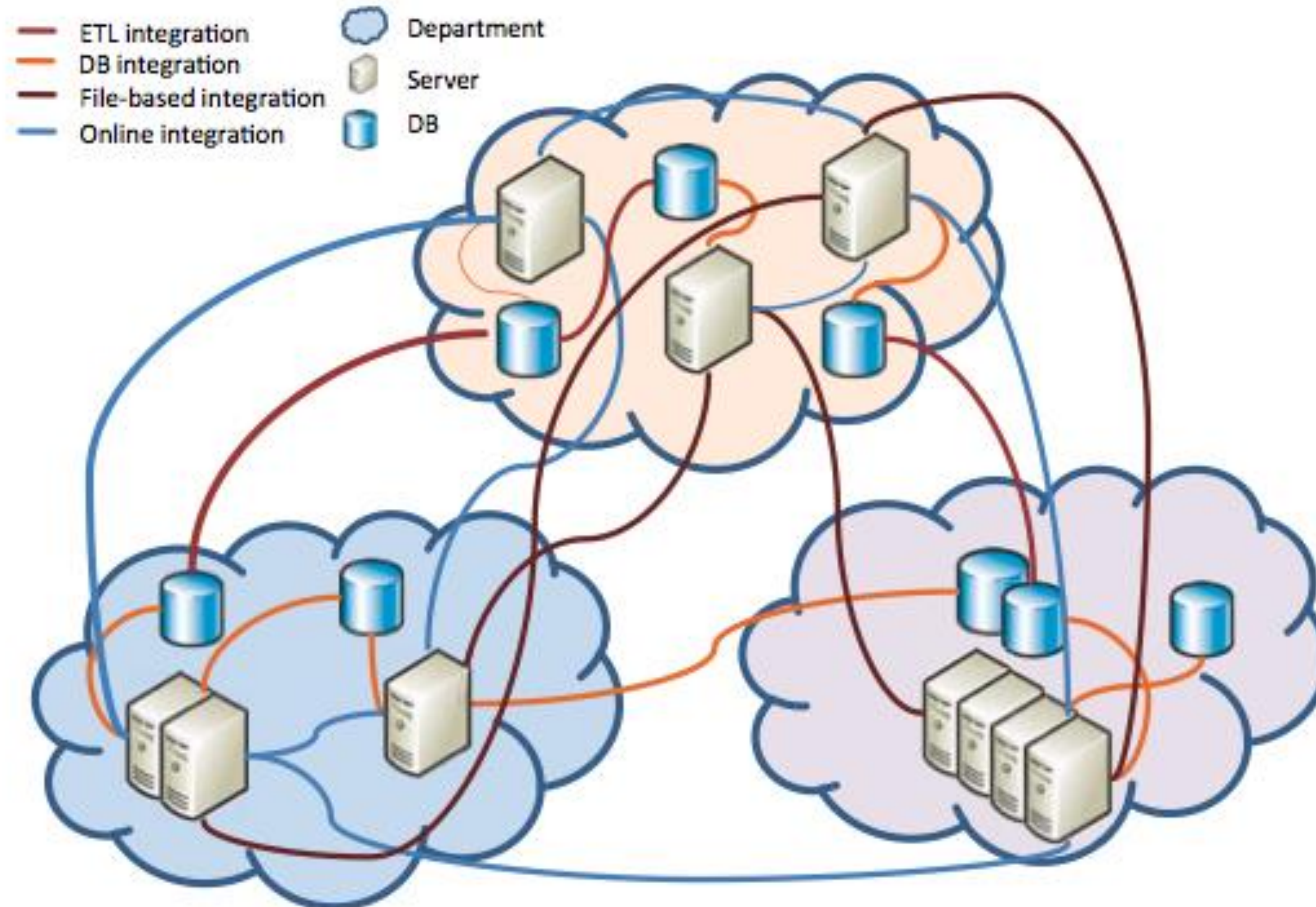


# Deduplication with graph algorithms

Tom Michiels

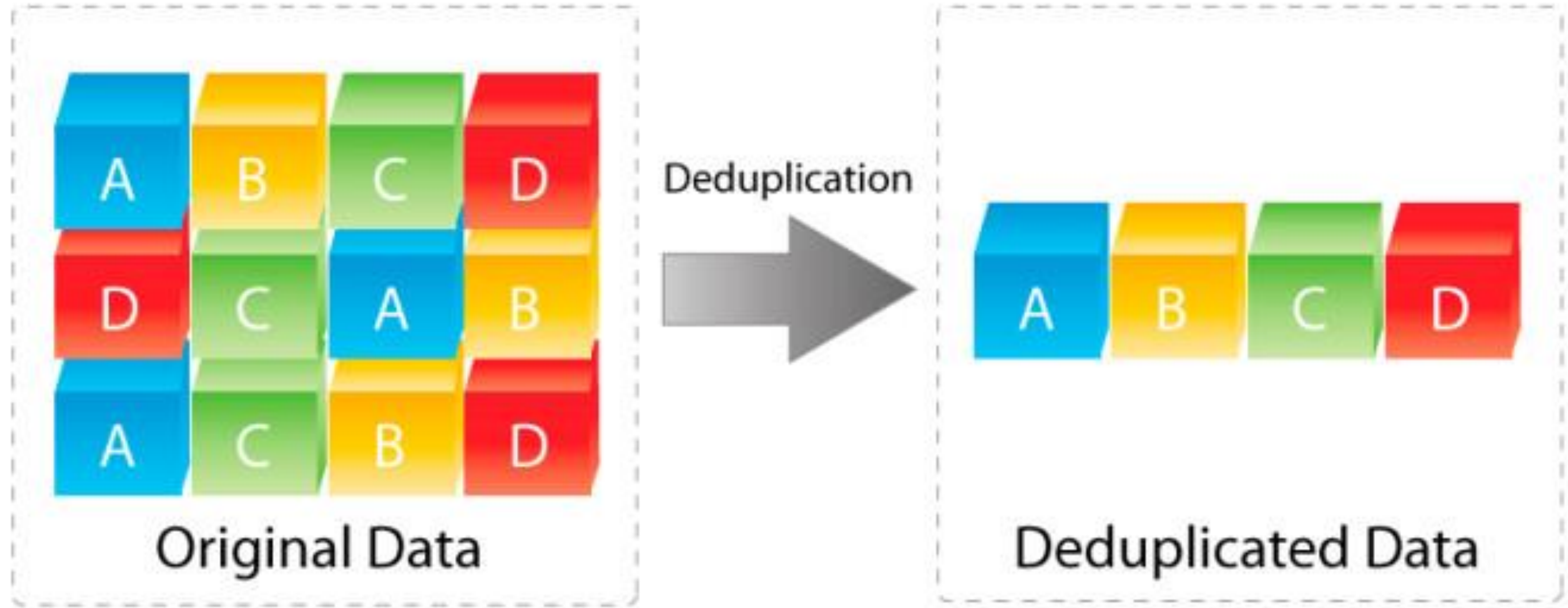
*Hoe kunnen we datascience en machine learning technieken (AI, ChatGPT...) inzetten voor de verbetering van bestaande en de productie van nieuwe openbare statistieken*

# Data Quality

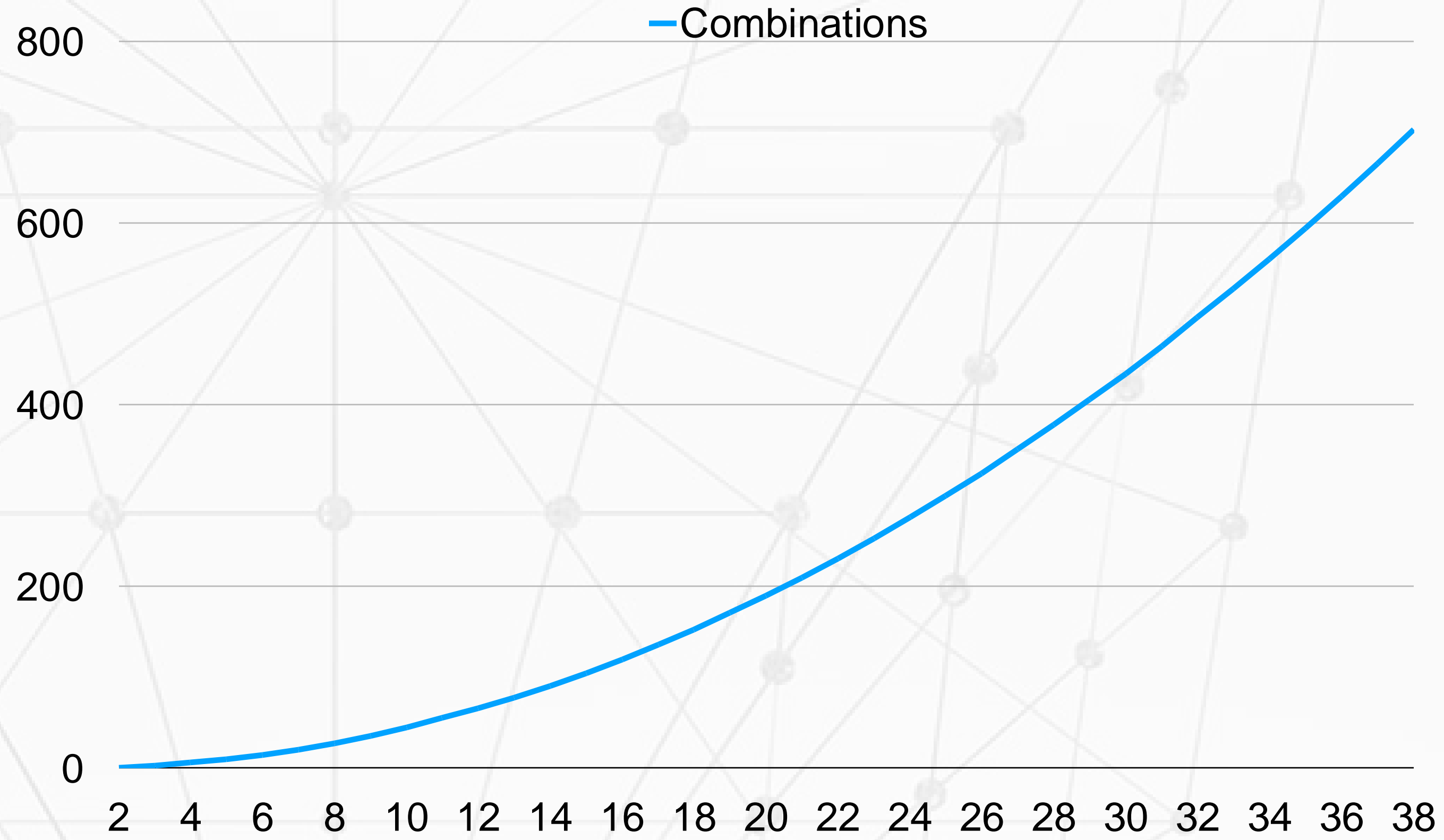




# Deduplication

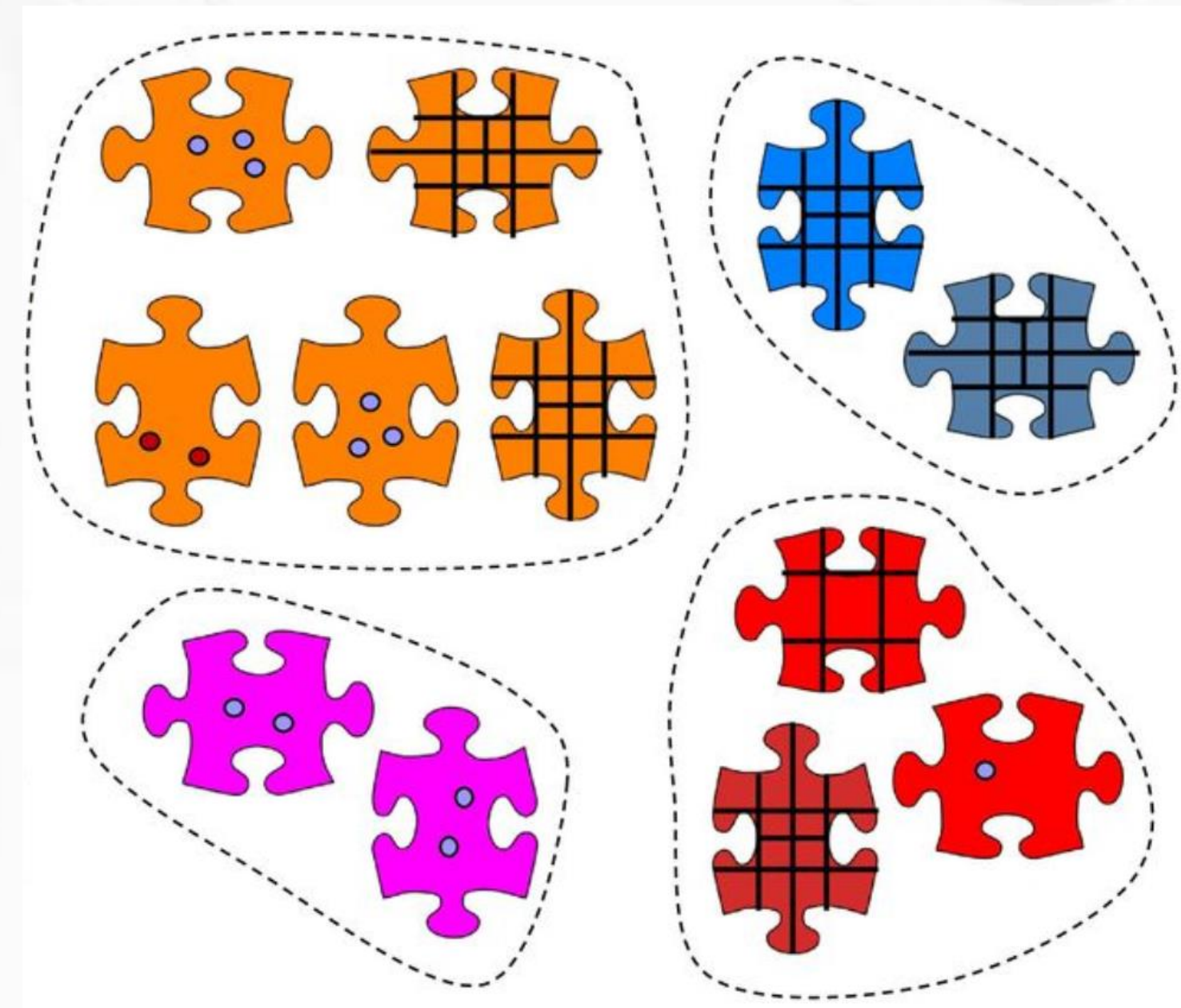
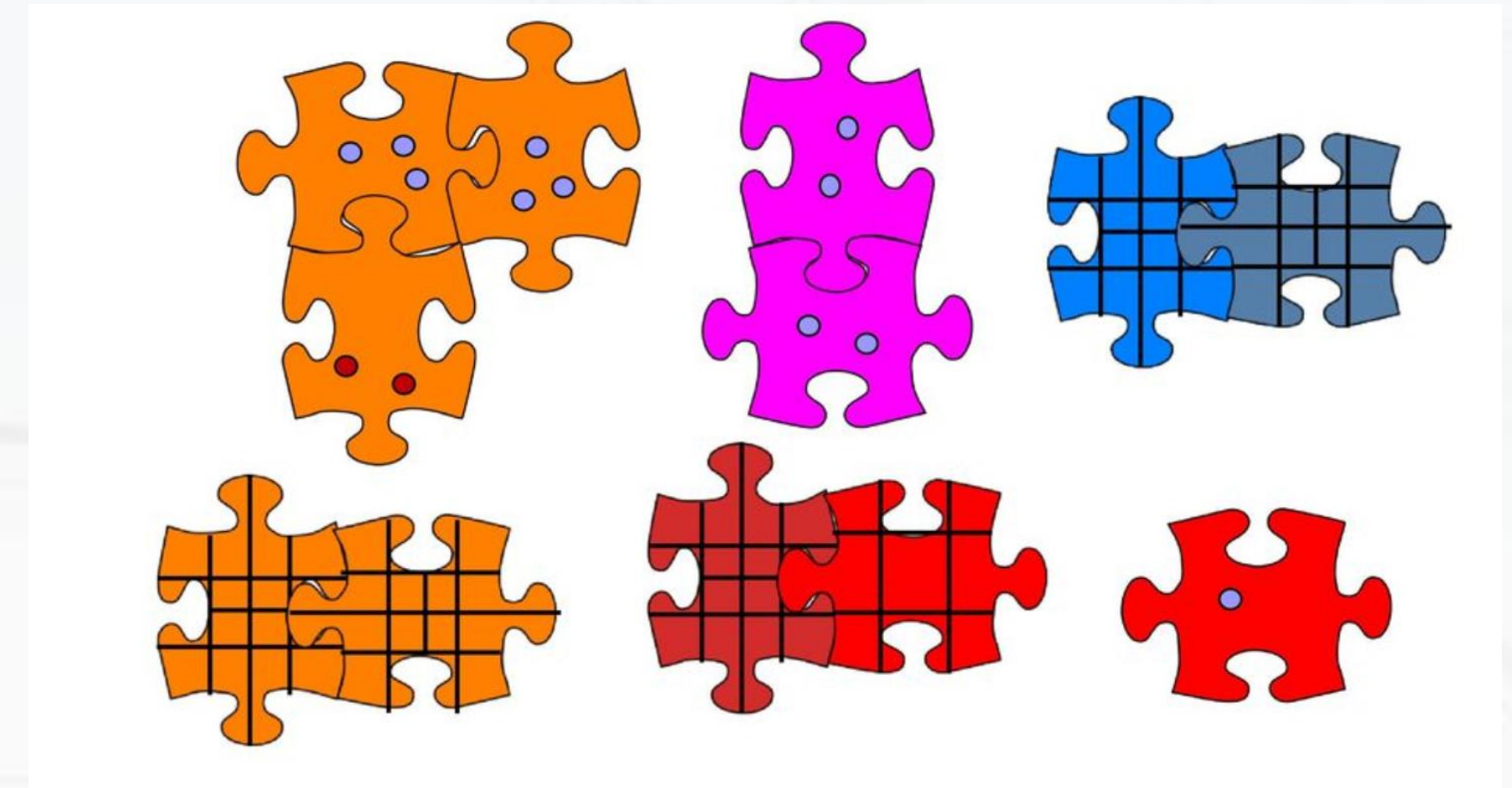
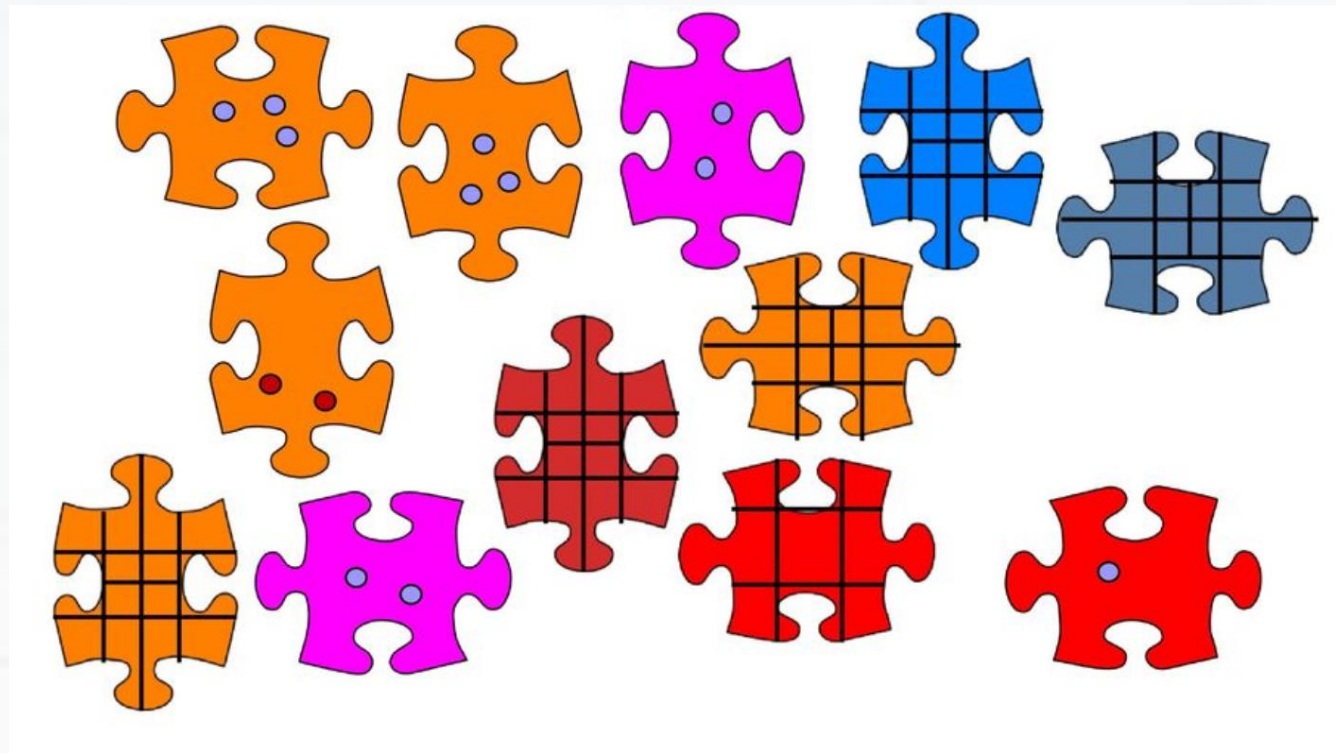


# Blocking



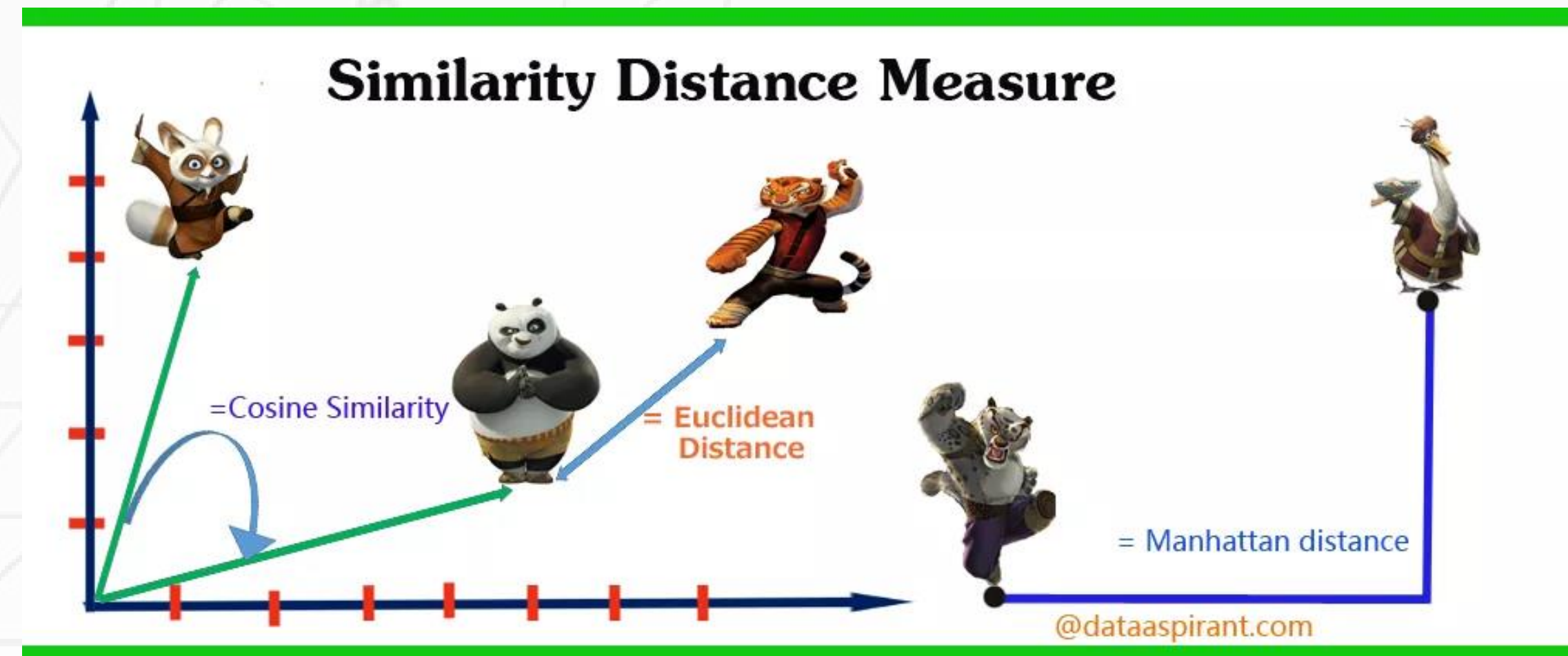
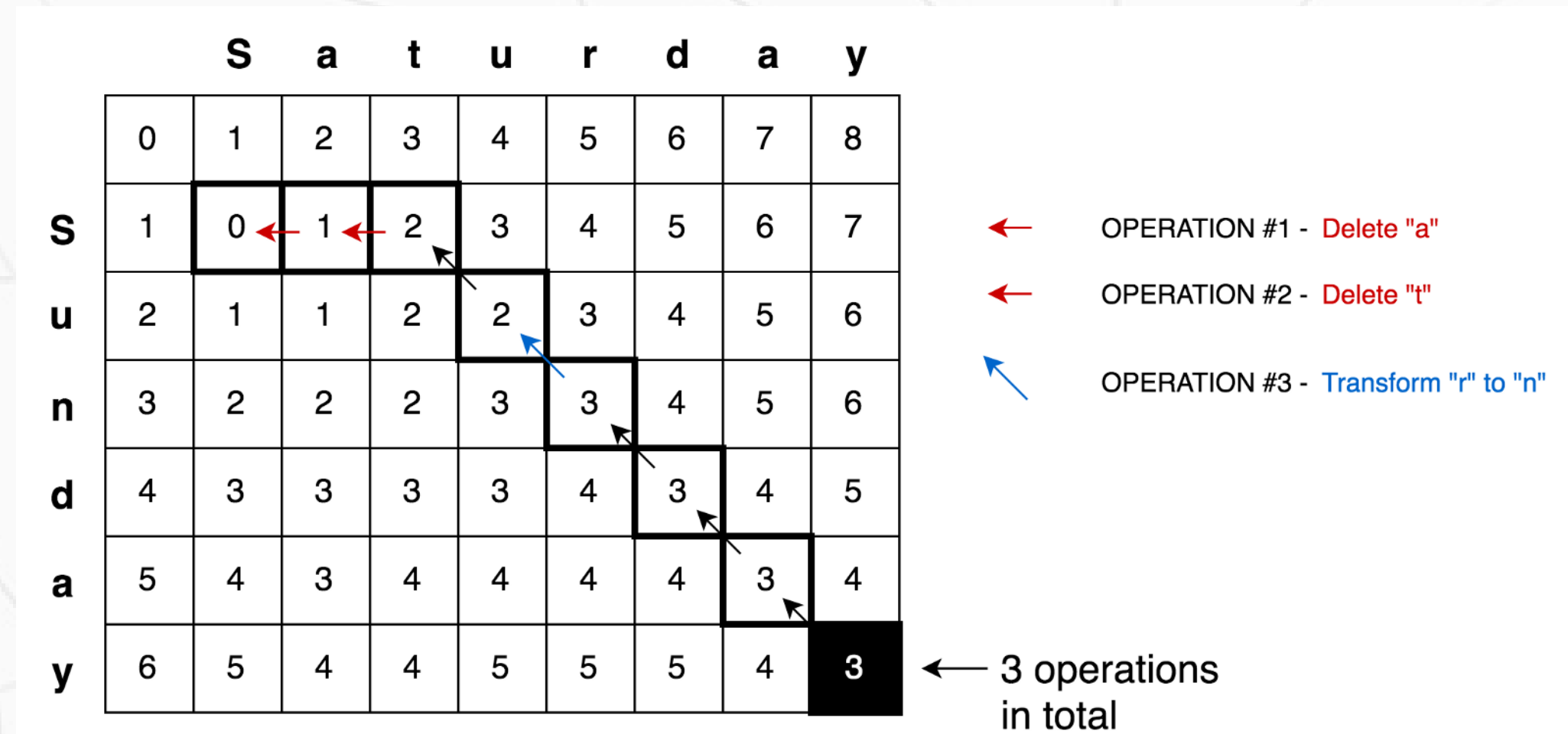


# Blocking



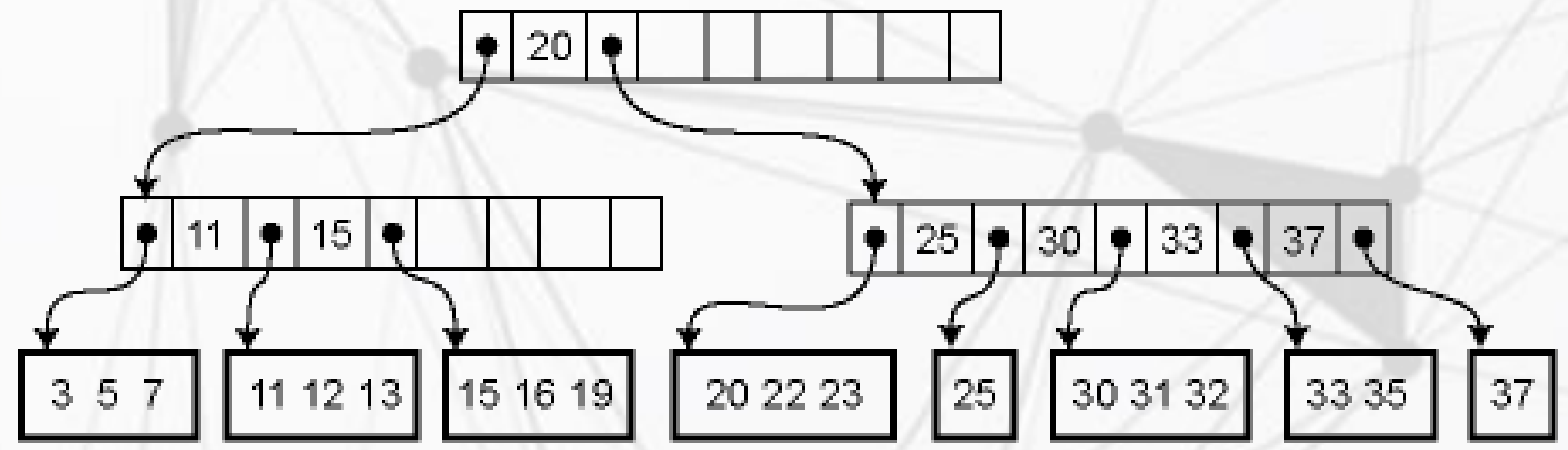


# Similarity functions



# Traditional Solution

▼ Index based



Insert(13)



# Pairwise Combination problem

Table 1

RecordID	Naam	Voornaam	Tel	Email	Twitter	Cookield
A	Michiels	T		<a href="mailto:tom@gmail.com">tom@gmail.com</a>	@TmMichiels	3445454
B	Michiels	T	027599608		@TmMichiels	
C	Moreau	M	027599608	<a href="mailto:t+m@gmail.com">t+m@gmail.com</a>		
D	Michiels	T		<a href="mailto:t+m@gmail.com">t+m@gmail.com</a>		
E	Michiels	T		<a href="mailto:t+m@gmail.com">t+m@gmail.com</a>		3556565

# Pairwise Combination problem

Table 1

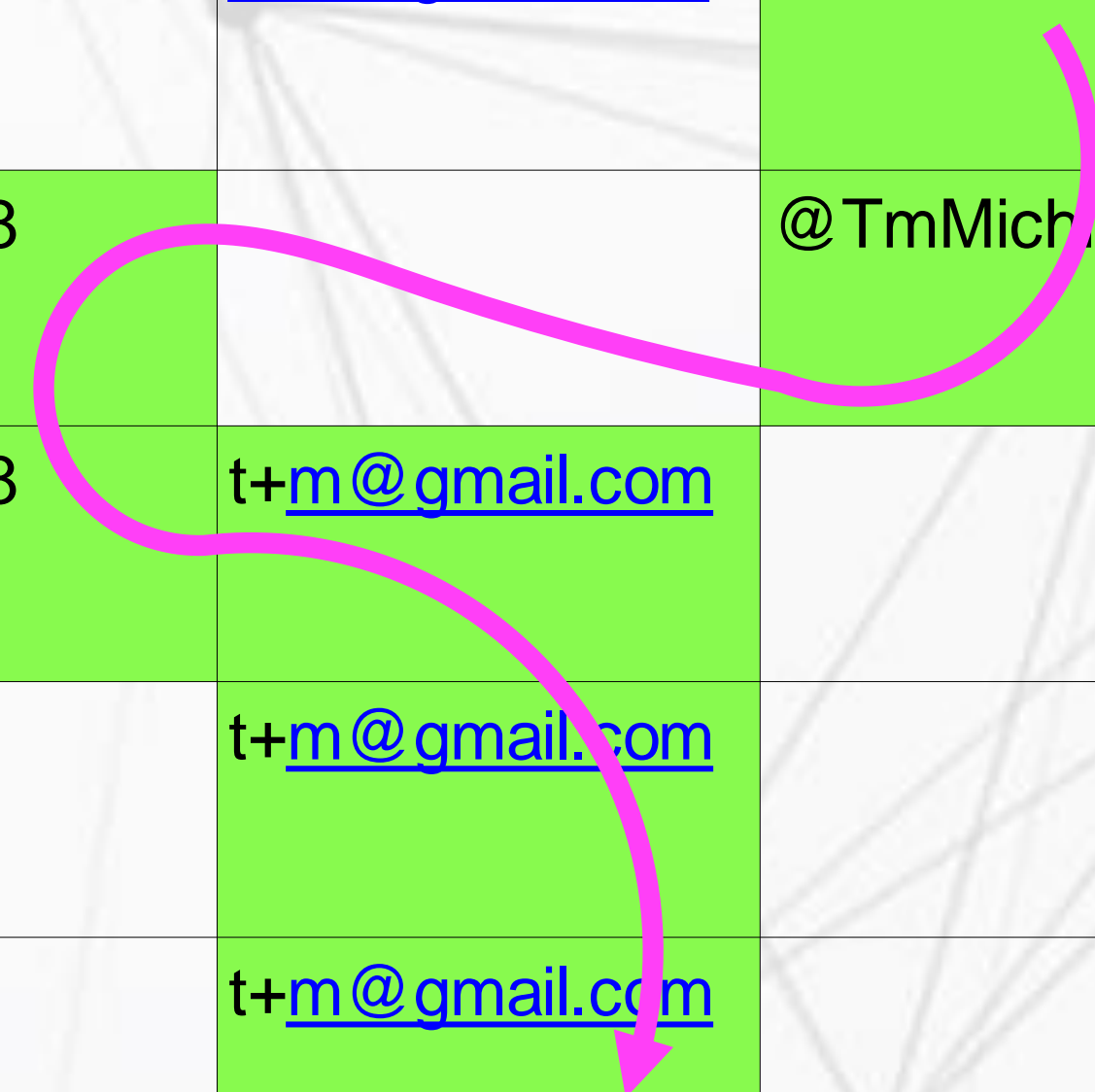
RecordID	Naam	Voornaam	Tel	Email	Twitter	Cookield
A	Michiels	T		<a href="mailto:tom@gmail.com">tom@gmail.com</a>	@TmMichiels	3445454
B	Michiels	T	027599608		@TmMichiels	
C	Moreau	M	027599608	<a href="mailto:t+m@gmail.com">t+m@gmail.com</a>		
D	Michiels	T		<a href="mailto:t+m@gmail.com">t+m@gmail.com</a>		
E	Michiels	T		<a href="mailto:t+m@gmail.com">t+m@gmail.com</a>		3556565



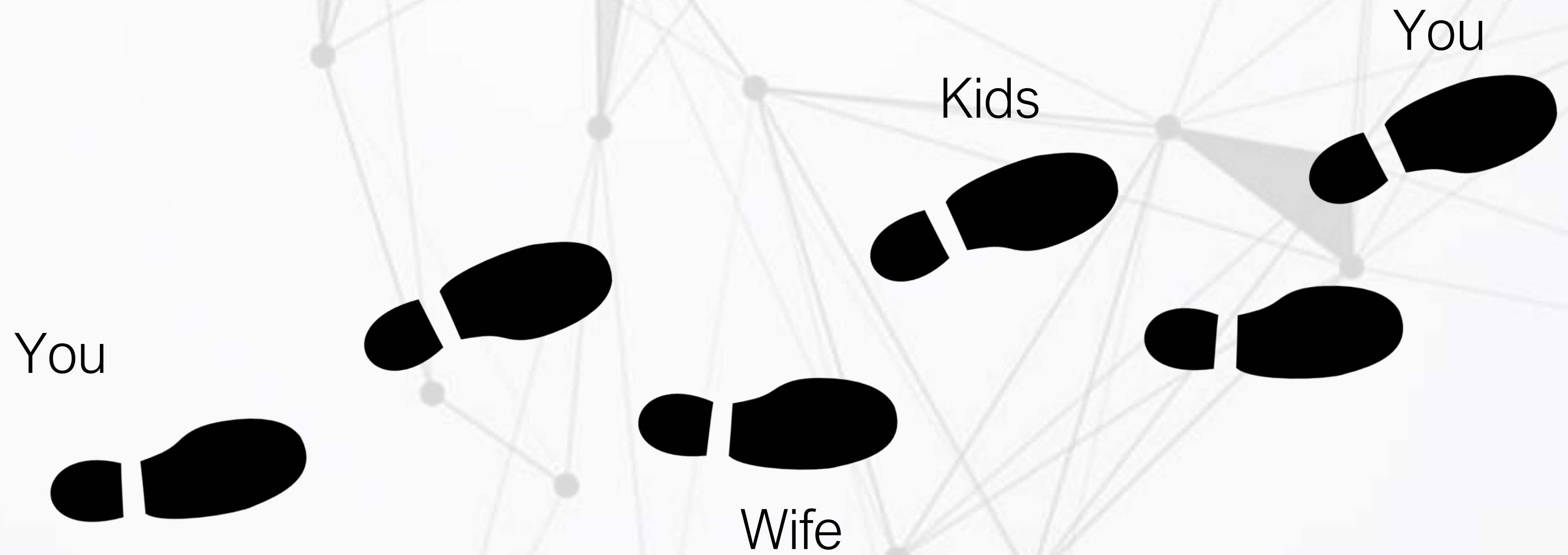
# Pairwise Combination problem

Table 1

RecordID	Naam	Voornaam	Tel	Email	Twitter	Cookield
A	Michiels	T		<a href="mailto:tom@gmail.com">tom@gmail.com</a>	@TmMichiels	3445454
B	Michiels	T	027599608		@TmMichiels	
C	Moreau	M	027599608	<a href="mailto:t+m@gmail.com">t+m@gmail.com</a>		
D	Michiels	T		<a href="mailto:t+m@gmail.com">t+m@gmail.com</a>		
E	Michiels	T		<a href="mailto:t+m@gmail.com">t+m@gmail.com</a>		3556565



# Graphs to the rescue



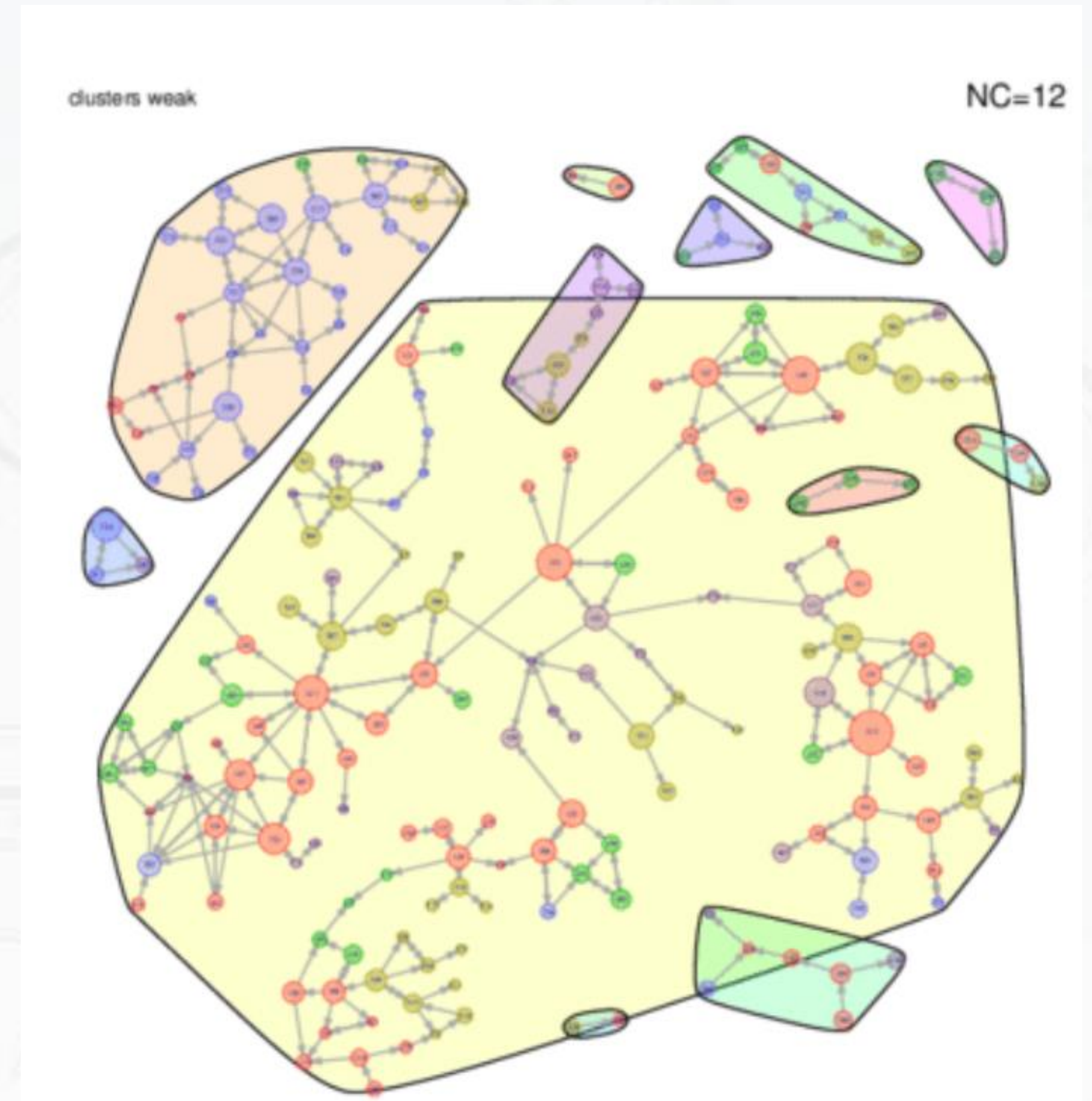
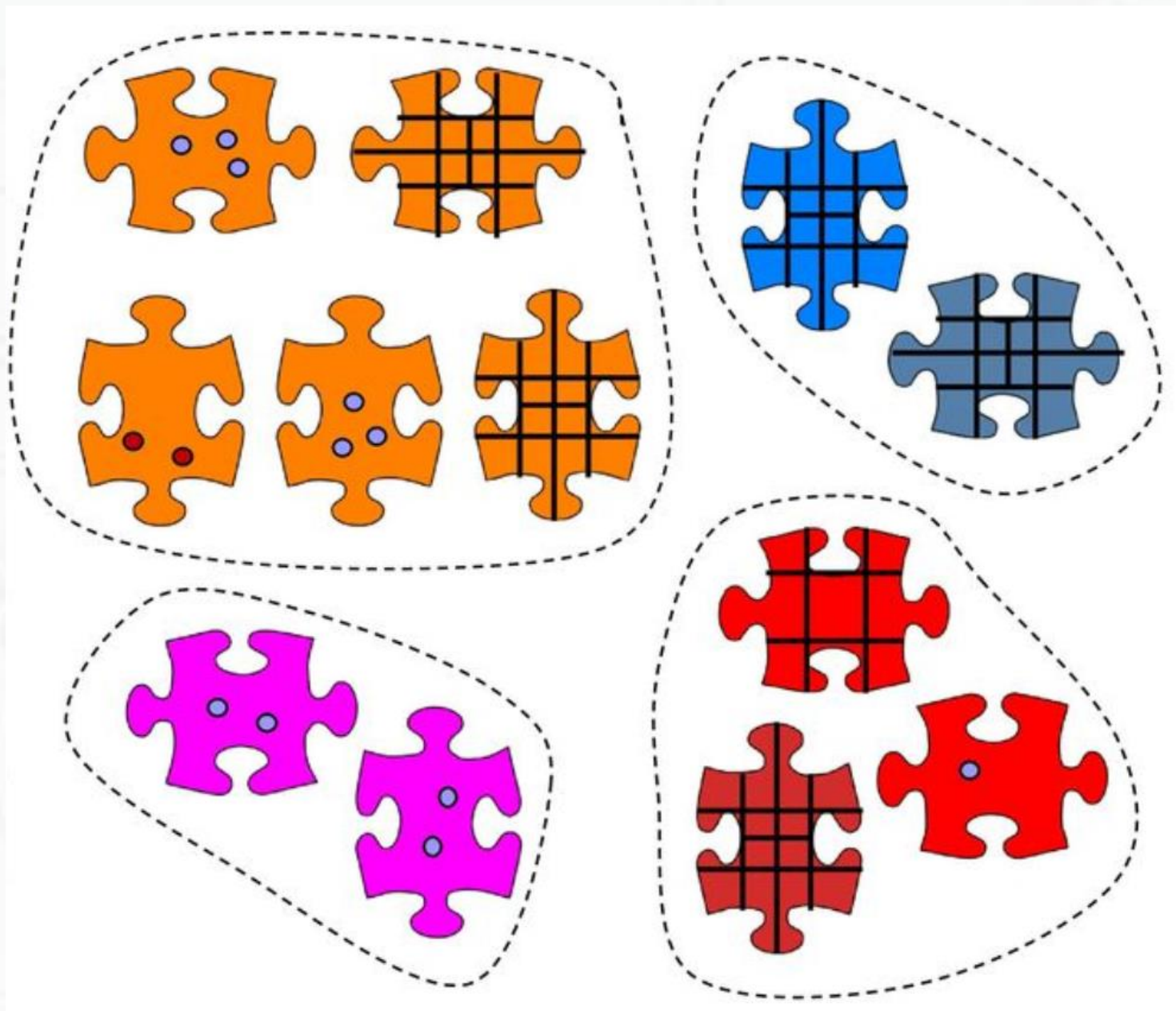


## Graphs to the rescue

- ▼ Graph locality
- ▼ Pointer hopping
- ▼ Graph algorithms



# Graphs to the rescue



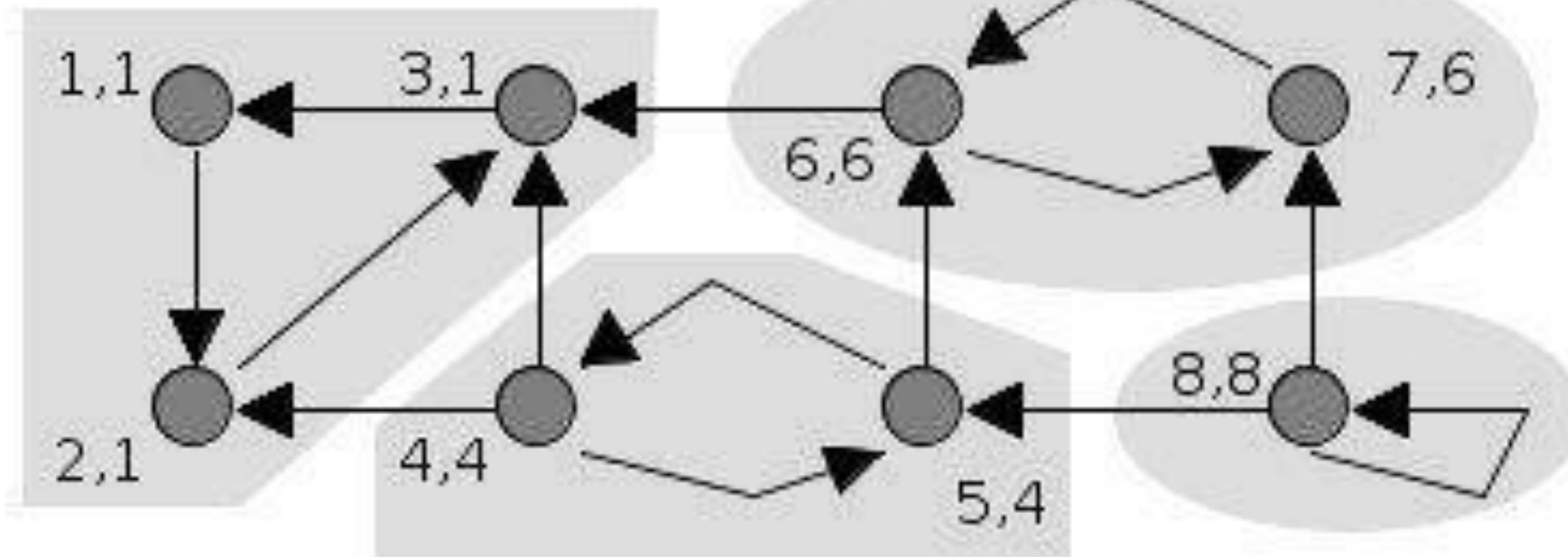


# Weakly connected component

Table 1

RecordID	Naam	Voornaam	Tel	Email	Twitter	Cookield
A	Michiels	T		<a href="mailto:tom@gmail.com">tom@gmail.com</a>	@TmMichiels	3445454
B	Michiels	T	027599608		@TmMichiels	
C	Moreau	M	027599608	<a href="mailto:t+m@gmail.com">t+m@gmail.com</a>		
D	Michiels	T		<a href="mailto:t+m@gmail.com">t+m@gmail.com</a>		
E	Michiels	T		<a href="mailto:t+m@gmail.com">t+m@gmail.com</a>		3556565

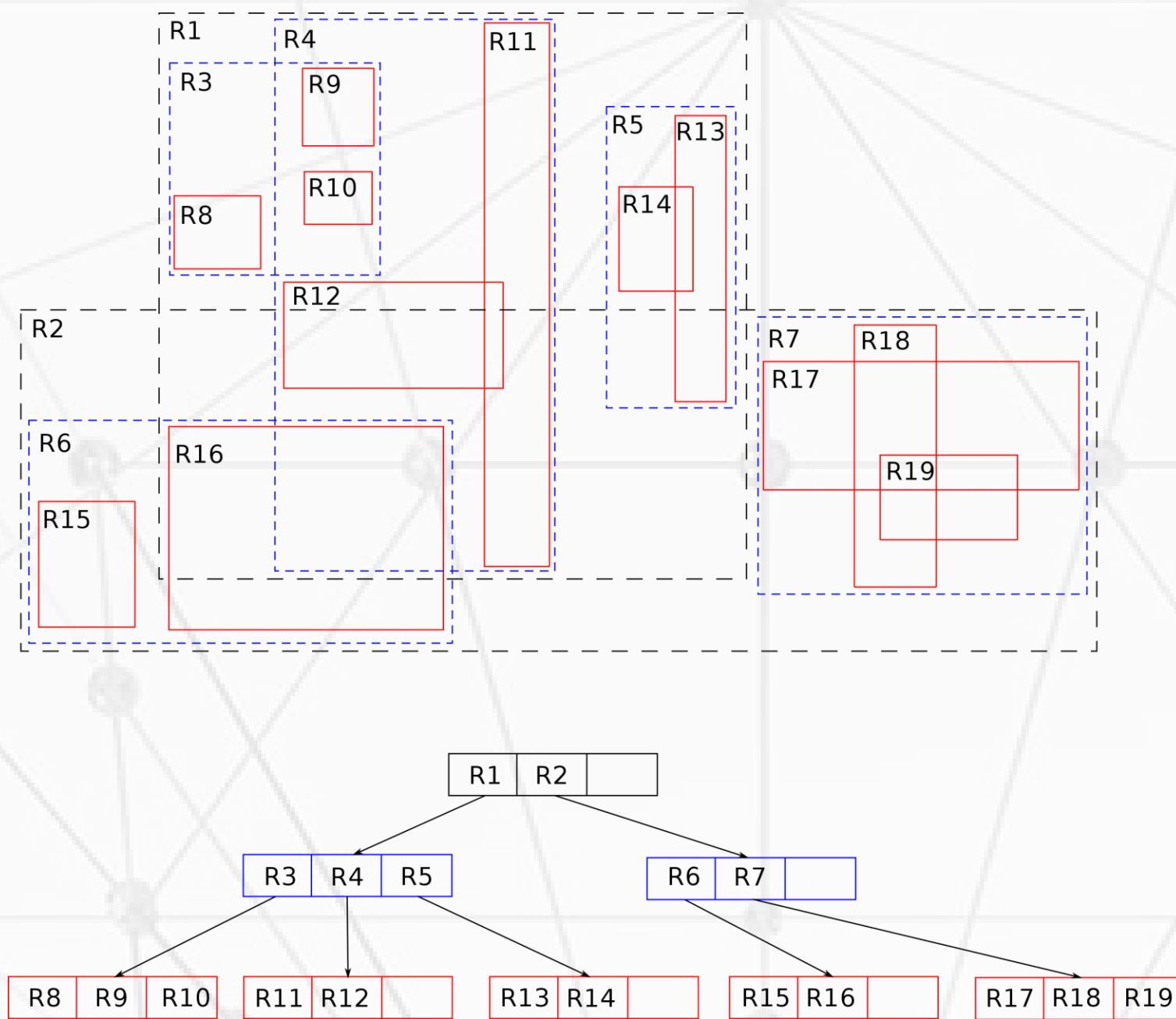
# Tarjan's connected component



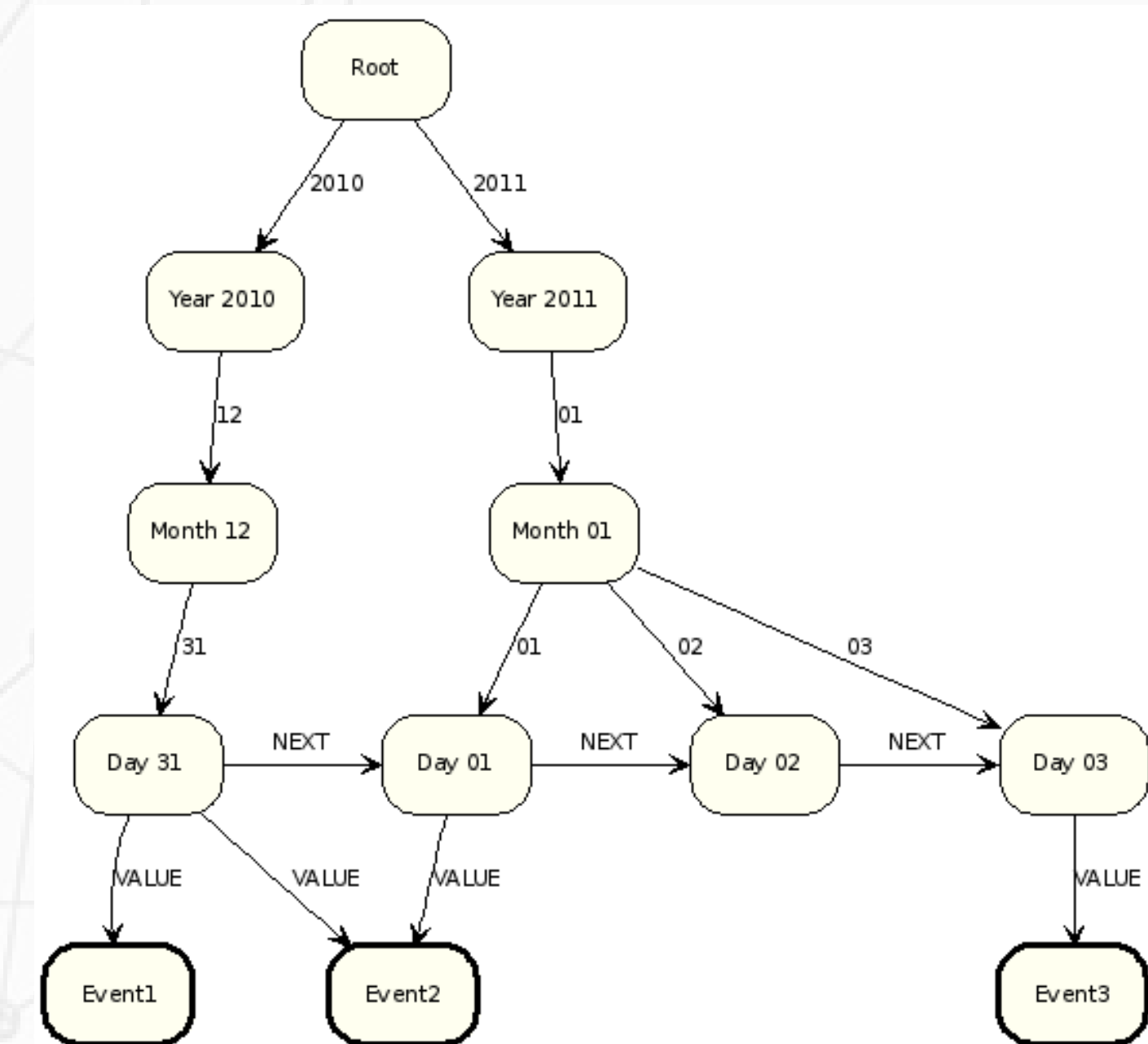


# Extension - Compare in Space and Time.

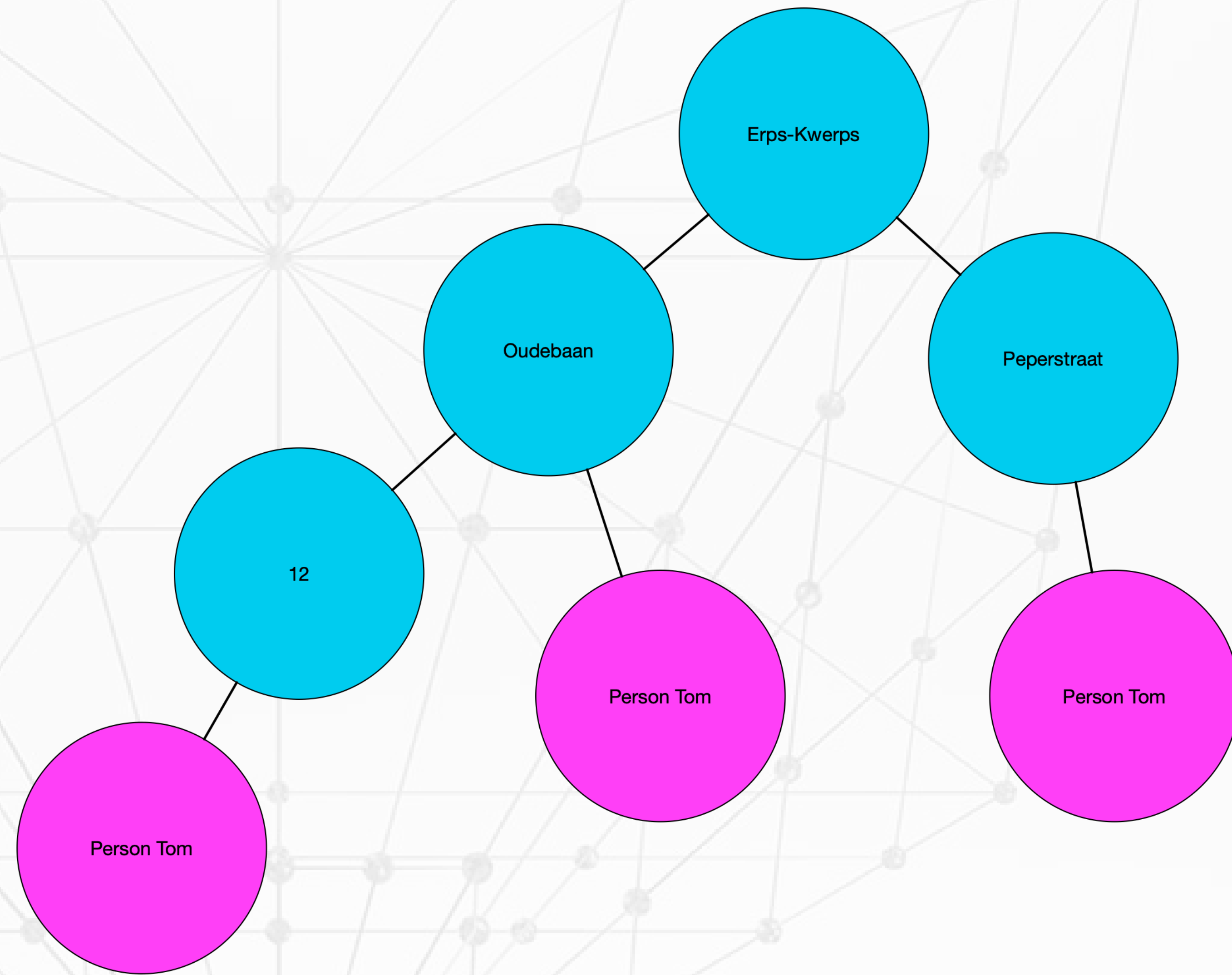
## R-Tree



## Timetree

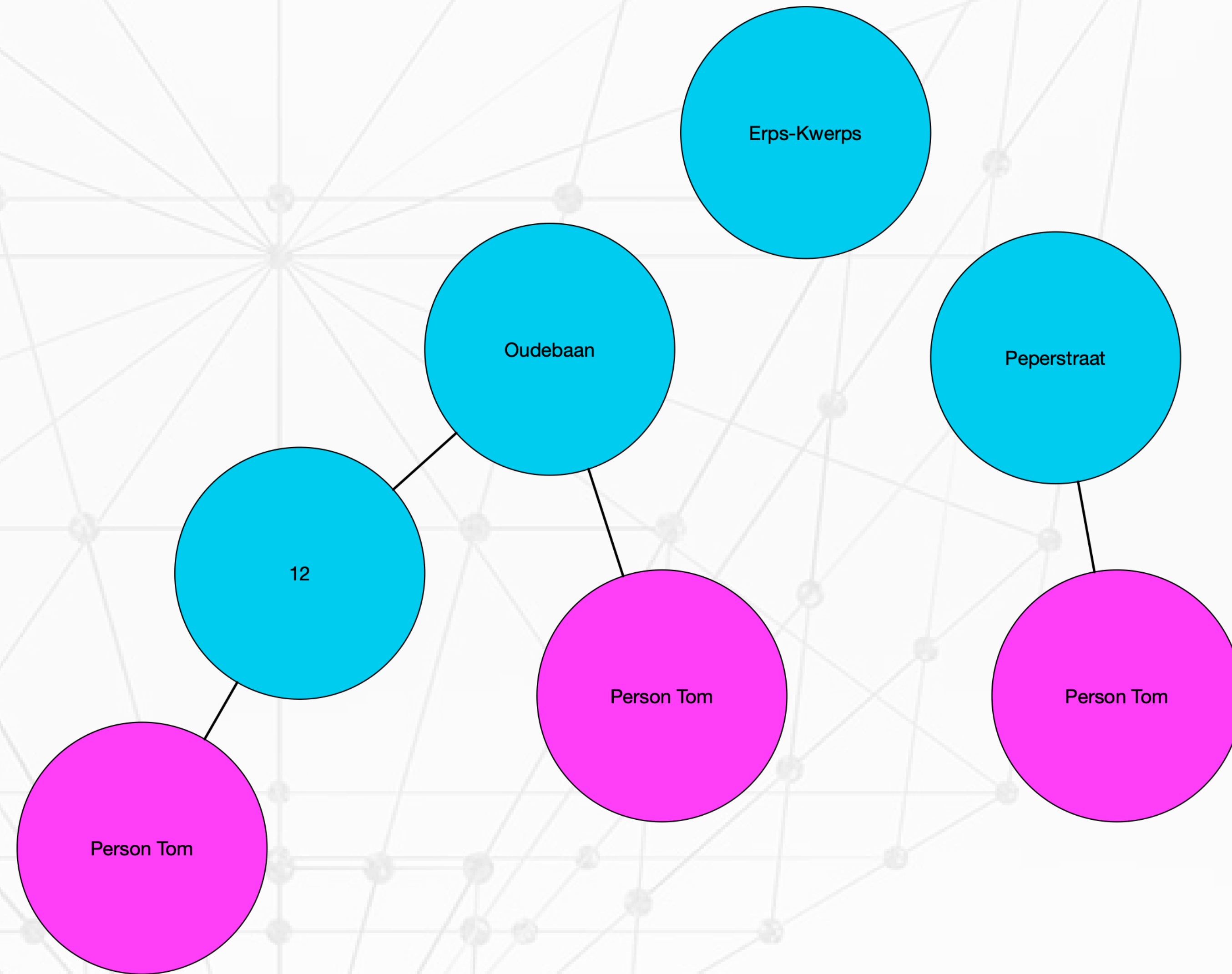


# Extension - Compare in Space and Time.

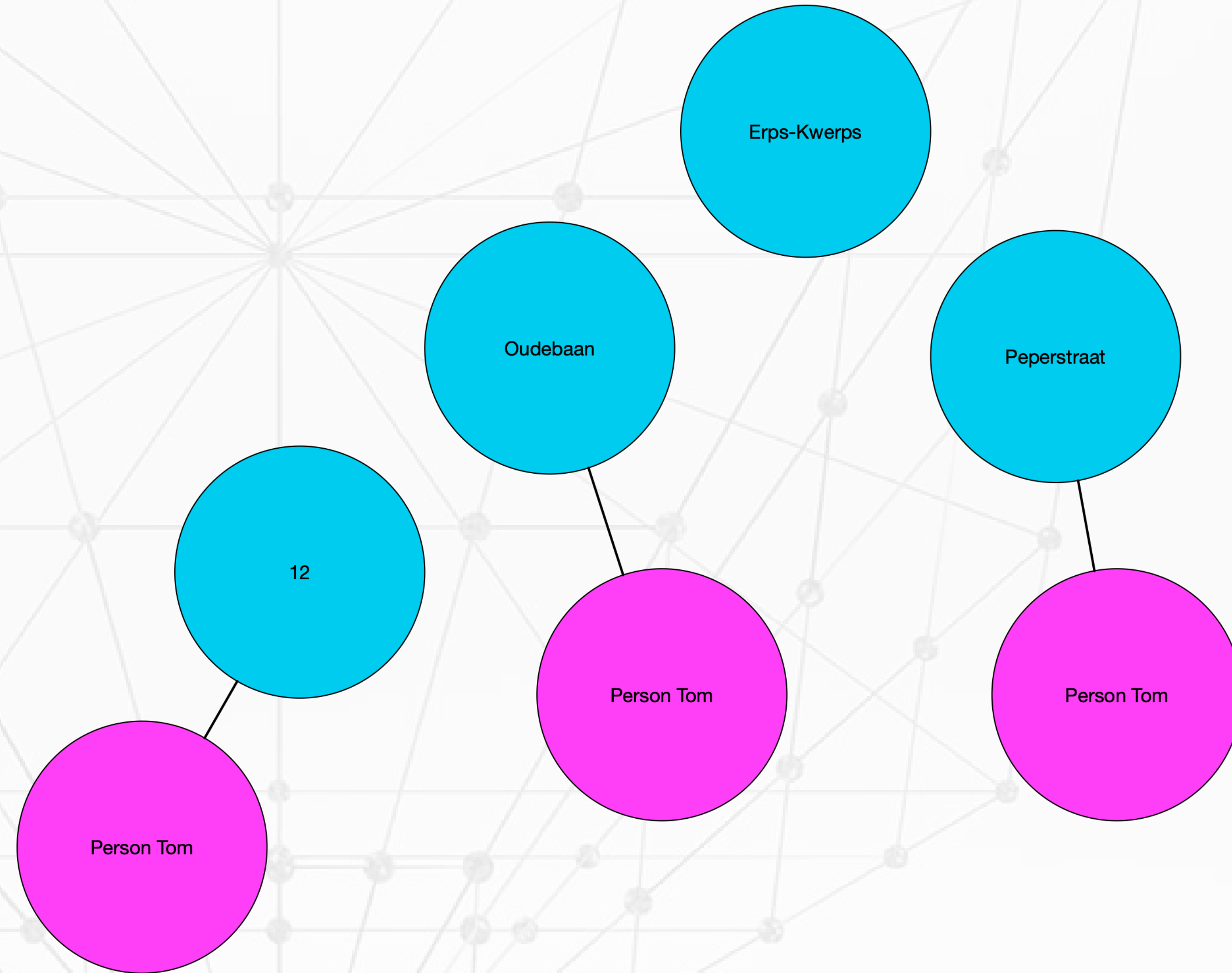




# Extension - Compare in Space and Time.

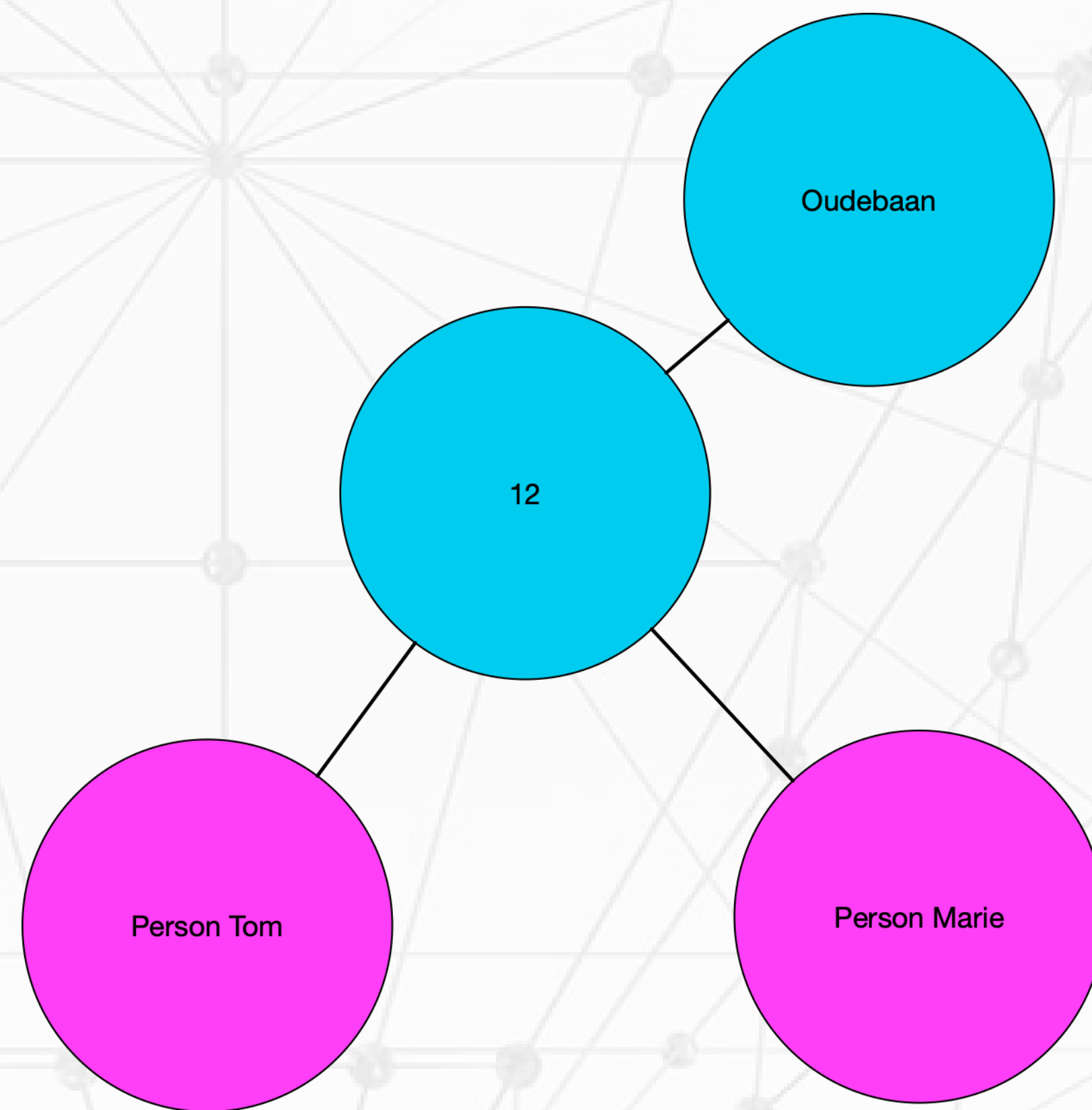


# Extension - Compare in Space and Time.





# Added Value Connected Component = Household

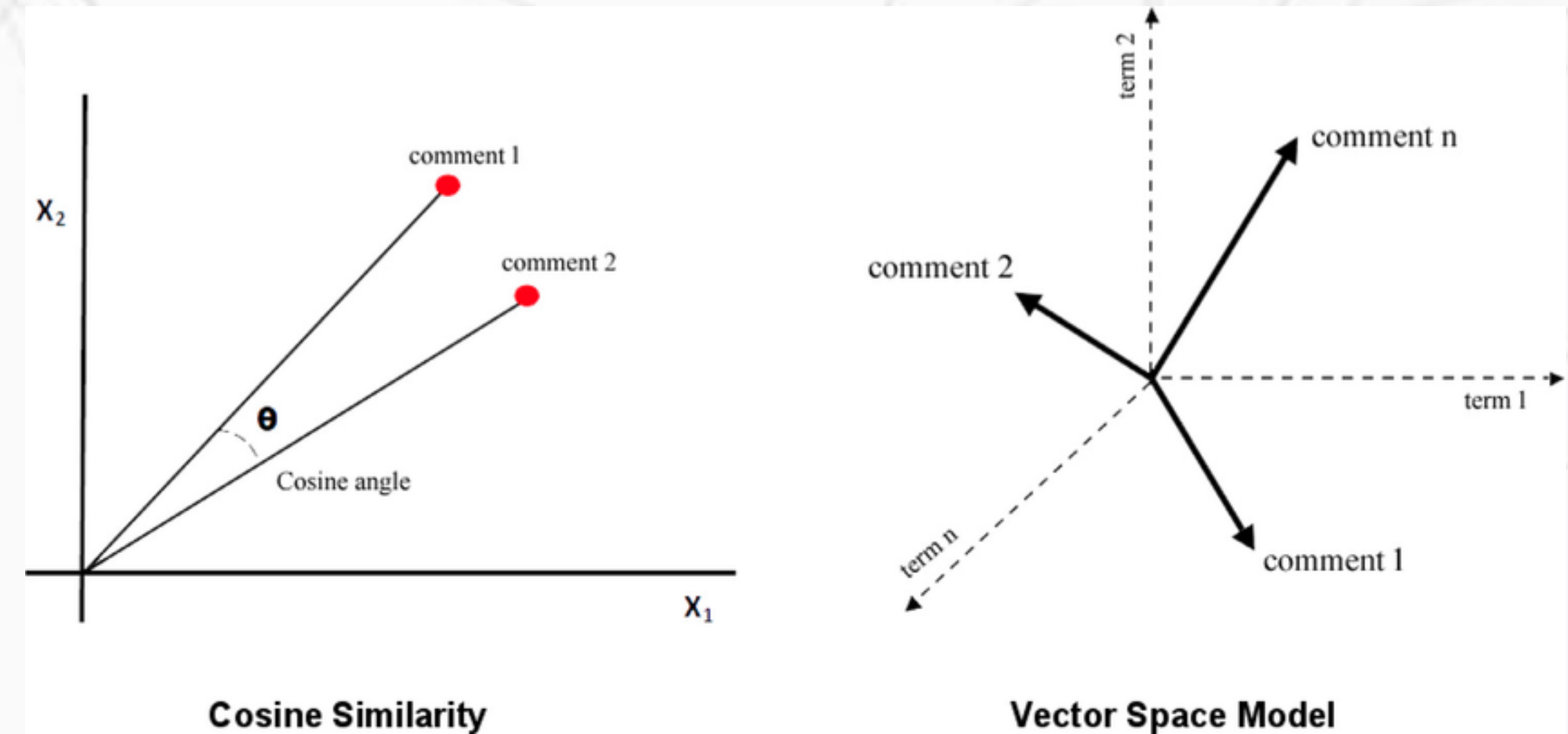
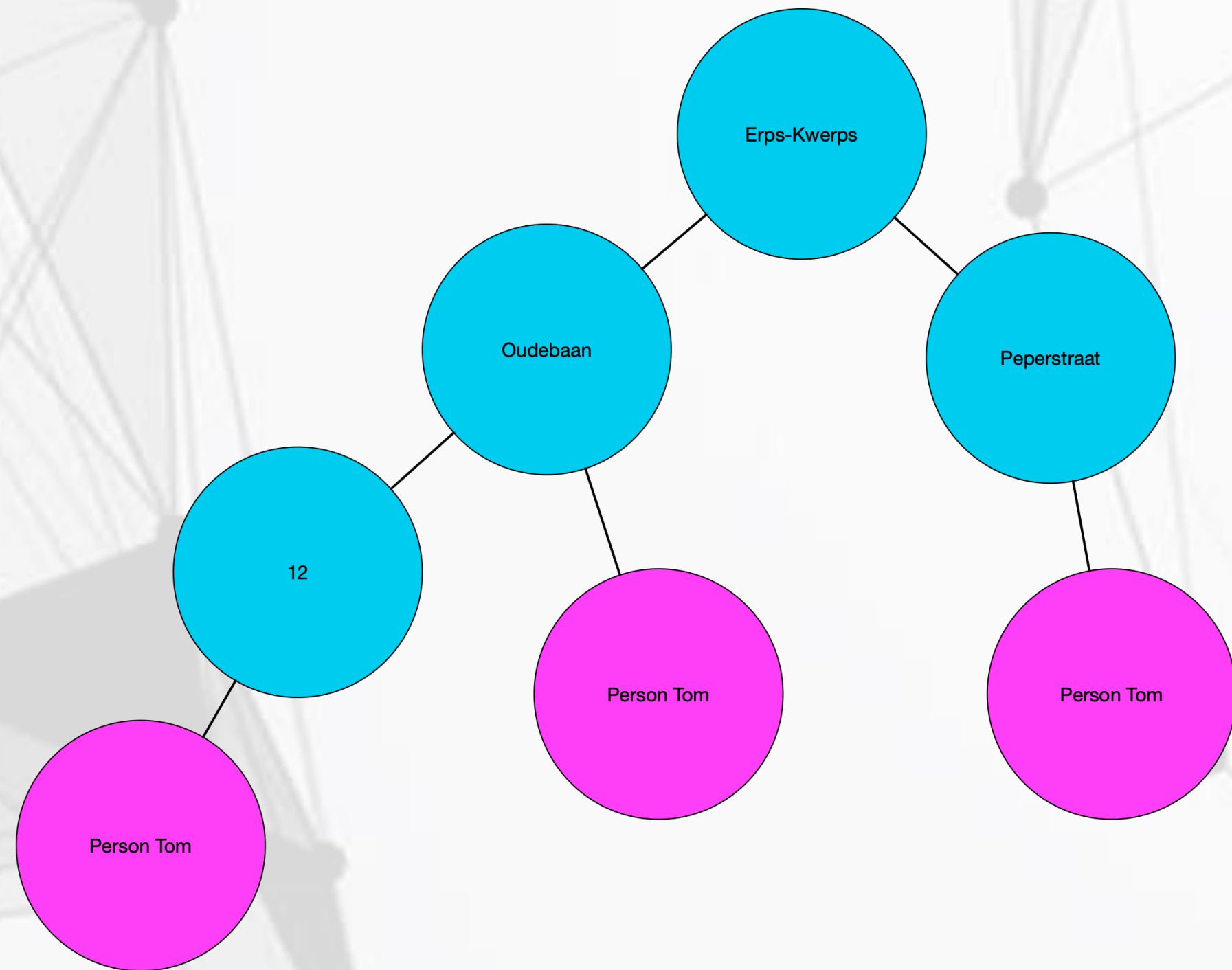


## Advantages

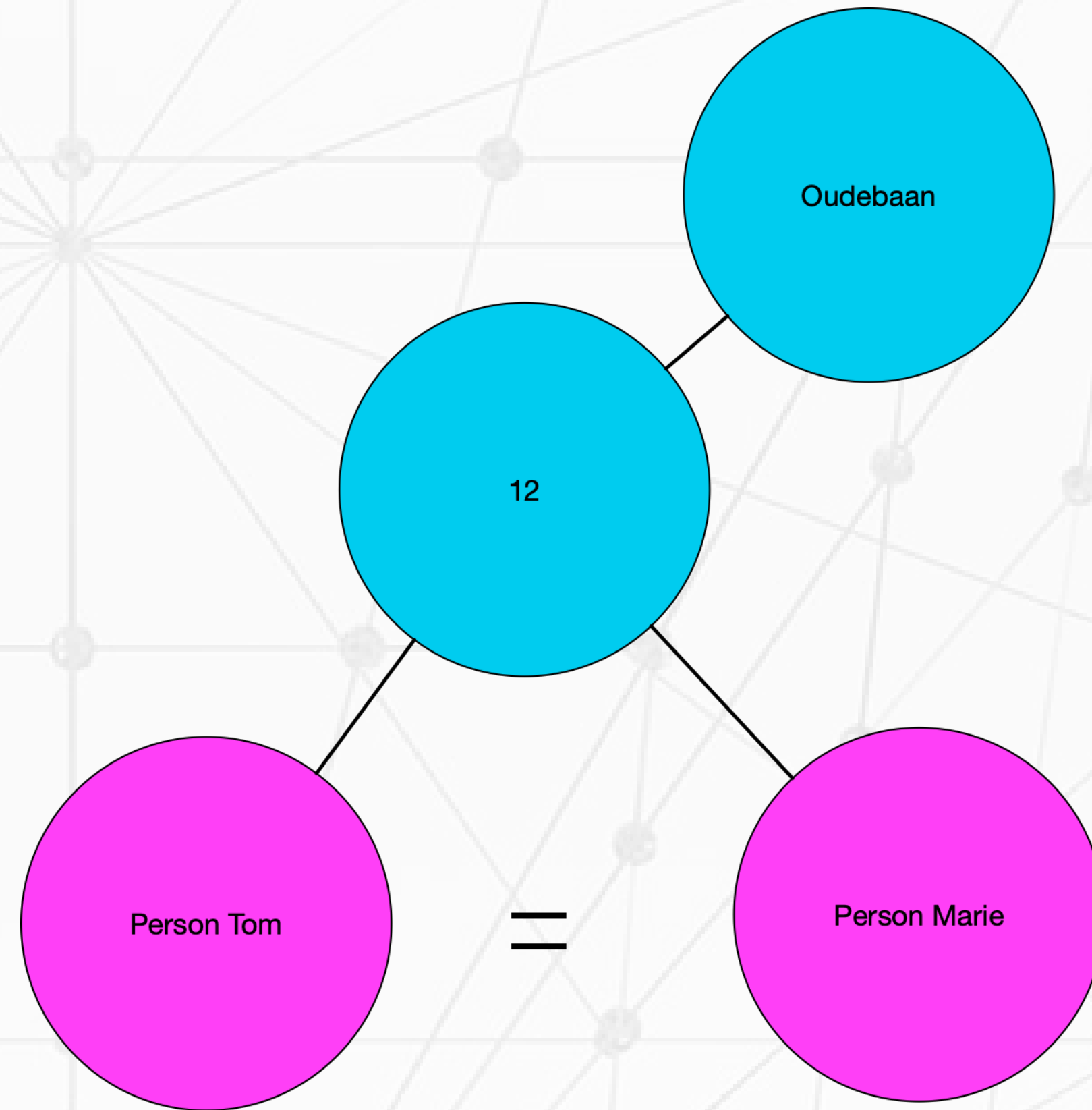
- ▶ Linking through other records (Paths)
- ▶ Finding candidates through trees (Time - Location)
- ▶ Deduplication and Cluster detection is 2 step process.



# Graph2Vec



# Comparison of pairs





# Comparison

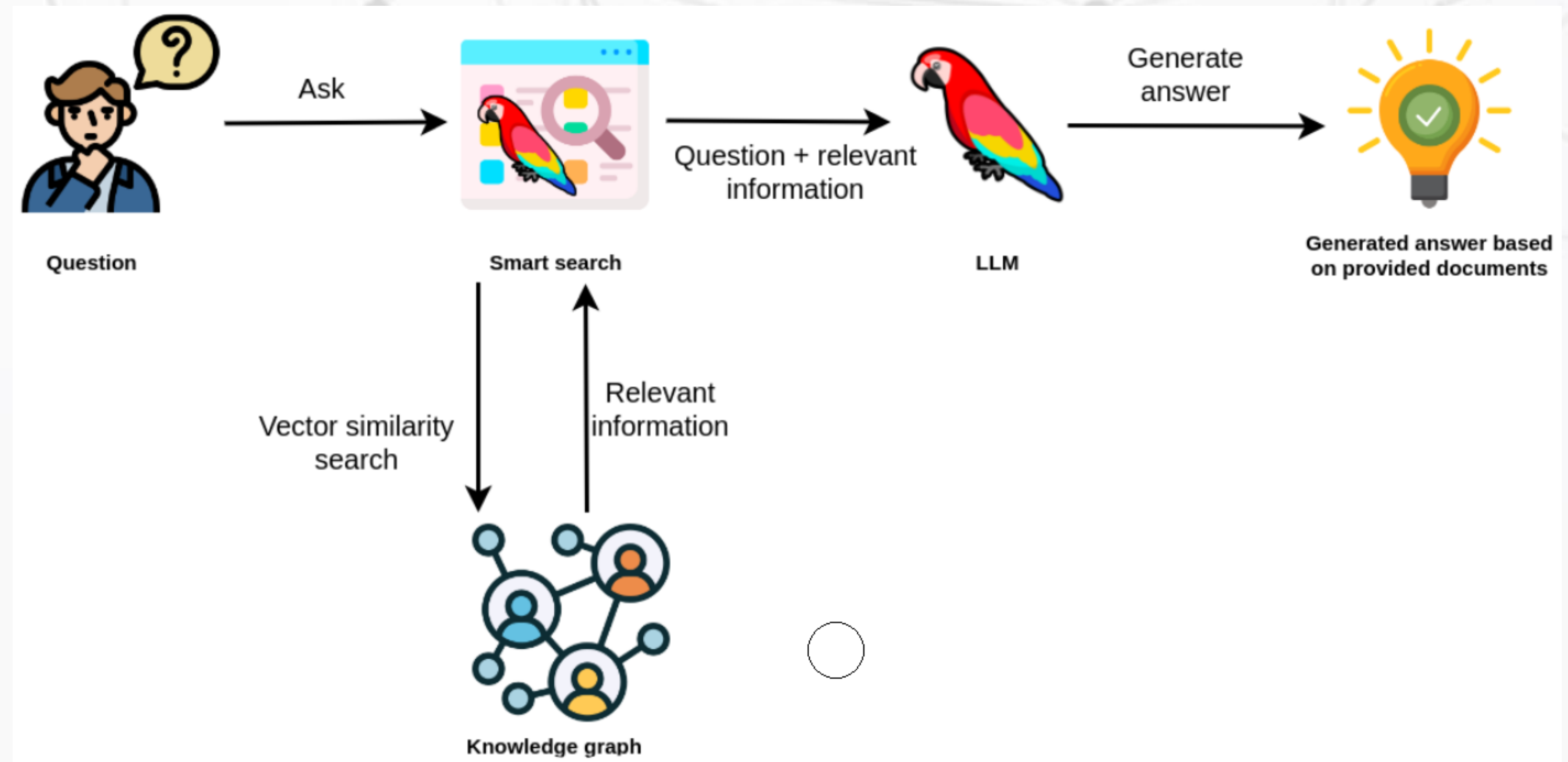
RecordID	Naam Metric1	Voornaam Metric2	Naam occurence	Voornaam occurence	Match
A - B	0,6	1	0,1	0,1	True
A - C	0,8	0,9	0,1	0,1	False

## Comparison

- ▼ Rule based (Deterministic)
- ▼ Machine Learning. Active Learning, Random forests, ...

## Other Possibilities

- Find Documents linked to each other by entities (Same principles as finding Pairs)





# Fraud Detection

- Find things that do not want to be found

## Modeling a fraud ring as a graph



## Conclusion

- ▼ Data Quality In improves data quality out.
- ▼ Graph algorithms can be used in other use cases like