

Approximation Techniques

Supplemental Reading
(papers posted on Piazza)

Query Approximation

- Summarize the data with a **Sketching** Algorithm.
 - Bloom Filters (Set Containment)
 - Flajolet & Martin Count-Distinct Sketch
 - Count Sketch (Frequent Items/Top-K)
- Sample the data and estimate the error
 - ... or better yet, keep generating samples!
 - Online Aggregation & Ripple Joins

Example: F&M Count Sketch

$\text{hash}(O_1) = 01011011$

$\text{hash}(O_2) = 00110111$

$\text{hash}(O_3) = 00111000$

$\text{hash}(O_4) = 10010010$

Example: F&M Count Sketch

$\text{hash}(O_1) = 0101101\textcircled{1}$	$\rho(O_1) = 0$
$\text{hash}(O_2) = 001101\textcircled{1}1$	$\rho(O_2) = 0$
$\text{hash}(O_3) = 0011\textcircled{1}000$	$\rho(O_3) = 3$
$\text{hash}(O_4) = 100100\textcircled{1}0$	$\rho(O_4) = 1$

Example: F&M Count Sketch

$\text{hash}(O_1) = 0101101\textcircled{1}$ $\rho(O_1) = 0 : 0001$

$\text{hash}(O_2) = 001101\textcircled{1}1$ $\rho(O_2) = 0 : 0001$

$\text{hash}(O_3) = 0011\textcircled{1}000$ $\rho(O_3) = 3 : 1000$

$\text{hash}(O_4) = 100100\textcircled{1}0$ $\rho(O_4) = 1 : 0010$

1011
←
 $R = \textcircled{2}10$

Example: F&M Count Sketch

$$\text{hash}(O_1) = 0101101\underset{\textcircled{1}}{1} \quad \rho(O_1) = 0 : 0001$$

$$\text{hash}(O_2) = 001101\underset{\textcircled{1}}{1}1 \quad \rho(O_2) = 0 : 0001$$

$$\text{hash}(O_3) = 0011\underset{\textcircled{1}}{1}000 \quad \rho(O_3) = 3 : 1000$$

$$\text{hash}(O_4) = 100100\underset{\textcircled{1}}{1}0 \quad \rho(O_4) = 1 : 0010$$

$$2^R/\varphi = 4 / 0.77351 = 5.2$$

$$R = \overleftarrow{210}$$

1011

Example: CountSketch

Find the 2 Most Frequent Objects

O₃ O₁ O₄ O₂ O₄ O₁ O₃ O₃ O₁

Example: CountSketch

Estimate the Count of Each Object

~~Find the 2 Most Frequent Objects~~

O₃ O₁ O₄ O₂ O₄ O₁ O₃ O₃ O₁

Example: CountSketch

Estimate the Count of Each Object

~~Find the 2 Most Frequent Objects~~

O₃ O₁ O₄ O₂ O₄ O₁ O₃ O₃ O₁

hash(O₁) = 01011011

hash(O₂) = 00110111

hash(O₃) = 00111000

hash(O₄) = 10010010

Example: CountSketch

Estimate the Count of Each Object

~~Find the 2 Most Frequent Objects~~

O₃ O₁ O₄ O₂ O₄ O₁ O₃ O₃ O₁

hash(O₁) = 1

hash(O₂) = 1

hash(O₃) = 0

hash(O₄) = 0

Example: CountSketch

Estimate the Count of Each Object

~~Find the 2 Most Frequent Objects~~

$O_3 \ O_1 \ O_4 \ O_2 \ O_4 \ O_1 \ O_3 \ O_3 \ O_1$

$-| +| -| +| -| +| -| -| +|$

$\text{hash}(O_1) = 1 +|$

$\text{hash}(O_2) = 1 +|$

$\text{hash}(O_3) = 0 -|$

$\text{hash}(O_4) = 0 -|$

Example: CountSketch

Estimate the Count of Each Object

~~Find the 2 Most Frequent Objects~~

$$\begin{array}{cccccccccc} O_3 & O_1 & O_4 & O_2 & O_4 & O_1 & O_3 & O_3 & O_1 \\ -| & +| & -| & +| & -| & +| & -| & -| & +| & = -| \end{array}$$

hash(O_1) =	1	+	$E[\text{Count}(O_1)] = - $
hash(O_2) =	1	+	$E[\text{Count}(O_2)] = - $
hash(O_3) =	0	-	$E[\text{Count}(O_3)] = $
hash(O_4) =	0	-	$E[\text{Count}(O_4)] = $

Example: CountSketch

Estimate the Count of Each Object

~~Find the 2 Most Frequent Objects~~

$O_3 \ O_1 \ O_4 \ O_2 \ O_4 \ O_1 \ O_3 \ O_3 \ O_1$

$\text{hash}(O_1) = 01011011$

$\text{hash}(O_2) = 00110111$

$\text{hash}(O_3) = 00111000$

$\text{hash}(O_4) = 10010010$

0	0	0	0
0	0	0	0

Example: CountSketch

Estimate the Count of Each Object

~~Find the 2 Most Frequent Objects~~

$O_3 \ O_1 \ O_4 \ O_2 \ O_4 \ O_1 \ O_3 \ O_3 \ O_1$

Each Pair: <Top/Bottom, +1/-1>

hash(O_1) = 01|011011
hash(O_2) = 00|110111
hash(O_3) = 00|111000
hash(O_4) = 10|010010

0	0	0	0
0	0	0	0

Example: CountSketch

Estimate the Count of Each Object

~~Find the 2 Most Frequent Objects~~

O₃ O₁ O₄ O₂ O₄ O₁ O₃ O₃ O₁

Each Pair: <Top/Bottom, +|-|>

hash(O₁) <B,+|>, <B,+|>, <T,-|>, <T,+|>

hash(O₂) <B,-|>, <T,+|>, <B,+|>, <T,+|>

hash(O₃) <B,-|>, <T,+|>, <T,-|>, <B,-|>

hash(O₄) <T,-|>, <B,+|>, <B,-|>, <T,-|>

0	0	0	0
0	0	0	0

Example: CountSketch

Estimate the Count of Each Object

~~Find the 2 Most Frequent Objects~~

O₃ O₁ O₄ O₂ O₄ O₁ O₃ O₃ O₁

hash(O₁) <B,+|>, <B,+|>, <T,-|>, <T,+|>
hash(O₂) <B,-|>, <T,+|>, <B,+|>, <T,+|>
hash(O₃) <B,-|>, <T,+|>, <T,-|>, <B,-|>
hash(O₄) <T,-|>, <B,+|>, <B,-|>, <T,-|>

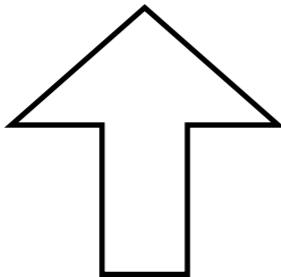
0	+1	-1	0
-1	0	0	-1

Example: CountSketch

Estimate the Count of Each Object

~~Find the 2 Most Frequent Objects~~

$O_3 \ O_1 \ O_4 \ O_2 \ O_4 \ O_1 \ O_3 \ O_3 \ O_1$



hash(O_1) $\langle B, +1 \rangle, \langle B, +1 \rangle, \langle T, -1 \rangle, \langle T, +1 \rangle$
hash(O_2) $\langle B, -1 \rangle, \langle T, +1 \rangle, \langle B, +1 \rangle, \langle T, +1 \rangle$
hash(O_3) $\langle B, -1 \rangle, \langle T, +1 \rangle, \langle T, -1 \rangle, \langle B, -1 \rangle$
hash(O_4) $\langle T, -1 \rangle, \langle B, +1 \rangle, \langle B, -1 \rangle, \langle T, -1 \rangle$

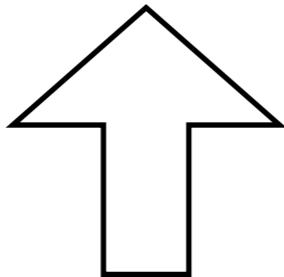
0	+1	-2	+1
0	+1	0	-1

Example: CountSketch

Estimate the Count of Each Object

~~Find the 2 Most Frequent Objects~~

$O_3 \ O_1 \ O_4 \ O_2 \ O_4 \ O_1 \ O_3 \ O_3 \ O_1$



$\text{hash}(O_1) \langle B, +1 \rangle, \langle B, +1 \rangle, \langle T, -1 \rangle, \langle T, +1 \rangle$

$\text{hash}(O_2) \langle B, -1 \rangle, \langle T, +1 \rangle, \langle B, +1 \rangle, \langle T, +1 \rangle$

$\text{hash}(O_3) \langle B, -1 \rangle, \langle T, +1 \rangle, \langle T, -1 \rangle, \langle B, -1 \rangle$

$\text{hash}(O_4) \langle T, -1 \rangle, \langle B, +1 \rangle, \langle B, -1 \rangle, \langle T, -1 \rangle$

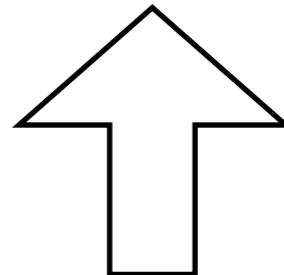
-1	+1	-2	0
0	+2	-1	-1

Example: CountSketch

Estimate the Count of Each Object

~~Find the 2 Most Frequent Objects~~

$O_3 \ O_1 \ O_4 \ O_2 \ O_4 \ O_1 \ O_3 \ O_3 \ O_1$



hash(O_1) $\langle B, +1 \rangle, \langle B, +1 \rangle, \langle T, -1 \rangle, \langle T, +1 \rangle$
hash(O_2) $\langle B, -1 \rangle, \langle T, +1 \rangle, \langle B, +1 \rangle, \langle T, +1 \rangle$
hash(O_3) $\langle B, -1 \rangle, \langle T, +1 \rangle, \langle T, -1 \rangle, \langle B, -1 \rangle$
hash(O_4) $\langle T, -1 \rangle, \langle B, +1 \rangle, \langle B, -1 \rangle, \langle T, -1 \rangle$

-2	+4	-6	+4
-1	+5	-1	+3

Example: CountSketch

Estimate the Count of Each Object

~~Find the 2 Most Frequent Objects~~

O₃ O₁ O₄ O₂ O₄ O₁ O₃ O₃ O₁

hash(O₁) <B,+|>, <B,+|>, <T,-|>, <T,+|>

hash(O₂) <B,-|>, <T,+|>, <B,+|>, <T,+|>

hash(O₃) <B,-|>, <T,+|>, <T,-|>, <B,-|>

hash(O₄) <T,-|>, <B,+|>, <B,-|>, <T,-|>

-2	+4	-6	+4
-1	+5	-1	+3

Example: CountSketch

Estimate the Count of Each Object

~~Find the 2 Most Frequent Objects~~

O₃ O₁ O₄ O₂ O₄ O₁ O₃ O₃ O₁

hash(O₁) <B,-|>, <B,+|>, <T,->, <T,+|>
hash(O₂) <B,-|>, <T,+|>, <B,+|>, <T,+|>
hash(O₃) <B,-|>, <T,+|>, <T,-|>, <B,-|>
hash(O₄) <T,-|>, <B,+|>, <B,-|>, <T,-|>

-2	+4	-6	+4
-1	+5	-1	+3

Example: CountSketch

Estimate the Count of Each Object

~~Find the 2 Most Frequent Objects~~

$O_3 \ O_1 \ O_4 \ O_2 \ O_4 \ O_1 \ O_3 \ O_3 \ O_1$

hash(O_1) $\langle B, -1 \rangle, \langle B, +1 \rangle, \langle T, -1 \rangle, \langle T, +1 \rangle$
hash(O_2) $\langle B, -1 \rangle, \langle T, +1 \rangle, \langle B, +1 \rangle, \langle T, +1 \rangle$
hash(O_3) $\langle B, -1 \rangle, \langle T, +1 \rangle, \langle T, -1 \rangle, \langle B, -1 \rangle$
hash(O_4) $\langle T, -1 \rangle, \langle B, +1 \rangle, \langle B, -1 \rangle, \langle T, -1 \rangle$

-2	+4	-6	+4
-1	+5	-1	+3
-1	+4	+5	+6

Example: CountSketch

Estimate the Count of Each Object

~~Find the 2 Most Frequent Objects~~

O₃ O₁ O₄ O₂ O₄ O₁ O₃ O₃ O₁

hash(O₁) <B,+|>, <B,+|>, <T,-|>, <T,+|>

hash(O₂) <B,-|>, <T,+|>, <B,+|>, <T,+|>

hash(O₃) <B,-|>, <T,+|>, <T,-|>, <B,-|>

hash(O₄) <T,-|>, <B,+|>, <B,-|>, <T,-|>

-2	+4	-6	+4
-1	+5	-1	+3

Example: CountSketch

Estimate the Count of Each Object

~~Find the 2 Most Frequent Objects~~

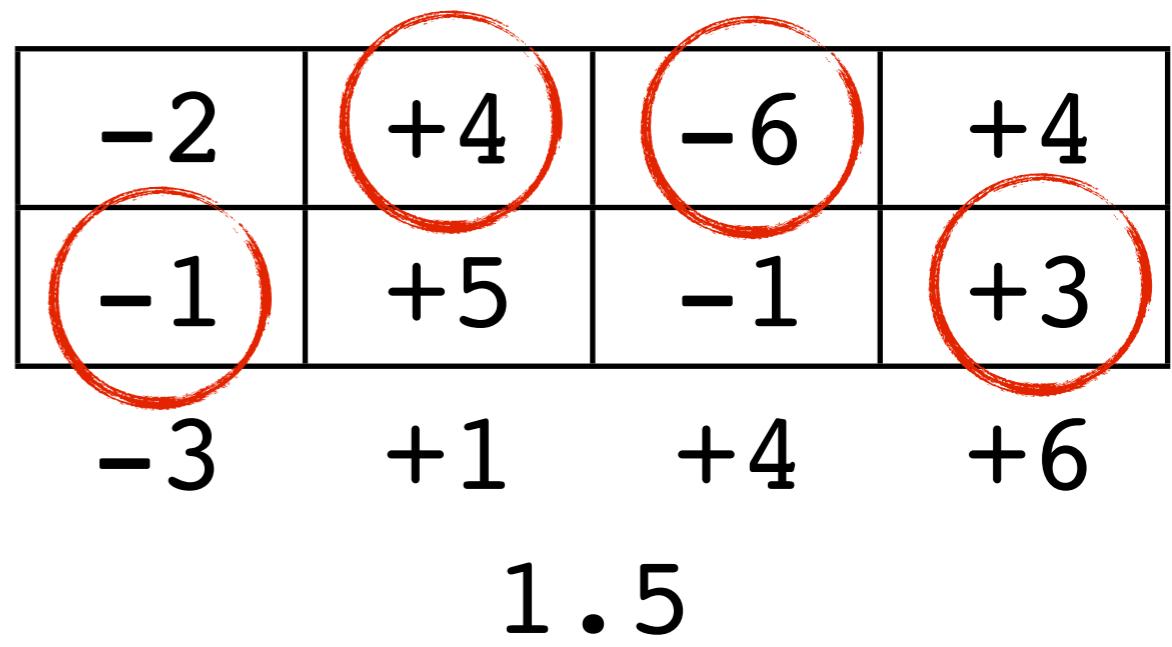
O₃ O₁ O₄ O₂ O₄ O₁ O₃ O₃ O₁

hash(O₁) <B,+|>, <B,+|>, <T,-|>, <T,+|>

hash(O₂) <B,-|>, <T,+|>, <B,+|>, <T,+|>

hash(O₃) <B,-|>, <T,+|>, <T,-|>, <B,-|>

hash(O₄) <T,-|>, <B,+|>, <B,-|>, <T,-|>



Wiggle Room

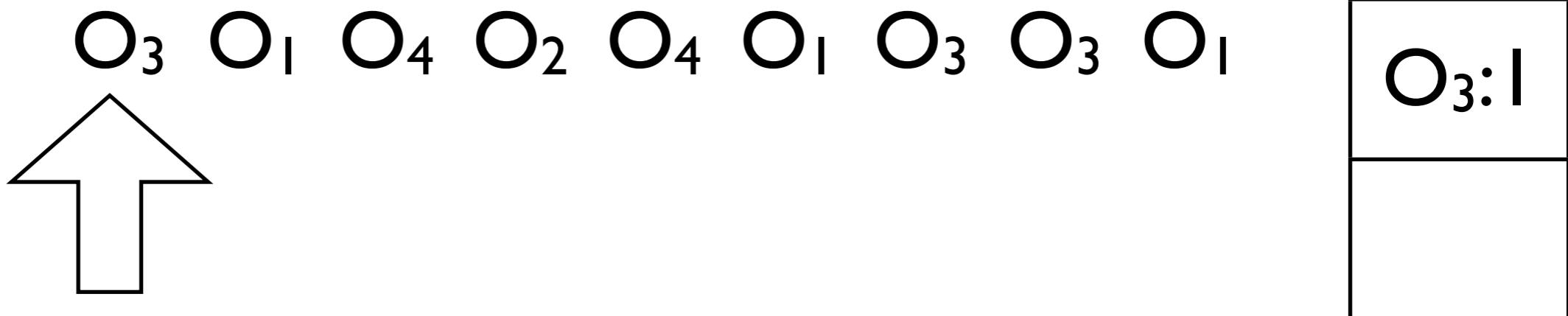
More Hashes

=

More Accurate Estimate

Example: CountSketch

Find the 2 Most Frequent Objects



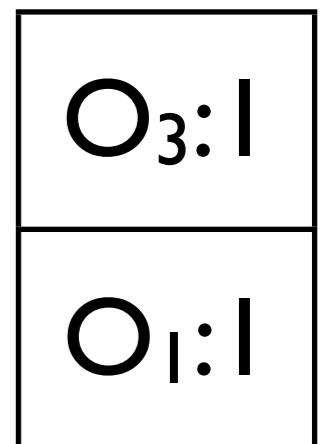
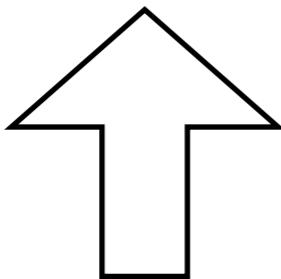
hash(O_1) $\langle B, +1 \rangle, \langle B, +1 \rangle, \langle T, -1 \rangle, \langle T, +1 \rangle$
hash(O_2) $\langle B, -1 \rangle, \langle T, +1 \rangle, \langle B, +1 \rangle, \langle T, +1 \rangle$
hash(O_3) $\langle B, -1 \rangle, \langle T, +1 \rangle, \langle T, -1 \rangle, \langle B, -1 \rangle$
hash(O_4) $\langle T, -1 \rangle, \langle B, +1 \rangle, \langle B, -1 \rangle, \langle T, -1 \rangle$

0	+1	-1	0
-1	0	0	-1

Example: CountSketch

Find the 2 Most Frequent Objects

$O_3 \ O_1 \ O_4 \ O_2 \ O_4 \ O_1 \ O_3 \ O_3 \ O_1$



$\text{hash}(O_1) \langle B, +1 \rangle, \langle B, +1 \rangle, \langle T, -1 \rangle, \langle T, +1 \rangle$

$\text{hash}(O_2) \langle B, -1 \rangle, \langle T, +1 \rangle, \langle B, +1 \rangle, \langle T, +1 \rangle$

$\text{hash}(O_3) \langle B, -1 \rangle, \langle T, +1 \rangle, \langle T, -1 \rangle, \langle B, -1 \rangle$

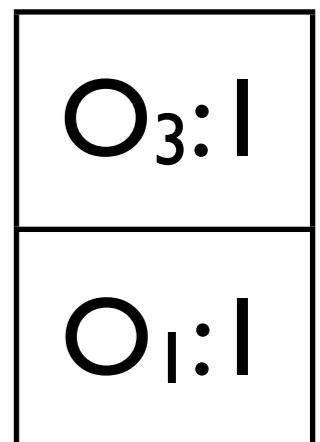
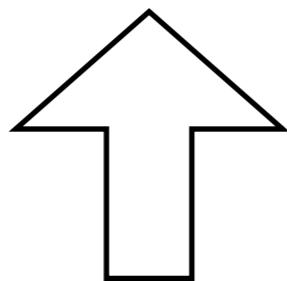
$\text{hash}(O_4) \langle T, -1 \rangle, \langle B, +1 \rangle, \langle B, -1 \rangle, \langle T, -1 \rangle$

0	+1	-2	1
0	+1	0	-1

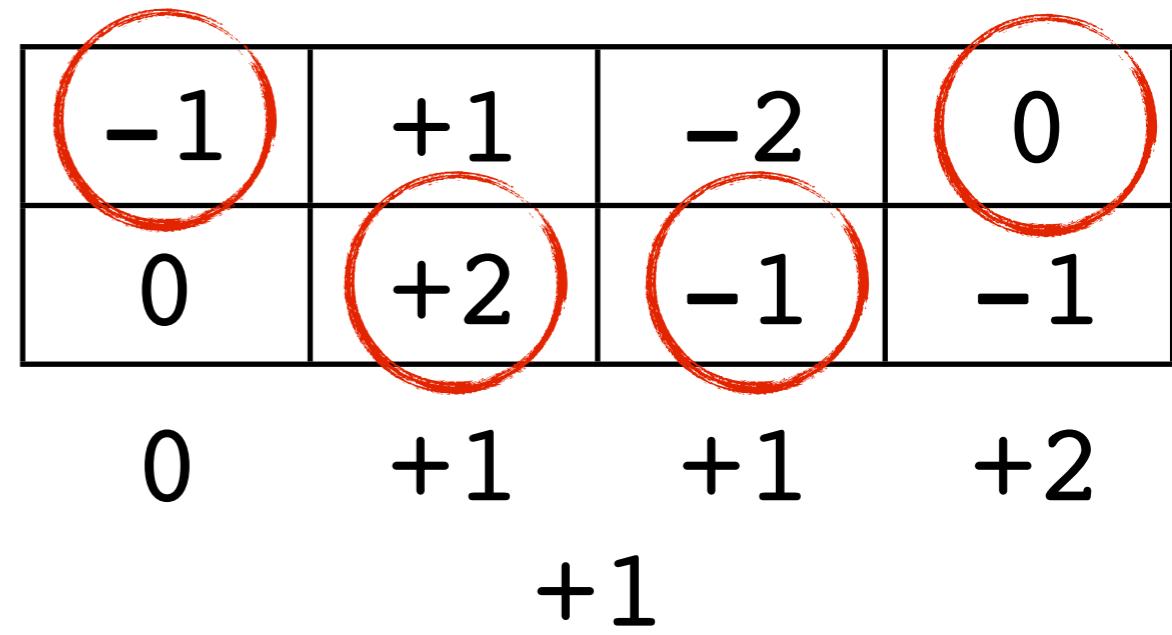
Example: CountSketch

Find the 2 Most Frequent Objects

$O_3 \ O_1 \ O_4 \ O_2 \ O_4 \ O_1 \ O_3 \ O_3 \ O_1$



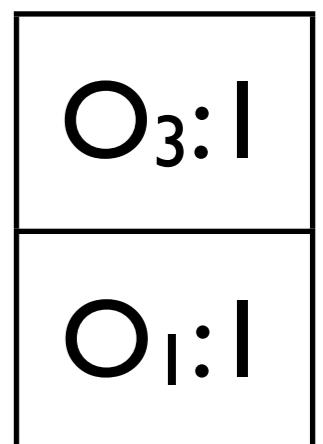
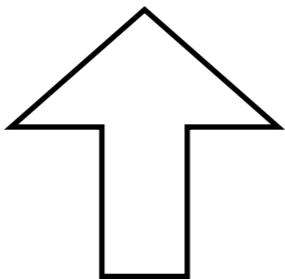
hash(O_1) $\langle B, +1 \rangle, \langle B, +1 \rangle, \langle T, -1 \rangle, \langle T, +1 \rangle$
hash(O_2) $\langle B, -1 \rangle, \langle T, +1 \rangle, \langle B, +1 \rangle, \langle T, +1 \rangle$
hash(O_3) $\langle B, -1 \rangle, \langle T, +1 \rangle, \langle T, -1 \rangle, \langle B, -1 \rangle$
hash(O_4) $\langle T, -1 \rangle, \langle B, +1 \rangle, \langle B, -1 \rangle, \langle T, -1 \rangle$



Example: CountSketch

Find the 2 Most Frequent Objects

$O_3 \ O_1 \ O_4 \ O_2 \ O_4 \ O_1 \ O_3 \ O_3 \ O_1$

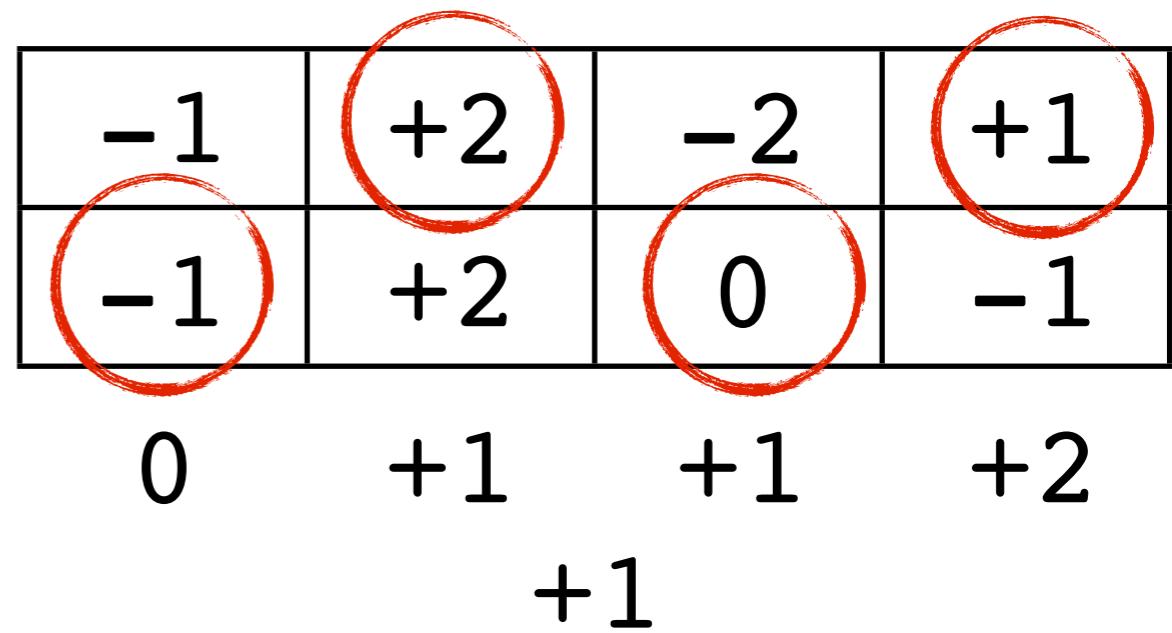


$\text{hash}(O_1) \langle B, +1 \rangle, \langle B, +1 \rangle, \langle T, -1 \rangle, \langle T, +1 \rangle$

$\text{hash}(O_2) \langle B, -1 \rangle, \langle T, +1 \rangle, \langle B, +1 \rangle, \langle T, +1 \rangle$

$\text{hash}(O_3) \langle B, -1 \rangle, \langle T, +1 \rangle, \langle T, -1 \rangle, \langle B, -1 \rangle$

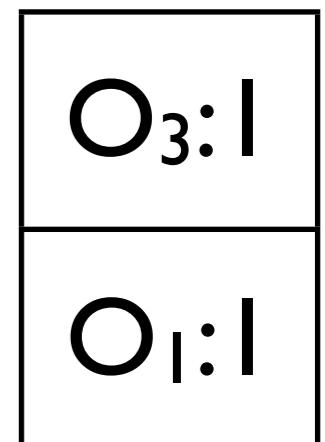
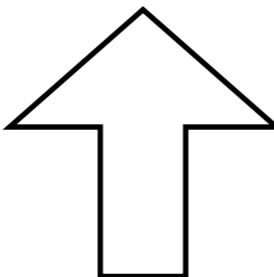
$\text{hash}(O_4) \langle T, -1 \rangle, \langle B, +1 \rangle, \langle B, -1 \rangle, \langle T, -1 \rangle$



Example: CountSketch

Find the 2 Most Frequent Objects

$O_3 \ O_1 \ O_4 \ O_2 \ O_4 \ O_1 \ O_3 \ O_3 \ O_1$

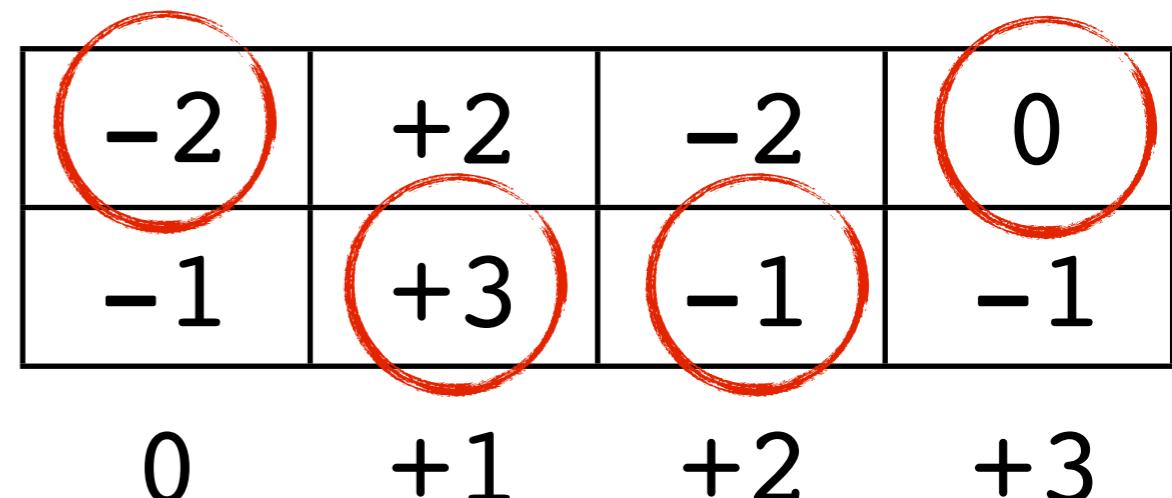


$\text{hash}(O_1) \langle B, +1 \rangle, \langle B, +1 \rangle, \langle T, -1 \rangle, \langle T, +1 \rangle$

$\text{hash}(O_2) \langle B, -1 \rangle, \langle T, +1 \rangle, \langle B, +1 \rangle, \langle T, +1 \rangle$

$\text{hash}(O_3) \langle B, -1 \rangle, \langle T, +1 \rangle, \langle T, -1 \rangle, \langle B, -1 \rangle$

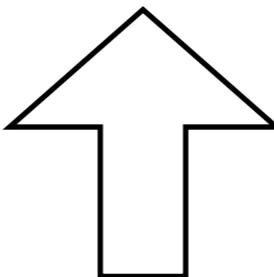
$\text{hash}(O_4) \langle T, -1 \rangle, \langle B, +1 \rangle, \langle B, -1 \rangle, \langle T, -1 \rangle$



Example: CountSketch

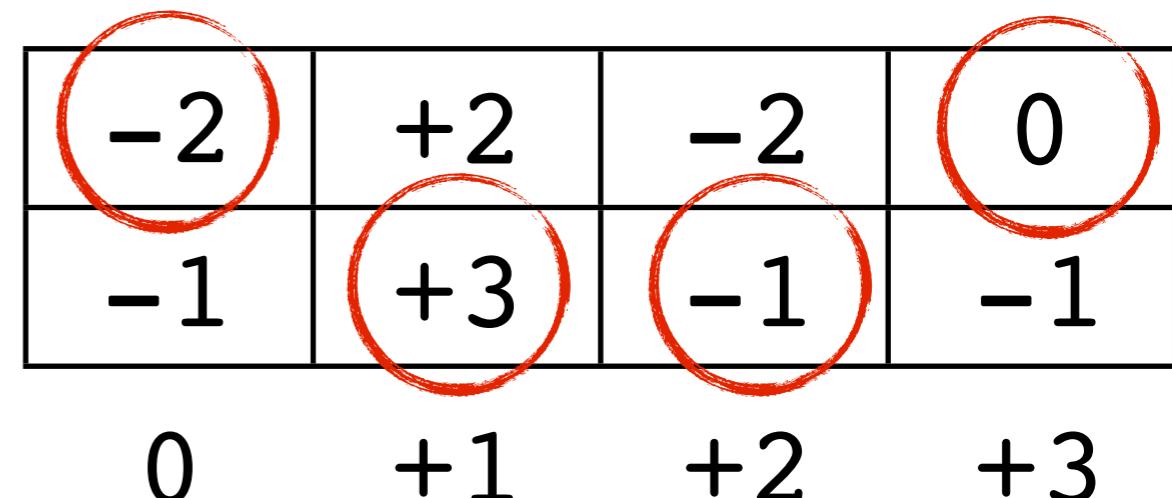
Find the 2 Most Frequent Objects

$O_3 \ O_1 \ O_4 \ O_2 \ O_4 \ O_1 \ O_3 \ O_3 \ O_1$



$O_4: 1.5$
$O_3: 1$

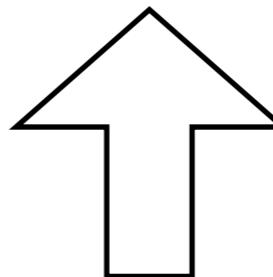
hash(O_1) $\langle B, +1 \rangle, \langle B, +1 \rangle, \langle T, -1 \rangle, \langle T, +1 \rangle$
hash(O_2) $\langle B, -1 \rangle, \langle T, +1 \rangle, \langle B, +1 \rangle, \langle T, +1 \rangle$
hash(O_3) $\langle B, -1 \rangle, \langle T, +1 \rangle, \langle T, -1 \rangle, \langle B, -1 \rangle$
hash(O_4) $\langle T, -1 \rangle, \langle B, +1 \rangle, \langle B, -1 \rangle, \langle T, -1 \rangle$



Example: CountSketch

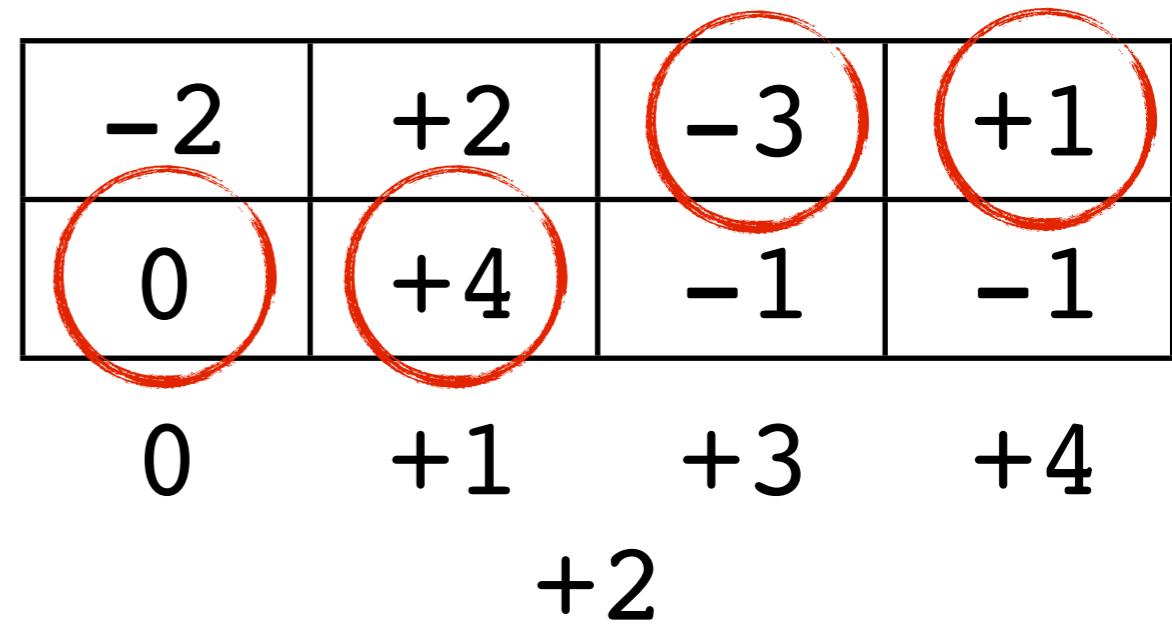
Find the 2 Most Frequent Objects

$O_3 \ O_1 \ O_4 \ O_2 \ O_4 \ O_1 \ O_3 \ O_3 \ O_1$



$O_4: 1.5$
$O_3: 1$

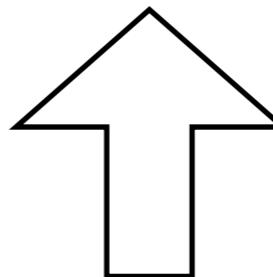
hash(O_1) $\langle B, +1 \rangle, \langle B, +1 \rangle, \langle T, -1 \rangle, \langle T, +1 \rangle$
hash(O_2) $\langle B, -1 \rangle, \langle T, +1 \rangle, \langle B, +1 \rangle, \langle T, +1 \rangle$
hash(O_3) $\langle B, -1 \rangle, \langle T, +1 \rangle, \langle T, -1 \rangle, \langle B, -1 \rangle$
hash(O_4) $\langle T, -1 \rangle, \langle B, +1 \rangle, \langle B, -1 \rangle, \langle T, -1 \rangle$



Example: CountSketch

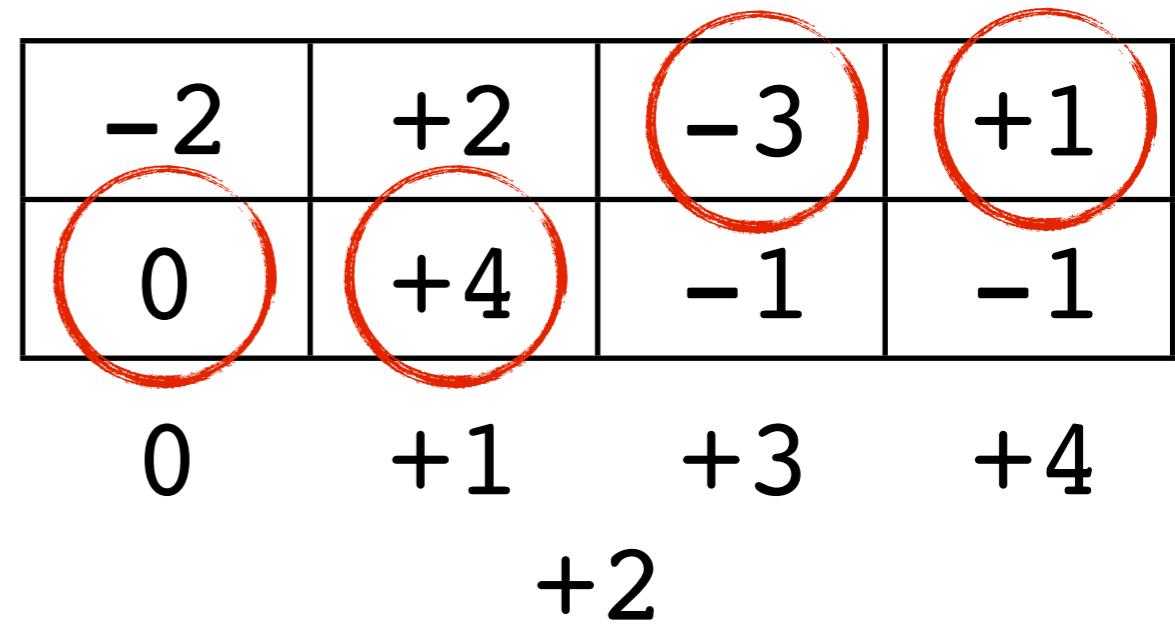
Find the 2 Most Frequent Objects

$O_3 \ O_1 \ O_4 \ O_2 \ O_4 \ O_1 \ O_3 \ O_3 \ O_1$



$O_1:2$
$O_4:1.5$

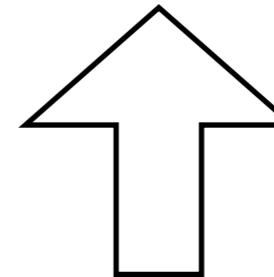
hash(O_1) $\langle B, +1 \rangle, \langle B, +1 \rangle, \langle T, -1 \rangle, \langle T, +1 \rangle$
hash(O_2) $\langle B, -1 \rangle, \langle T, +1 \rangle, \langle B, +1 \rangle, \langle T, +1 \rangle$
hash(O_3) $\langle B, -1 \rangle, \langle T, +1 \rangle, \langle T, -1 \rangle, \langle B, -1 \rangle$
hash(O_4) $\langle T, -1 \rangle, \langle B, +1 \rangle, \langle B, -1 \rangle, \langle T, -1 \rangle$



Example: CountSketch

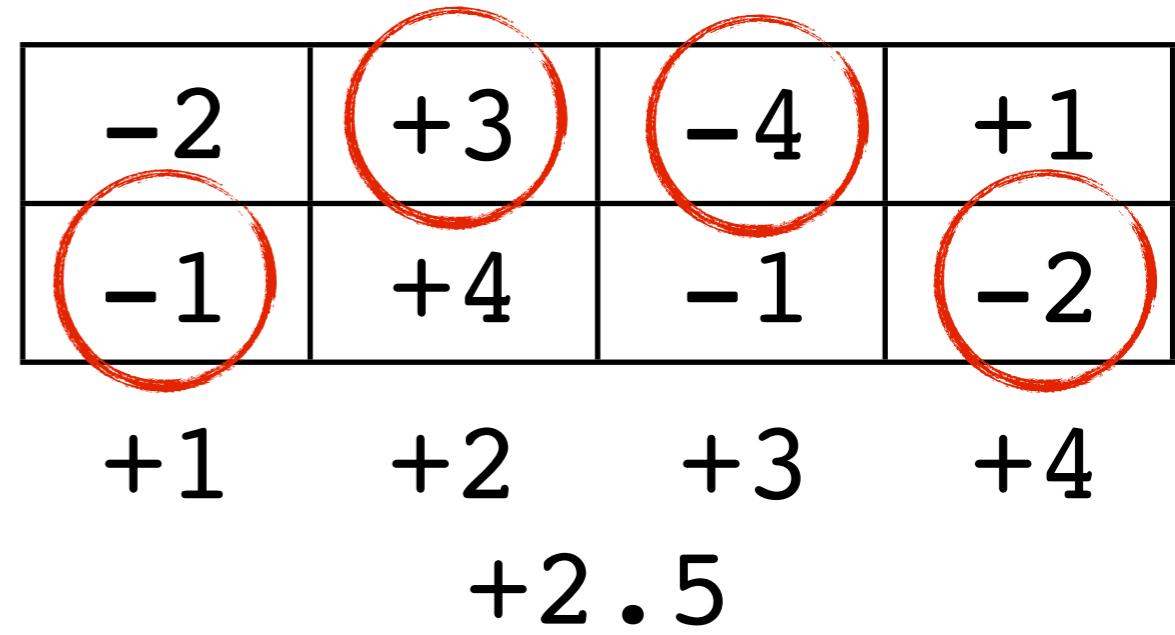
Find the 2 Most Frequent Objects

$O_3 \ O_1 \ O_4 \ O_2 \ O_4 \ O_1 \ O_3 \ O_3 \ O_1$



$O_1:2$
$O_4:1.5$

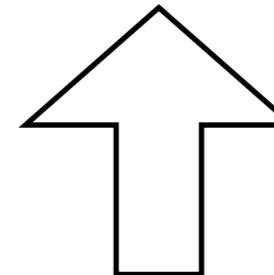
$\text{hash}(O_1) \langle B, +1 \rangle, \langle B, +1 \rangle, \langle T, -1 \rangle, \langle T, +1 \rangle$
 $\text{hash}(O_2) \langle B, -1 \rangle, \langle T, +1 \rangle, \langle B, +1 \rangle, \langle T, +1 \rangle$
 $\text{hash}(O_3) \langle B, -1 \rangle, \langle T, +1 \rangle, \langle T, -1 \rangle, \langle B, -1 \rangle$
 $\text{hash}(O_4) \langle T, -1 \rangle, \langle B, +1 \rangle, \langle B, -1 \rangle, \langle T, -1 \rangle$



Example: CountSketch

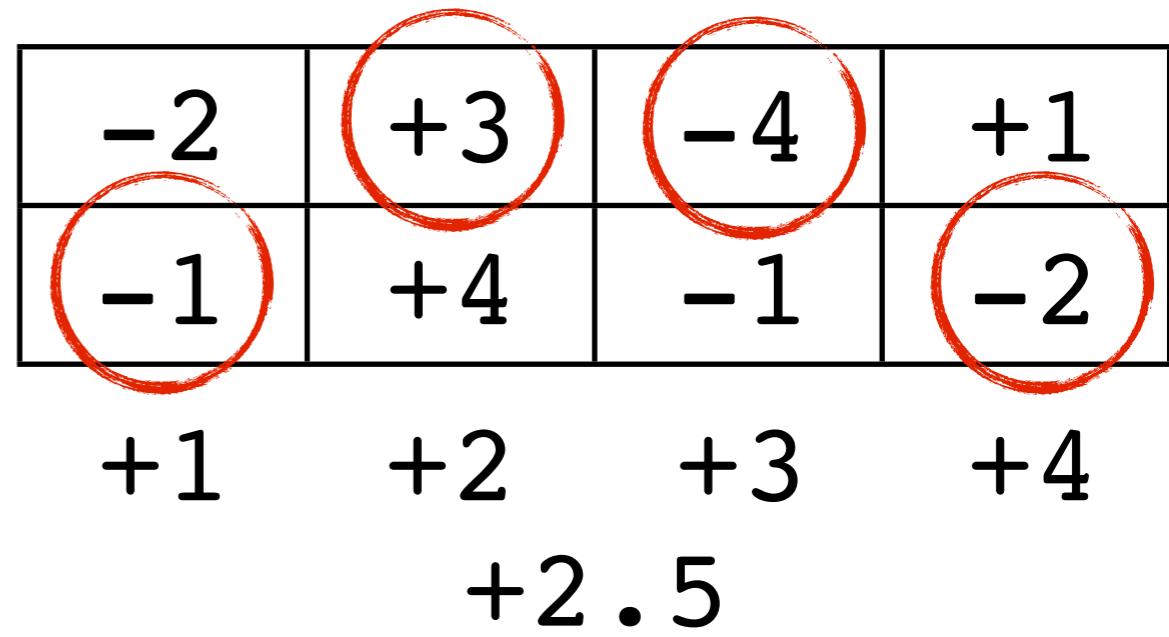
Find the 2 Most Frequent Objects

$O_3 \ O_1 \ O_4 \ O_2 \ O_4 \ O_1 \ O_3 \ O_3 \ O_1$



$O_3:2.5$
$O_1:2$

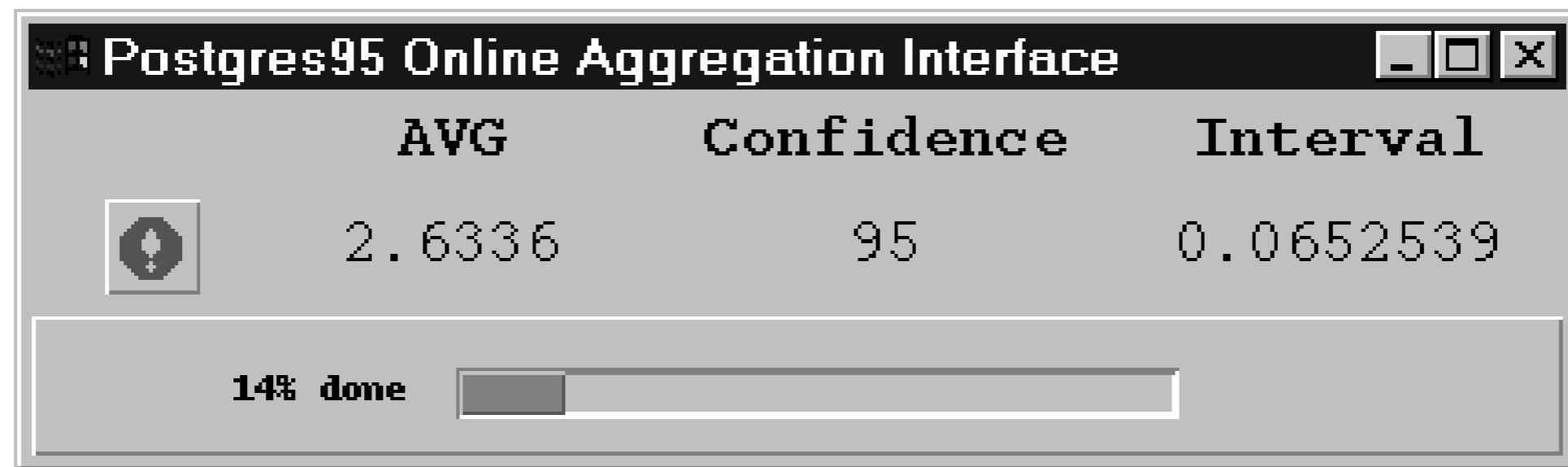
$\text{hash}(O_1) \langle B, +1 \rangle, \langle B, +1 \rangle, \langle T, -1 \rangle, \langle T, +1 \rangle$
 $\text{hash}(O_2) \langle B, -1 \rangle, \langle T, +1 \rangle, \langle B, +1 \rangle, \langle T, +1 \rangle$
 $\text{hash}(O_3) \langle B, -1 \rangle, \langle T, +1 \rangle, \langle T, -1 \rangle, \langle B, -1 \rangle$
 $\text{hash}(O_4) \langle T, -1 \rangle, \langle B, +1 \rangle, \langle B, -1 \rangle, \langle T, -1 \rangle$



Online Aggregation

- Give (some) results immediately.
 - (Imprecise) results obtained by sampling.
- More time, more samples, more accuracy.
 - Query takes longer to complete.
 - ... but result estimate available sooner.

Online Aggregation



OA Challenges

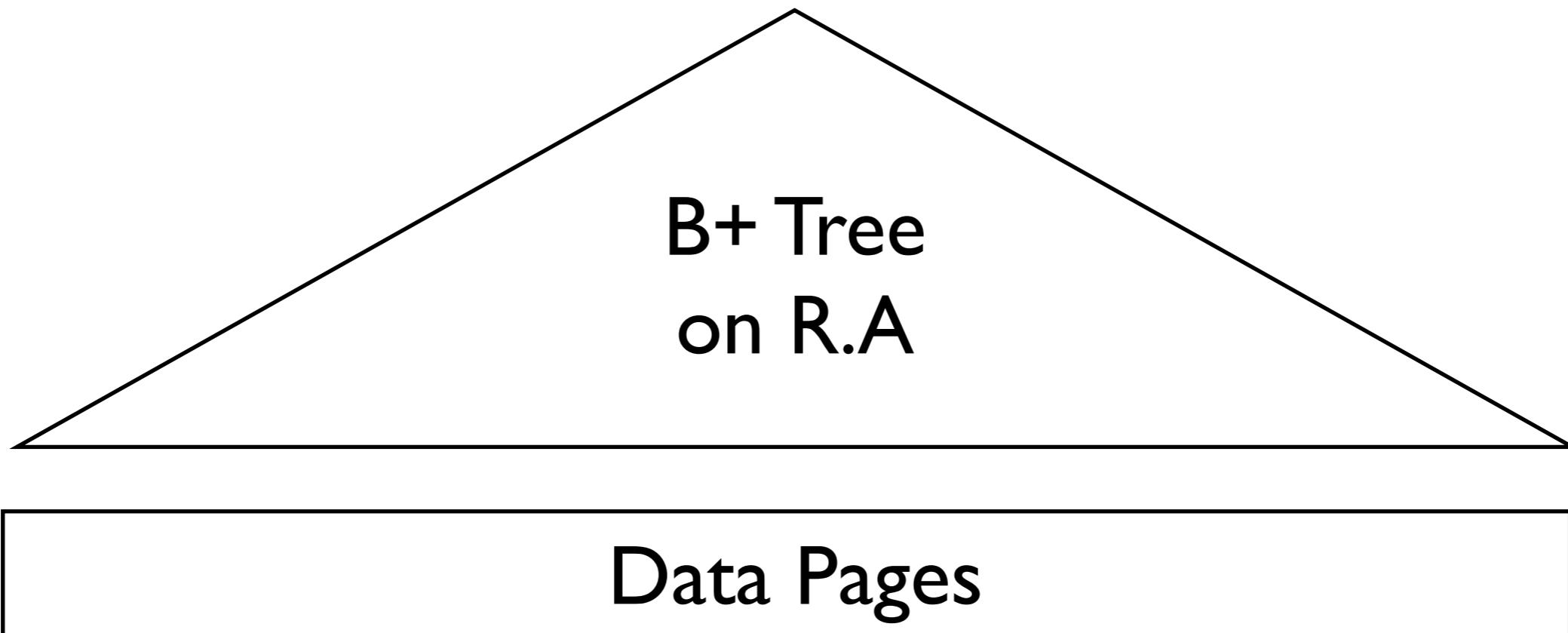
- Sampling: Need Random Access to Data.
 - Heap (Unsorted) Files, Flash Drives
- Fairness: Sampling For “Rare” Group-By Columns.
 - Index Striding
- Blocking: Joins Must Be Streamed.
 - Ripple Join

Fairness

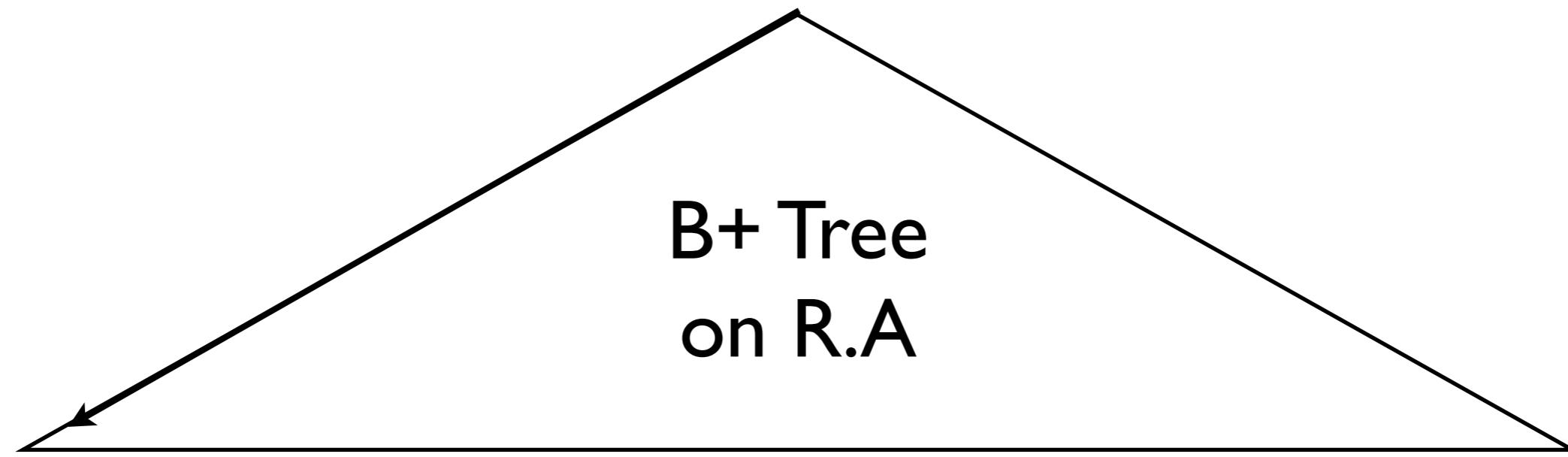
	major	Avg	Confidence	Interval
	1	2.27216	95	0.160417
	2	2.56146	95	0.160417
	3	2.66702	95	0.160417
	4	2.86235	95	0.160417
	5	3.12048	95	0.160417
	9	2.89645	95	0.160417

Goal: Keep all Confidences / Intervals in Synch.

Index Striding

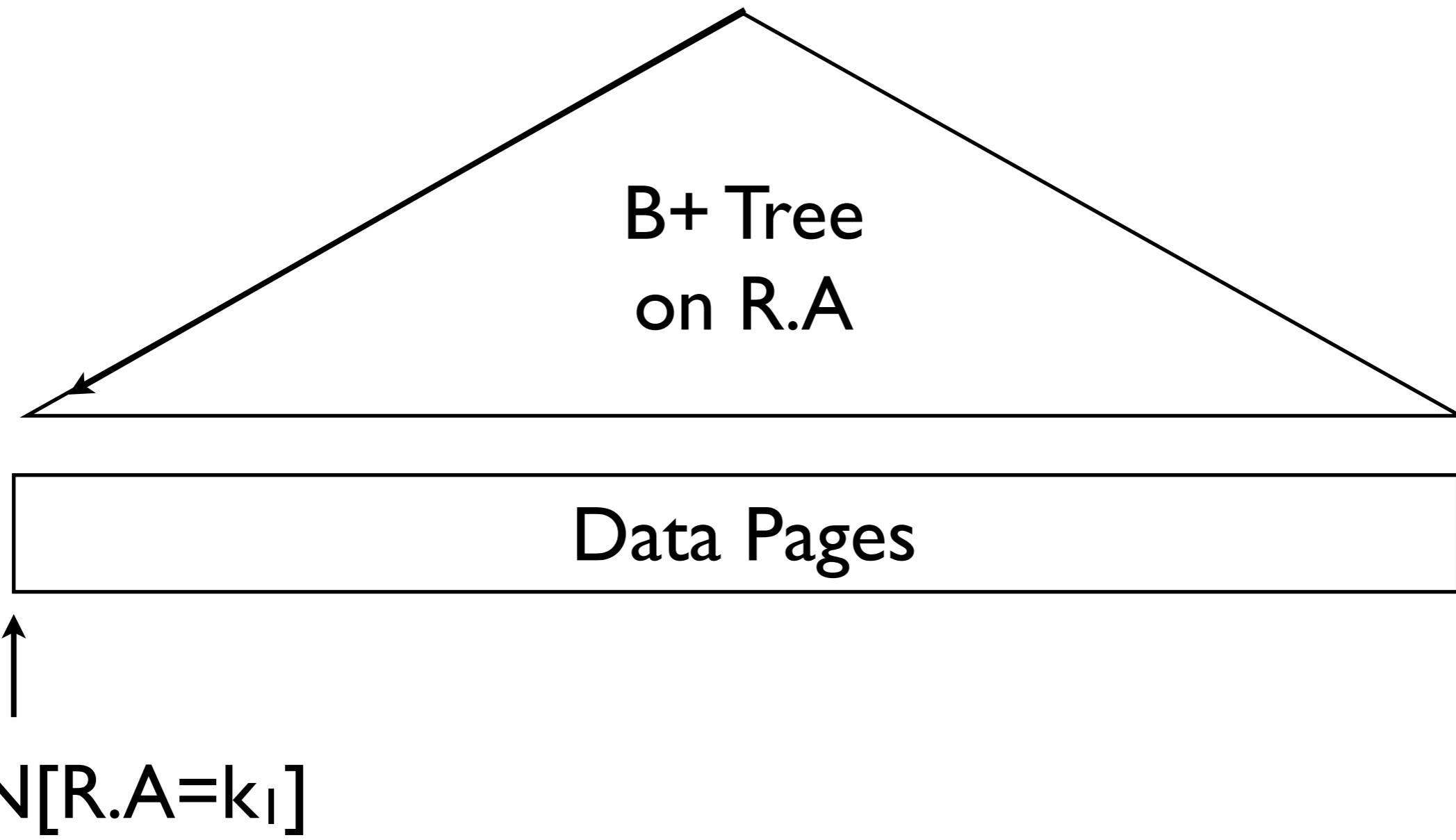


Index Striding

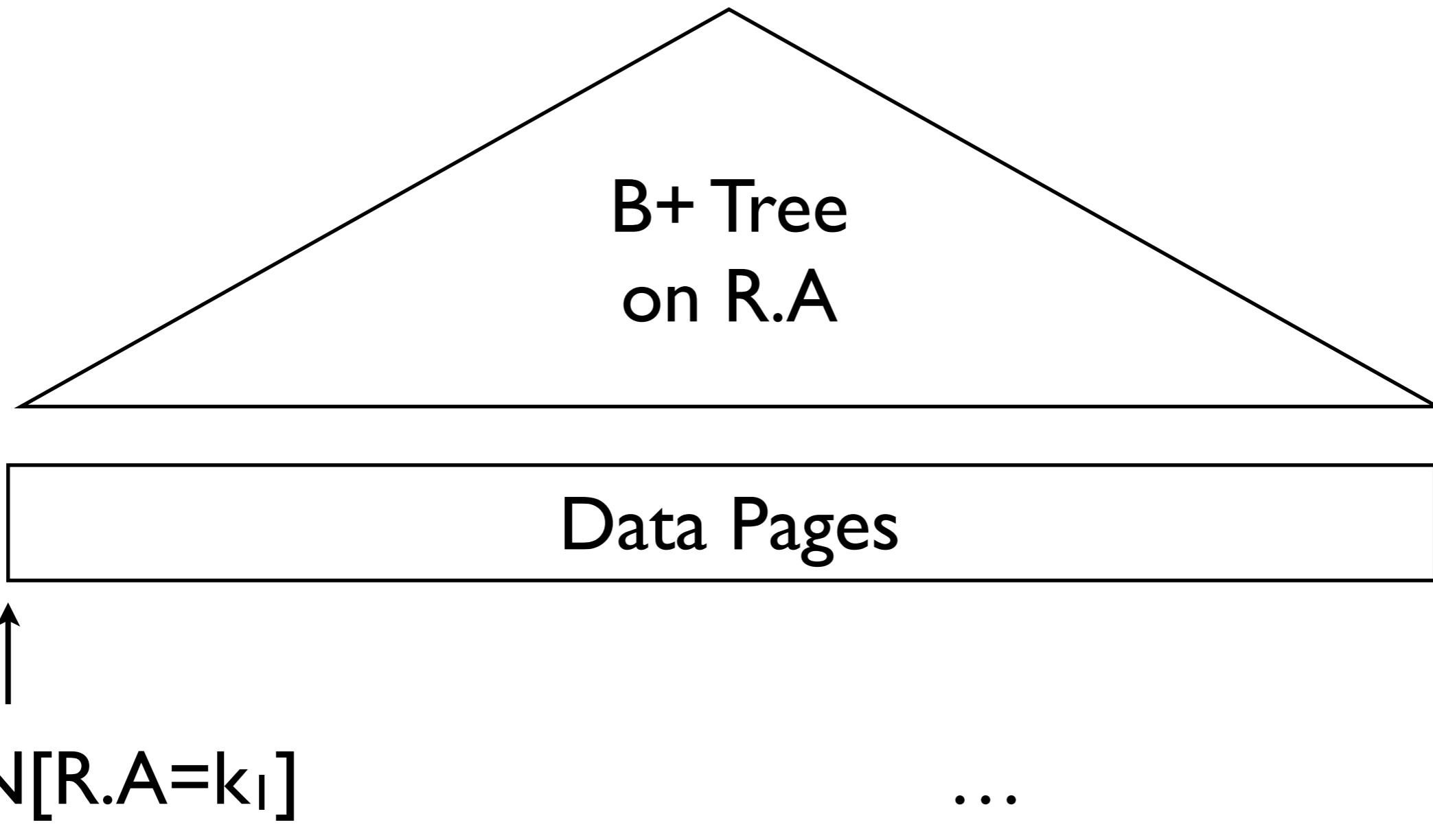


↑
SCAN[*]

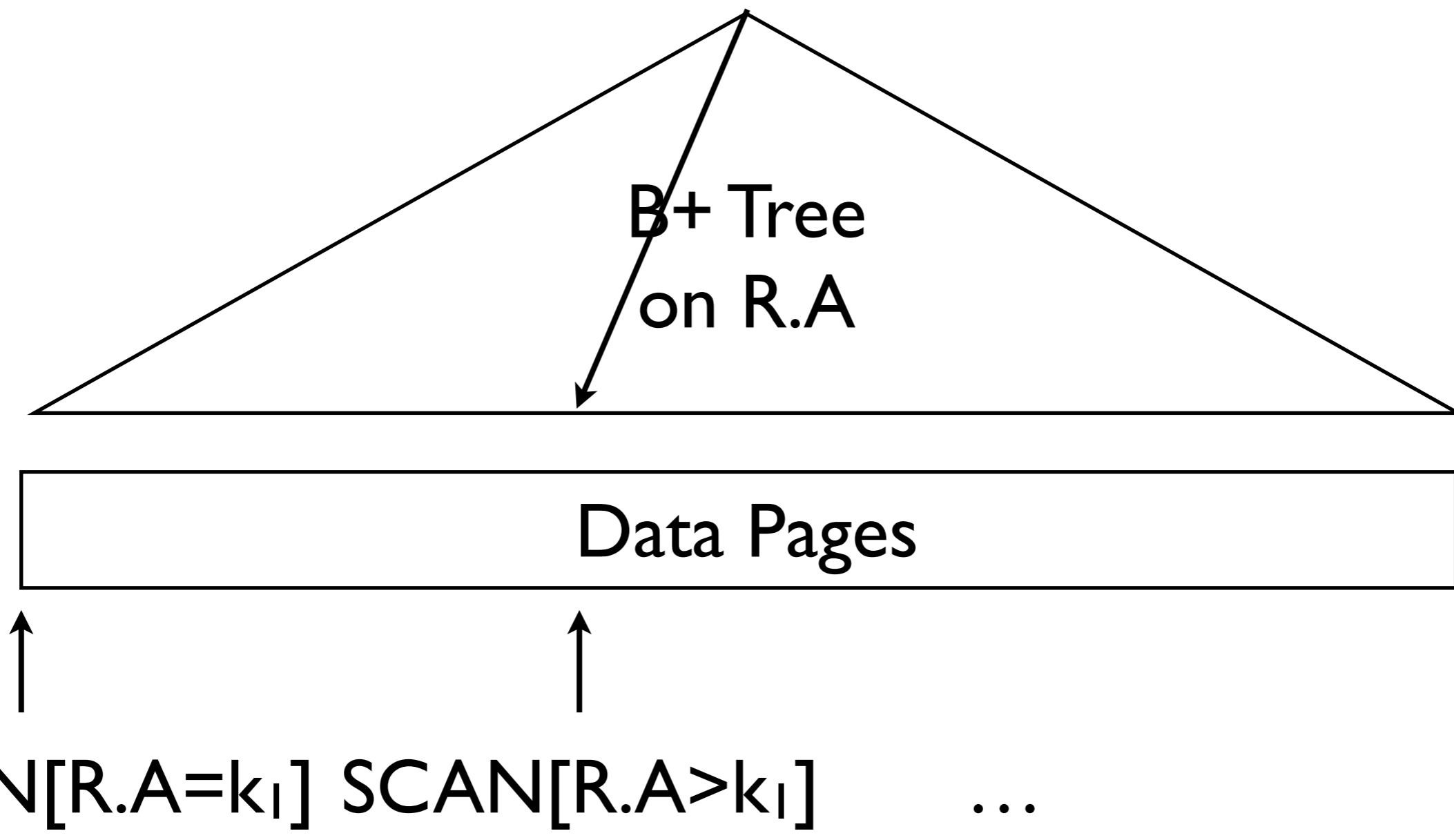
Index Striding



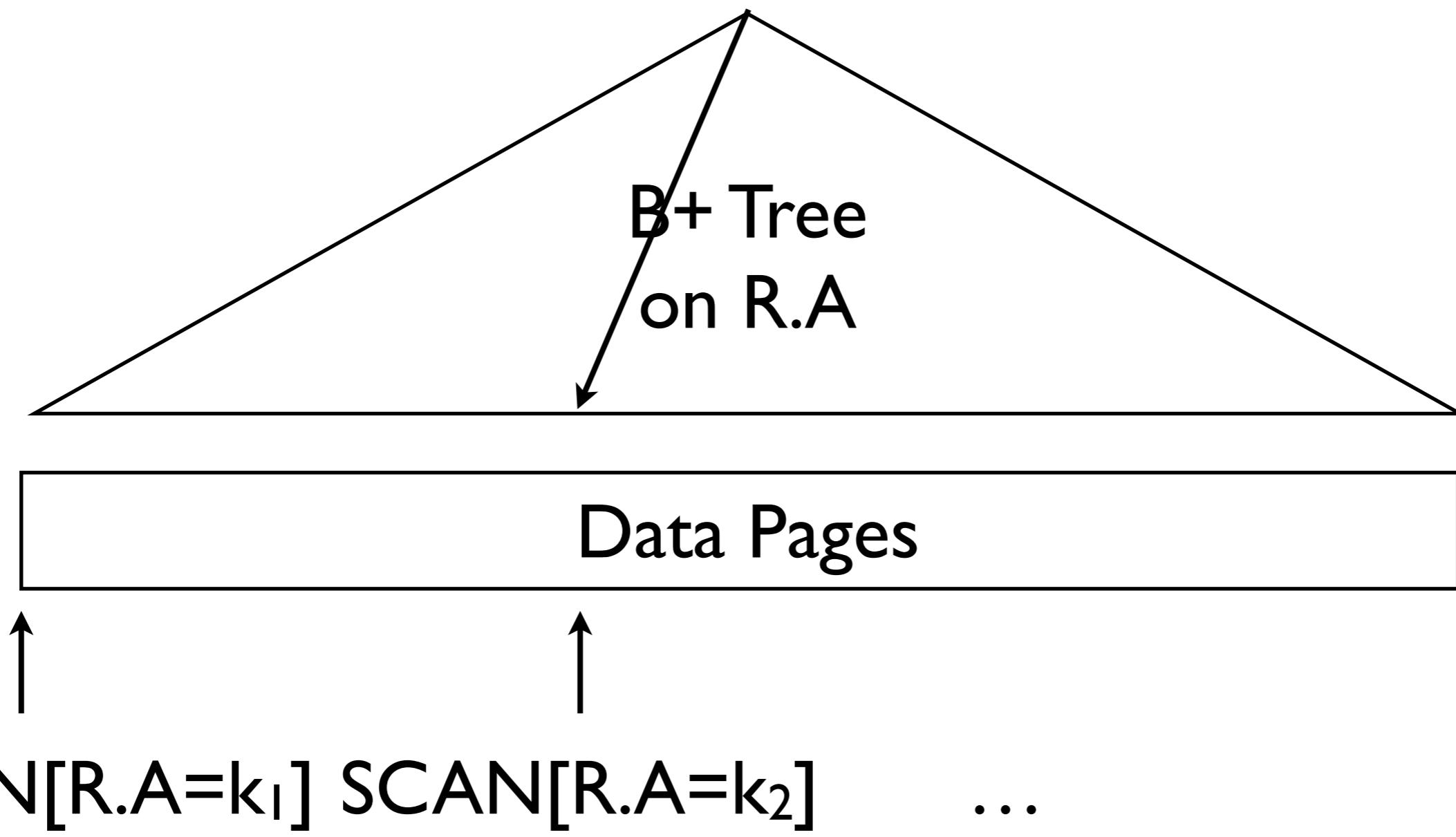
Index Striding



Index Striding



Index Striding



Index Striding

SCAN[R.A=k₁]

One Scan for Each GB Key

SCAN[R.A=k₂]

Each Scan is now a Heap Scan

SCAN[R.A=k₃]

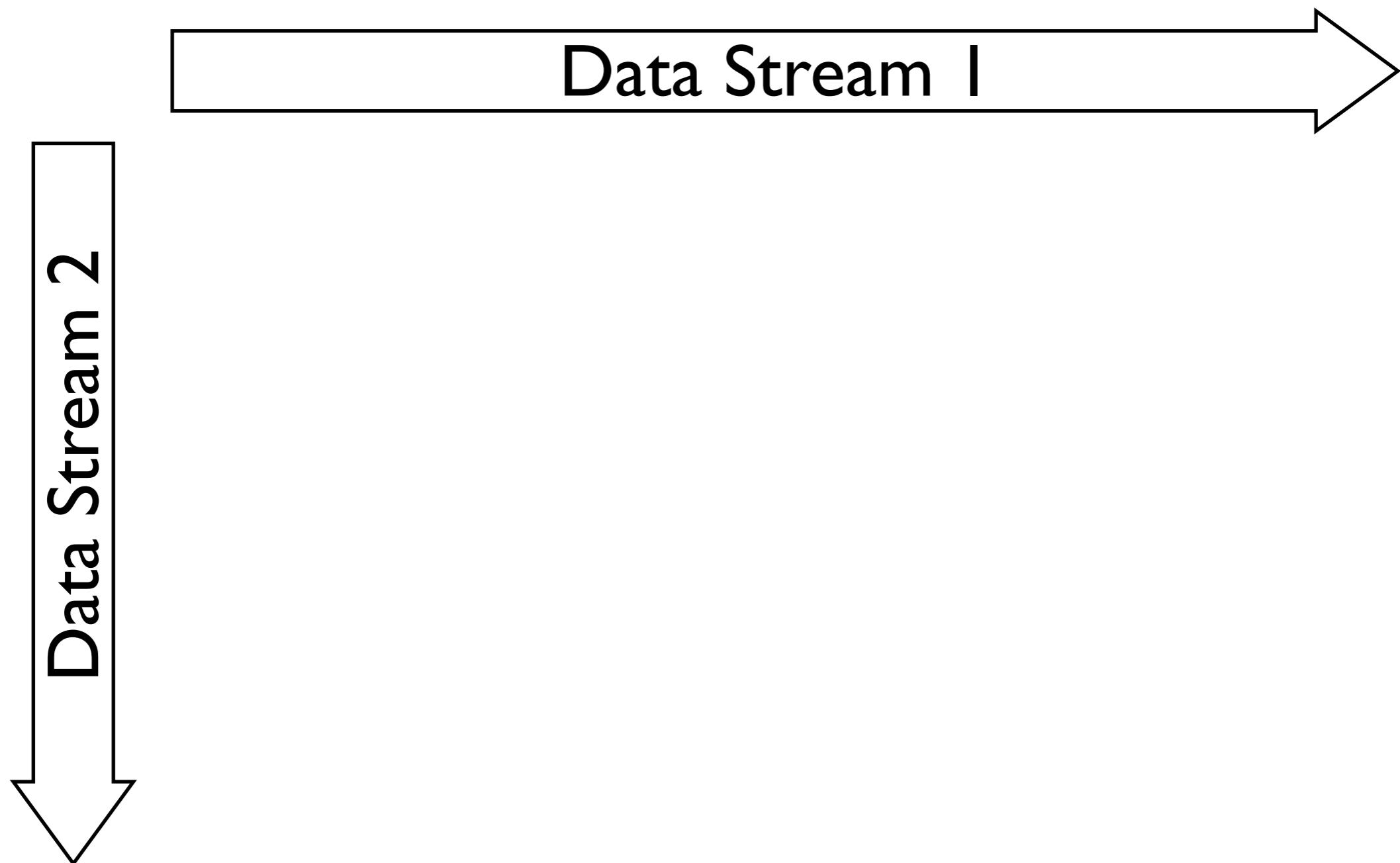
Split resources evenly between
each scan created.

...

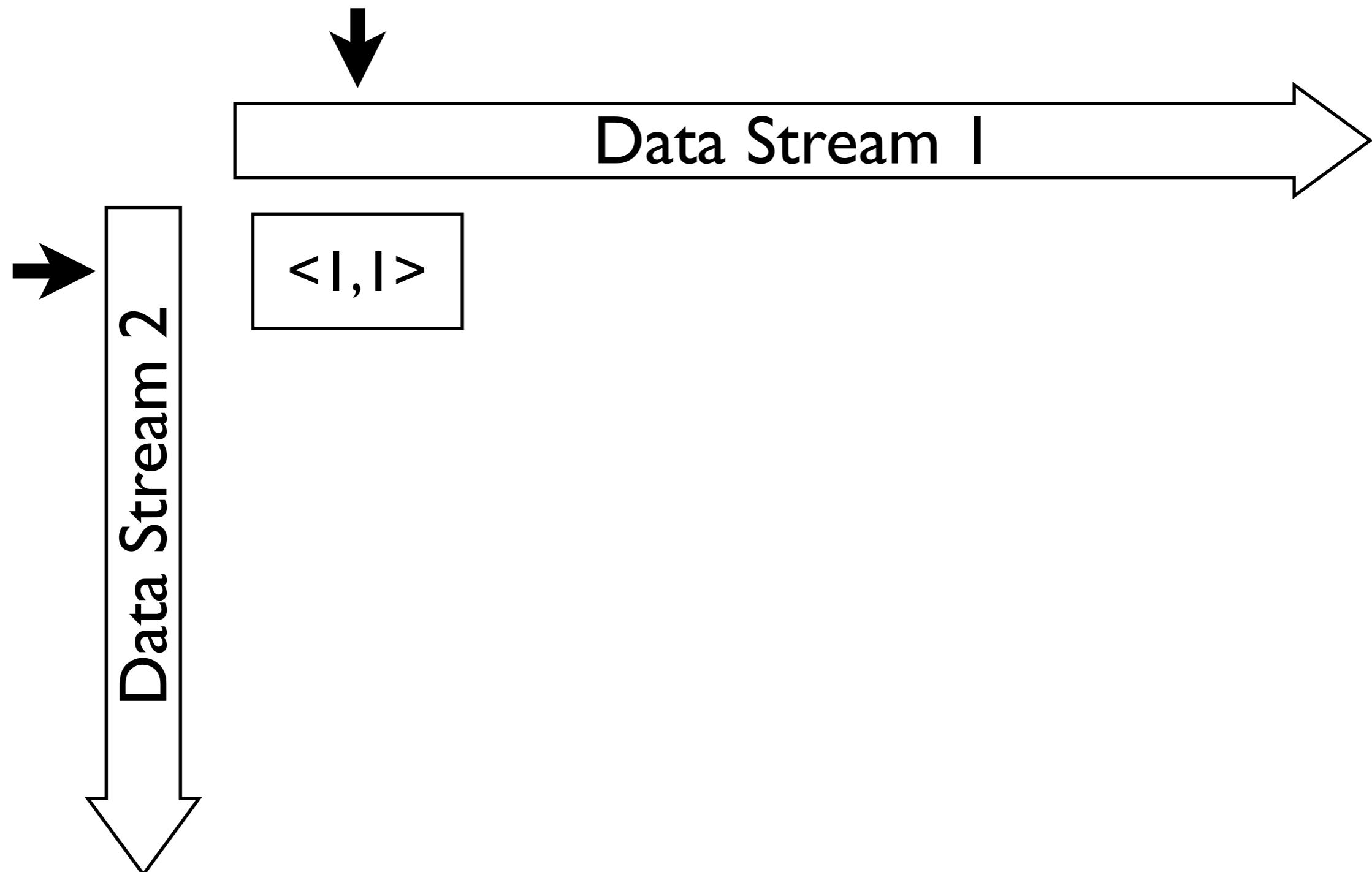
Non-Blocking Joins

- Sort/Merge Join
 - We want the data unsorted
- Index-Nested Loop Join
 - Could work if only few tuples matched.
- Hybrid Hash Join
 - Could work if one table is small.

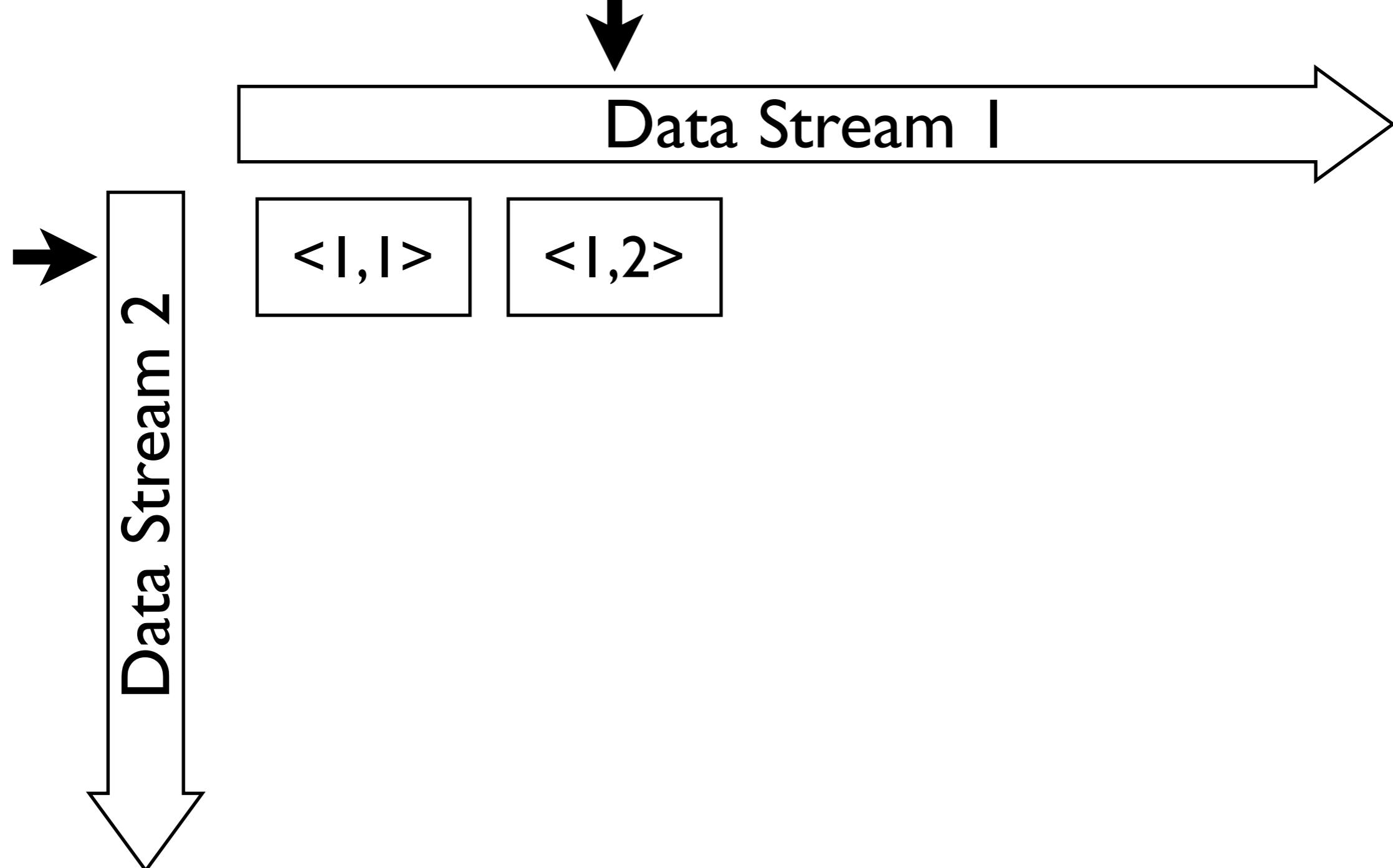
Ripple Join



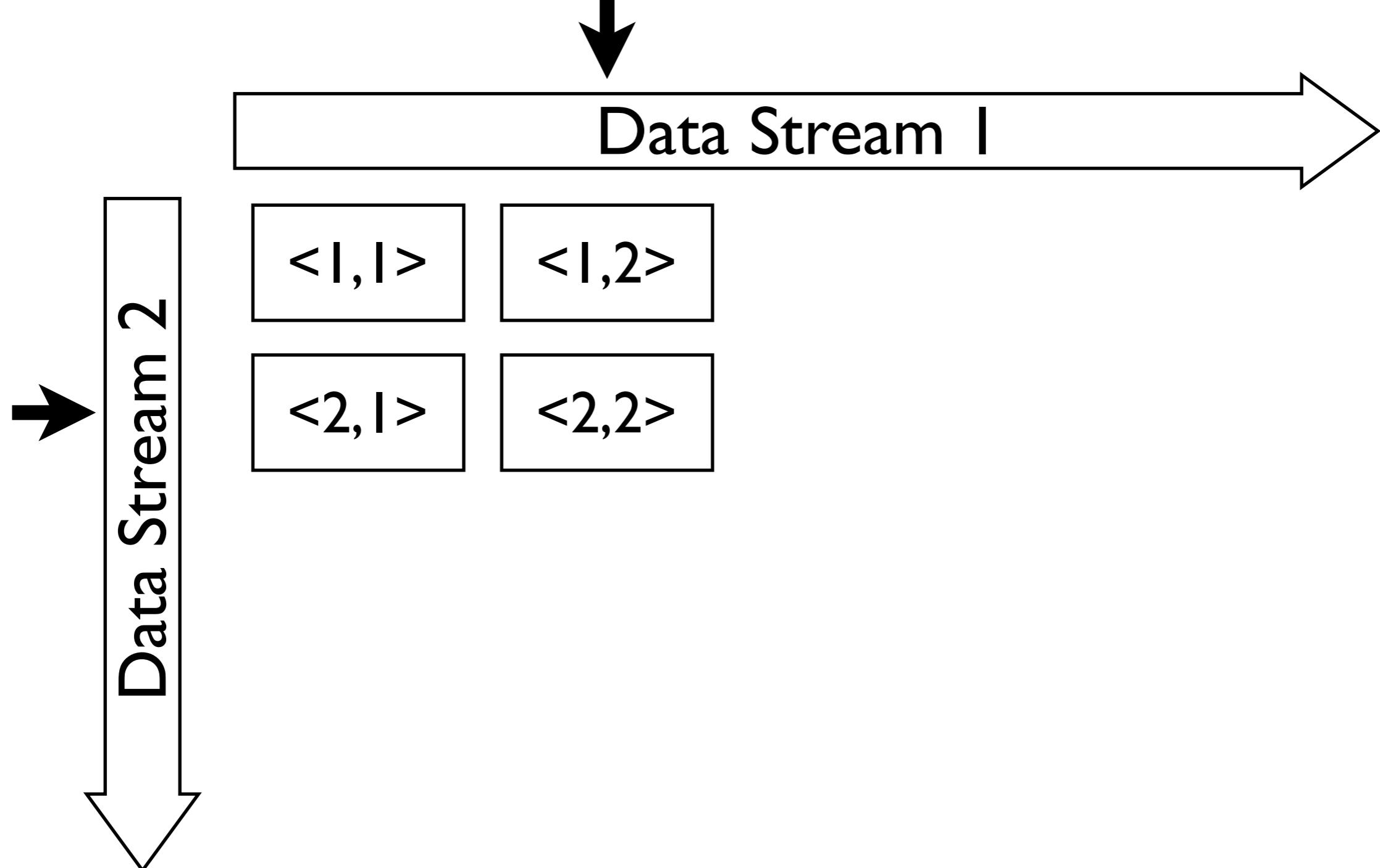
Ripple Join



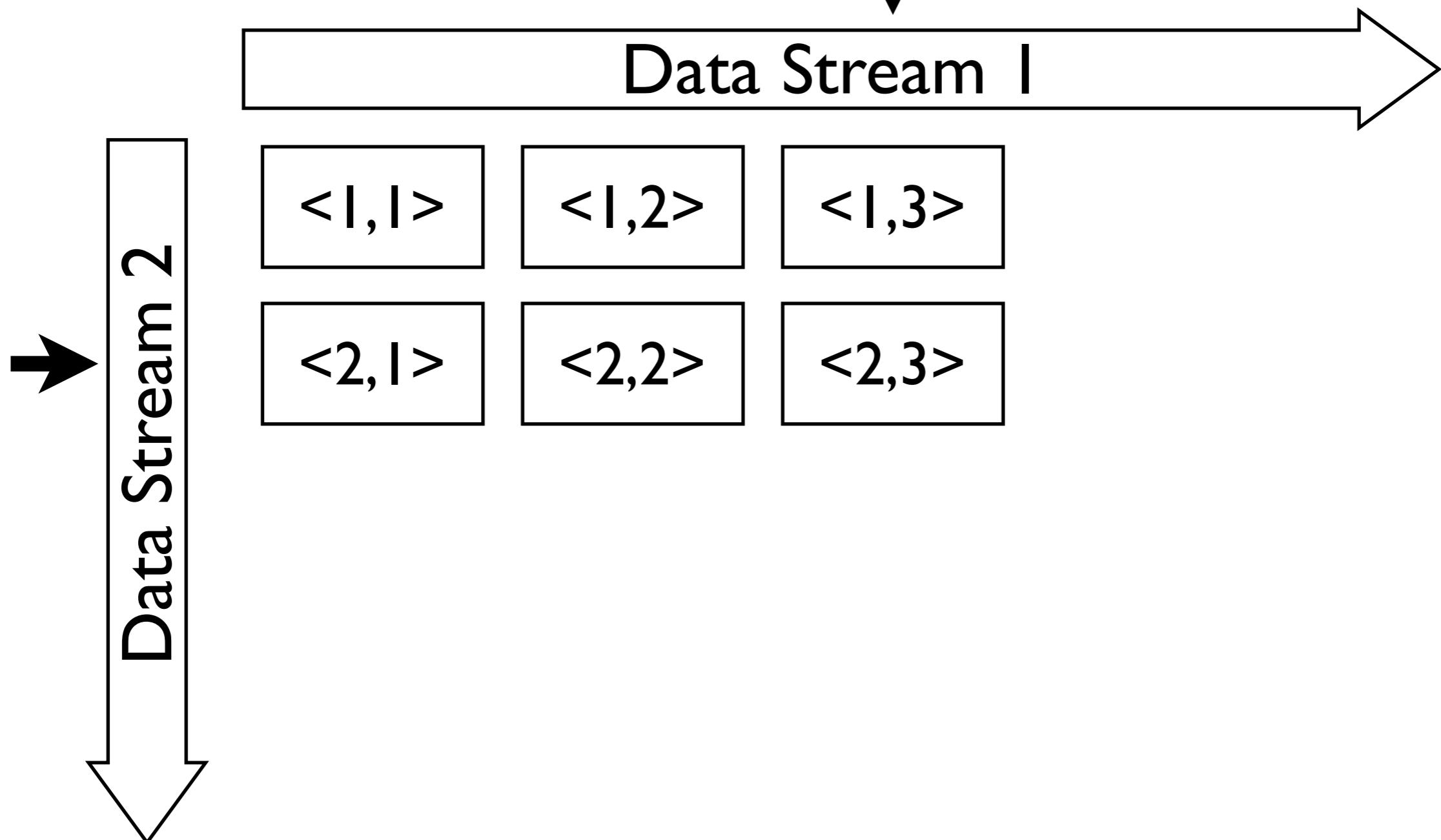
Ripple Join



Ripple Join



Ripple Join



Ripple Join



Ripple Join

