

# An in-memory Graph System for Scalable and Consistent Legacy System Integration



Bilal Arshad (PhD Researcher) – UK Systems Research Workshop 30<sup>th</sup> Nov -2<sup>nd</sup> Dec 2021

# Agenda

- Introduction – Automotive Data Integration
- Dealer Management Systems – DMS
- Graph-based data integration
- Entity resolution in data integration
- Evaluation
- Conclusion

# Automotive Dealerships - UK

- There are over 4K automotive dealerships within UK
- Each dealership could be part of a Franchise or be independent (multiple dealerships within a franchise)
- Data could be coming from various data sources with varying data types and formats (financial data-sets, telephony, sales, services etc)
- Each dealership could have its own format for data storage
- Unknown datasets coming from multiple data sources that require data transformation – Black Box
- Initially manual mapping is required to extract and transform data to store in a data warehouse

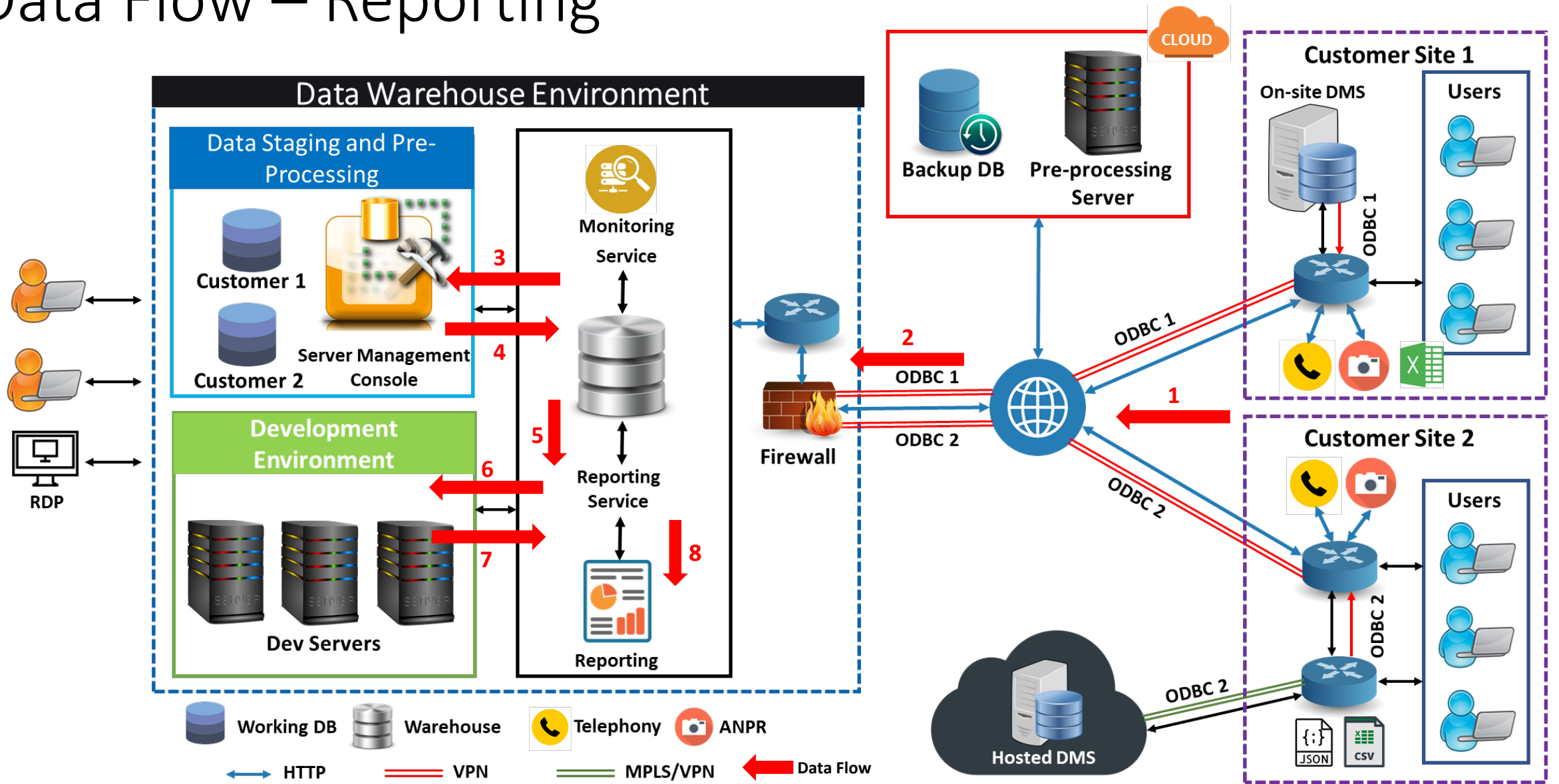
# Automotive Data Integration

- In order to provide a 360-view of a dealership's performance, data from multiple sources is integrated to provide a complete picture
- Management Information System (MIS) allows data from multiple sources to be brought together to provide a comprehensive real-time reporting dashboard with advanced analytics capabilities.
- Advantages of Integration
  - Feeds from management systems, telephony, account packages, sales tracking systems etc
  - Real-time data delivered as usable information
  - Acquire data from legacy systems – Dealer Management Systems (DMS)

# Dealer Management Systems - DMS

- These dealerships have deployed Dealer Management Systems (DMS) to manage:
  - Vehicle sales stock
  - Customer leads
  - Service appointments
  - Online advertising appointment
- DMS are Proprietary Software provided by limited market leaders for automotive dealerships – closed source software
- Legacy Systems ~ approx. four decades old
  - Some of these DMS are quite old but limited choice forces dealerships to continue usage
- In order to extract data Windows based license is required – ODBC
  - Currently no support for Linux based licenses thus limitation to design systems around this requirement

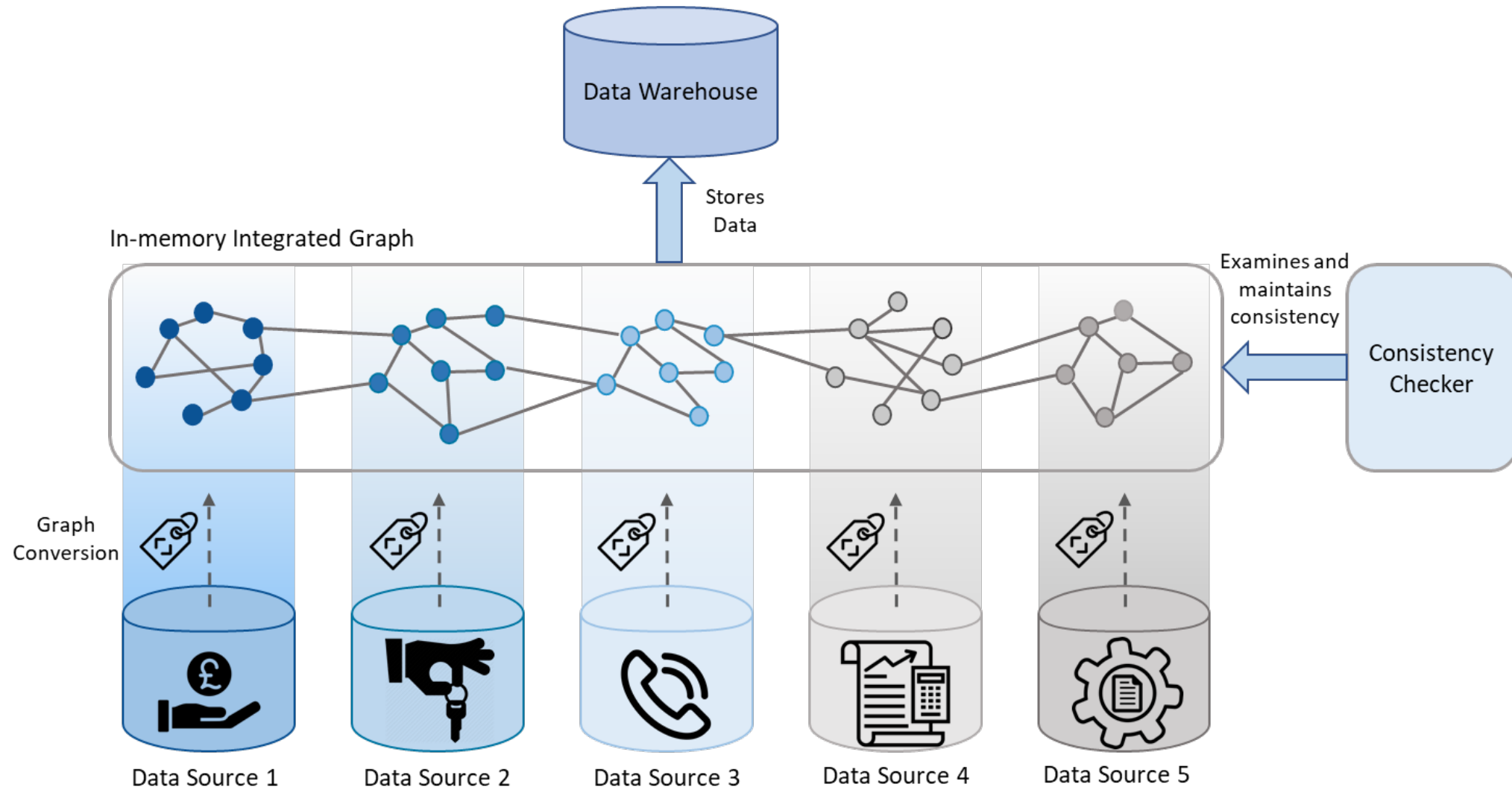
# Data Flow – Reporting



# Issues in data integration

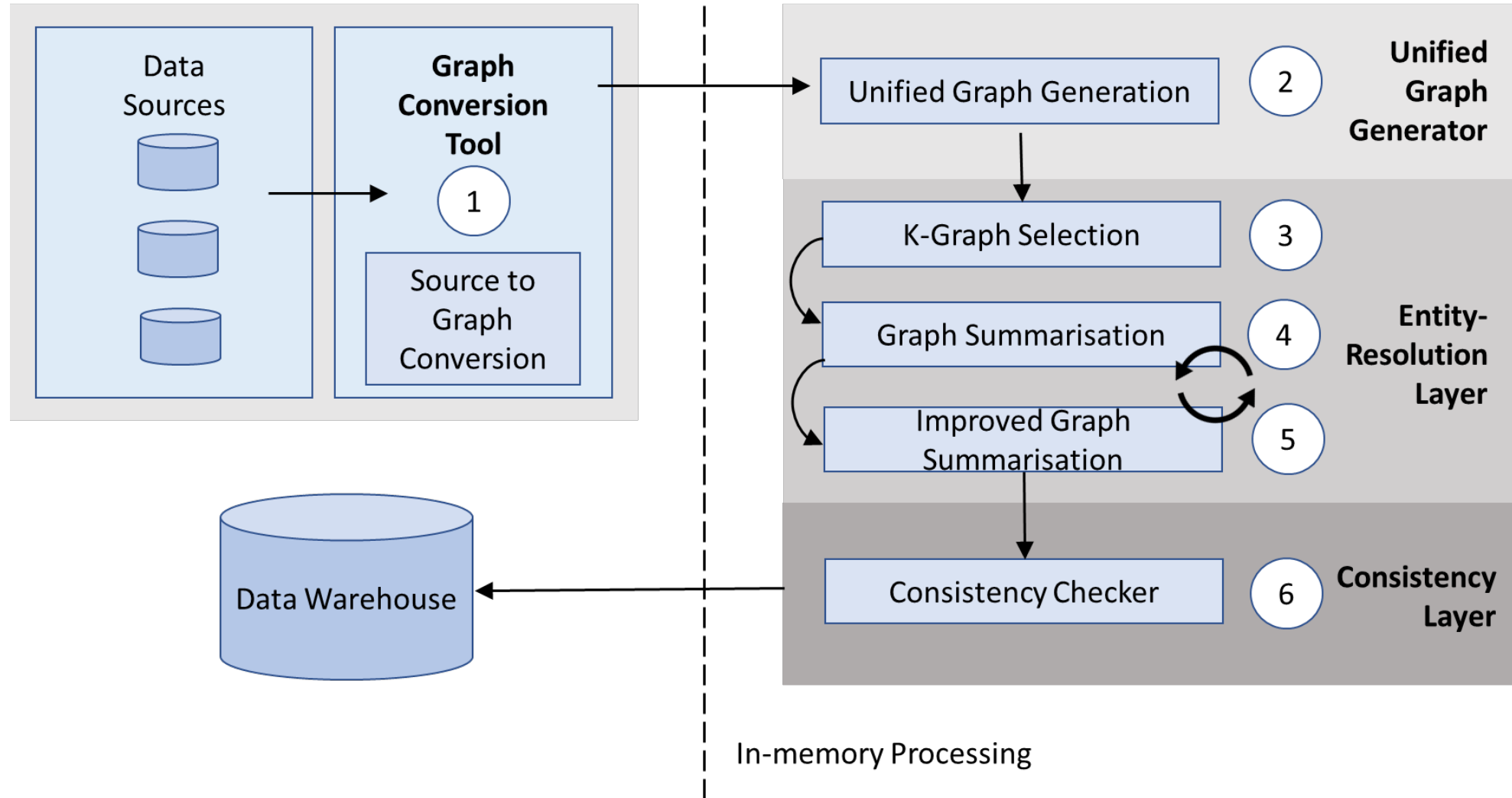
- Dealerships require integration of data from DMS and other data sources (telephony, ANPR, edge devices etc.)
- VPN links – cost
  - Solution is to set up an edge device that cuts down the VPN link's cost and send updated data only (lower bandwidth and processing cost)
- Issues:
  - Volume and velocity of change
  - Data consistency as these sources evolve

# Graph-based Data Integration





# Summarised Architecture



# Data-sets for Evaluation

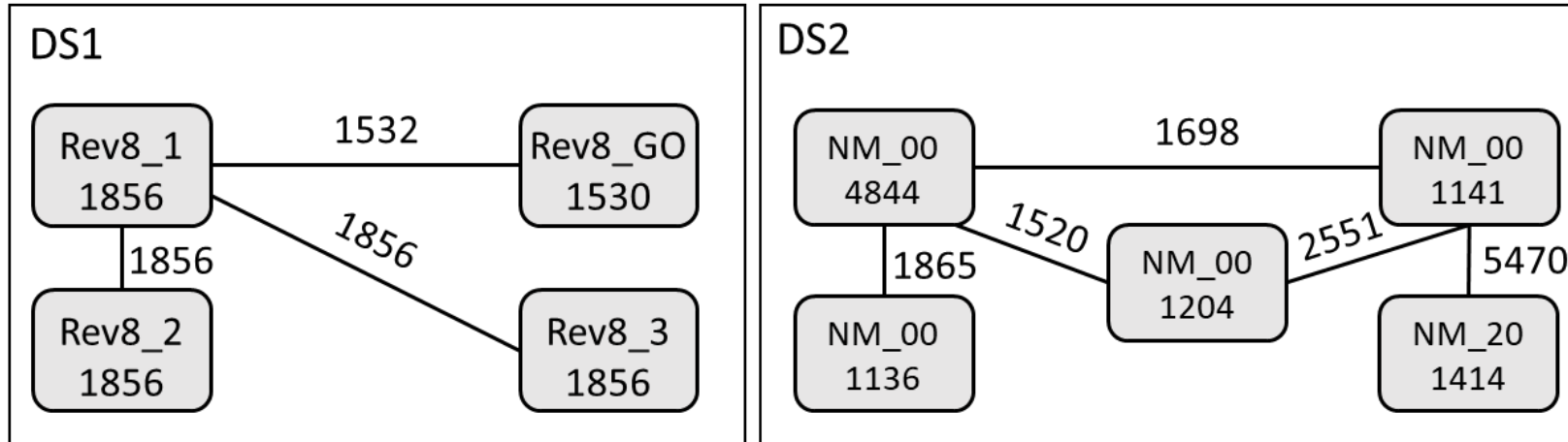
Datasets	Type	No. of Vertices	No. of Edges
Data-set 1	Real World	350	2875
Data-set 2	Real World	11600	65425
Data-set 3	Real World	25767	98598
Data-set 4	Real World	42494	109271
Data-set 5	Synthetic	65536	1048576
Data-set 6	Synthetic	131072	2097152

- We have collated data from three major DMS systems, Drive, Rev8 and Pinnacle
- The synthetic datasets are generated for scales 16 and 17 with average degree of 14 per vertex.
- The structures within these synthetic graphs are similar to the ones present in automotive data sets to ensure uniformity across the testbeds and results.

\*DataSynth - Arasu, A., Kaushik, R. and Li, J., 2011. DataSynth: Generating synthetic data using declarative constraints. Proceedings of the VLDB Endowment, 4(12), pp.1418-1421.

\*Graph500 RMat - Murphy, R.C., Wheeler, K.B., Barrett, B.W. and Ang, J.A., 2010. Introducing the graph 500. Cray Users Group (CUG), 19, pp.45-74.

# Entity Resolution Evaluation



Structures within DS1 and DS2 along with the number of links and entities

- A subset of the previously mentioned data-sets is presented for evaluation
- We begin by splitting these datasets into two named as DS1 gathered from Rev8 DMS and DS2 gathered from the Drive DMS.

# Cluster Sizes in DI Phases

Cluster sizes in integration phases for evaluation dataset - DS1

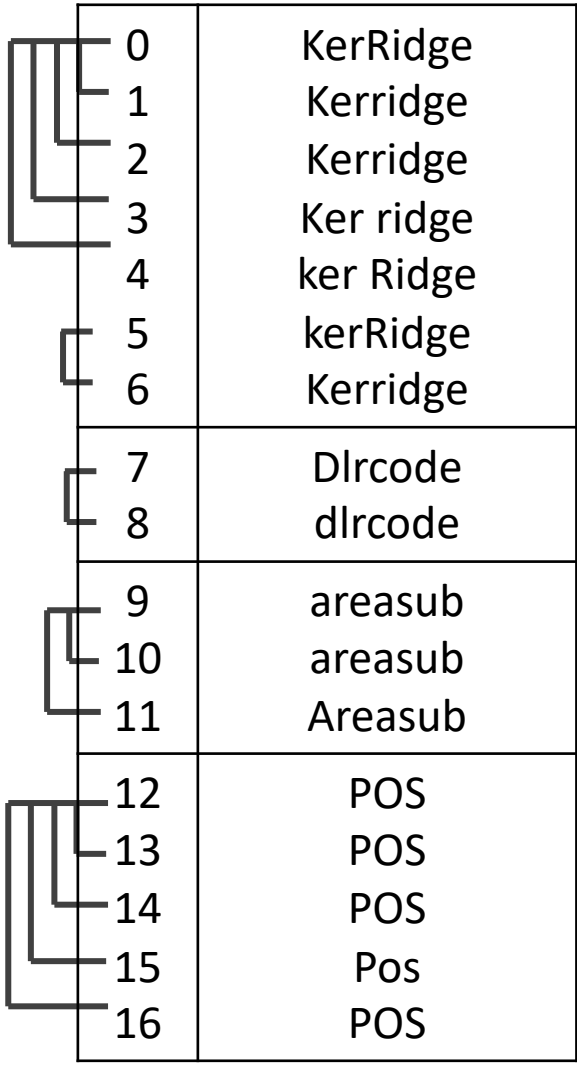
Cluster Size	Initial Clustering	Decomposition		Cluster Merge
		Type-based	Sim-based merge	
1	-	55	186	169
2	-	120	167	164
3	130	185	232	233
4	1690	1720	1604	1608

Cluster Size	Initial Clustering	Decomposition		Cluster Merge
		Type-based	Sim-based merge	
1	-	44	353	327
2	720	796	841	833
3	756	771	818	805
4	1071	1013	439	1085
5	574	570	423	436

Cluster sizes in integration phases for evaluation dataset – DS2

# Sample entities from datasets to express Entity Clustering

id	label	source	type
0	KerRidge	Drive	DMS
1	Kerridge	Rev8	DMS
2	Kerridge	Pinnacle	DMS
3	Ker ridge	Drive	DMS
4	ker Ridge	Pinnacle	-
5	kerRidge	Rev8	DMS
6	Kerridge	Rev8	DMS
7	Dlrcode	Drive	DLR
8	dlrcode	Rev8	
9	areasub	Pinnacle	Sub-category
10	areasub	Rev8	Sub-category
11	Areasub	Rev8	Sub-category
12	POS	-	TabCode
13	POS	Drive	TabCode
14	POS	Rev8	TabCode
15	Pos	Pinnacle	TabCode
16	POS	Drive	TabCode



# SplitMerge for Entity Resolution in Graphs

---

## Algorithm 5: SplitMerge Clustering Algorithm

---

**Input:** Set of entities  $e$  from  $n$  sources, edge set  $d$ , simFun  $f_{sim}$ , thresholds

$t_s, t_m$

**Output:** Set of clusters  $C$

$C \leftarrow \emptyset$

$e, l \leftarrow$

preprocessing( $e, l, f_{sim}$ ) /\*preprocessing\*/

$C_{init} \leftarrow$  computeConnectedComponents( $e, l$ )

$l_c \leftarrow$  computeLinksSim( $C_{init}, f_{sim}$ )

$C - int \leftarrow$  refineConnectedComponents( $C_{init}, L_c$ ) /\*initial clustering\*/

**foreach**  $c \in C_{int}$  **do** /\*cluster decomposition\*/

$C_{split} \leftarrow$  groupByType( $c, L_c$ )

$C_{split} \leftarrow$  simBasedRefinement( $C_{split}, L_c, t_s$ )

$C_{split} \leftarrow$  createRepresentatives( $C_{split}$ )

$C \leftarrow C \cup C_{split}$

$CM \leftarrow$  computeClusterSim ( $C, f_{sim}, t_m$ ) /\*create cluster mapping CM \*/

**while**  $CM \neq \emptyset$  **do**

$(c_1, c_2) \leftarrow$  CM.getBestMatch()

$c_m \leftarrow$  mergeClusters( $c_1, c_2$ )  $C \leftarrow C - c_1, c_2 \cup c_m$  /\*cluster merge\*/

$CM \leftarrow$  adaptMapping( $CM, C, c_m, c_1, c_2, f_{sim}, t_m$ )

**return**  $C$

---

SplitMerge Phases:

**Note:**  $f_{sim}$  = *Similarity function*;  $t_s, t_m$  = *similarity threshold*

### 1. Preprocessing

- property values required for similarity computation are normalized

### 2. Initial Clustering

- Connected components
- In order to phase out deduplicated entities and the refineConnectedComponents (Line 5) connected components is used on one entity per source.

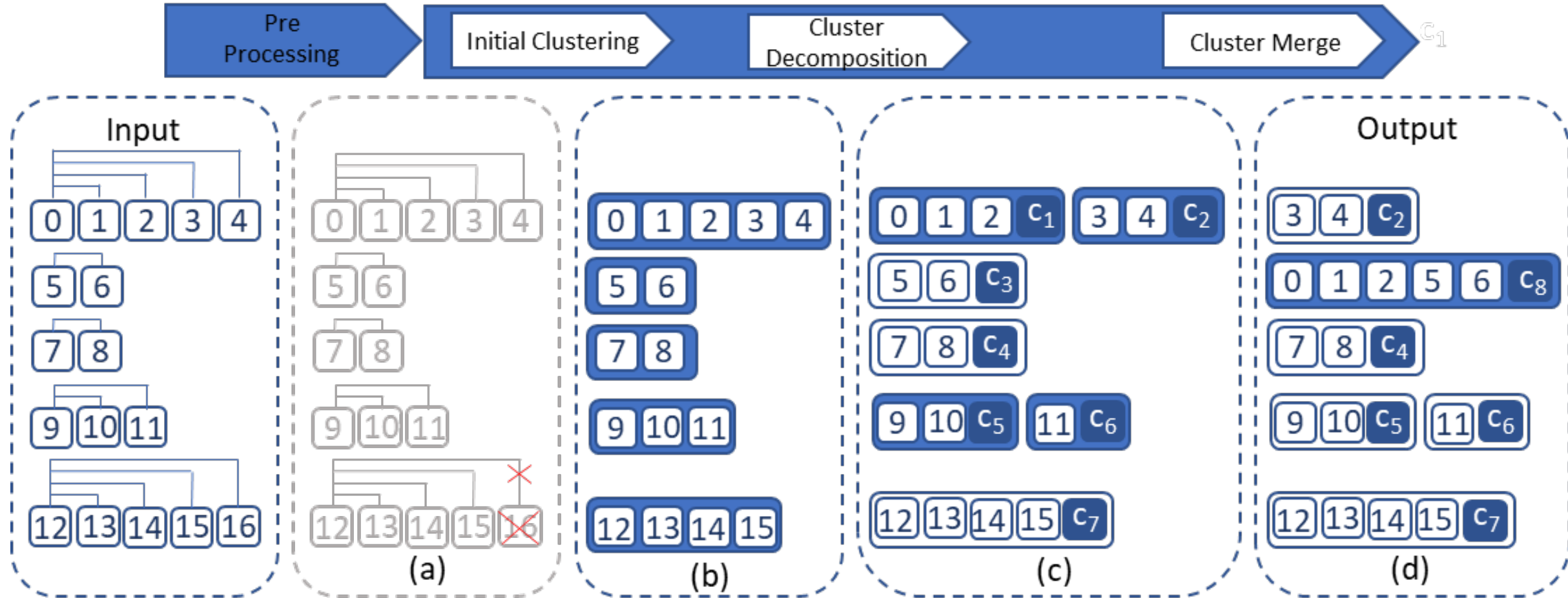
### 3. Cluster Decomposition – two main approaches

- Split clusters based on inconsistent semantic types
- clusters containing inadequate similarity to other cluster members are split up

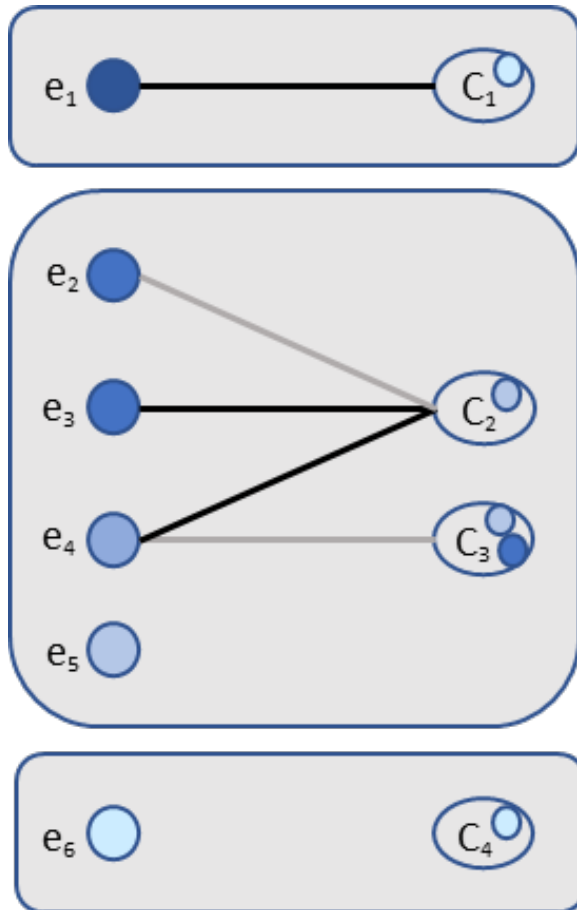
### 4. Cluster Merge

- merge clusters that range below the maximally possible cluster size  $k$

# SplitMerge Example continued...



# Incremental Clustering Approach



(a) Root Approach

---

**Algorithm 8:** Set-based incremental entity clustering - root approach

---

**Input:** Existing clusters  $C_{exist}$ , edge set  $l_c$ , split threshold  $t_s$

**Output:** Set of clusters  $C_{result}$

$C_{new} \leftarrow \text{createInitialClusters}(E_{new}, f_{blocking})$

$C_{exist} \leftarrow \text{addBlockingInfo}(C_{exist}, f_{blocking})$

**for** block  $i$  **in Parallel** **do**

$L_i \leftarrow \text{getClusterCandidates}(C_{exist}, C_{new}, f_{sim}, t_{min})$

$L_{sorted} \leftarrow \text{sortLinkSim}(L_i)$

**foreach**  $(c_{new}, c_{exist}, sim) \in L_{sorted}$  **do**

**if**  $c_{exist} \notin c_{new}$  **then**

continue()

**if** isSrcConsistent( $c_{new}, c_{exist}$ ) **then**

$c_{exist}.add(c_{new})$

$C_{new}.remove(c_{new})$

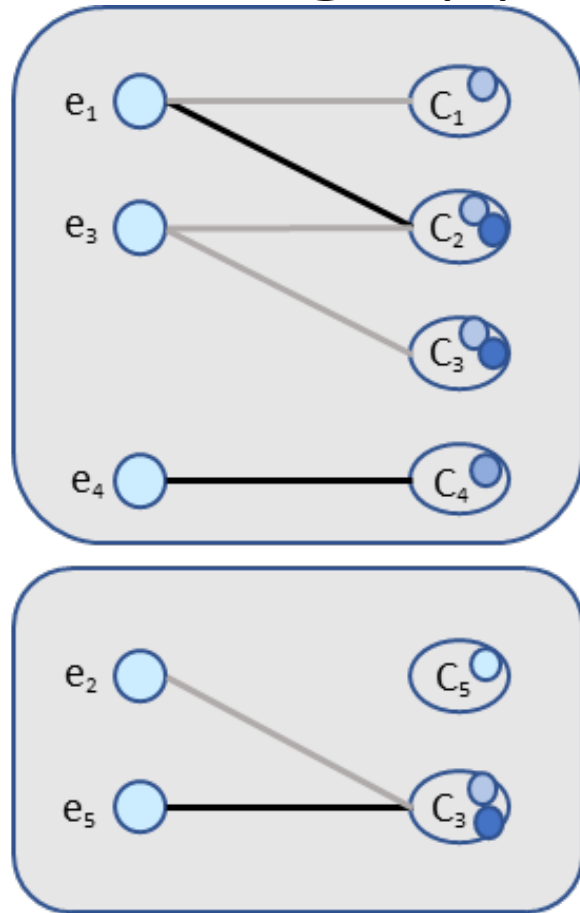
**return**  $C_{exist} \cup C_{new}$

---

- Two different scenarios for cluster generation:
  - root approach and
  - source-specific approach
- (colours signify various data sources, it is assumed that all the links exceed the minimal similarity threshold)



# Incremental Clustering Approach



(b) Source-specific Approach

---

## Algorithm 5: SplitMerge Clustering Algorithm

---

**Input:** Set of entities  $e$  from  $n$  sources, edge set  $d$ , simFun  $f_{sim}$ , thresholds  $t_s, t_m$

**Output:** Set of clusters  $C$

$C \leftarrow \emptyset$

$e, l \leftarrow$

preprocessing( $e, l, f_{sim}$ ) /\*initial clustering\*/

$C_{init} \leftarrow \text{computeConnectedComponents}(e, l)$

$l_c \leftarrow \text{computeLinksSim}(C_{init}, f_{sim})$

$C - int \leftarrow \text{refineConnectedComponents}(C_{init}, L_c)$  /\*initial clustering\*/

**foreach**  $c \in C_{int}$  **do** /\*cluster decomposition\*/

$C_{split} \leftarrow \text{groupByType}(c, L_c)$

$C_{split} \leftarrow \text{simBasedRefinement}(C_{split}, L_c, t_s)$

$C_{split} \leftarrow \text{createRepresentatives}(C_{split})$

$C \leftarrow C \cup C_{split}$

$CM \leftarrow \text{computeClusterSim}(C, f_{sim}, t_m)$  /\*create cluster mapping CM\*/

**while**  $CM \neq \emptyset$  **do**

$(c_1, c_2) \leftarrow CM.\text{getBestMatch}()$

$c_m \leftarrow \text{mergeClusters}(c_1, c_2)$   $C \leftarrow C - c_1, c_2 \cup c_m$  /\*cluster merge\*/

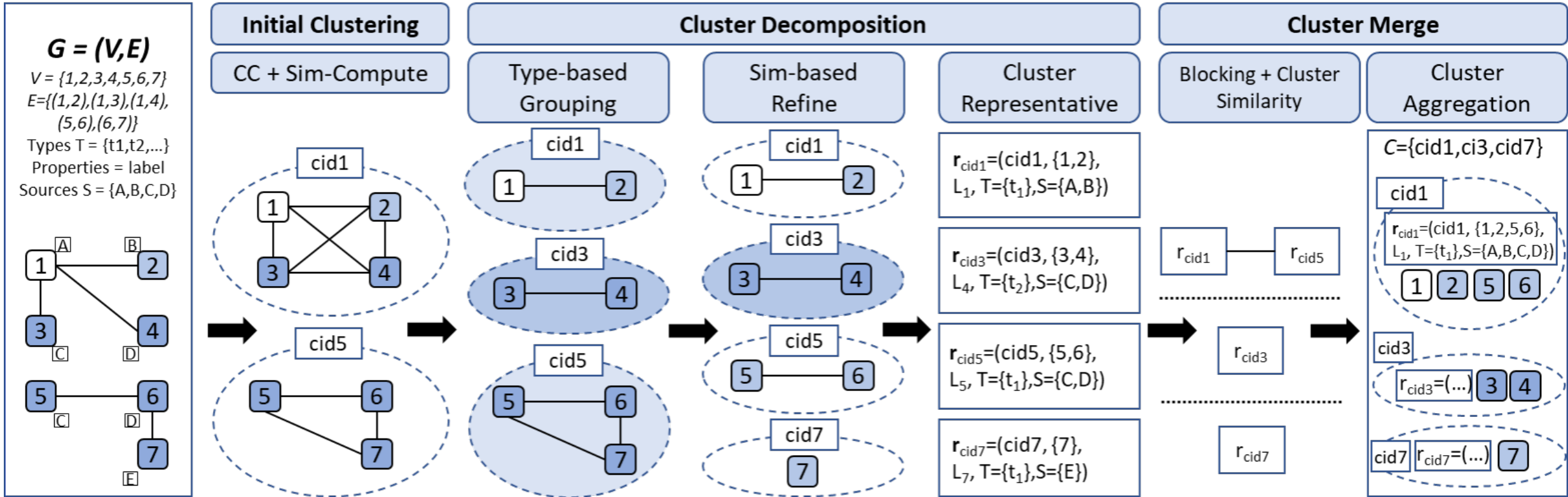
$CM \leftarrow \text{adaptMapping}(CM, C, c_m, c_1, c_2, f_{sim}, t_m)$

**return**  $C$

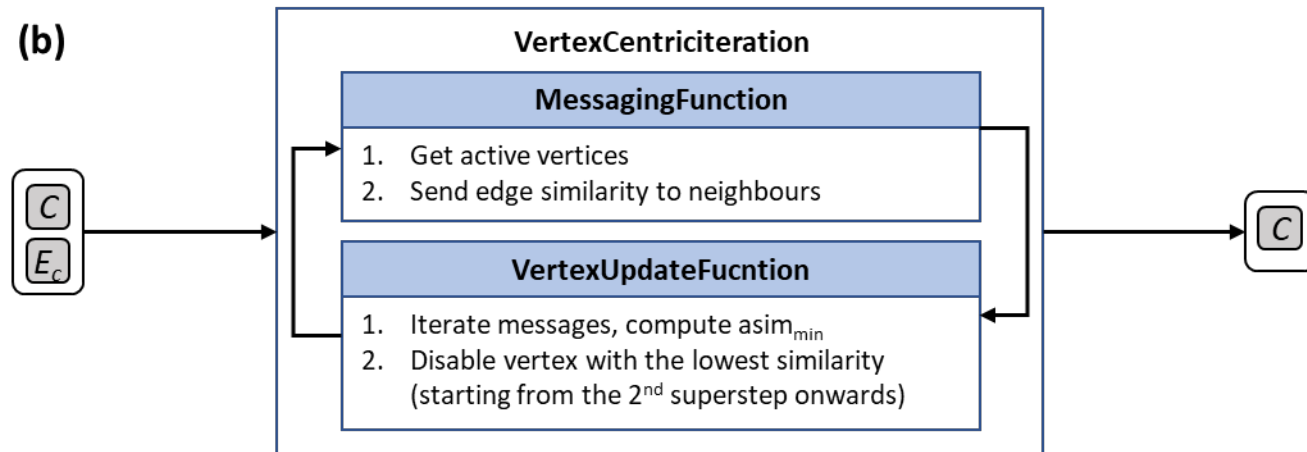
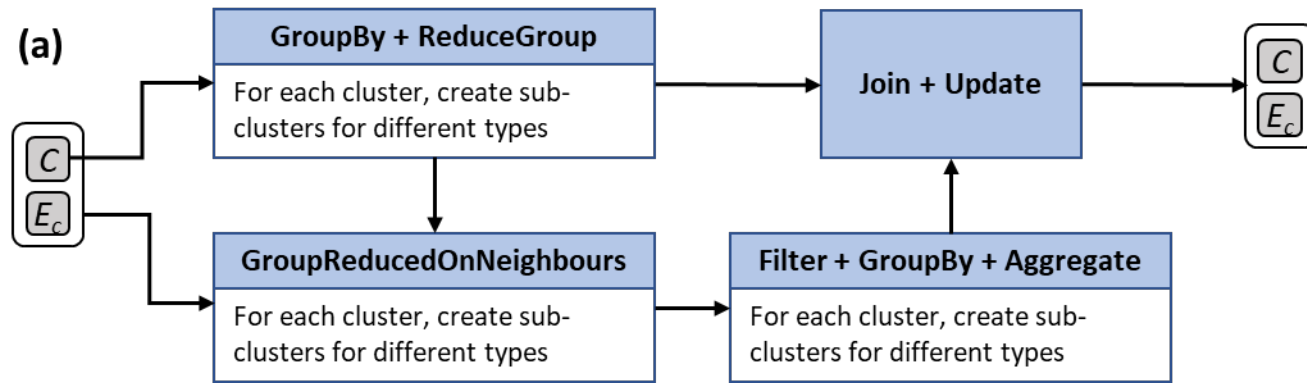
---

The algorithm resolves source-consistent candidate links between the newer entities and existing set of clusters in parallel with partitioned blocks.

# Distributed Clustering

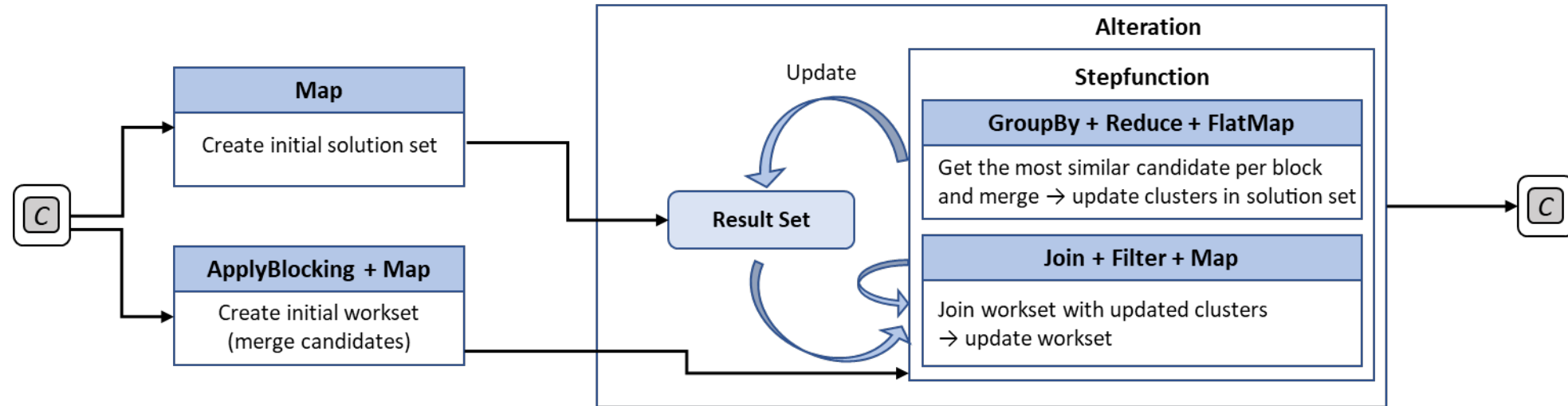


# Cluster Decomposition - Distributed ER



- Sub-workflows for type-based grouping shown in (a) and similarity-based refinement in (b)
- The first step in decomposition of the cluster is to break the clusters into sub-clusters based on the compatibility of property types (Type-based grouping).
- The second step is to decompose these clusters using non-similar entities from clusters based on the step Similarity-based refinement.

# Cluster Merge – Distributed ER



- Spark's iterative operative in addition to user-defined functions are used to address the final merge stage.
- Clusters with high similarity which as usually small are aggregated iteratively into large clusters.
- The creation of representatives for each of the cluster enables to reduce the number of potential entities for the merge step.

# Evaluation

- Evaluation of these data-sets are based on dynamic graph queries. For each data-set we grouped the queries in five sets (i.e. ten queries per set): each set is homogeneous with respect to its complexity of the queries (e.g. number of connected components, number of results and so on.).
- For instance, referring to integrated Rev8 data-sets, the first set of queries searches information about services while the second set of queries seeks information about sales.
- For each set, we ran the queries ten times and measured the average response time.

# Cluster Sizes and Configuration Results

	Data-set	Node Properties	No. of Nodes	No. of Sources
DS1-A1	Drive	DOC	5,079	4
DS2-C1	Rev8	POS	11,600	5
DS2-C2		DlrCode	42,949	5
DS3-N1	Pinnacle	PinCode	131,072	5
DS3-N2		CustRef	500,000	10

- Evaluation data-set details

	Perfect Result		Best configuration -results		
	# of clusters	# of edges	conf(t_min,bk)	# of correct edges	F-measure
DS1-A1	790	6497	conf(0.4,1)	6,207	0.981
DS2-C1	5000	10,340	conf(0.5,1)	9,589	0.953
DS2-C2	20000	39,321	conf(0.7,1)	36,956	0.846
DS3-N1	110,440	101,843	conf(0.7,6)	100,057	0.804
DS3-N2	350,960	619,528	conf(0.7,6)	513,975	0.795

# Evaluation of Static vs Dynamic Clustering

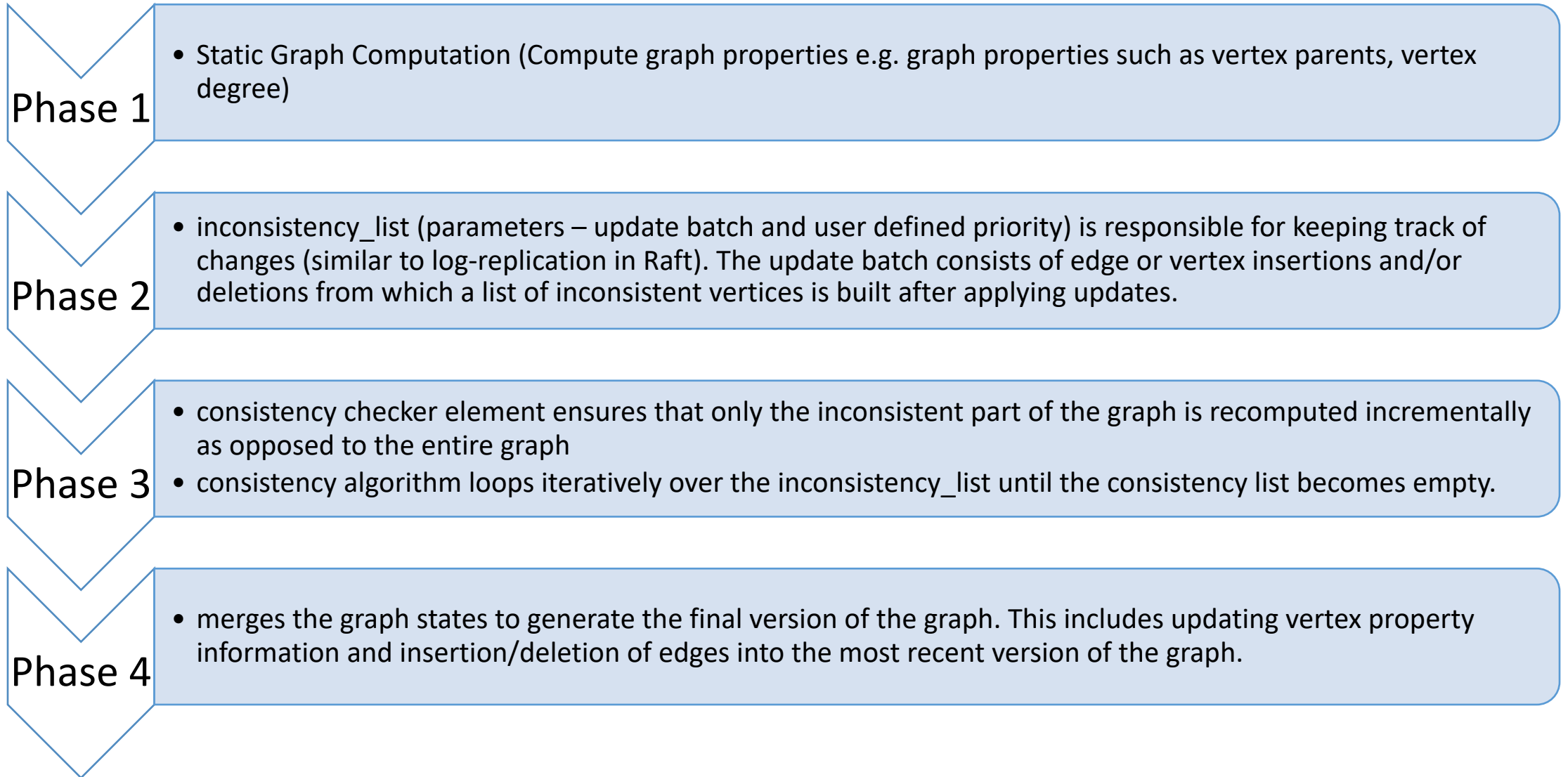
DS2 - C2	Incremental		Static	
	Root	Source	CLIP	SplitMerge
run time (sec)	4210	1052	1859 + 72	1859 + 732
Precision	0.765	0.897	0.868	0.848
recall	0.865	0.839	0.819	0.833
F-measure	0.812	0.879	0.855	0.845

• Data Quality and Run time for:

- DS2-C2
- DS3-N1

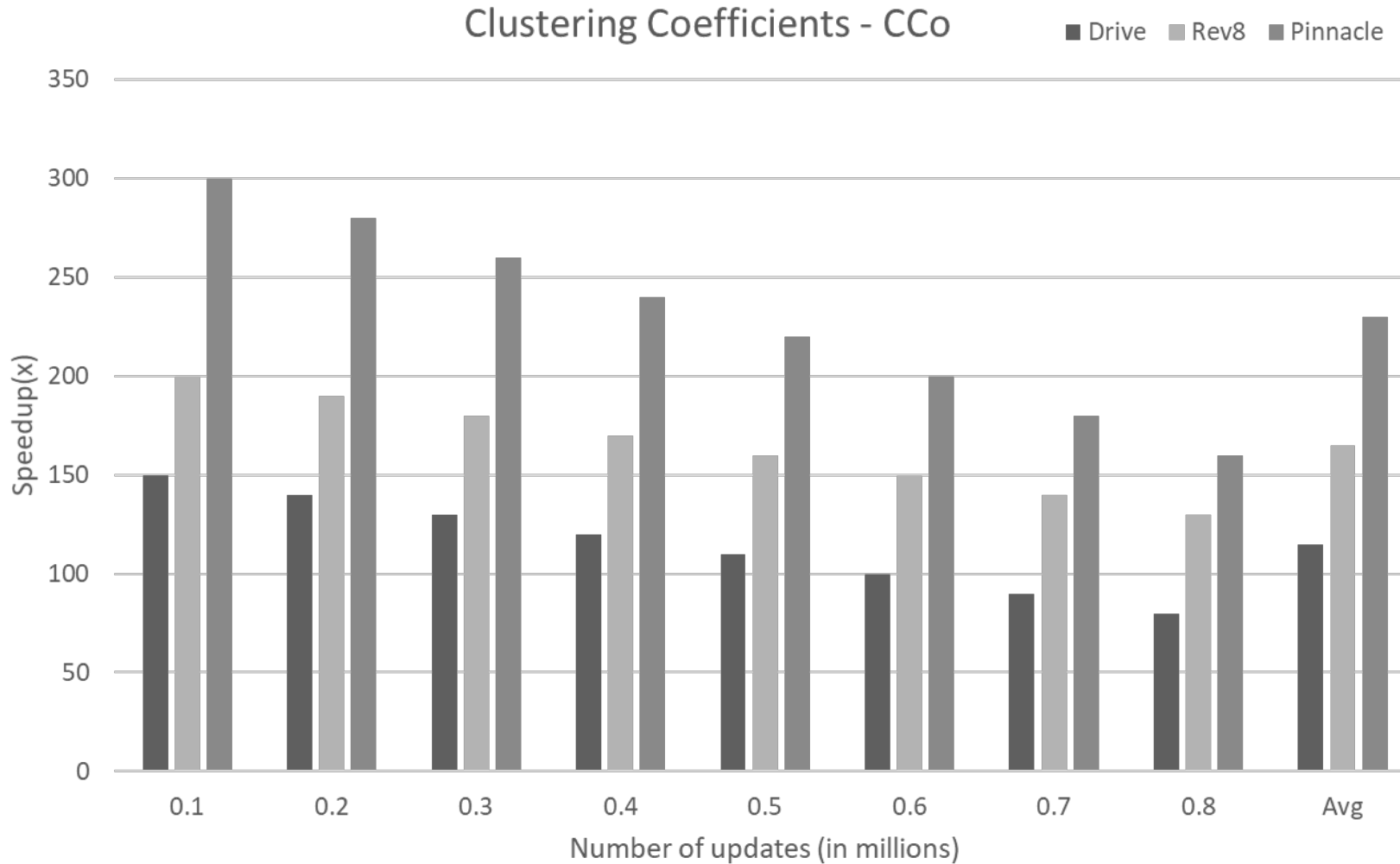
DS3 – N1	Incremental		Static	
	Root	Source	CLIP	SplitMerge
run time (sec)	642	221	110+105	110+763
Precision	0.565	0.817	0.860	0.789
recall	0.844	0.821	0.819	0.862
F-measure	0.676	0.819	0.846	0.832

# Consistency Phases



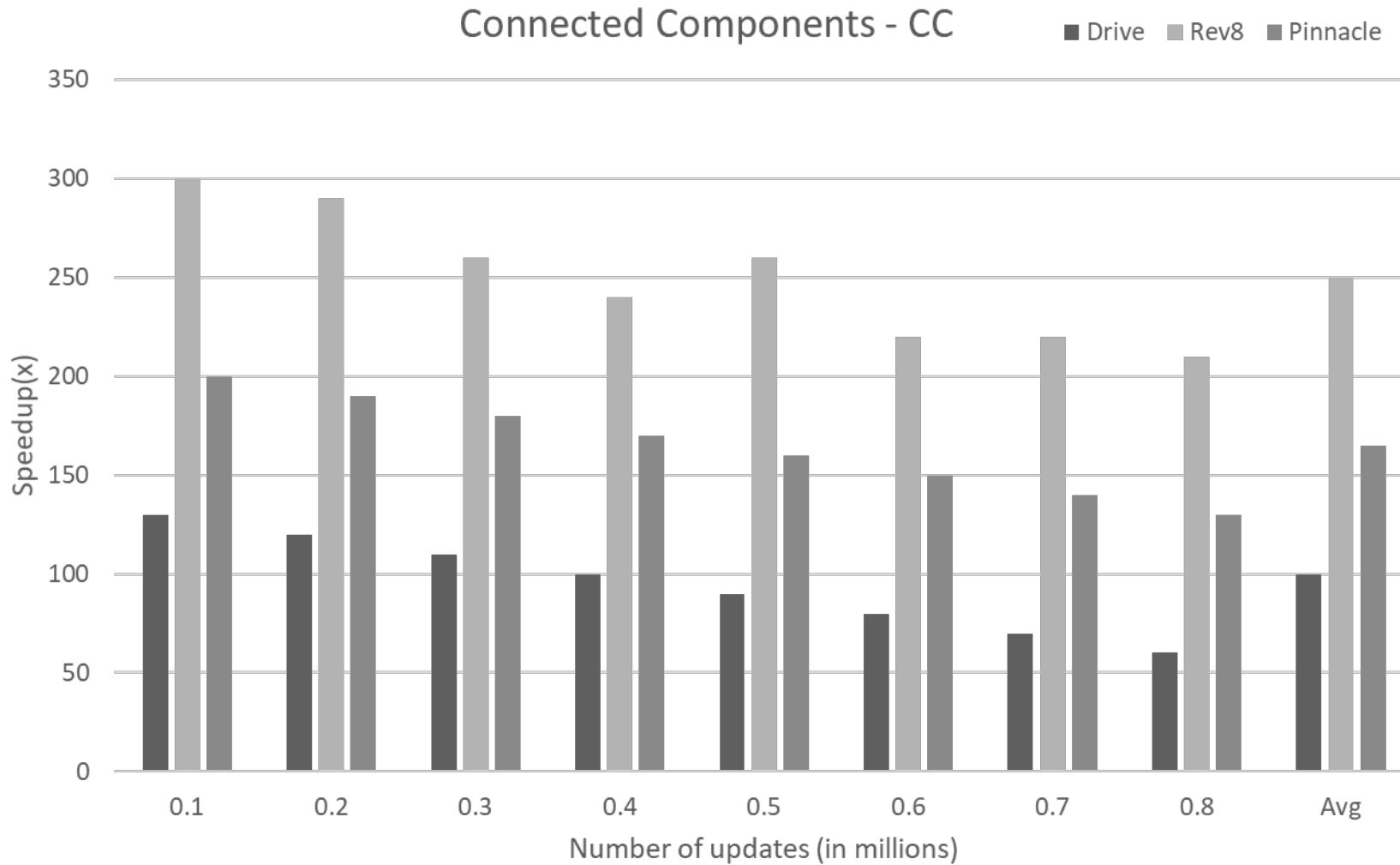


# Evaluation – static vs incremental computation (CCo)



Incremental Speedup over static execution verses the update batch size for Clustering Coefficient

# Evaluation – static vs incremental computation (CC)



- Incremental Speedup over static execution verses the update batch size for Clustering Coefficient

# Conclusion

- Entity Resolution techniques combined with graphs result in quicker and scalable data integration
- Test the efficacy of the approach on other domains – currently tested approach on clinical datasets (limitations due to limited public data availability)
- Further enhance the solution to provide performance and scalability guarantees
- Employ ML and AI techniques to automate the report generation process within a cloud-based environment

Thank you! Questions?

 b.arshad@derby.ac.uk

 @bilalarshadali