



Similarity Join over Array Data

Weijie Zhao¹, Florin Rusu^{1,2}, Bin Dong², Kesheng Wu²

University of California, Merced¹
Lawrence Berkeley National Laboratory²

June 30, 2016

- 1 Cross-Matching Astronomical Catalogs
- 2 Array Similarity Join
- 3 Query Optimization
- 4 Experimental Evaluation
- 5 Conclusions

- Sloan Digital Sky Survey (SDSS), Palomar Transient Factory (PTF), etc.
- Transient identification: supernovae detection, exoplanet search



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

- Image $\xrightarrow{\text{Image processing}}$ Objects \rightarrow 3-D array [ra, dec, time]
- Array is chunked and stored in a shared-nothing architecture

- Image $\xrightarrow{\text{Image processing}}$ Objects \rightarrow 3-D array [ra, dec, time]
- Array is chunked and stored in a shared-nothing architecture
- Array is sparse and skewed

- Detect an object: Does it exist in a catalog?
- Coordinates for the same object in different catalogs may not be EXACTLY the same
- Find object's "neighbors" in the catalog
- → Join each cell with "neighbor cells" in catalog

- Formal array similarity join definition

- Formal array similarity join definition
- Array similarity join operator – data transfer, network, load balancing

- Formal array similarity join definition
- Array similarity join operator – data transfer, network, load balancing
- Modeling data transfer as vertex cover problem

- Formal array similarity join definition
- Array similarity join operator – data transfer, network, load balancing
- Modeling data transfer as vertex cover problem
- Experimentally compare against existing solutions on PTF catalog and LinkedGeoData

- Relational similarity join: ϵ – k dB tree, EGO family, etc.
 - Proposed operator does not require repartitioning and heavy preprocessing
- MapReduce similarity join: MRSimJoin, ClusterJoin, MR-DSJ, PHiDJ, etc.
 - Proposed operator takes advantage of array chunking to schedule join processing and reduce data transfer
- Array equi-join: structural join (Soroush 2011), shuffle join (Duggan 2015)
 - Proposed operator overlaps communication and computation
 - More complicated MIP optimization formulation
 - Graph formulation and heuristic solution
 - Generalization to more general joins

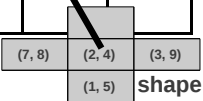
- 1 Cross-Matching Astronomical Catalogs
- 2 Array Similarity Join**
- 3 Query Optimization
- 4 Experimental Evaluation
- 5 Conclusions

i \ j	[1]	[2]	[3]	[4]	[5]	[6]	[7]	[8]
[1]	(3, 4)			(3, 3)			(8, 8)	(3, 1)
[2]		(9, 8)	(4, 4)	(9, 3)				(2, 4)
[3]	(2, 3)	(6, 6)		(1, 0)	(1, 1)			(4, 3)
[4]	(1, 5)		(9, 9)	(4, 5)	(7, 8)	(2, 4)	(3, 9)	(2, 8)
[5]			(8, 5)	(3, 8)		(1, 5)		
[6]			(5, 5)		(2, 5)			
[7]	(3, 4)	(7, 1)	(2, 2)			(6, 3)		
[8]	(7, 8)	(3, 2)		(8, 1)	(7, 5)	(7, 7)		

server X server Y server Z

APPLY operator (Baumann 1998, Marathe 2002)

i \ j	[1]	[2]	[3]	[4]	[5]	[6]	[7]	[8]
[1]	7			18			20	26
[2]		37	37	27				17
[3]	23	34		24	18			23
[4]	11		40	56	32	39	28	29
[5]			52	33		12		
[6]			27		7			
[7]	30	24	22			23		
[8]	27	28		21	35	35		



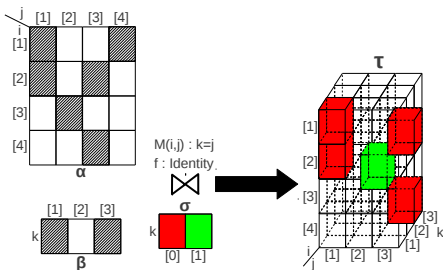
- SELECT f INTO τ FROM α SIMILARITY JOIN β ON \mathcal{M}
WITH SHAPE σ

- SELECT f INTO τ FROM α SIMILARITY JOIN β ON \mathcal{M}
WITH SHAPE σ

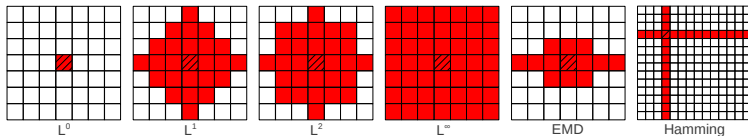
f : computation function, τ : result array, α, β : input arrays,
 \mathcal{M} : array mapping function, σ : similarity shape array

- SELECT f 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

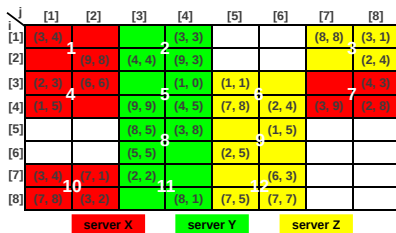


Similarity Shape Array for Popular Distance Metrics



- 1 Cross-Matching Astronomical Catalogs
- 2 Array Similarity Join
- 3 Query Optimization**
- 4 Experimental Evaluation
- 5 Conclusions

- Find all chunk pairs that have to be joined
- For each chunk pair (Υ, Ψ) , $\Upsilon \in \alpha$, $\Psi \in \beta$, Υ and Ψ must be transferred to the same node
- Compute f for each chunk pair transferred to the same node
- Overlap I/O (network and disk) and processing through multi-threading



$$\min \left\{ \begin{aligned} & \max_{ijk} \{ y_{ijkt} \cdot t \cdot T_{ntwk} \}, \\ & \max_k \left\{ \sum_{tij} y_{ijk(t-1)} \sum_{i' \in \Psi_i} (1 - z_{i'j(t-1)}) \cdot T_{disk} \right\} \end{aligned} \right\}$$

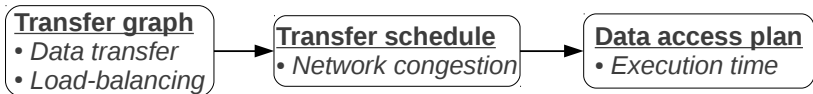
$$C_1 : x_{ijk} + x_{i'kj} = 1; \forall i, \forall i' \in \Psi_i, p(i) = j, p(i') = k$$

$$C_2 : \sum_t y_{ijkt} = x_{ijk}; \forall i, \forall j, \forall k \quad C_3 : \sum_{ijk} y_{ijkt} = 1; \forall t$$

$$C_4 : \sum_{ik} y_{ijkt} = 1; \forall t, \forall j \quad C_5 : \sum_{ij} y_{ijkt} = 1; \forall t, \forall k$$

$$C_6 : \sum_i z_{ijt} \leq B; \forall t, \forall j$$

$$C_7 : z_{i'jt} \leq z_{i'j(t-1)} + y_{ikjt}; \forall t, \forall k, \forall j, \forall i' \in \Psi_i$$



- Decompose the optimization formulation into 3 independent stages
- Each stage takes the optimization output of the previous stage as a pre-condition

- Consider chunk pair (Υ, Ψ) , Υ on Node X, Ψ on Node Y with $\text{size}(\Upsilon) = 3$, $\text{size}(\Psi) = 2$

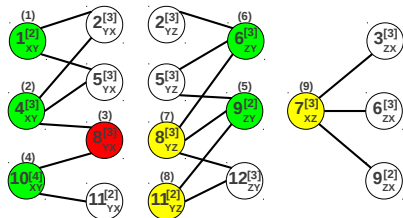
- Consider chunk pair (Υ, Ψ) , Υ on Node X , Ψ on Node Y with $\text{size}(\Upsilon) = 3$, $\text{size}(\Psi) = 2$
- Transfer the smaller chunk to the larger chunk's node?
($\Psi : Y \rightarrow X$)

- Consider chunk pair (Υ, Ψ) , Υ on Node X , Ψ on Node Y with $\text{size}(\Upsilon) = 3$, $\text{size}(\Psi) = 2$
- Transfer the smaller chunk to the larger chunk's node?
($\Psi : Y \rightarrow X$)
- No! If we have another chunk pair (Υ, Ψ') , Ψ' on Node Y with $\text{size}(\Psi') = 2$. $\Upsilon : X \rightarrow Y$ (cost 3) is better than $\Psi, \Psi' : Y \rightarrow X$ (cost $2 + 2 = 4$)

- Consider chunk pair (Υ, Ψ) , Υ on Node X, Ψ on Node Y with $\text{size}(\Upsilon) = 3$, $\text{size}(\Psi) = 2$
- Transfer the smaller chunk to the larger chunk's node? ($\Psi : Y \rightarrow X$)
- No! If we have another chunk pair (Υ, Ψ') , Ψ' on Node Y with $\text{size}(\Psi') = 2$. $\Upsilon : X \rightarrow Y$ (cost 3) is better than $\Psi, \Psi' : Y \rightarrow X$ (cost $2 + 2 = 4$)

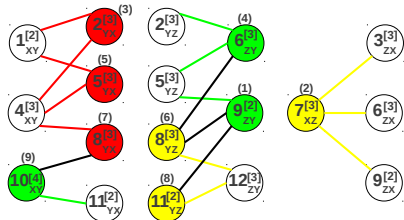
i \ j	[1]	[2]	[3]	[4]	[5]	[6]	[7]	[8]
[1]	(3, 4) 1		(3, 3) 2				(8, 8) 3	(3, 1)
[2]	(6, 6) 4	(4, 4) 5	(9, 3)				(2, 4)	
[3]	(2, 3)	(6, 6) 5	(1, 0)	(1, 1)			(4, 3) 7	
[4]	(1, 5) 4	(8, 9) 5	(4, 5)	(7, 8)	(2, 4) 6	(3, 9)	(2, 8)	
[5]			(8, 5) 8	(3, 8)		(1, 5)		
[6]			(5, 5) 8		(2, 5)			
[7]	(3, 4)	(7, 1)	(2, 2)			(6, 3)		
[8]	(7, 8) 10	(3, 2)	(8, 1) 11	(7, 5)	(7, 7)			

server X server Y server Z



- Consider chunk pair (Υ, Ψ) , Υ on Node X, Ψ on Node Y with $\text{size}(\Upsilon) = 3$, $\text{size}(\Psi) = 2$
- Transfer the smaller chunk to the larger chunk's node? ($\Psi : Y \rightarrow X$)
- No! If we have another chunk pair (Υ, Ψ') , Ψ' on Node Y with $\text{size}(\Psi') = 2$. $\Upsilon : X \rightarrow Y$ (cost 3) is better than $\Psi, \Psi' : Y \rightarrow X$ (cost $2 + 2 = 4$)

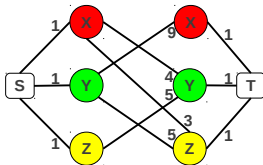
i \ j	[1]	[2]	[3]	[4]	[5]	[6]	[7]	[8]
[1]	(3, 4)		(3, 3)			(8, 8)	(3, 1)	
[2]	1	(6, 6)	2	(9, 3)			3	(2, 4)
[3]	(2, 3)	4	(6, 6)	5	(1, 0)	(1, 1)		7
[4]	(1, 5)		(8, 9)	5	(4, 5)	(7, 8)	(2, 4)	(3, 9)
[5]			(8, 5)	8	(3, 8)		(1, 5)	
[6]			(5, 5)		(2, 5)			
[7]	(3, 4)	10	(7, 1)			(6, 3)		
[8]	(7, 8)		(3, 2)	11	(8, 1)	(7, 5)	(7, 7)	
	server X			server Y			server Z	



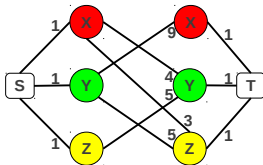
- Avoid many nodes talking to same receiver node at the same time

Transfer Schedule Optimization

- Avoid many nodes talking to same receiver node at the same time
- Node order: enforce strict node-to-node communication at any instant in time in order to minimize network congestion

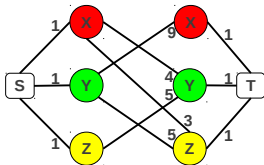


- Avoid many nodes talking to same receiver node at the same time
- Node order: enforce strict node-to-node communication at any instant in time in order to minimize network congestion



- In addition, even though we know two nodes will communicate in a period of time, what chunks in what order should be transferred?

- Avoid many nodes talking to same receiver node at the same time
- Node order: enforce strict node-to-node communication at any instant in time in order to minimize network congestion



- In addition, even though we know two nodes will communicate in a period of time, what chunks in what order should be transferred?
- Chunk order: choose chunk from the same sender node that shares the largest number of local chunks with the receiver

- Which chunk should be kicked out when memory budget is full?

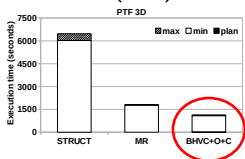
- Which chunk should be kicked out when memory budget is full?
- The one that will be used in the furthest future?

- Which chunk should be kicked out when memory budget is full?
- The one that will be used in the furthest future?
- According to chunk access plan, we know the future!

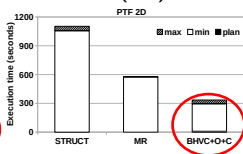
- 1 Cross-Matching Astronomical Catalogs
- 2 Array Similarity Join
- 3 Query Optimization
- 4 Experimental Evaluation**
- 5 Conclusions

Execution time:

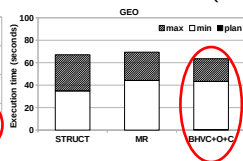
PTF 3D (L^∞)



PTF 2D (L^1)



LinkedGeoData (L^2)

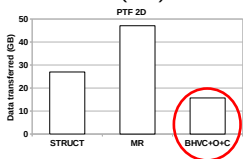


Data transferred:

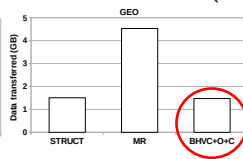
PTF 3D (L^∞)



PTF 2D (L^1)



LinkedGeoData (L^2)



- 1 Cross-Matching Astronomical Catalogs
- 2 Array Similarity Join
- 3 Query Optimization
- 4 Experimental Evaluation
- 5 Conclusions**

- We define formally array similarity join with a shape array

- We define formally array similarity join with a shape array
- We introduce the first array similarity join operator that minimizes the overall data transfer and network congestion while providing load-balancing in a multi-thread pipelined architecture

- We define formally array similarity join with a shape array
- We introduce the first array similarity join operator that minimizes the overall data transfer and network congestion while providing load-balancing in a multi-thread pipelined architecture
- We model the query optimization of array similarity join as a vertex cover problem and introduce efficient algorithms to find optimal execution plans

- We define formally array similarity join with a shape array
- We introduce the first array similarity join operator that minimizes the overall data transfer and network congestion while providing load-balancing in a multi-thread pipelined architecture
- We model the query optimization of array similarity join as a vertex cover problem and introduce efficient algorithms to find optimal execution plans
- The proposed array similarity join operator is running inside the PTF pipeline at Lawrence Berkeley National Lab's NERSC

Thank you!
Questions?