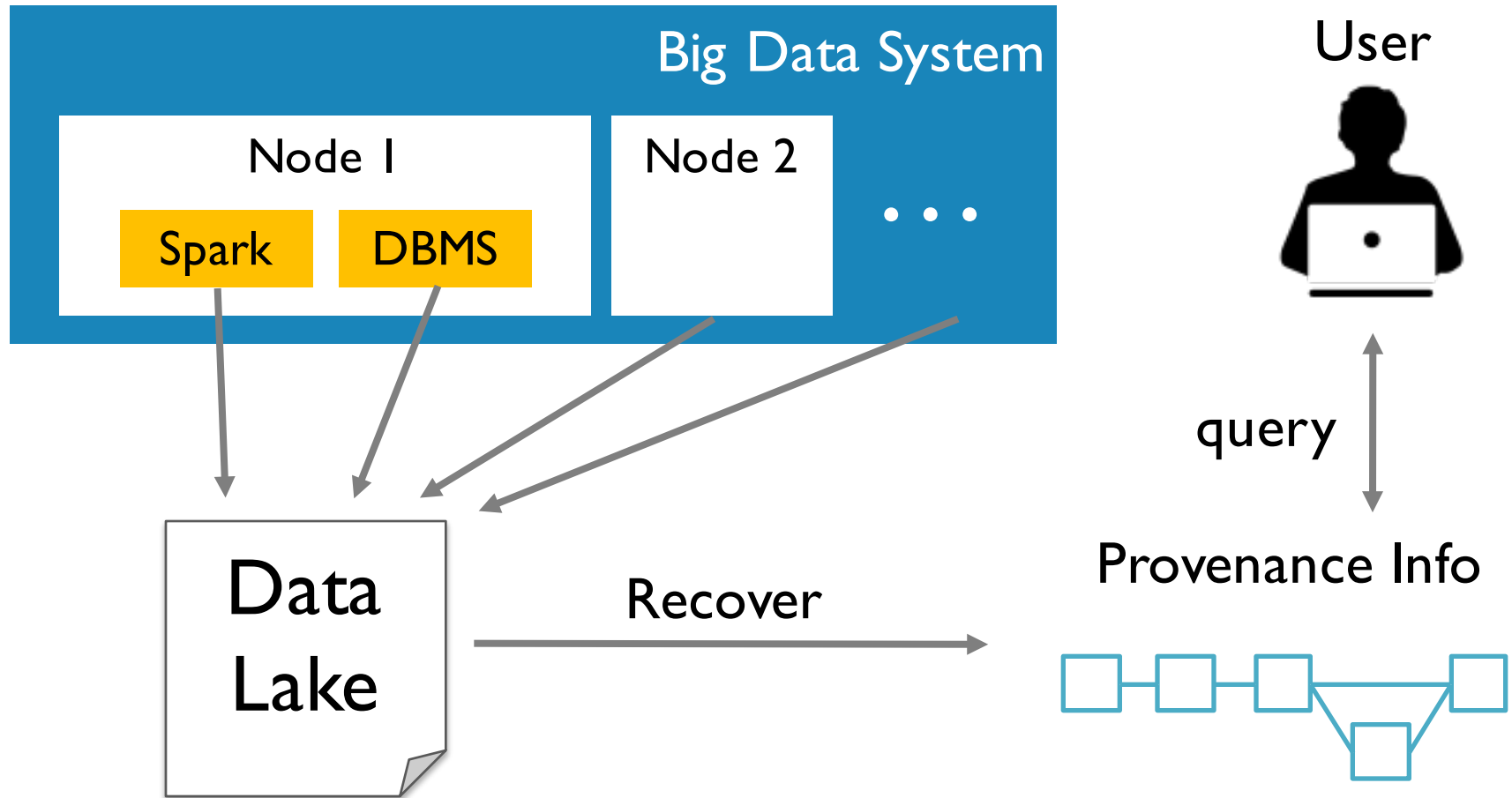
Projects@Columbia

**Open Problems**

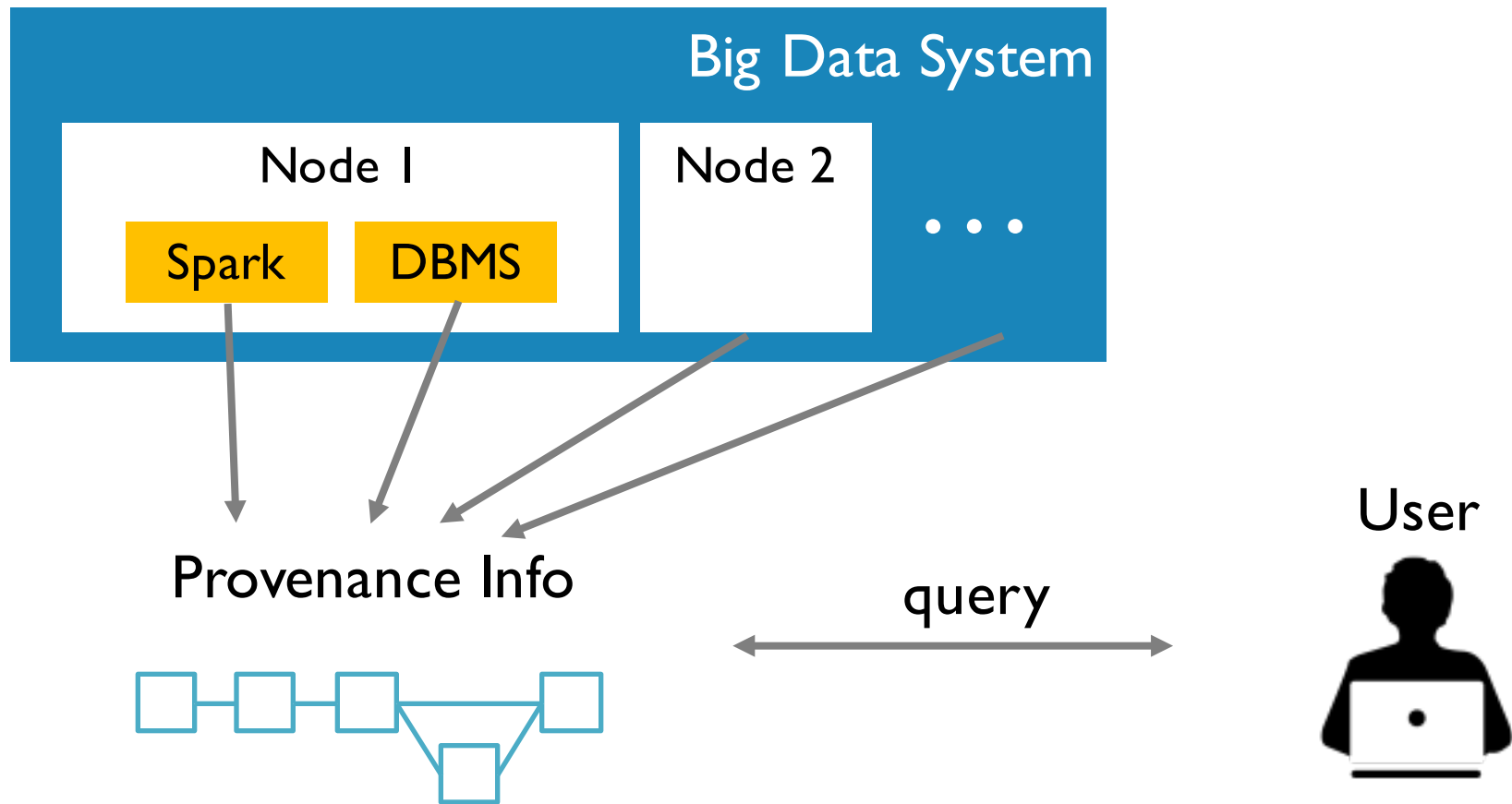
# Recovering Provenance



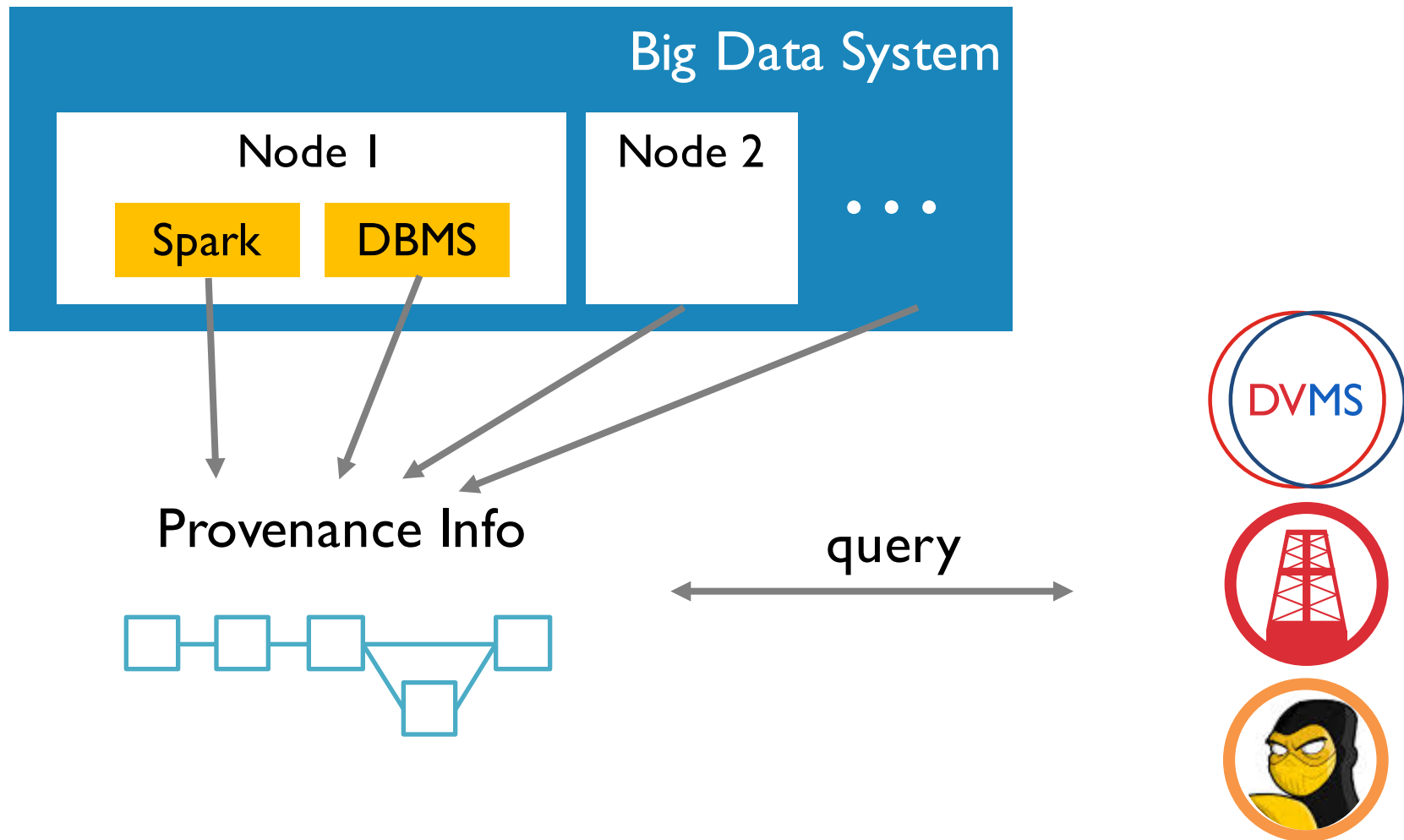
# Recovering Provenance



# Query-based Instrumentation



# Provenance-based Analysis



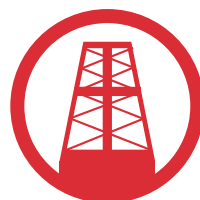
# eugenewu.net

ewu@cs.columbia.edu

## Explanation



## Cleaning



## Visualization



**CUDBG**

 **COLUMBIA UNIVERSITY**  
IN THE CITY OF NEW YORK