Distributed Caching for Processing Raw Arrays

Weijie Zhao¹, **Florin Rusu**^{1,2}, Bin Dong², Kesheng Wu², Anna Y. Q. Ho³, and Peter Nugent²

 $^1 \text{University}$ of California Merced, $^2 \text{Lawrence}$ Berkeley National Laboratory, $^3 \text{CalTech}$

July 2018





Palomar Transient Factory (PTF)



- 2,000-4,000 images 2,048x4,096 pixels per night
- 60-100 GB per night
- Complete processing of an image set within 30-45 minutes of aquisition

- 4 同 ト 4 ヨ ト 4 ヨ ト

Automatic Transient Identification & Classification



- 1-1.5 million candidates extracted per night; 10,000 every 45 minutes; 30-150 are real
- Real-or-bogus classification: Is the candidate real?
- What is the transient type of a real candidate: VarStar? SN/Nova? circumnuclear event? asteroid?

Array Construction



• Image $\xrightarrow{Image \ processing}$ Objects \rightarrow 3-D array [ra, dec, time]

э

- Array is chunked and stored in a shared-nothing array database
- Array is sparse and skewed

1 Motivation, Contributions, and Related Work

- 2 Arrays and Shape-Based Similarity Join
- 8 Raw Array Distributed Caching
- 4 Results and Conclusions

/□ ▶ ◀ ⋽ ▶ ◀

Raw PTF Data (FITS, HDF5, CSV)





(日)

- A file for each night; file ranges overlap
- High object cardinality variance across files
- Files stored across a distributed cluster
- Load and partition inside array database before query

- Design query-driven distributed caching framework for raw arrays
 - Identify cells to be cached locally and incrementally from each input file
 - Assign cells to nodes such that dependent cells are collocated
- Design an evolving R-tree index that refines the chunking of a sparse array to efficiently find the cells contained in a given subarray query
- Propose efficient cost-based algorithms for distributed cache eviction and placement that consider a historical query workload
- Evaluate experimentally distributed caching over raw arrays on two real sparse arrays with more than 1 billion cells stored in three file formats—FITS, HDF5, and CSV

・ 同 ト ・ ヨ ト ・ ヨ ト

- Raw or in-situ data processing over CSV and JSON files
 - NoDB, SCANRAW, OLA-RAW, VIDa, SLALOM, Proteus
- Raw array processing in SciDB (HDF5)
 - ArrayBridge
- Caching of raw semi-structured data (JSON, Parquet)
 - NoDB, ReCache, vertical partitioning, invisible loading
- Distributed caching in relational databases and web systems
 - LRU-K, DBMIN, LRFU, LRU-Threshold, Lowest-Latency-First, and Greedy-Dual-Size

・ 同 ト ・ ヨ ト ・ ヨ

• Distributed caching in Hadoop and Spark



2 Arrays and Shape-Based Similarity Join

- 3 Raw Array Distributed Caching
- 4 Results and Conclusions

/□ ▶ ◀ ⋽ ▶ ◀

Raw Array Data Model



- Raw array A<n:char,f:int>[i=1,6;j=1,8] consisting of 7 files distributed over 3 nodes
- Attributes n and f correspond to the node and the file id on the node, e.g., <X,2> is a cell in the second file on node X
- The cell color also shows the server on which the file is stored

• The dashed rectangle specifies a query

Array Similarity Join

- SELECT <code>f</code> INTO τ FROM α SIMILARITY JOIN β ON $\mathcal M$ WITH SHAPE σ
- f : computation function, τ : result array, α,β : input arrays,
- \mathcal{M} : array mapping function, σ : similarity shape array



Similarity Shape Array for Popular Distance Metrics



Motivation, Contributions, and Related Work

- Arrays and Shape-Based Similarity Join
- 3 Raw Array Distributed Caching
- 4 Results and Conclusions

/□ ▶ ◀ ⋽ ▶ ◀

Raw Array Distributed Caching Approach



- Build in-memory chunks incrementally based on the query (workload-driven chunking)
- Create higher granularity chunks for the entire file
- Design a global cost-based caching mechanism
- Combine query processing with cache replacement

・ 同 ト ・ ヨ ト ・ ヨ ト

Distributed Caching Architecture



W. Zhao, F. Rusu, B. Dong, K. Wu, A. Ho, and P. Nugent Distribute

Distributed Caching for Processing Raw Arrays

э

Raw Array Chunking

• Original raw array and resulting chunked array



• Query-driven chunking



Chunk Split Algorithm

- **Input:** Chunk α with bounding box \mathcal{BB}_{α} that intersects query subarray Q; Minimum number of cells threshold *MinC* **Output:** Chunks β and γ after splitting α
- 1: if (cells in $\alpha < MinC$) and (\exists cell in $\alpha \in Q$) then return
- 2: $\min_{v} = +\infty$
- 3: for each boundary $b \in Q$ that intersects with \mathcal{BB}_{lpha} do
- 4: $(\beta_b, \gamma_b) \leftarrow$ split cells in α into two sets by boundary b
- 5: **if** $\operatorname{vol}(\beta_b) + \operatorname{vol}(\gamma_b) < \min_{v}$ then
- 6: $\min_{v} \operatorname{vol}(\beta_b) + \operatorname{vol}(\gamma_b)$

7:
$$\beta \leftarrow \text{bounding}_box(\beta_b), \gamma \leftarrow \text{bounding}_box(\gamma_b)$$

- 8: end if
- 9: end for
 - Build an evolving R-tree incrementally based on queries
 - When to split? Cell number threshold or no cells at all

▲ □ ▶ ▲ □ ▶ ▲ □ ▶

• How to split? Always split into two chunks



- Find optimal caching plan
 - Centralized setting: minimize cache misses
 - Distributed setting: minimize network misses
- Weigh chunks based on how they are accessed in the query workload

・ 同 ト ・ ヨ ト ・ ヨ ト

- Access frequency: cache eviction
- Co-locality: cache placement

Cache Eviction

• Scan a file entirely even if one accessed chunk is not cached

• $cost_{evict}(Q_l, f_i, \{C_j\}) = w_{Q_l} \cdot \frac{size(f_i)}{\sum size(uncached C_j)}$

Input: Cache state as set of triples $S = \{(Q_l, f_i, \{C_j\})\}$ consisting of the chunks *j* accessed from file *i* at query *l*; Set of pairs $(f_i, \{C_j\})$ consisting of the chunks *j* accessed from file *i* at the current query Q_{l+1} ; Cumulated cache budget $B = \sum_k B_k$ consisting of local budgets B_k at each node

Output: Updated cache state $S' = \{(Q_{l+1}, f_i, \{C_j\})\}$

1:
$$S' \leftarrow \{(Q_{l+1}, f_i, \{C_j\})\}$$

- 2: while budget(S') < B do
- 3: Extract triple $t = (Q_I, f_i, \{C_j\})$ from S such that $budget(t \cup S') \le B$ and cost(t) is maximum
- 4: $S' \leftarrow S' \cup t$
- 5: $S \leftarrow S \setminus t$
- 6: Increase cost of triples $t' \in S$ that contain chunks in t
- 7: end while

Cache Placement

- Piggyback on the replication induced by query execution
- Preserve a single copy of every chunk
- Maximize co-locality of correlated chunks
- $cost_{place}(C_i, n, P', W) = \sum_{Q \in W} w_Q \cdot |C_j \in P'_n \land (C_i, C_j) \in Q|$

Input: Set $W = \{(Q_l, \{(C_i, C_j)\})\}$ consisting of chunk pairs (i, j) that join at query *I*; Set of locations $P = \{(C_i, \{N_k\})\}$ specifying all the nodes *k* that have a copy of cached chunk *i* at current query Q_{l+1} ; Cache budget B_k at node *k* **Output:** Updated locations $P' = \{(C_i, N_k)\}$

- 1: $P' \leftarrow \{p \in P\}$, where p has no replicas
- 2: for each $p = (C_i, \{N_k\}) \in P$ with multiple replicas do
- 3: Select node $n \in \{N_k\}$ such that $budget_n(C_i) \leq B_n$ and $cost(C_i, n, P', W)$ is maximum

< ロ > < 同 > < 三 > < 三 >

- 4: $P' \leftarrow P' \cup (C_i, n)$
- 5: end for

Motivation, Contributions, and Related Work

- Arrays and Shape-Based Similarity Join
- 3 Raw Array Distributed Caching



Experiments: Query Execution Time



W. Zhao, F. Rusu, B. Dong, K. Wu, A. Ho, and P. Nugent Distributed Caching for Processing Raw Arrays

-

First-Ever Neutron Star Merger Observation



 Produces gravitational waves and turns out to be the origin of heavy elements, including gold

(日)

• **Science** article: *Illuminating gravitational waves: A* concordant picture of photons from a neutron star merger

- Design query-driven distributed caching framework for raw arrays
- Design an evolving R-tree index that refines the chunking of a sparse array to efficiently find the cells contained in a given subarray query
- Propose efficient cost-based algorithms for distributed cache eviction and placement that consider a historical query workload

Thank you!

Questions?

W. Zhao, F. Rusu, B. Dong, K. Wu, A. Ho, and P. Nugent Distributed Caching for Processing Raw Arrays

э