



# NOCAP: Near-Optimal Correlation-Aware Partitioning for Joins

*Zichen Zhu* Xiao Hu Manos Athanassoulis





**Skew Optimization** 

#### Skew Optimization in DHH



Skew optimization reduces the number of I/Os when the matching exhibits skew

#### Problems



Q2: When should we trigger skew optimization?

Input *S* 

# OCAP (Optimal Correlation-Aware Partitioning)

$$\arg\min_{P,m}\sum_{j=2}^{m+1} Cost(||R_j||, ||S_j||)$$

s.t. 
$$\forall i \in [n], \sum_{j=1}^{m+1} P_{i,j} = 1$$

$$||R_1|| + m + 2 \le B$$

$$P_{i,j} \in \{0,1\}, \forall i \in [n], \forall j \in [m+1]$$

 $||R_{j}|| = \sum_{i=1}^{n} P_{i,j} / b_{R}$  $||S_{j}|| = \sum_{i=1}^{n} P_{i,j} \cdot CT[i] / b_{S}$ 

Correlation Table (CT)	
Index i	Frequency in S
1	1
n-1	77
n	100

Define an  $n \times (m + 1)$  Boolean matrix P to represent the partitioning assignment

Notation	Meaning	
$n\left(n_{R} ight)$	The number of tuples in relation R	
m	The number of partitions on disk	
$ = \begin{bmatrix} 0 & \cdots & 1 \\ \vdots & \ddots & \vdots \\ 1 & \cdots & 0 \end{bmatrix}_{n \times (m+1)} $	A Boolean matrix P where $P_{i,j} = 1$ represents the $i^{th}$ record belongs to the $j^{th}$ partition	
$R_1$	A partition cached in memory	
$Cost(  R_j  ,   S_j  )$	The write cost of $R_j$ , $S_j$ plus the join cost between $R_j$ , $S_j$	
Exponential searching space !		
	↓ ↓	
$O(n^2 \cdot \log n)$	gB/B) for PK-FK joins	

# Practical Challenges for OCAP

1. We cannot have the whole CT in practice

Index <i>i</i>	Frequency in S
1	1
10M	1000

2. Partitioning assignment also occupies memory 
$$P = \begin{bmatrix} 0 & \cdots & 1 \\ \vdots & \ddots & \vdots \\ 1 & \cdots & 0 \end{bmatrix}_{n \times (m+1)}$$

### Observations from OCAP

01: MCVs can be prioritized in *two* ways: build an in-memory hash table (if B is large) or **assign them into a small partition on disk (if B is small)** 

02: We should ensure on-disk partitions fully fill  $z \cdot (B - 2)$  pages ( $z \in Z^+$ )



# Prioritizing MCVs with Constrained Memory



#### Total Available Memory NOCAP *HT<sub>mem</sub>* useless Call DHH with $m_r$ pages Partitioning Phase of R 9 **Disk Partitions for R** $P_7$ $P_4$ ..... P<sub>1</sub> $P_2$ YES $P_7$ $P_4$ $r \in HS$ ? Split Hash Function h<sub>split</sub> NO Partitioning Workflow: $P_1$ $P_2$ NO 1 YES Input **R** $r \in f$ ? **OCAP** for top-k' frequent keys $P_7$ 9 **Disk Partitions for S** Input S **DHH** to partition the rest 0 0 0 0 0 $P_1$ $P_2$ Split Hash Function $h_{split}$ P<sub>4</sub> 🕨 NO YES $P_1 \overline{P_2}$ $r \in f$ ? NO Partitioning Phase of S Probe Hash Function $h_{probe}$ **Disk-resident Partition** YES **Staged Partition** Input Page 12 Join Output Page Probe In-memory Hash Table

# Experiment Setup

**Storage**: PCIe NVM SSD (15 μs for reading a 4KB page) **Measured read/write symmetry**:

random\_write\_latency/sequential\_read\_latency = 3.3

sequential\_write\_latency/sequential\_read\_latency = 3.2

**PK-FK join input size**: 1M #records join with 8M #records

Record size: 1KB per record

Page size: 4KB





#### Selected Experimental Results



Correlation-aware joins (**DHH and NOCAP**) can **adaptively** reduce I/O cost when it comes to a **skew** distribution.

Note: OCAP only represents a lower bound, not a practical algorithm <sup>14</sup>

Varying skew



While DHH helps reduce #I/Os, **NOCAP** can better exploit the correlation skew to **achieve even lower I/O cost**.

#### Summary of NOCAP

NOCAP join outperforms DHH by up to 30%, and the textbook GHJ by up to 4X. Even for uniform distribution, NOCAP outperforms DHH by up to 10%!







# Thanks

Q&A