

Dirk Habich

# Scalable Data Management (SDM)

Modern Database System Architecture

# Previous focus

1

## Foundations

- Relational Data and SQL
- Traditional Database Architecture
- Transactions

2

## Parallelism

- Data-oriented Architecture
- Coarse- and Fine-grained parallelism
- Data Fragmentation

3

## Transactions

- Data Replication
- Distributed Transactions
- CAP-Theorem

4

## Additional

- Map/Reduce and Hadoop
- Spark and its Ecosystem
- Database in the cloud

# New focus – Modern Database Systems

1

## Foundations

- Relational Data and SQL
- Traditional Database Architecture
- Transactions

2

## Parallelism

- Data-oriented Architecture
- Coarse- and Fine-grained parallelism
- Data Fragmentation

3

## Transactions

- Data Replication
- Distributed Transactions
- CAP-Theorem

4

## Additional

- Map/Reduce and Hadoop
- Spark and its Ecosystem
- Database in the cloud

# Importance of Complex Data Analyses



**Industry 4.0**



**Electronic Patient Records**



**Business Intelligence**

**Data is the new oil!**



Gerd

# Trends to Face

## Application trends

Complex analytical queries

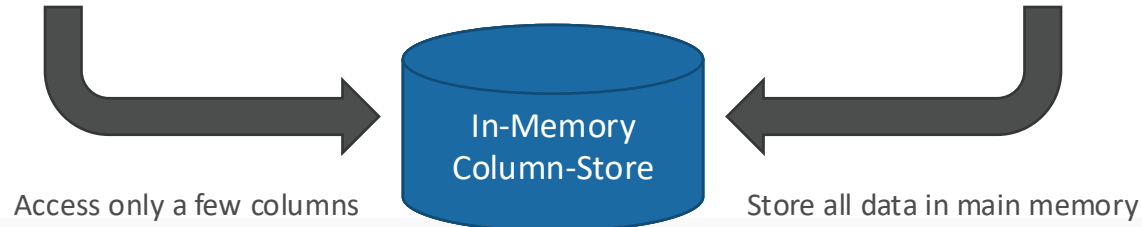
### Example: Star Schema Benchmark (SSB) Query 3.1

```
SELECT c_nation, s_nation, d_year, SUM(lo_revenue)
FROM customer, lineorder, supplier, date
WHERE lo_custkey = c_custkey
      AND lo_suppkey = s_suppkey
      AND lo_orderdate = d_datekey
      AND c_region = 'ASIA' AND s_region = 'ASIA'
      AND d_year >= 1992 AND d_year <= 1997
GROUP BY c_nation, s_nation, d_year
ORDER BY d_year ASC, revenue DESC;
```

## Hardware trends

Main  
Memory

Increasingly large capacities



# Recap: Traditional Row-based Storage

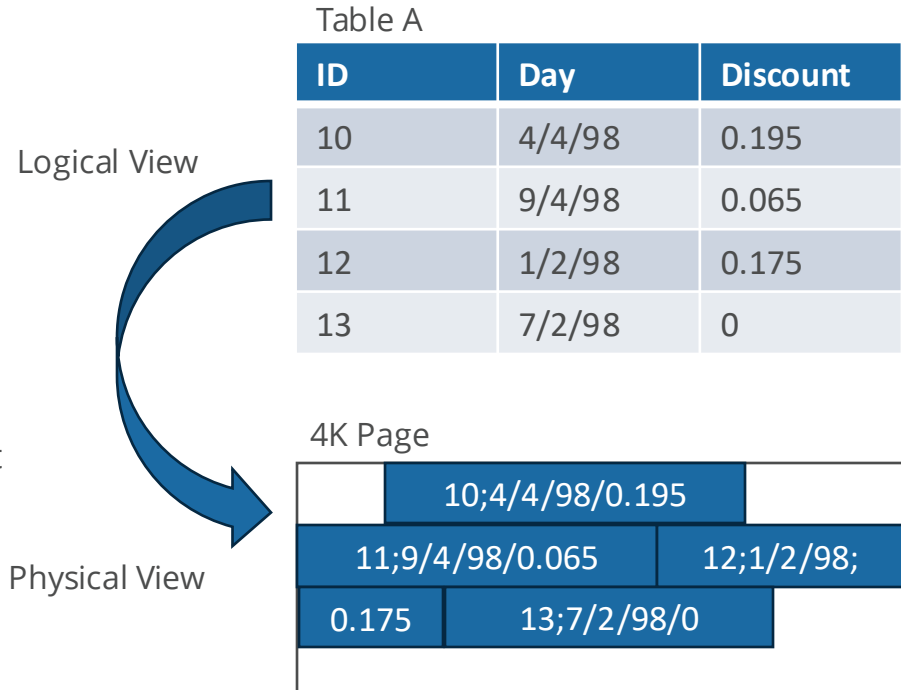
## Physical Data Storage

- Properties: The unit of a tuple with all attributes is preserved

## Pro (+) / Cons (-)

- + Easy to add/modify a tuple
- Even if only a few attributes are needed, all attributes have to be read  
→ unnecessary I/O-cost
- Compression is hard since consecutive values are from different domains (e.g., name, age, street addr, zip code, etc).

## Illustration



# Modern Column-based Storage

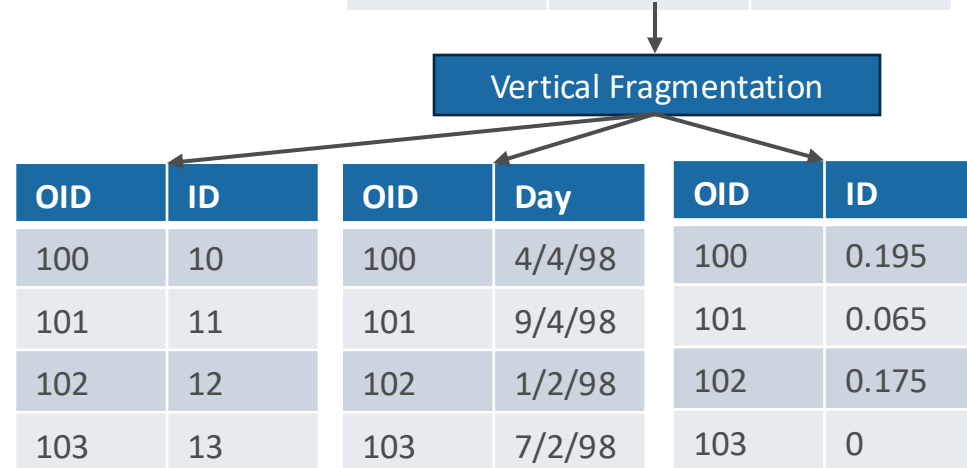
## Physical Data Storage

- Storage model based on vertical fragmentation
- Instead of storing all attributes of each relational tuple in one record, each column is stored in a separate table called BAT
- BAT = Binary Association Tuple
  - The left column is called head and is required to reconstruct the relational tuples
  - The right column is called tail and is holding the actual attribute

## Illustration

Table A

ID	Day	Discount
10	4/4/98	0.195
11	9/4/98	0.065
12	1/2/98	0.175
13	7/2/98	0



## Physical Data Storage

- Every relational table is internally represented as a collection of BATs
- For a table of k attributes, there exists k BATs and each BAT stores the attribute as (OID, value) pairs
- System generated OID values identifies the relational tuple that the attribute value belongs to, i.e., all attribute values of a single tuple are assigned the same OID.
- BATs are stored in regular 4K-pages → each BAT considers as single table with two attributes

## Illustration

OID	ID	OID	Day	OID	ID
100	10	100	4/4/98	100	0.195
101	11	101	9/4/98	101	0.065
102	12	102	1/2/98	102	0.175
103	13	103	7/2/98	103	0

4K Page

	100/10	101/11
102/12		



# Modern Column-based Storage/3

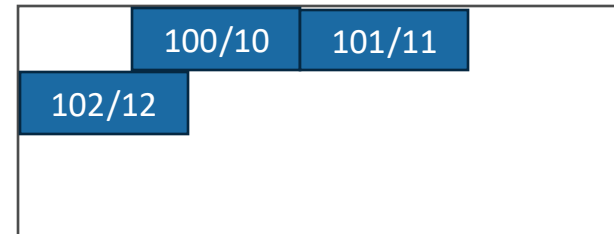
## Pro / Cons

- + Only need to read in relevant data, more efficient scan operators
- tuple writes require multiple accesses
- Reading all columns of a single row requires expensive row-reconstruction (1:1 join)
- High storage overhead due to additional OID attribute

## Illustration

OID	ID	OID	Day	OID	ID
100	10	100	4/4/98	100	0.195
101	11	101	9/4/98	101	0.065
102	12	102	1/2/98	102	0.175
103	13	103	7/2/98	103	0

4K Page





# Columnar Storage Optimization

## (1) Elimination of OIDs

- Implicit instead of explicit OIDs
- The order of the tuples is now crucial

## Advantages

- Each attribute is now a simple array without any additional information
- Each array has a unique data type
- OID could be on-the-fly generated based on the position within the array

## (2) Compression

- Compression easily possible for each single array

OID	ID	OID	Day	OID	ID
100	10	100	4/4/98	100	0.195
101	11	101	9/4/98	101	0.065
102	12	102	1/2/98	102	0.175
103	13	103	7/2/98	103	0



ID	Day	ID
10	4/4/98	0.195
11	9/4/98	0.065
12	1/2/98	0.175
13	7/2/98	0

# Compression Overview

## Storing Relational Data Column-wise

- Each column as a sequence of values

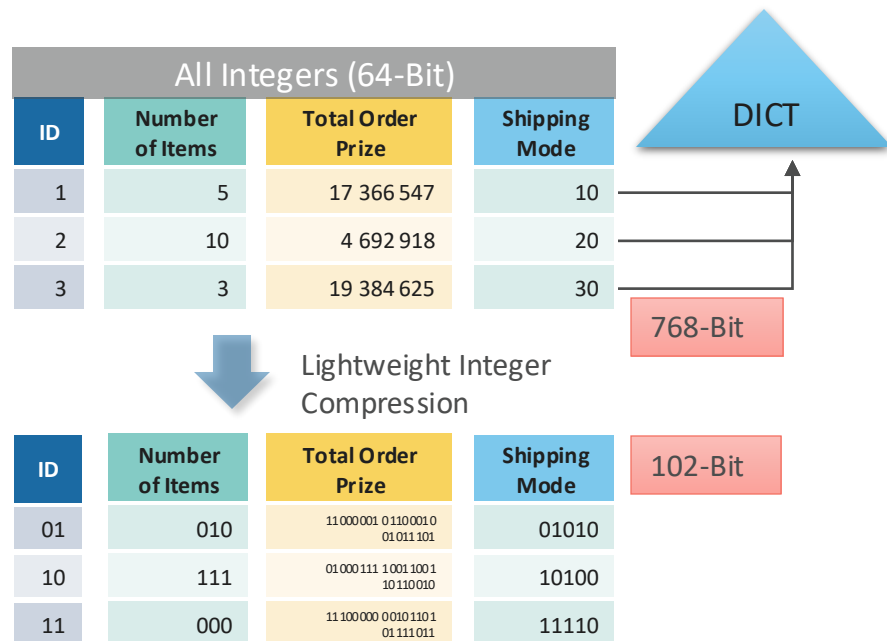
Integers		Fixed Point	Strings
ID	Number of Items	Total Order Prize	Shipping Mode
1	5	173 665.47	"Mail"
2	10	46 929.18	"Air"
3	3	193 846.25	"Rail"

✓ Only need to read in relevant data attributes during query processing

– Tuple writes require multiple accesses

## Reducing Memory Footprint

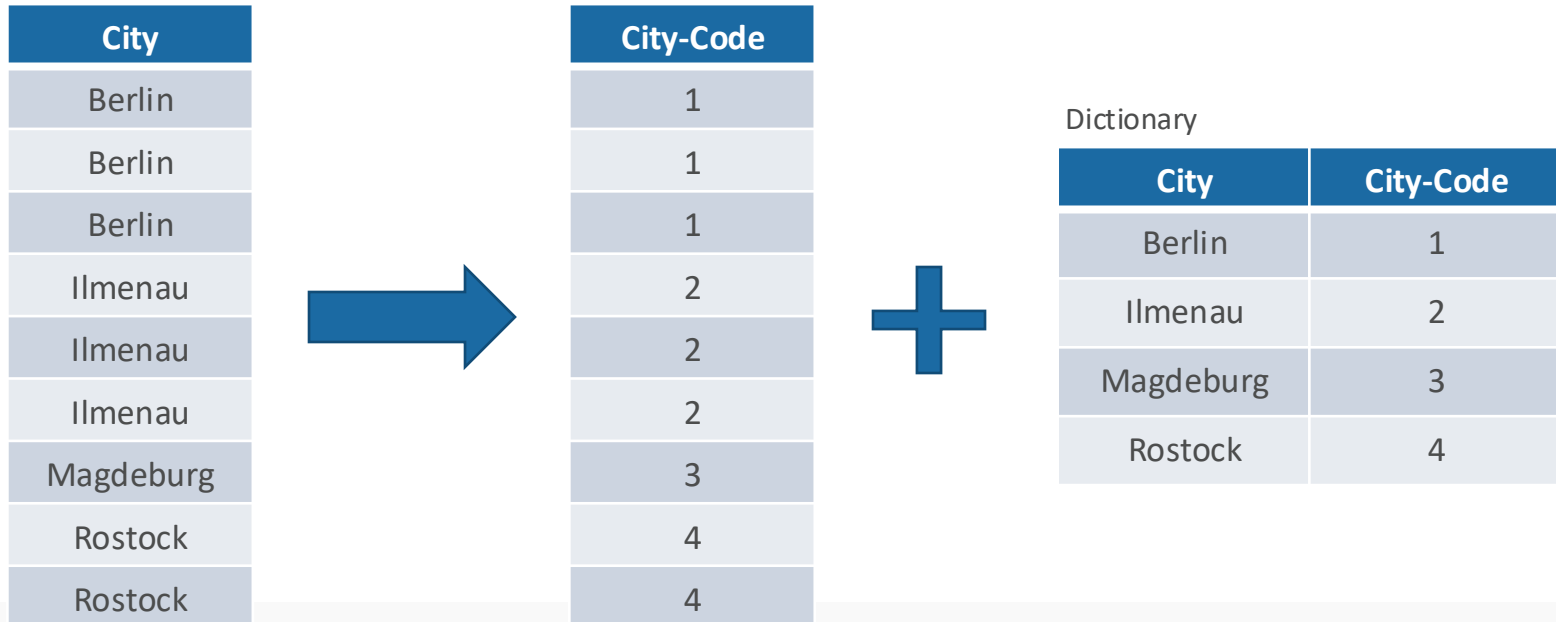
- All logical data is mapped to integer data on the physical level



# Preprocessing Technique: DICT

## Dictionary Coding – DICT

- replaces each (string-)value by its unique key in a dictionary
- achieve small integer values



# Preprocessing Technique: FOR

## Frame-of-Reference – FOR

- represents each value as the difference to a certain given reference value (FOR)
- Achieve smaller integer values

Preis
45
54
48
55
51
53
40
50
49

FOR = 40




Preis (FOR)
5
14
8
15
11
13
0
10
9

Reference values (FOR) has to be stored to be able to decompress data

# Preprocessing Technique: DELTA

## DELTA Coding – DELTA

- represents each value as the difference to its predecessor value
- achieve smaller integer values
- Sorting helpful for optimized compression ratio

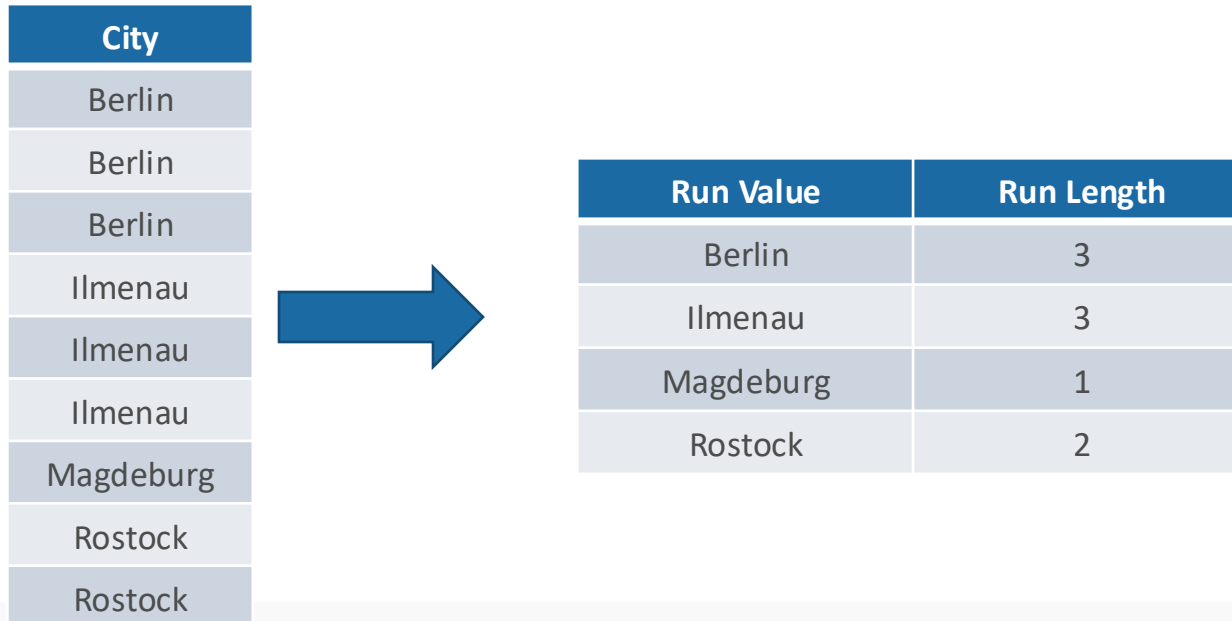


Preis	DELTA
10	10
15	5
20	5
21	1
23	2
25	2
30	5

# Compression Technique: RLE

## Run Length Encoding – RLE

- tackles uninterrupted sequences of occurrences of the same value, so called runs
- each run is represented by its value and length

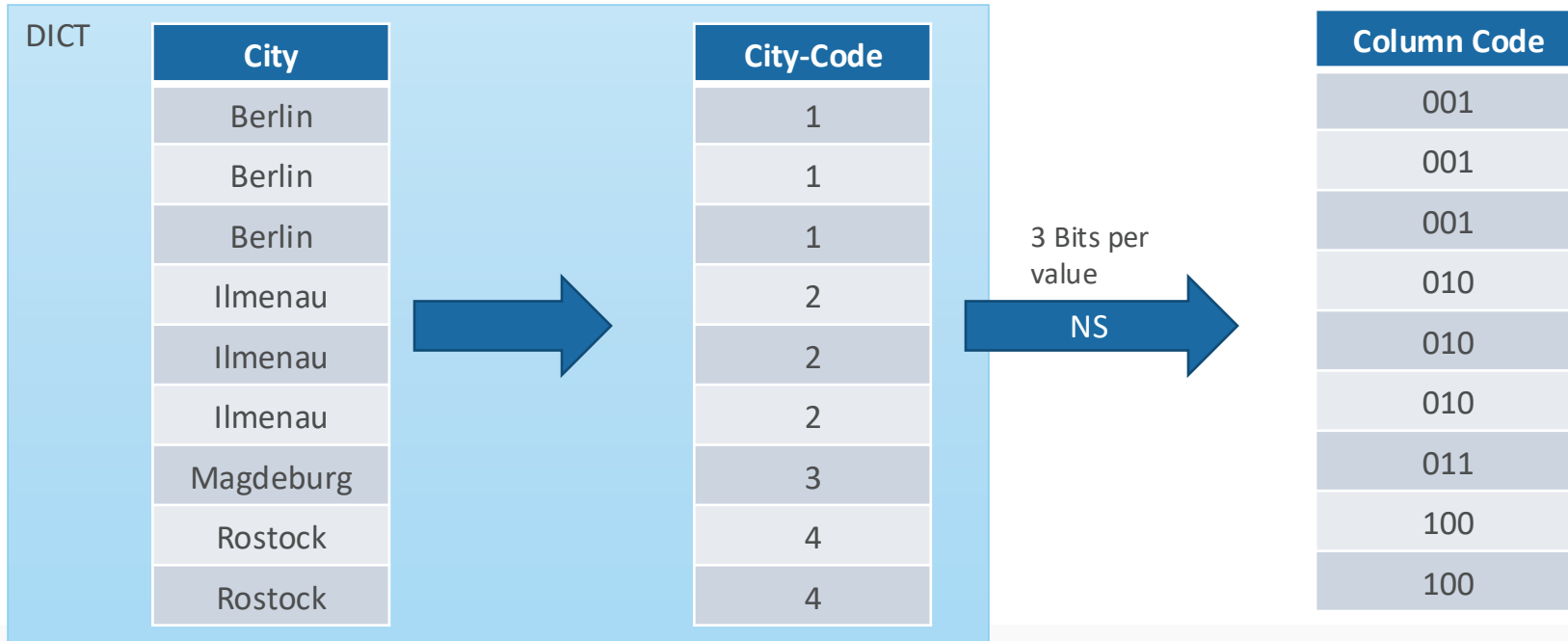




# Compression Technique: NS

## Null Suppression – NS

- addresses the physical level of bits or bytes to reduce the number of bits per value



# Null Suppression Algorithms

## Bit-Aligned

- compress an integer value with a minimal number of bits

## Byte-Aligned

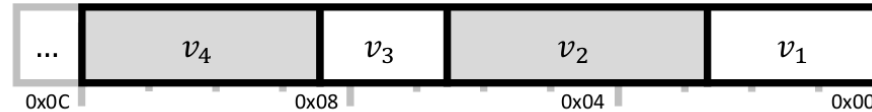
- Compress an integer value with a minimal number of bytes

## Word-Aligned

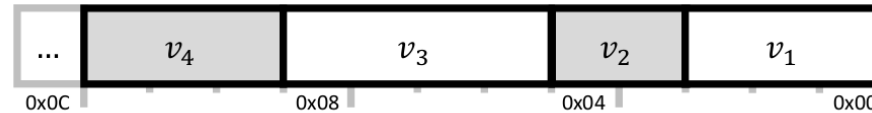
- Encode as many integer values as possible into 32-bit or 64-bit words

### Alignment Examples

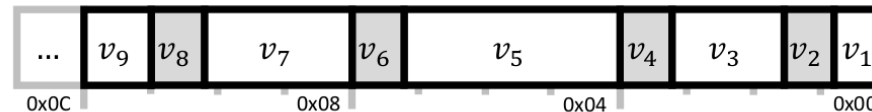
Bit-  
Aligned



Byte-  
Aligned



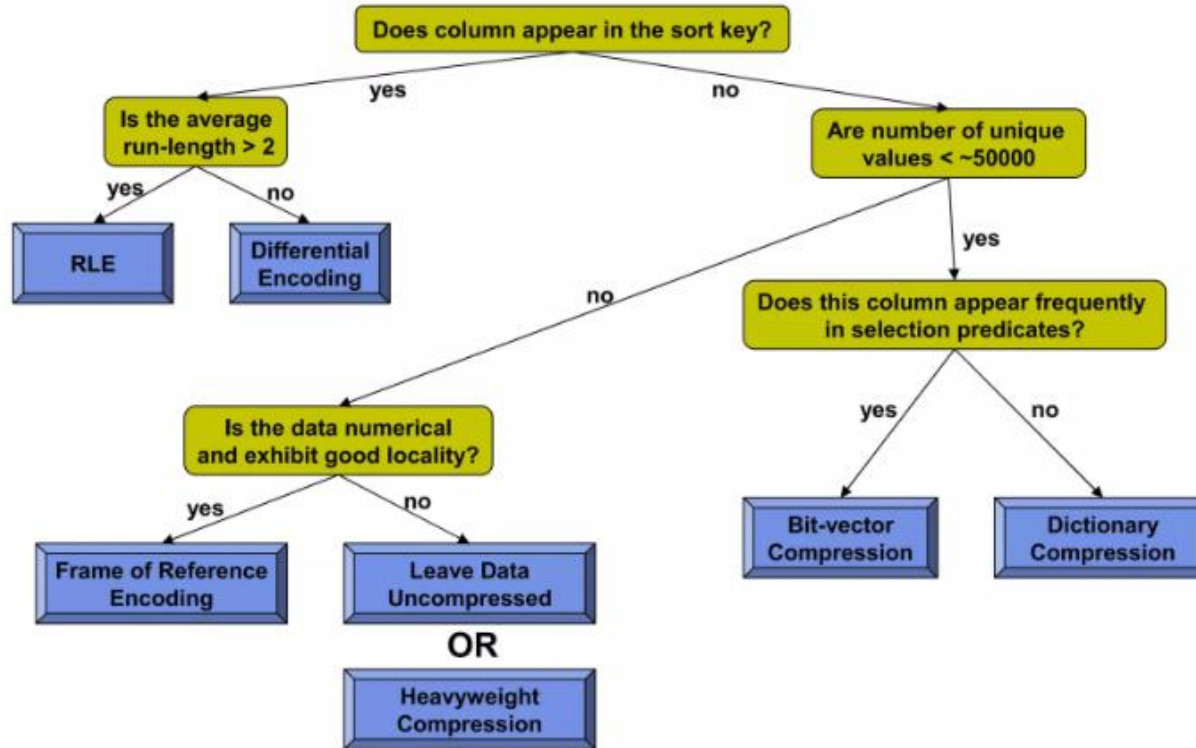
Word-  
Aligned



# Compression Techniques Summary

RLE	DELTA	FOR	DICT	NS
<b>Run Length Encoding</b>  Replace run by value & length	<b>Delta Coding</b>  Replace data elem. by difference to predecessor	<b>Frame-of-Reference</b>  Replace data elem. by difference to reference value	<b>Dictionary Coding</b>  Replace data elem. by key in dictionary	<b>Null Suppression</b>  Eliminate leading zero-bits
Repetitions	Sorting	Value range	# distinct values	Small integers
<b>Logical</b> natural numbers <i>preprocessing</i>				<b>Physical</b> bits and bytes <i>actual compression</i>
Level data role				

# What Compression Scheme To Use?



# Storage Summary

## Most Important Properties

- Every column is encoded as a sequence of integer values
- Every column is stored in a contiguous memory area
- Heavily use of integer compression
  - to fit more data in main memory
  - To increase effective bandwidth

## Impact on Query Processing

- Only relevant columns have to be accessed
- Query processing is mostly done on integer sequences



# Processing Model

# Processing Model

## A processing model

- Defines how the system execute a query plan and moves data from one operator to the next

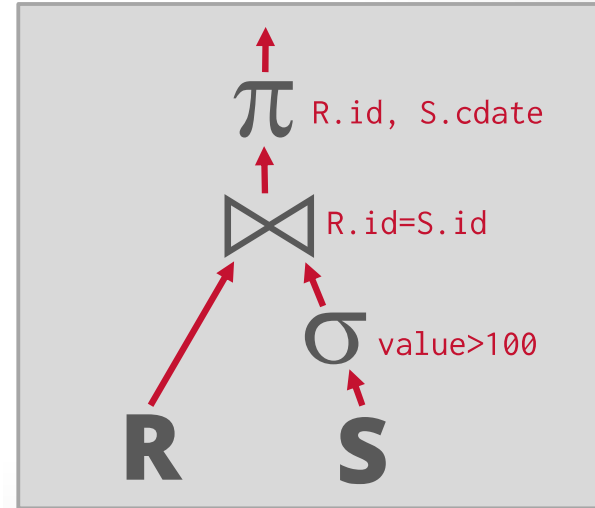
## Each processing model

- Is comprised of two types of execution paths
  - Control Flow: How the system invokes an operator
  - Data Flow: How an operator send its result

## The output

- of an operator can be either whole tuples (row-storage) or subset of columns (columnar storage)

```
SELECT R.id, S.cdate  
FROM R JOIN S  
      ON R.id = S.id  
WHERE S.value > 100
```



# Evolution of Processing Models





## Each query plan operator implements a Next() function.

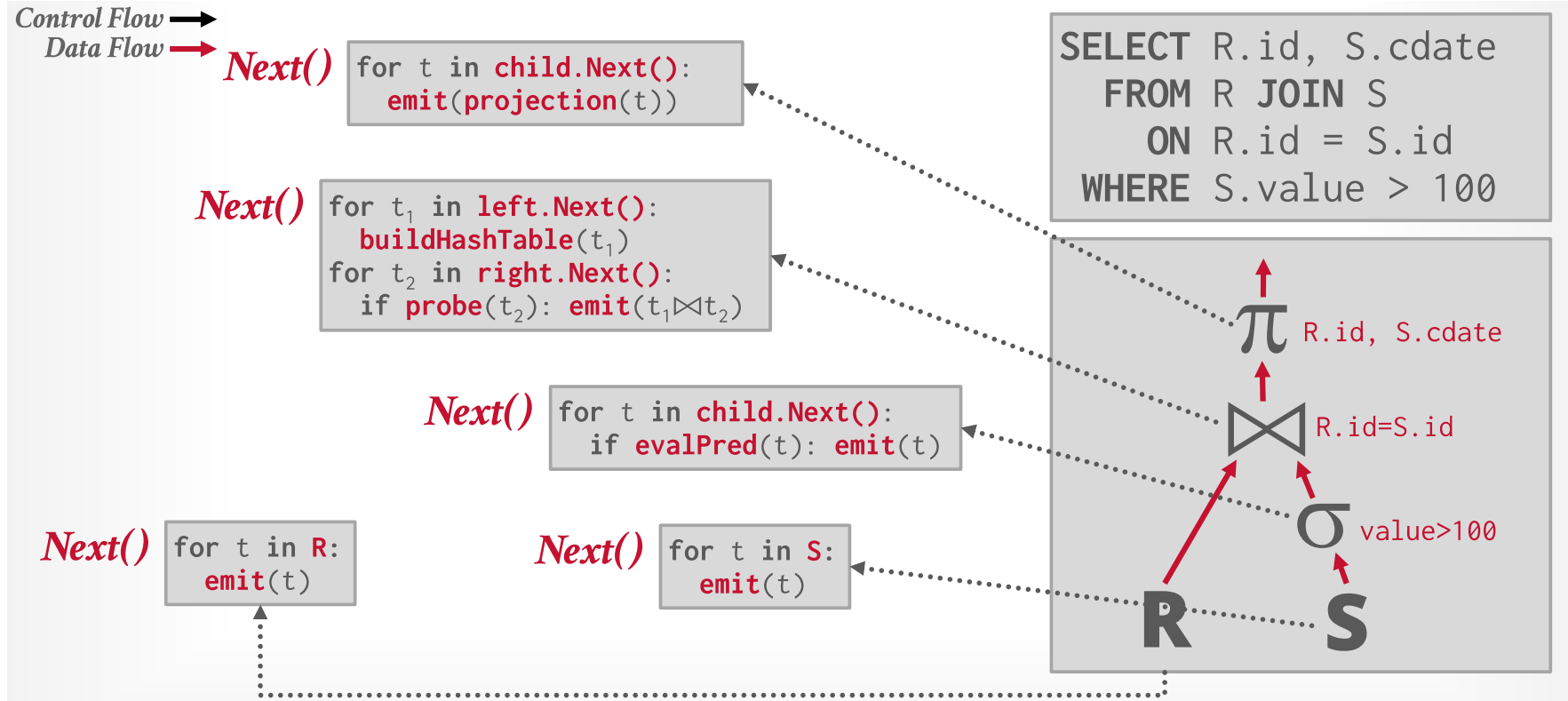
- On each invocation, the operator returns either a single tuple or a EOF marker if there are no more tuples.
- The operator implements a loop that calls next on its children to retrieve their tuples and then process them.

## Each operator implementation

- also has Open() and Close() functions.
- Analogous to constructors/destructors, but for operators.

## Also called Volcano or Pipeline Model

# Iterator Model - Example



Control Flow  $\rightarrow$   
Data Flow  $\rightarrow$

**1**

```
for t in child.Next():  
    emit(projection(t))
```

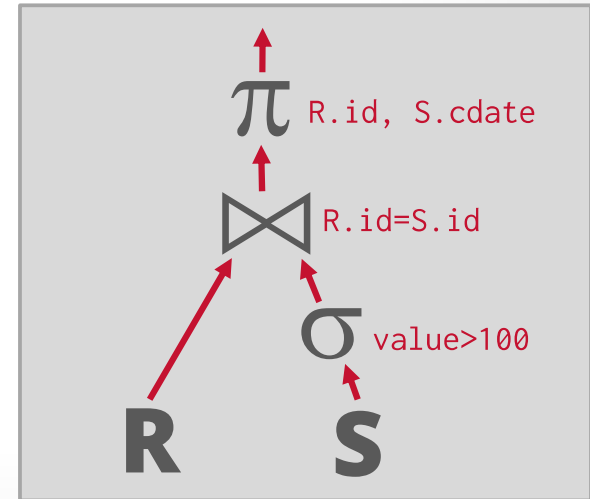
```
for t1 in left.Next():  
    buildHashTable(t1)  
for t2 in right.Next():  
    if probe(t2): emit(t1 ⋈ t2)
```

```
for t in child.Next():  
    if evalPred(t): emit(t)
```

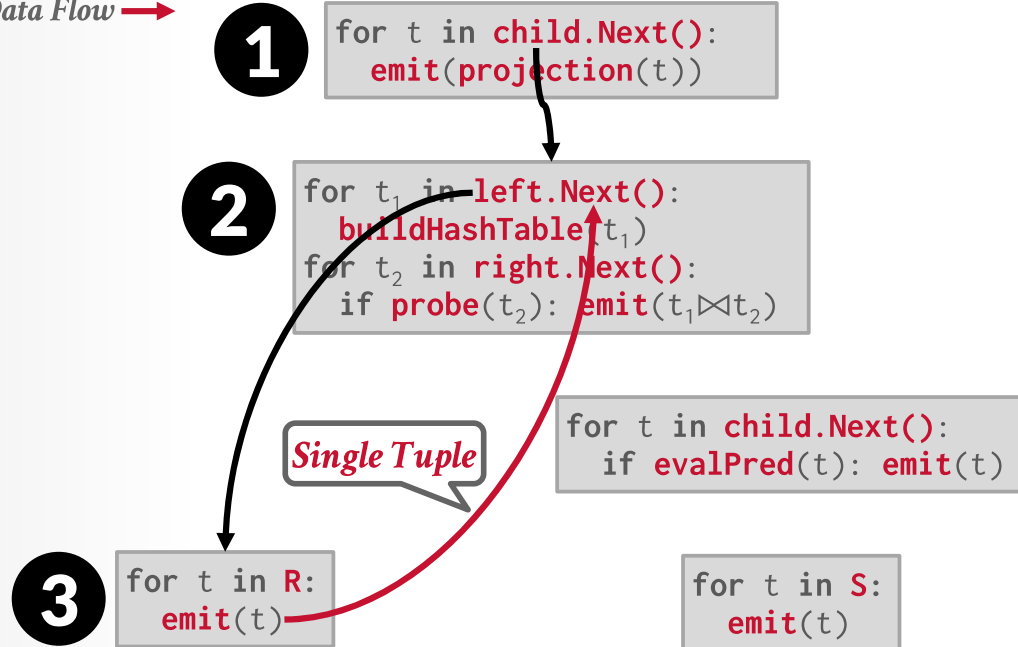
```
for t in R:  
    emit(t)
```

```
for t in S:  
    emit(t)
```

```
SELECT R.id, S.cdate  
FROM R JOIN S  
ON R.id = S.id  
WHERE S.value > 100
```

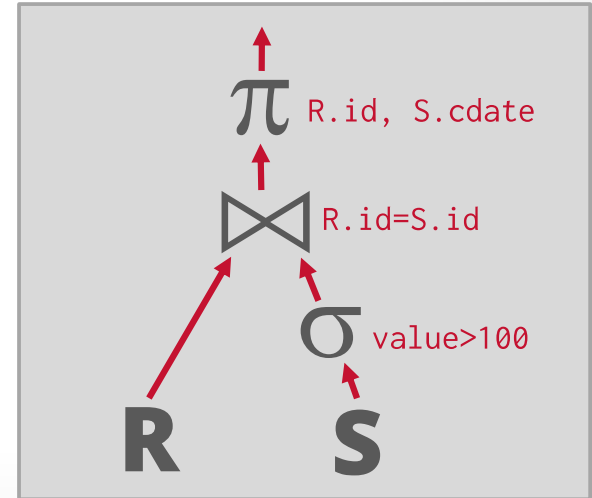


Control Flow  $\blackrightarrow$   
Data Flow  $\color{red}\rightarrow$

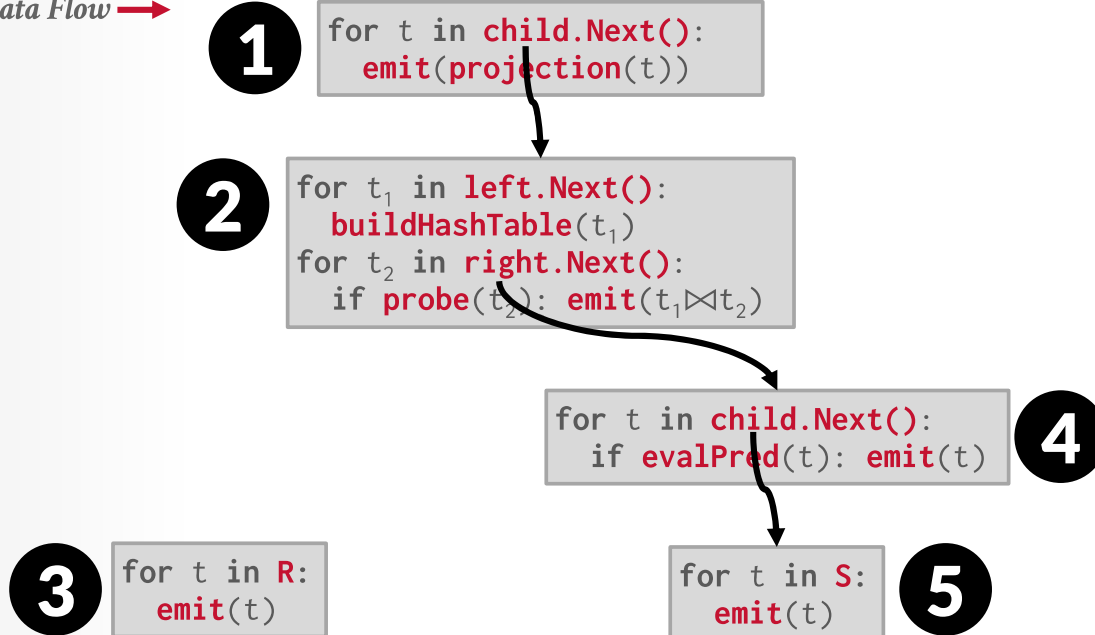


```

SELECT R.id, S.cdate
FROM R JOIN S
ON R.id = S.id
WHERE S.value > 100
  
```

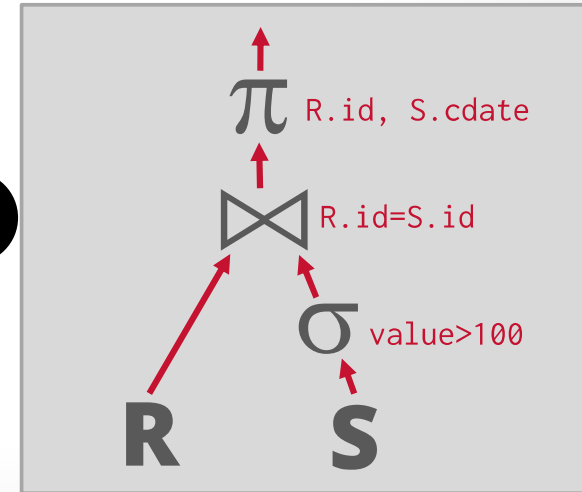


Control Flow →  
Data Flow →

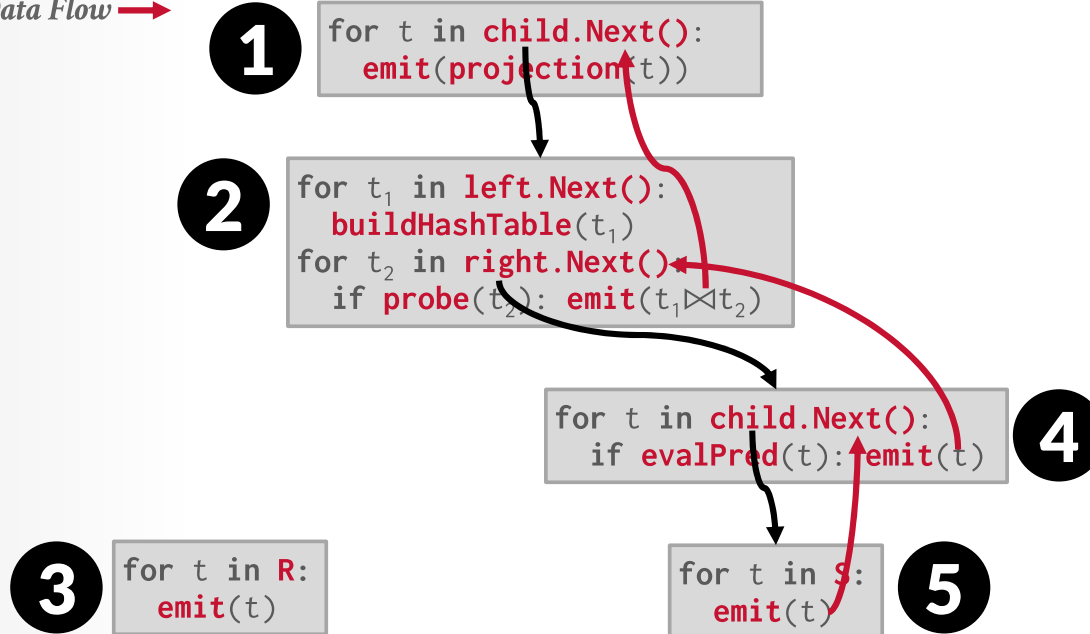


```

SELECT R.id, S.cdate
FROM R JOIN S
      ON R.id = S.id
WHERE S.value > 100
  
```

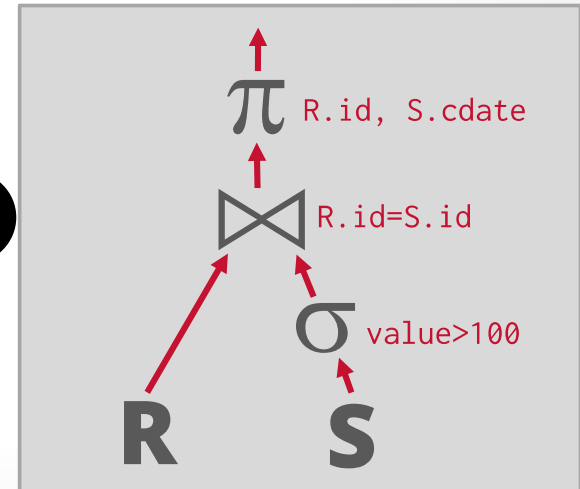


Control Flow  $\rightarrow$   
Data Flow  $\rightarrow$

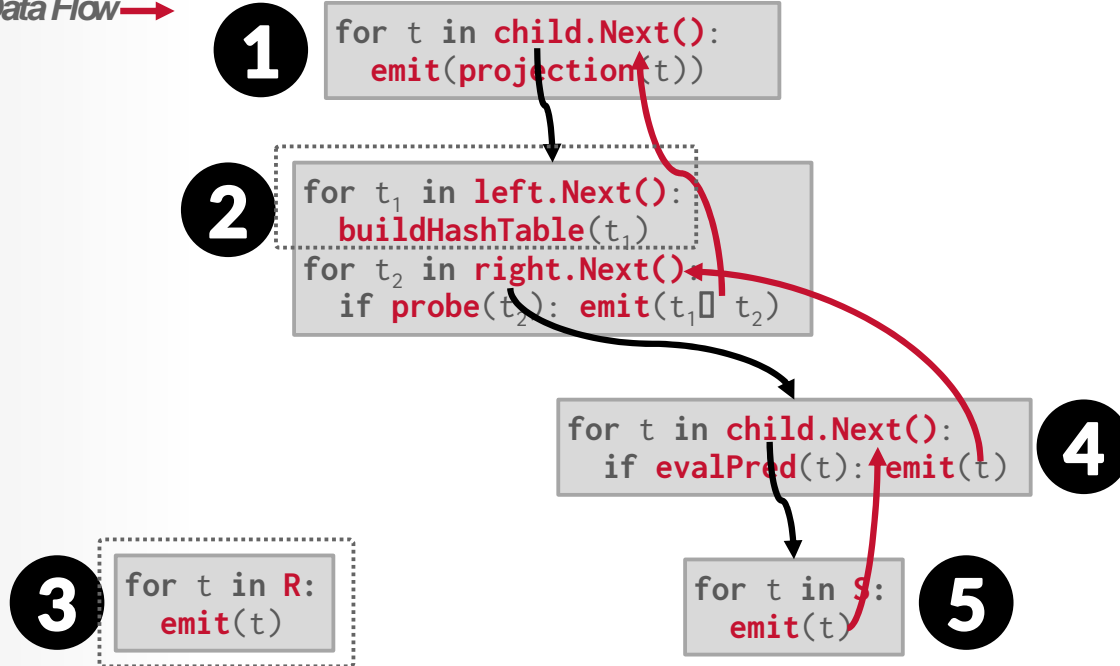


```

SELECT R.id, S.cdate
FROM R JOIN S
      ON R.id = S.id
WHERE S.value > 100
  
```

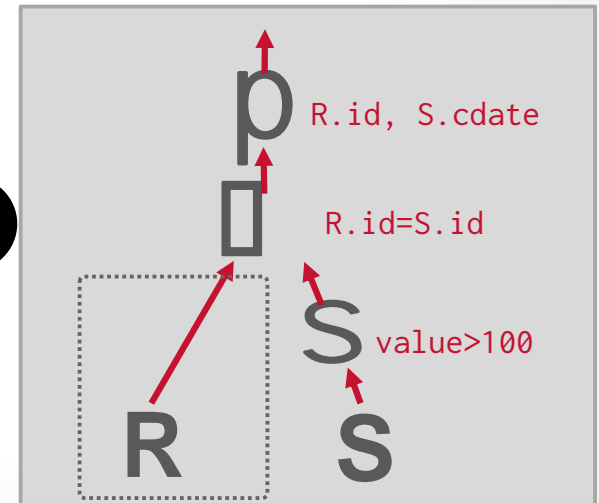


Control Flow →  
Data Flow →



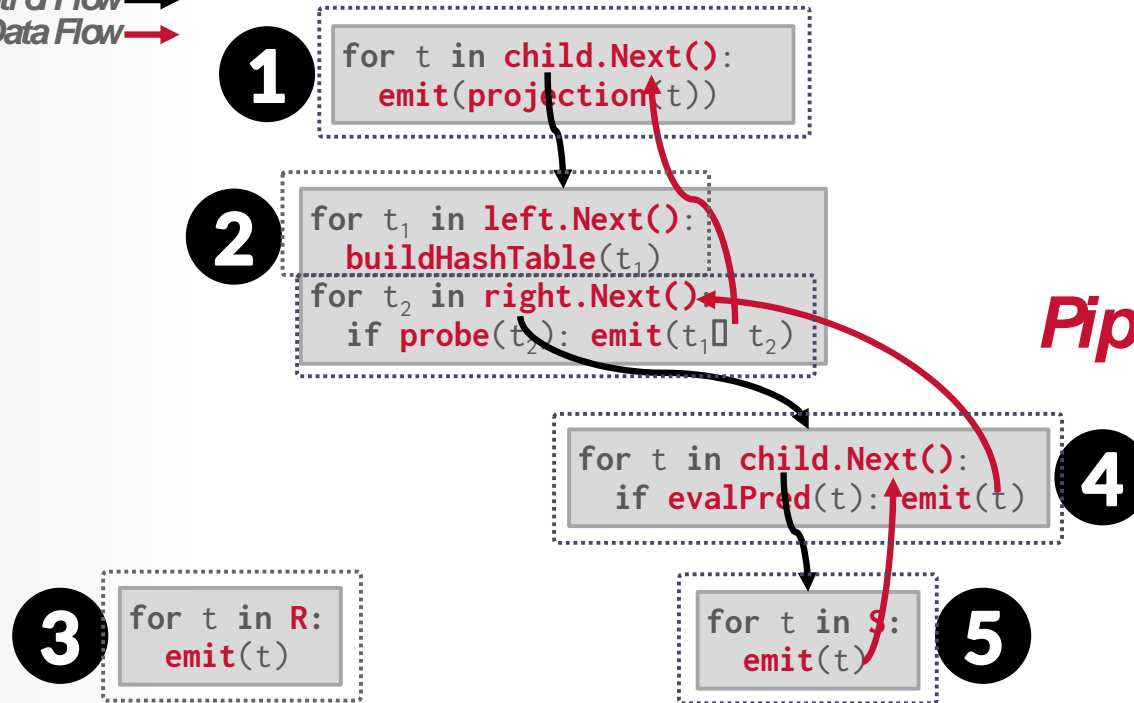
```

SELECT R.id, S.cdate
FROM R JOIN S
ON R.id = S.id
WHERE S.value > 100
  
```



Pipeline #1

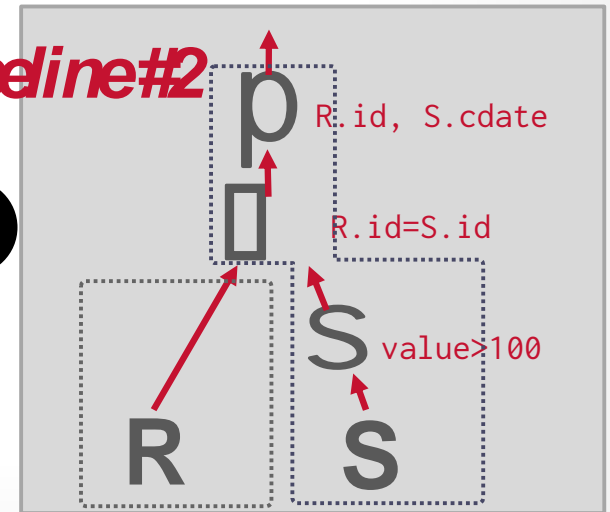
Control Flow →  
Data Flow →



```

SELECT R.id, S.cdate
FROM R JOIN S
ON R.id = S.id
WHERE S.value > 100
  
```

*Pipeline#2*





# Iterator Model

**The Iterator model is used in almost every DBMS.**

- Easy to implement / debug.
- Output control works easily with this approach.

**Allows for pipelining**

- where the DBMS tries to process each tuple through as many operators as possible before retrieving the next tuple.

**A pipeline breaker**

- is an operator that cannot finish until all its children emit all their tuples.
- → Joins (Build Side), Subqueries, Order By

**Systems**



# Evolution of Processing Models



## Description

- Each operator processes its input all at once and then emits its output all at once.
- The operator "materializes" its output as a single result.
- Can send either a materialized row or a single column.

## Origin

- Originally developed in MonetDB in the 1990s to process entire columns at a time instead of single tuples.
- MonetDB = 1st pure database system with a columnar storage approach

# Example

Control Flow   
Data Flow 

1

```
out = [ ]  
for t in child.Output():  
    out.add(projection(t))  
return out
```

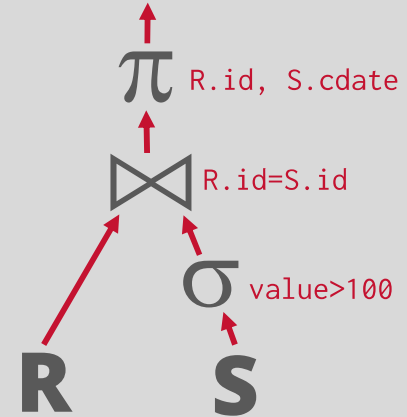
```
out = [ ]  
for t1 in left.Output():  
    buildHashTable(t1)  
for t2 in right.Output():  
    if probe(t2): out.add(t1 ⋈ t2)  
return out
```

```
out = [ ]  
for t in child.Output():  
    if evalPred(t): out.add(t)  
return out
```

```
out = [ ]  
for t in R:  
    out.add(t)  
return out
```

```
out = [ ]  
for t in S:  
    out.add(t)  
return out
```

```
SELECT R.id, S.cdate  
FROM R JOIN S  
ON R.id = S.id  
WHERE S.value > 100
```



# Example

Control Flow  $\rightarrow$   
Data Flow  $\rightarrow$

1

```
out = [ ]  
for t in child.Output():  
    out.add(projection(t))  
return out
```

2

```
out = [ ]  
for t1 in left.Output():  
    buildHashTable(t1)  
for t2 in right.Output():  
    if probe(t2): out.add(t1 ⋈ t2)  
return out
```

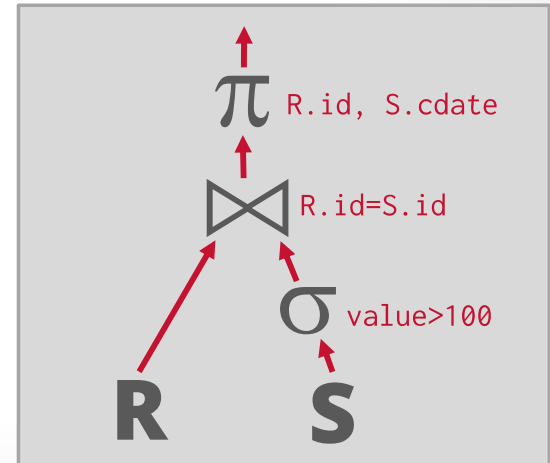
3

```
out = [ ]  
for t in R:  
    out.add(t)  
return out
```

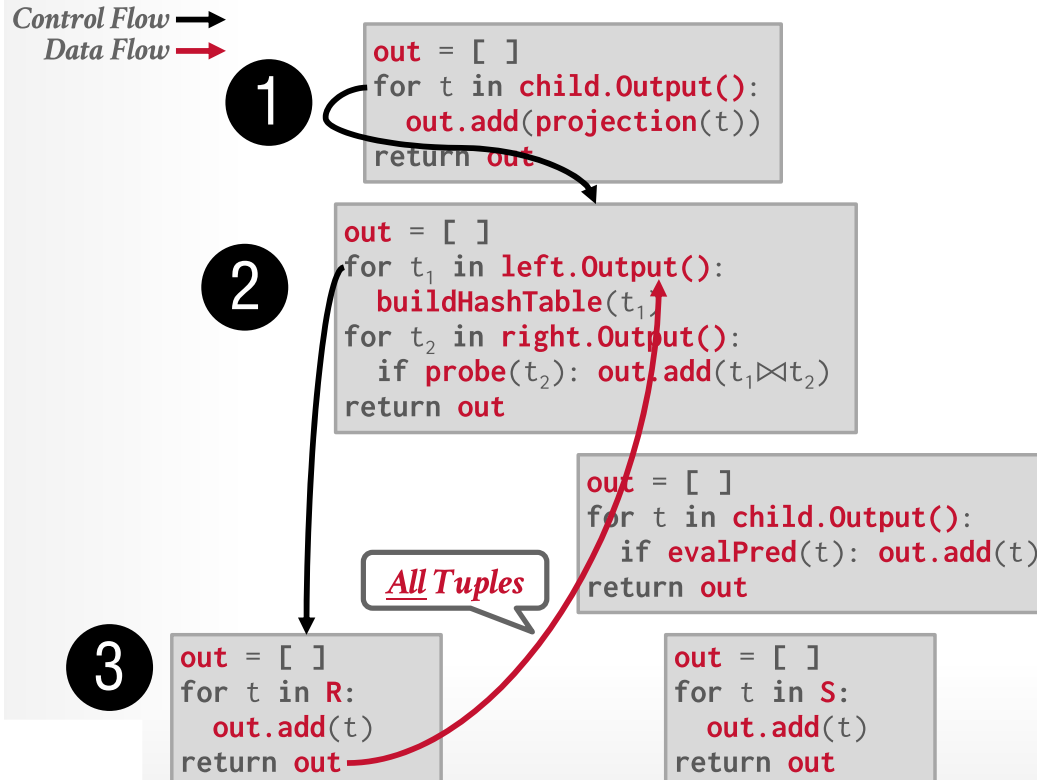
```
out = [ ]  
for t in child.Output():  
    if evalPred(t): out.add(t)  
return out
```

```
out = [ ]  
for t in S:  
    out.add(t)  
return out
```

```
SELECT R.id, S.cdate  
FROM R JOIN S  
ON R.id = S.id  
WHERE S.value > 100
```

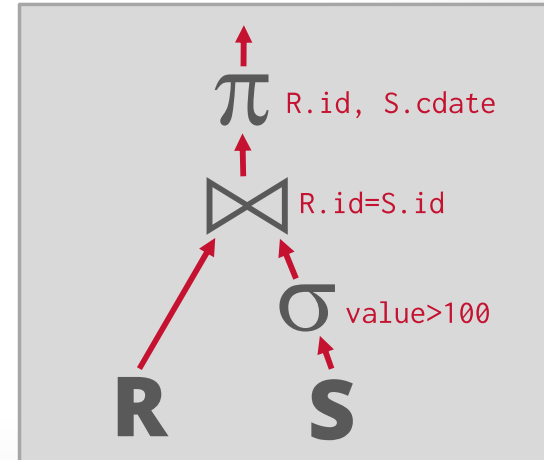


# Example



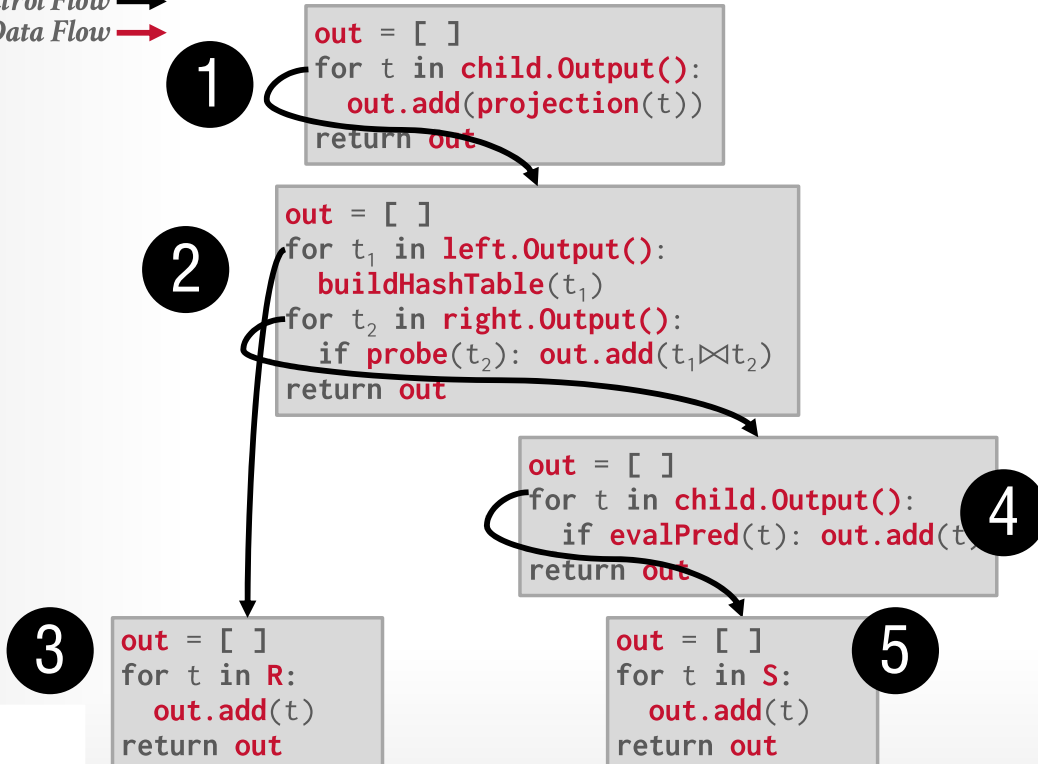
```

SELECT R.id, S.cdate
FROM R JOIN S
ON R.id = S.id
WHERE S.value > 100
    
```

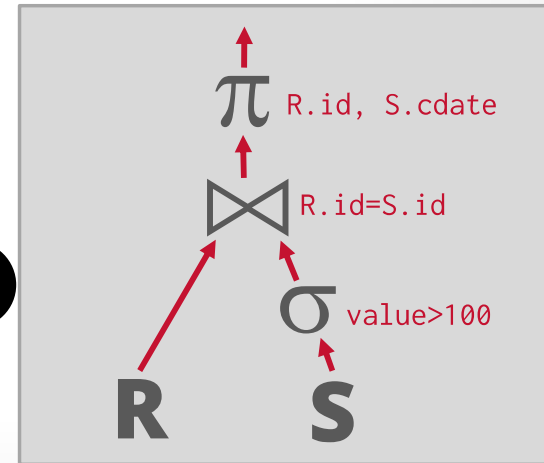


# Example

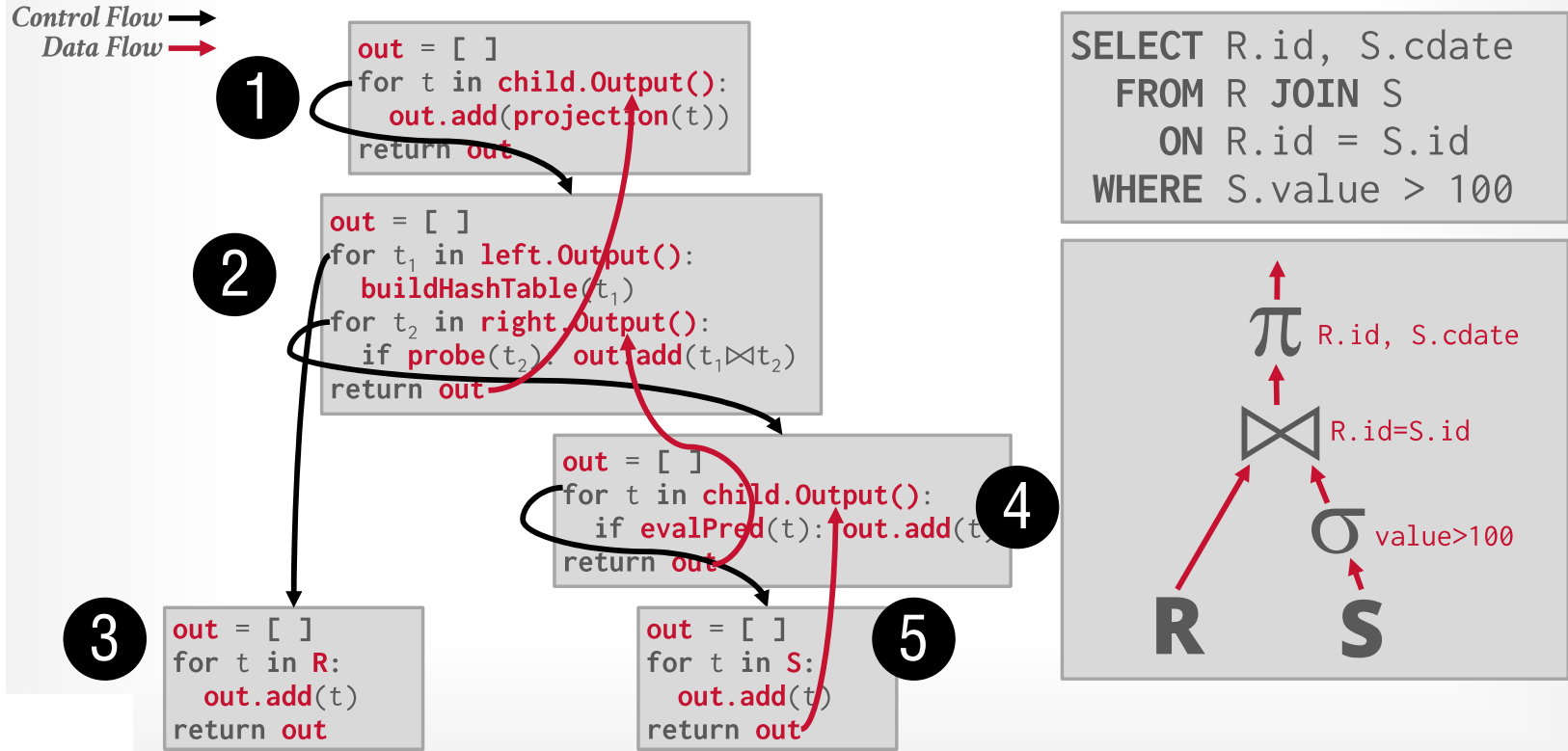
Control Flow  $\rightarrow$   
Data Flow  $\rightarrow$



```
SELECT R.id, S.cdate
FROM R JOIN S
ON R.id = S.id
WHERE S.value > 100
```

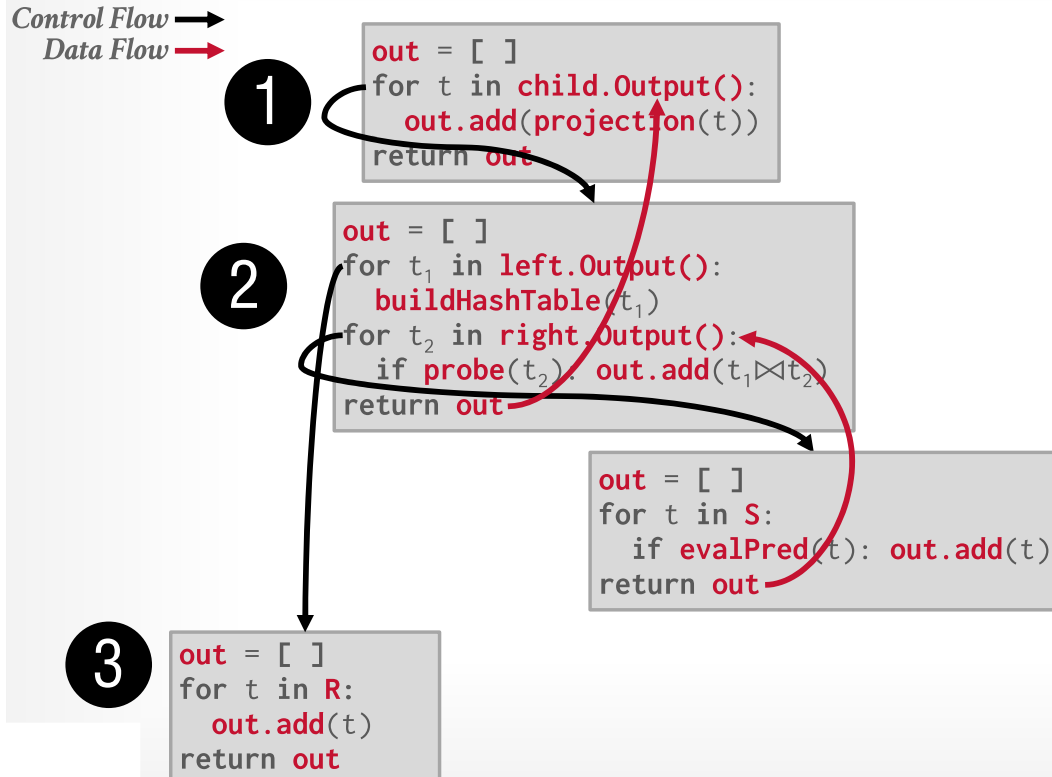


# Example

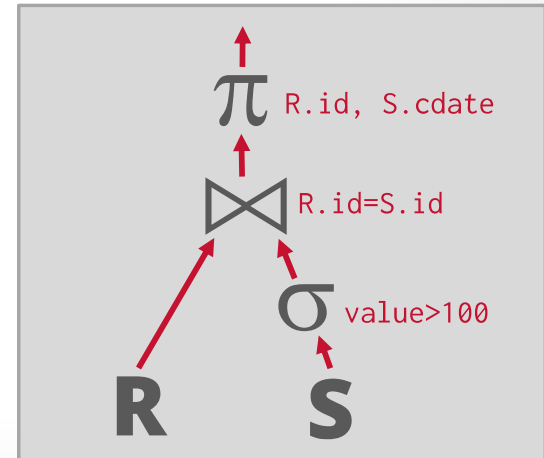




# Example



```
SELECT R.id, S.cdate
FROM R JOIN S
ON R.id = S.id
WHERE S.value > 100
```



# Materialization Model

## Example: Star Schema Benchmark (SSB) Query 3.1

```
SELECT c_nation, s_nation, d_year, SUM(lo_revenue)
FROM customer, lineorder, supplier, date
WHERE lo_custkey = c_custkey
      AND lo_suppkey = s_suppkey
      AND lo_orderdate = d_datekey
      AND c_region = 'ASIA' AND s_region = 'ASIA'
      AND d_year >= 1992 AND d_year <= 1997
GROUP BY c_nation, s_nation, d_year
ORDER BY d_year ASC, revenue DESC;
```

### Number of occurrences

33x  operator

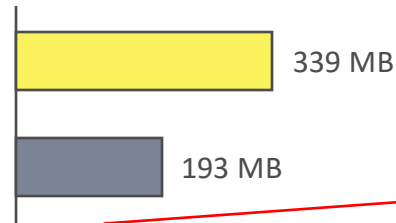
35x  intermediate

12x  base column

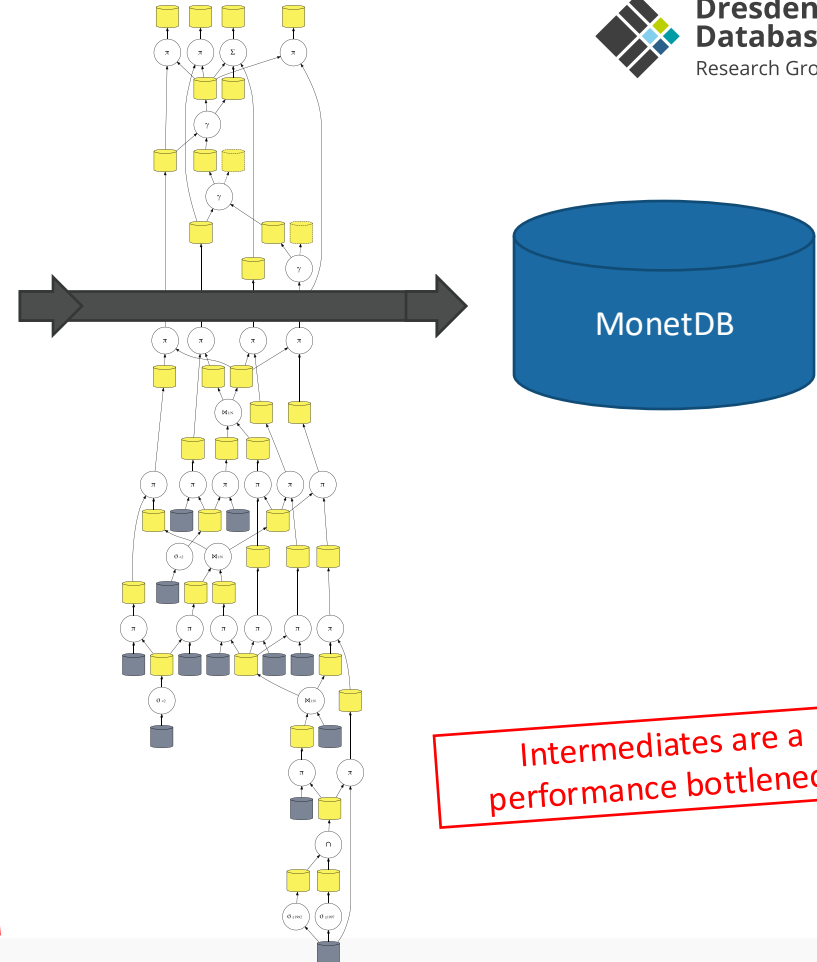
Many Intermediates

### Total data size @sf 1

(64-bit per data element)



Large total size



Intermediates are a performance bottleneck

# Evolution of Processing Models



# Batch Model

**Like the Iterator Model where each operator implements a Next() function, but...**

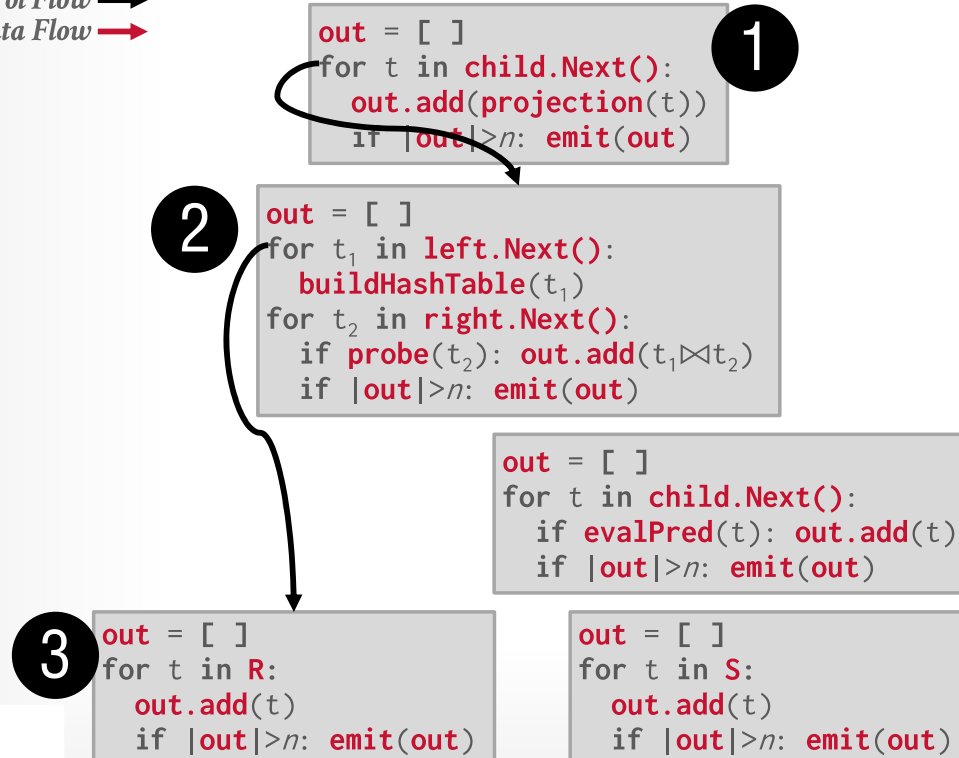
- Each operator emits a batch of tuples/columns instead of a single tuple.
- The operator's internal loop processes multiple tuples/column values at a time.
- The size of the batch can vary based on hardware or query properties.

## Iterator Model

- Best-of-both- Worlds: Combination of Iterator Model and Materialization Model

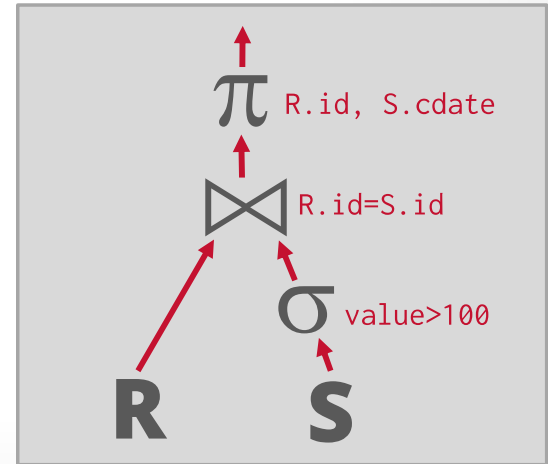
# Example

Control Flow  $\rightarrow$   
Data Flow  $\rightarrow$



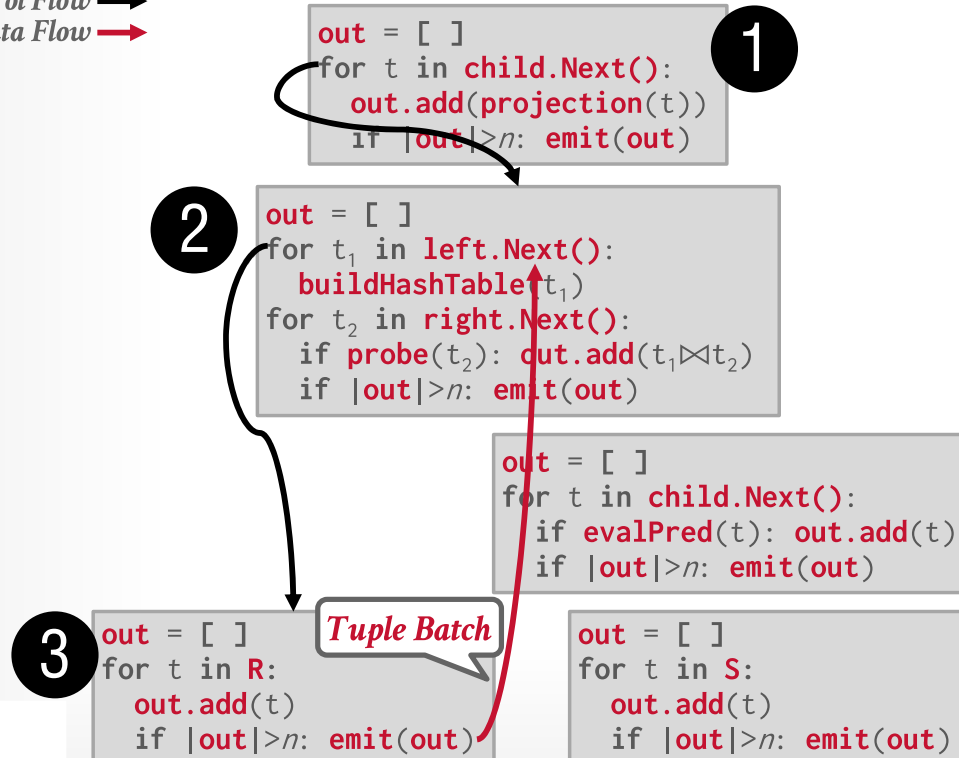
```

SELECT R.id, S.cdate
FROM R JOIN S
      ON R.id = S.id
WHERE S.value > 100
    
```



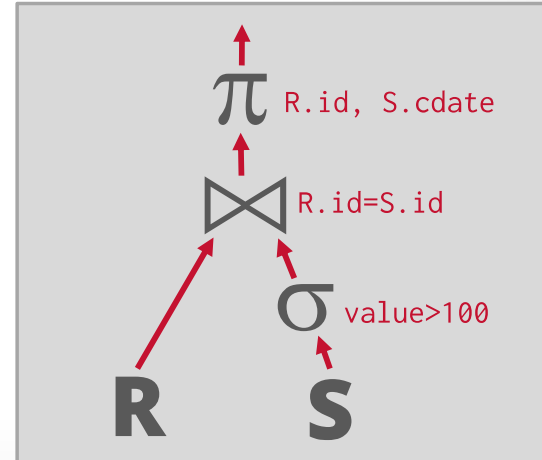
# Example

Control Flow →  
Data Flow →



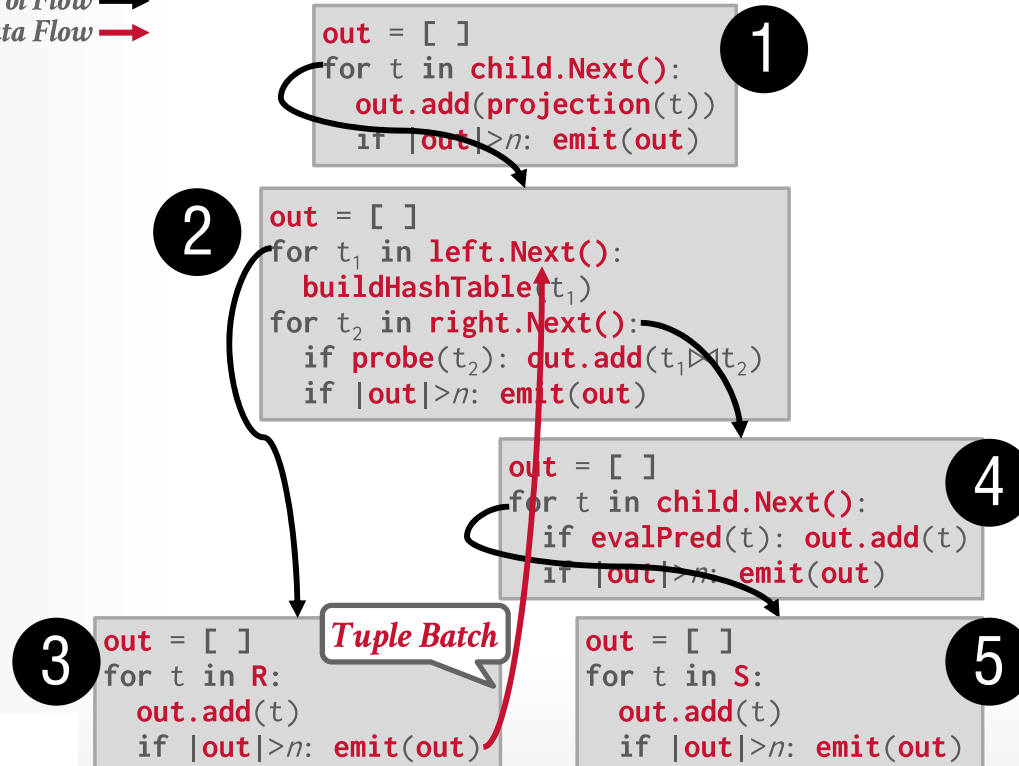
```

SELECT R.id, S.cdate
FROM R JOIN S
ON R.id = S.id
WHERE S.value > 100
    
```

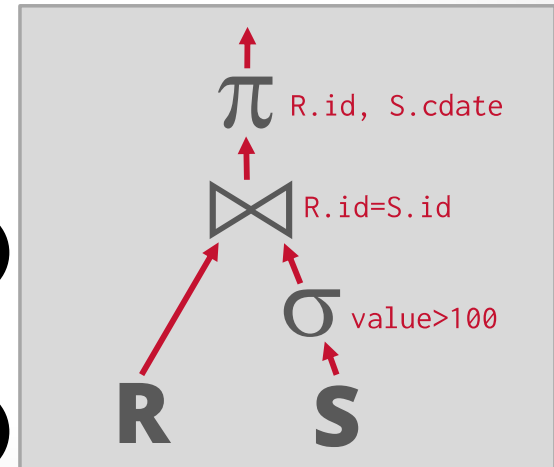


# Example

Control Flow  $\blackrightarrow$   
Data Flow  $\color{red}\rightarrow$

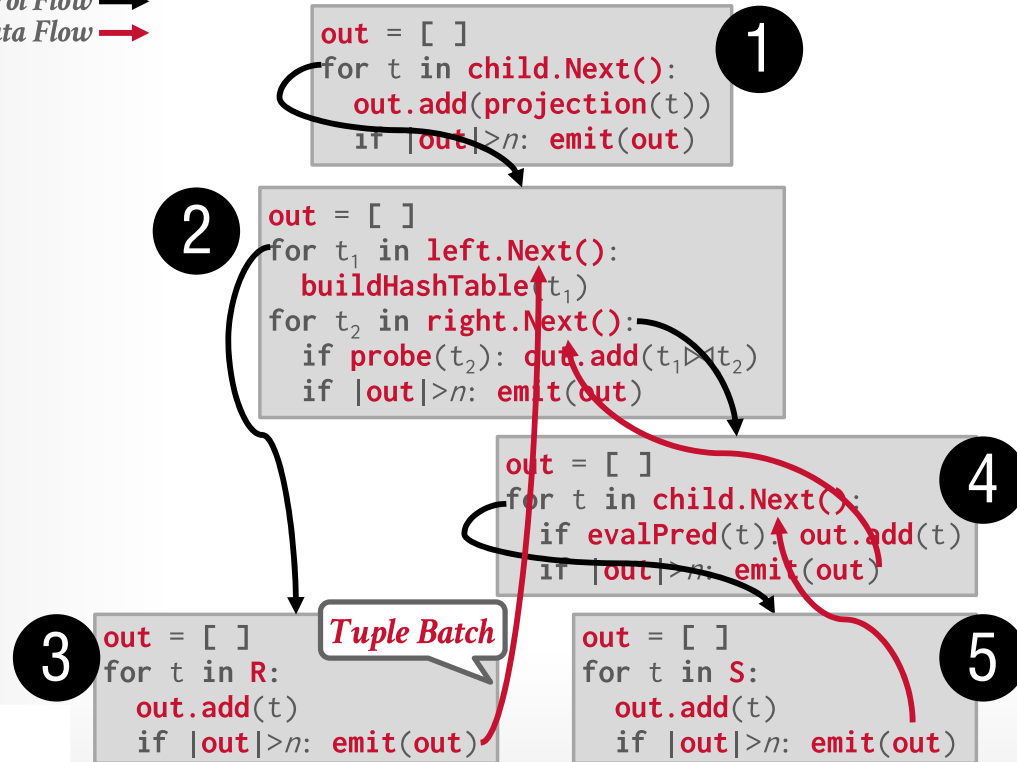


```
SELECT R.id, S.cdate
FROM R JOIN S
ON R.id = S.id
WHERE S.value > 100
```



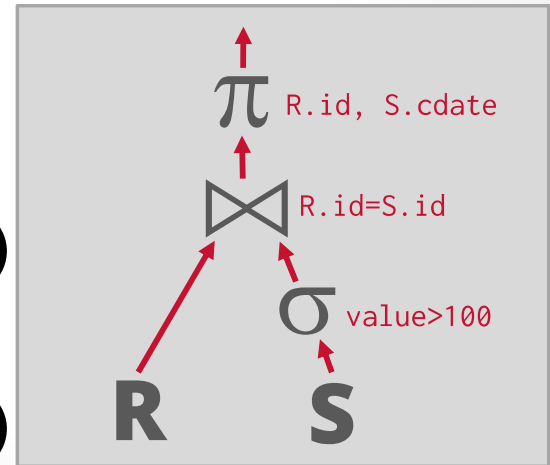
# Example

Control Flow  $\rightarrow$   
Data Flow  $\rightarrow$



```

SELECT R.id, S.cdate
FROM R JOIN S
      ON R.id = S.id
WHERE S.value > 100
    
```





# Batch Model

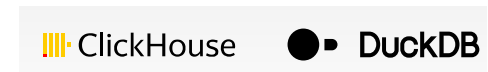
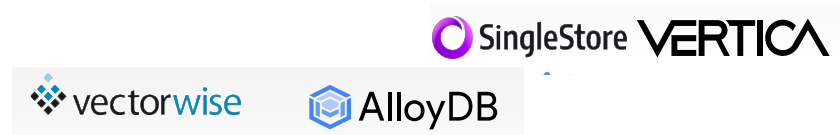
## Ideal for OLAP queries

- because it greatly reduces the number of invocations per operator.

## Allows an out-of-order CPU

- to efficiently execute operators over batches of tuples.
- Operators perform work in tight for-loops over arrays, which compilers know how to optimize / vectorize.
- No data or control dependencies.
- Hot instruction cache.

## Systems



## Observation

- In the previous examples, the DBMS starts executing a query by invoking Next() at the root of the query plan and pulling data up from leaf operators.
- This is the how most DBMSs implement their execution engine.

## Approach #1: Top-to-Bottom (Pull)

- Start with the root and "pull" data up from its children.
- Tuples are always passed between operators using function calls (unless it's a pipeline breaker).

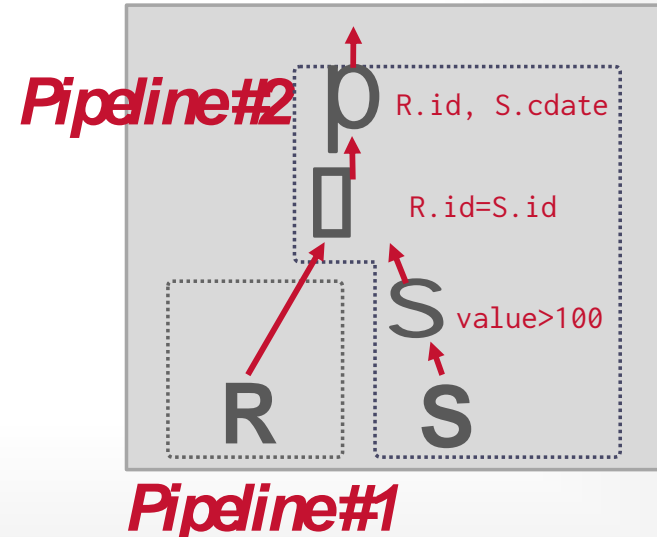
## Approach #2: Bottom-to-Top (Push)

- Start with leaf nodes and "push" data to their parents.
- Can "fuse" operators together within a for-loop to minimize intermediate result staging.

# Push-based Processing Direction

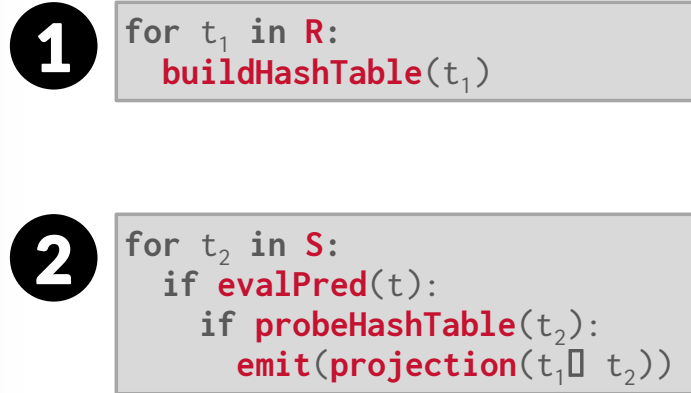
Control Flow →  
Data Flow →

```
SELECT R.id, S.cdate  
FROM R JOIN S  
ON R.id = S.id  
WHERE S.value > 100
```



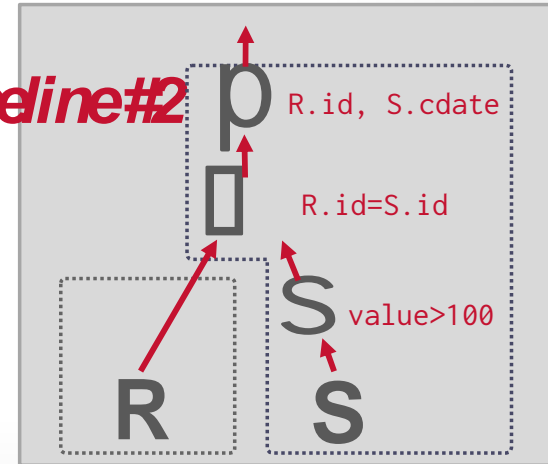
# Push-based Processing Direction

Control Flow →  
Data Flow →



```
SELECT R.id, S.cdate
FROM R JOIN S
      ON R.id = S.id
WHERE S.value > 100
```

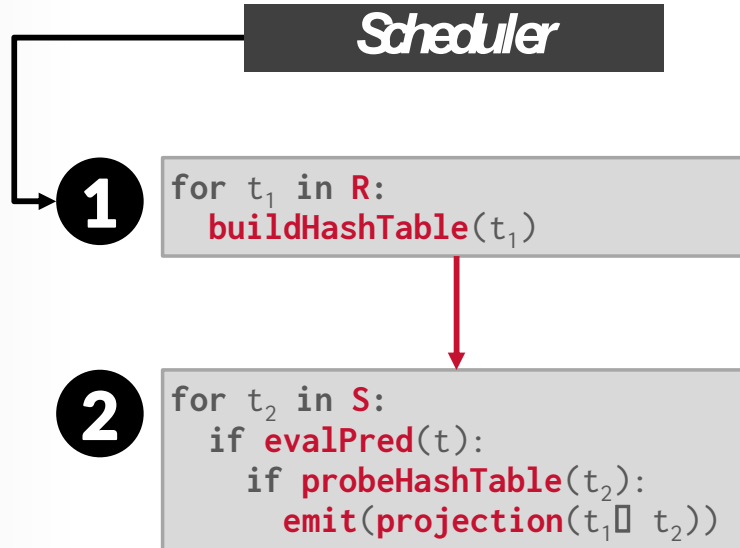
*Pipeline#2*



*Pipeline#1*

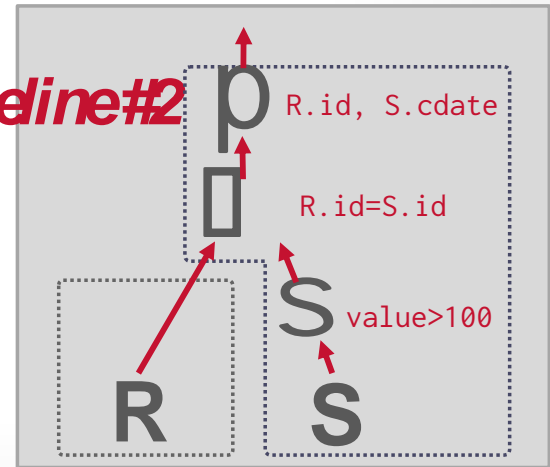
# Push-based Processing Direction

Control Flow →  
Data Flow →



```
SELECT R.id, S.cdate  
FROM R JOIN S  
ON R.id = S.id  
WHERE S.value > 100
```

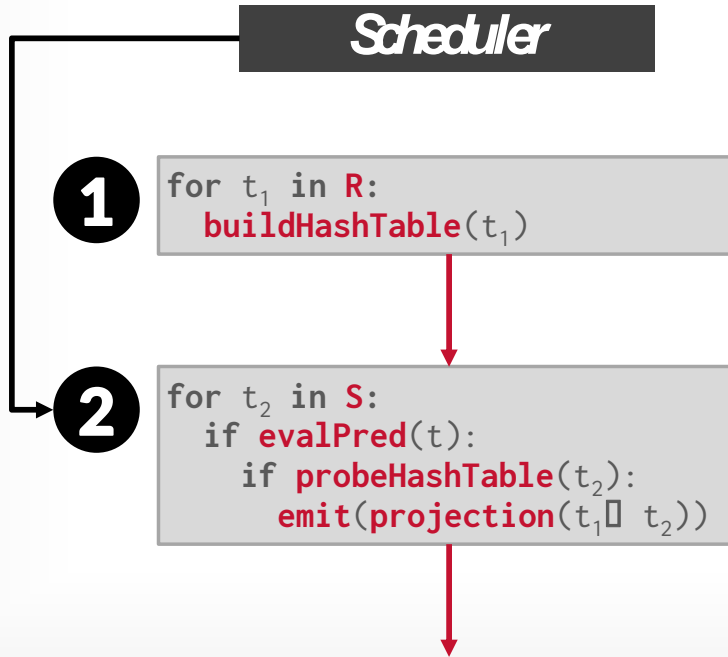
*Pipeline#2*



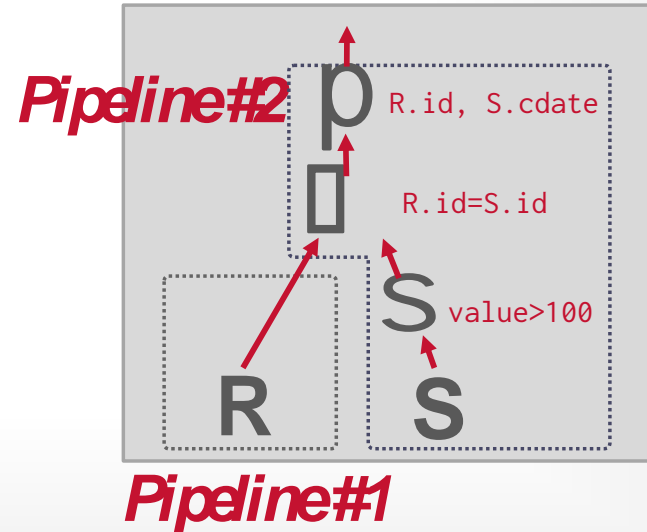
*Pipeline#1*

# Push-based Processing Direction

Control Flow  $\rightarrow$   
Data Flow  $\rightarrow$



```
SELECT R.id, S.cdate
FROM R JOIN S
      ON R.id = S.id
WHERE S.value > 100
```



# Summary

## Modern In-Memory Database Systems

- Optimized for analytical queries
- Columnar Data Organization
- Batch-based Processing Model

## Homework

- Please read the following paper (available on OPAL: 04-MonetDB.pdf)

DOI:10.1145/1409360.1409380

# Breaking the Memory Wall in MonetDB

By Peter A. Boncz, Martin L. Kersten, and Stefan Manegold