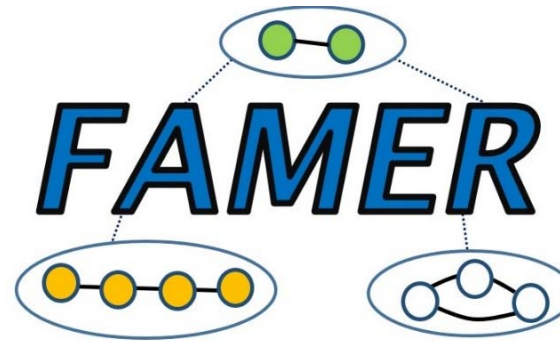




UNIVERSITÄT
LEIPZIG

ScaDS 
DRESDEN LEIPZIG



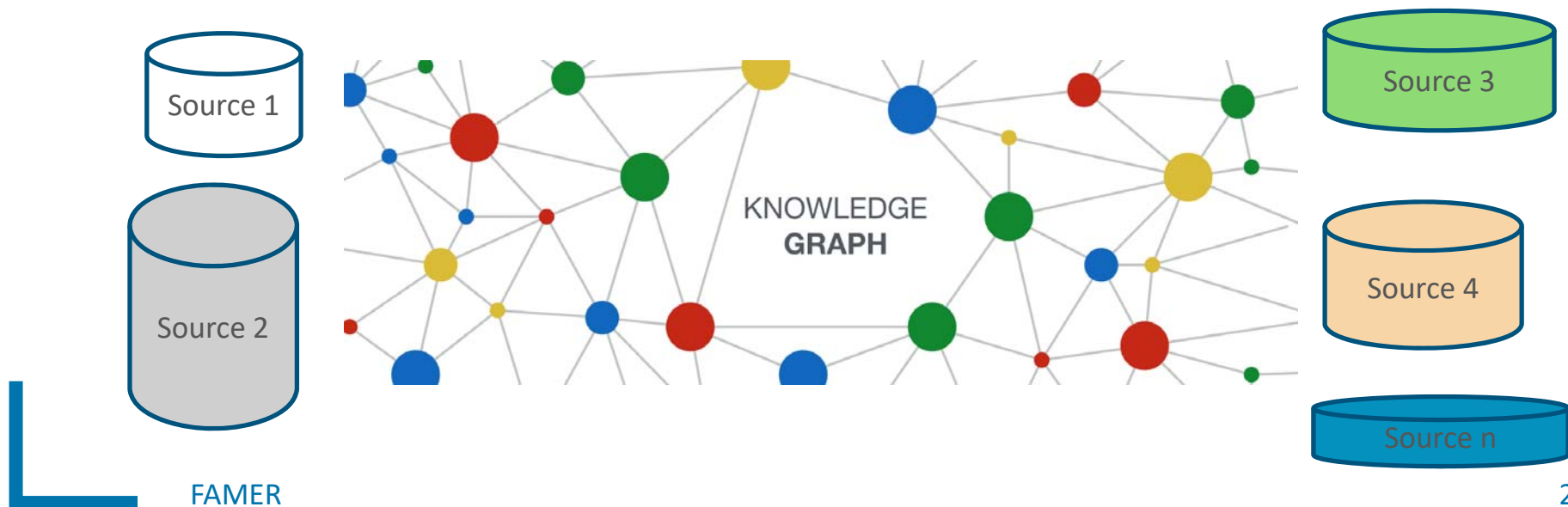
FAst Multi-source Entity Resolution System

Alieh Saeedi, Eric Peukert, Erhard Rahm

www.scads.de

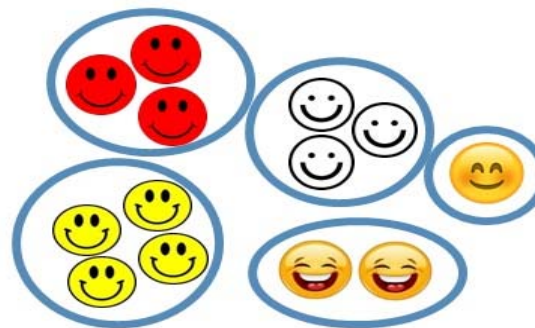
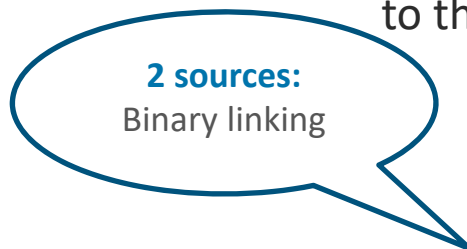


- Physical data integration
 - Knowledge graph: Store data from multiple sources in a graph-like structure





- Automatic construction & maintenance of KG: **data quality**
- Challenges for data quality
 - **Entity Resolution**: The task of **identifying** and **linking** entities that refer to the **same** real-world entity





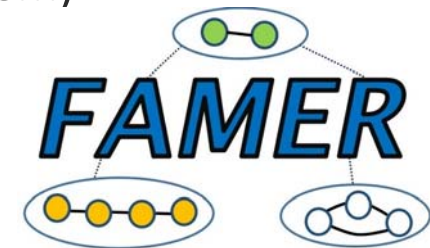
- **FAMER (FASt Multi-source Entity Resolution system)**

- Scalable ER approaches for big data

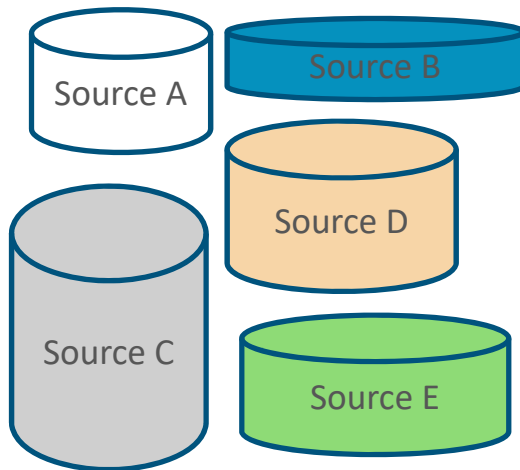
- Multiple data sources
- Large volumes of data

- Built on top of the distributed data flow framework Apache Flink and Gradoop

- High scalability
- Large amounts of data
- Many machines

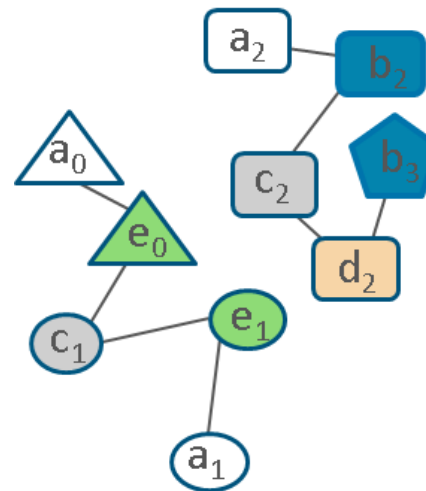


Input

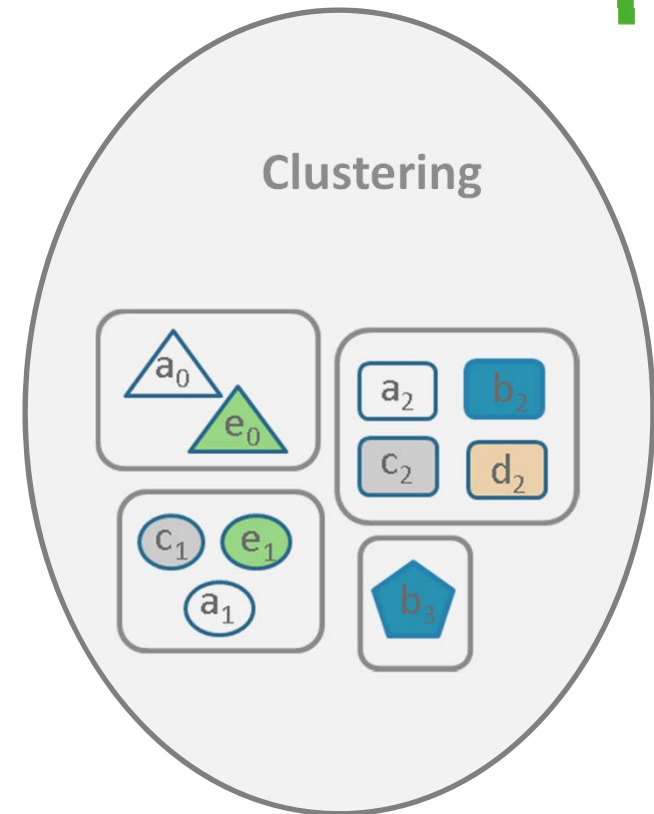


FAMER

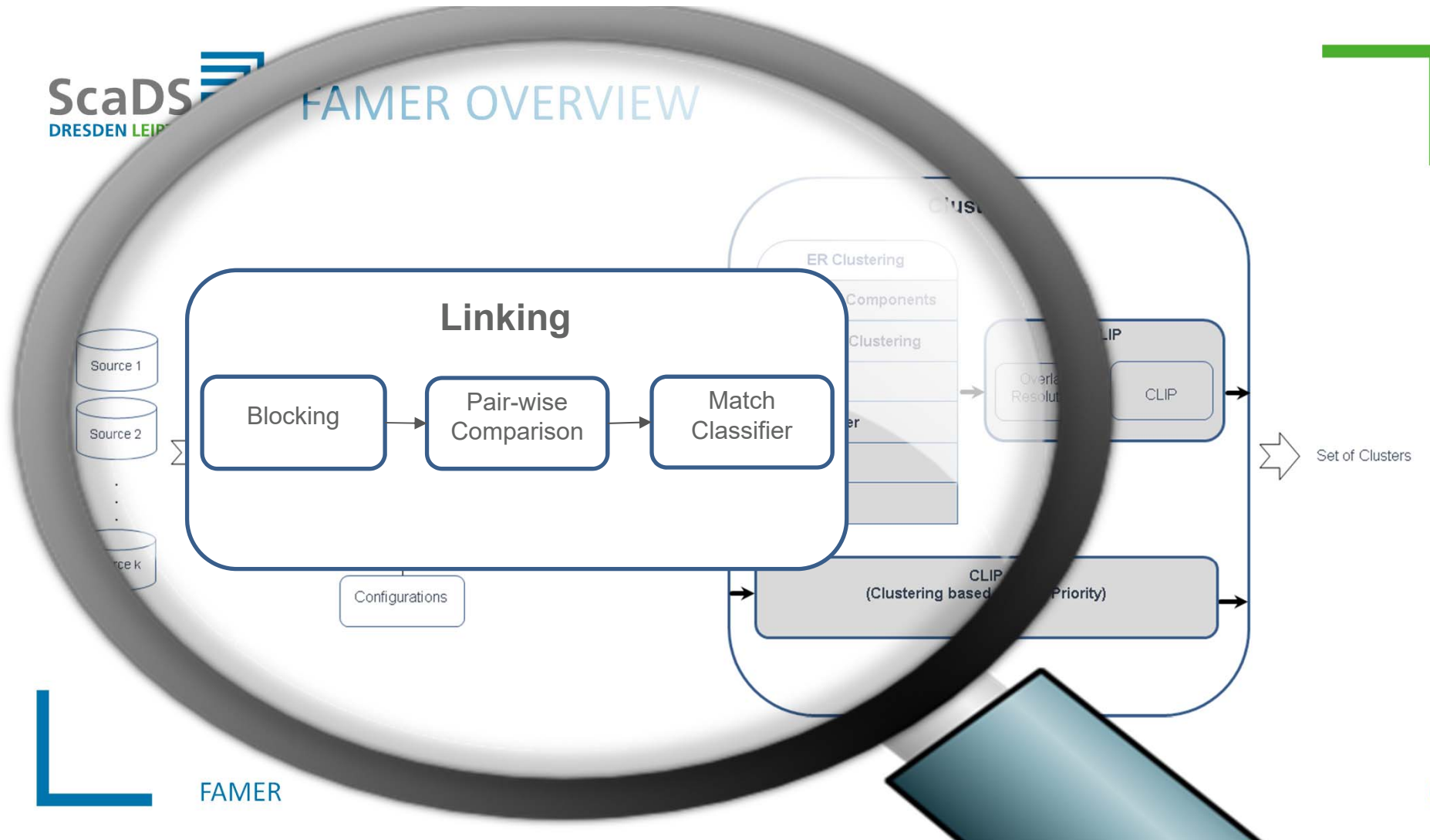
Linking: Similarity Graph



Clustering



FAMER OVERVIEW





Id	Name	Surname	Suburb	Post code	SourceId
a ₀	geOrge	Walker	winston salem	271o6	Src A
b ₀	George	Alker	winstom salem	27106	Src B
c ₀	George	Walker	Winstons	27106	Src C
d ₀	Geoahge	Waker	Winston	271oo	Src D
a ₁	Bernie	Davis	pink hill	28572	Src A
b ₁	Bernie	Daviis	Pinkeba	2787z	Src B
c ₁	Bernii	Davs	pink hill	28571	Src C
a ₂	Bertha	Summercille	Charlotte	28282	Src A
b ₂	Bertha	Summeahville	Charlotte	2822	Src B
d ₂	Brtha	Summerville	Charlotte	28222	Src D
c ₃	Bereni	dan'lel	Pinkeba	27840	Src C
d ₃	Bereni	Dasniel	Pinkeba	2788o	Src D

Id key

a₀ wa
c₀ wa
d₀ wa

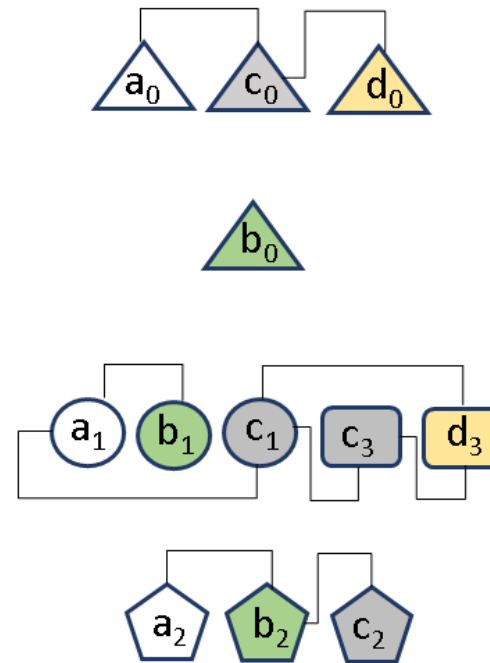
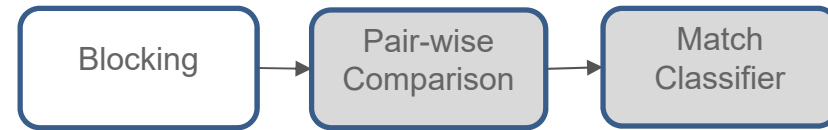
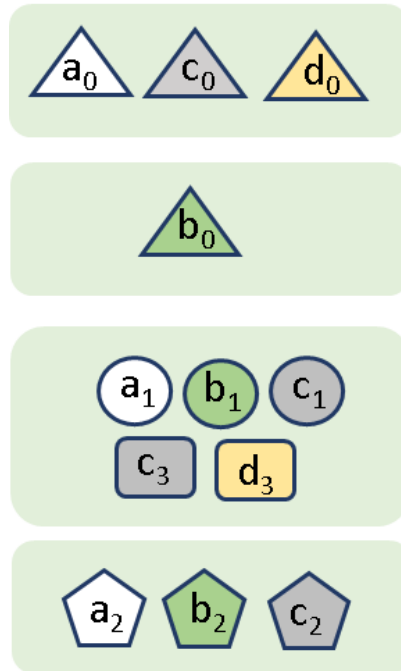
b₀ al

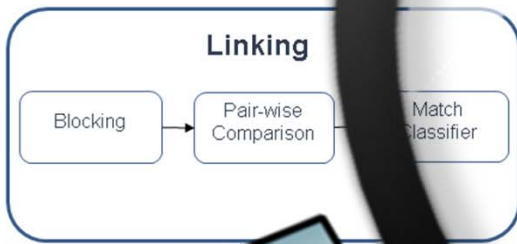
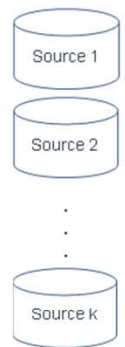
a₁ da
b₁ da
c₁ da
c₃ da
d₃ da

a₂ su
b₂ su
d₂ su

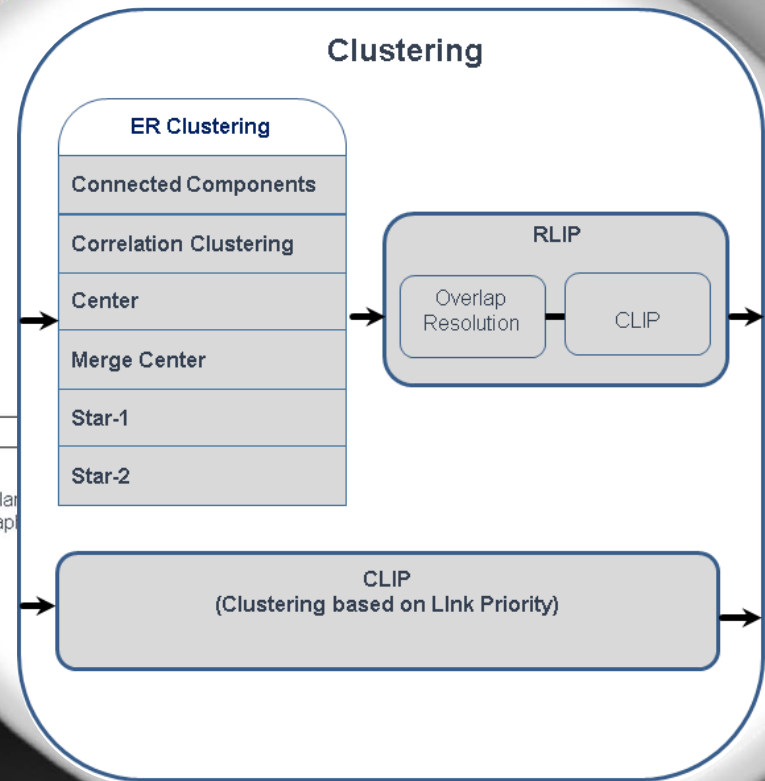
FAMER LINKING

Blocks

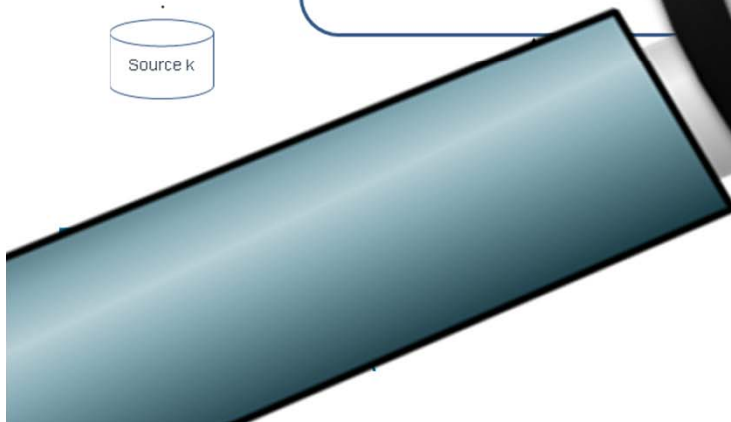




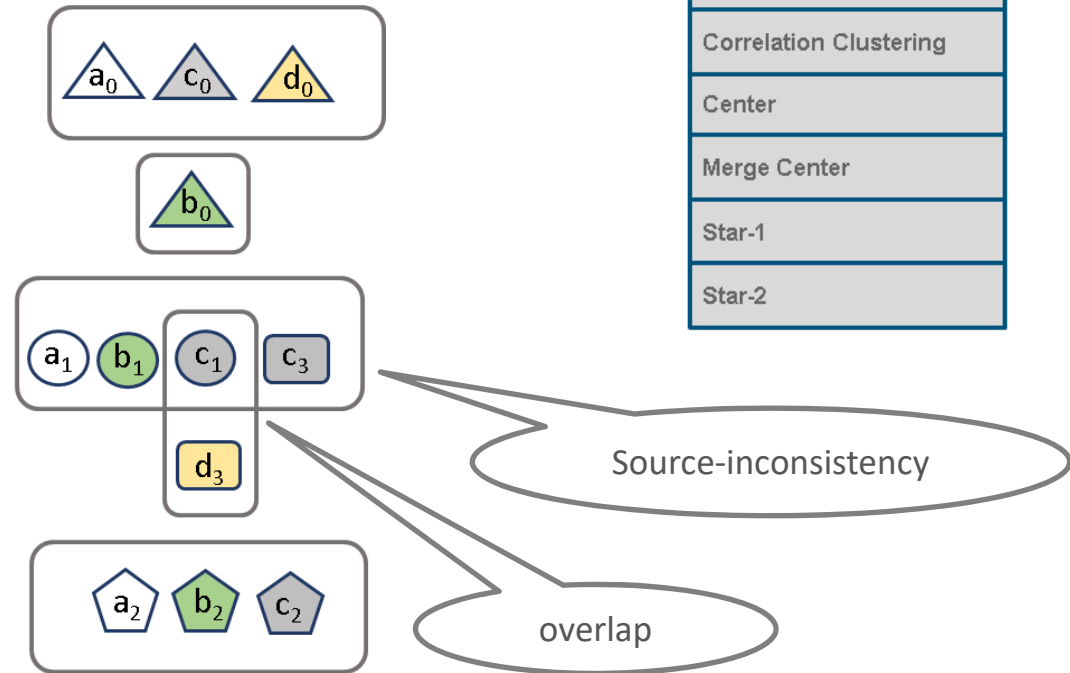
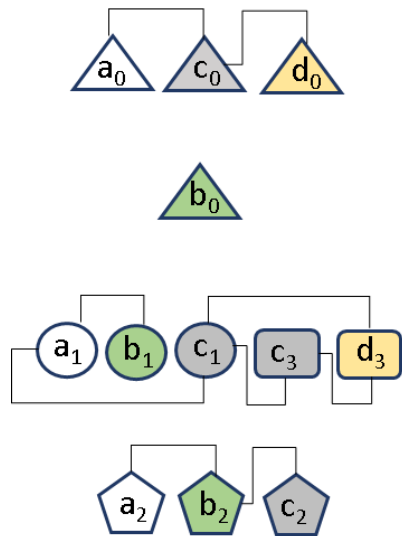
Similar Graphs



Set of Clusters

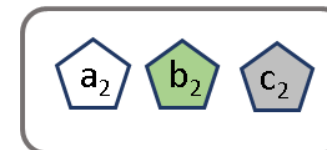
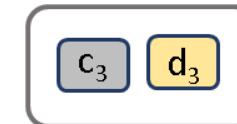
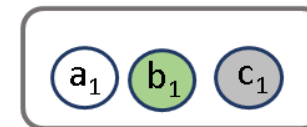
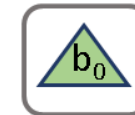
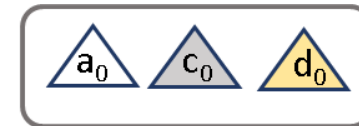
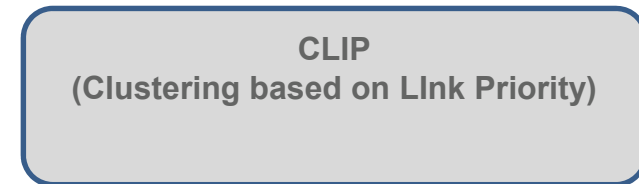


Similarity graph



ER Clustering
Connected Components
Correlation Clustering
Center
Merge Center
Star-1
Star-2

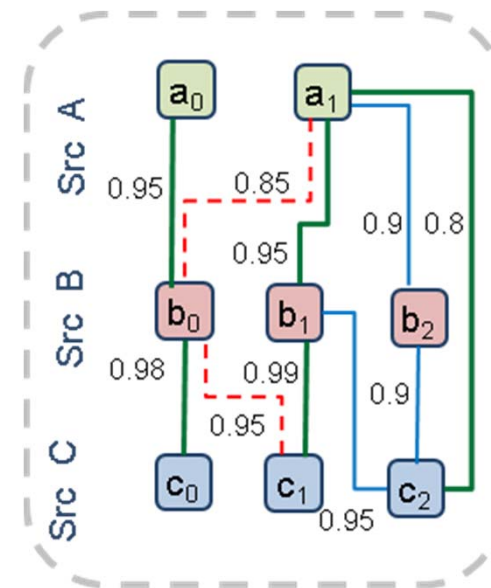
- Prioritize links based on
 - Link strength
 - Strong, Normal, Weak
 - Link degree
 - Similarity value
- produces
 - Source-consistent clusters
 - No overlap

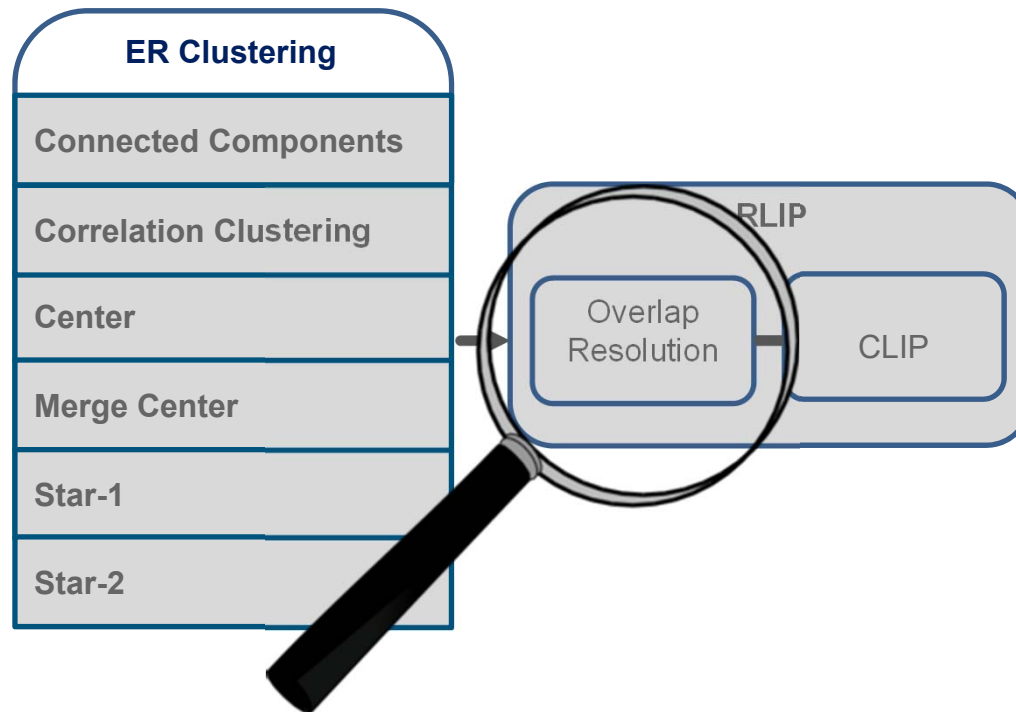




– Link Strength

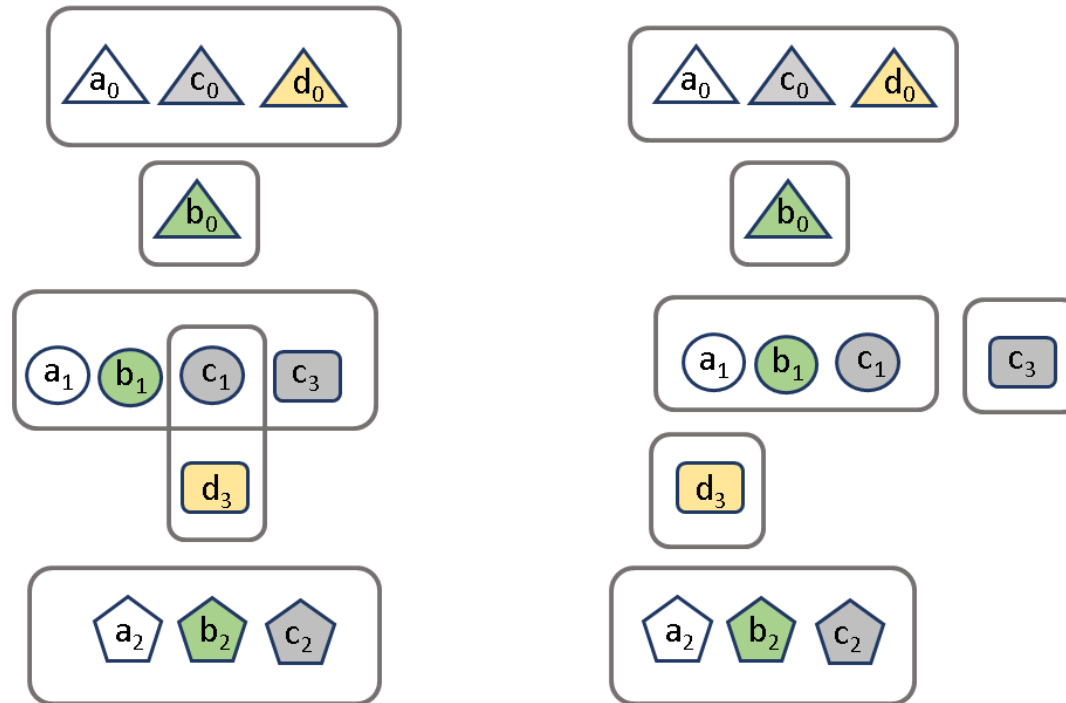
- Strong
- Normal
- Weak





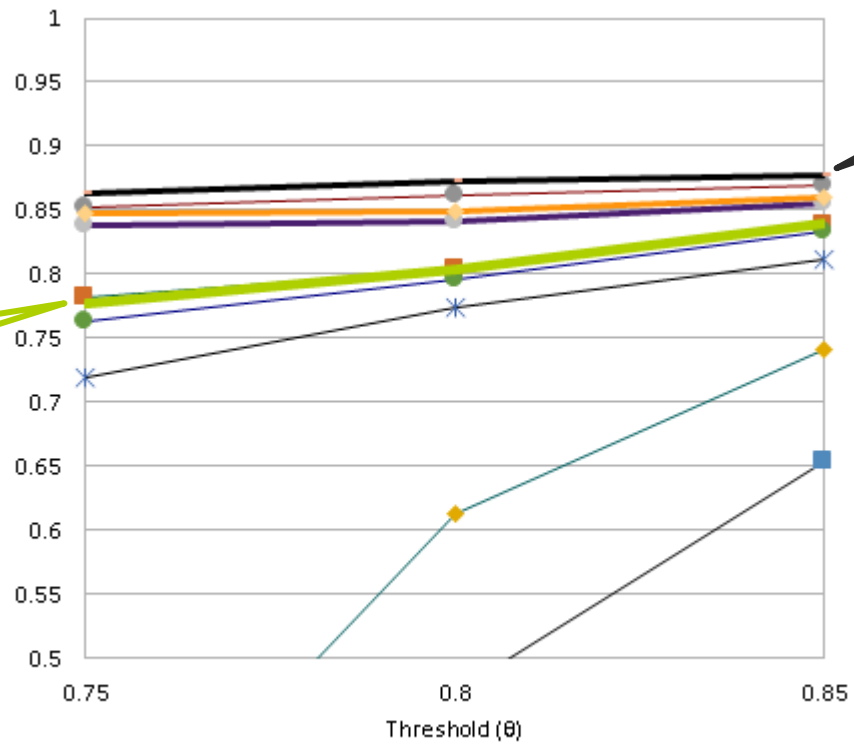


ER Clustering output



- Geographical domain
- 4 sources
- F-Measure

Similarity Graph

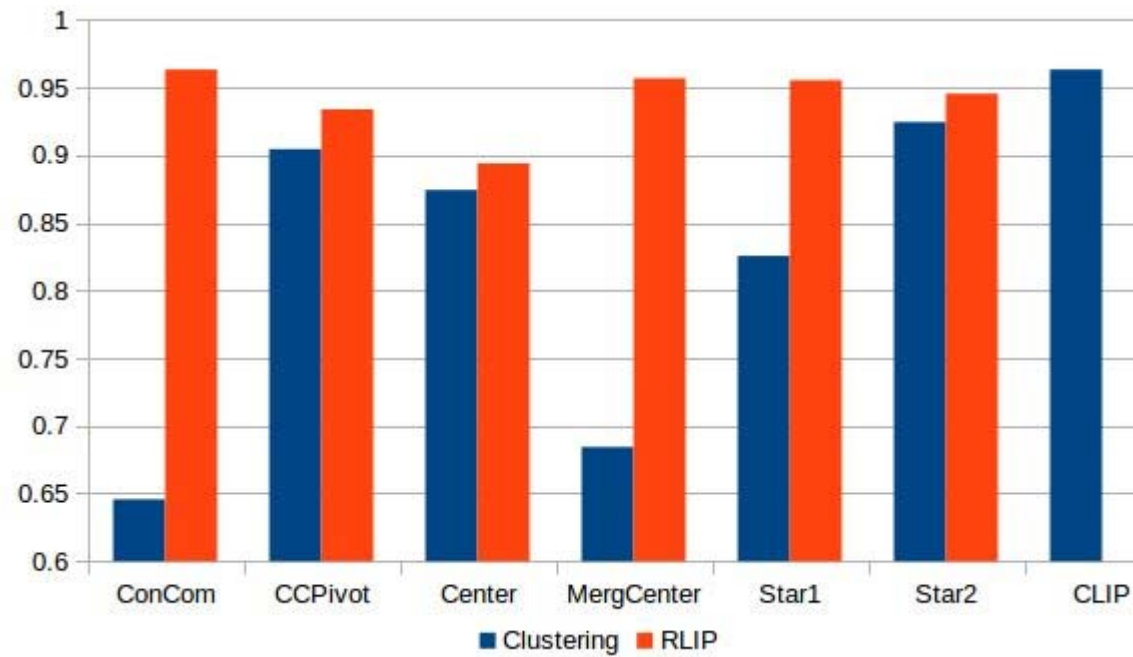


CLIP





■ Average F-Measure



- 5 sources
- 5,000,000 entities
- Flink cluster of 16 workers
- CLIP
 - 69s (1 min) for 16 workers
 - 10 sources/
10 million entities: 228 s

	#workers		
clustering	4	8	16
ConCom	51	57	55
CCPiv	1530	10008	688
Center	390	208	117
Merge Center	640	349	194
Star-1	288	149	85
Star-2	214	124	67
CLIP	190	101	69



4th International Summer School for Big Data and Machine Learning

Jun. 30th - Jul. 6th

- **Hackathon:** June 30th- July 1st 2018
- **Summer School:** July 2nd- July 6th 2018
- **Location:** University of Leipzig
- **Bonus:** LSWT attendants obtain early registration rates

- Topics & Highlights
 - Deep Learning
 - Data Mining in Stream Data
 - Machine Learning for the Sciences
 - Apache SystemML
 - Secure Cloud Databases
 - Visual Analytics for Big Data
 - Graph Analytics on Spark
 - Web-scale information extraction
 - Big Data Integration
 - Big Data and HPC

