



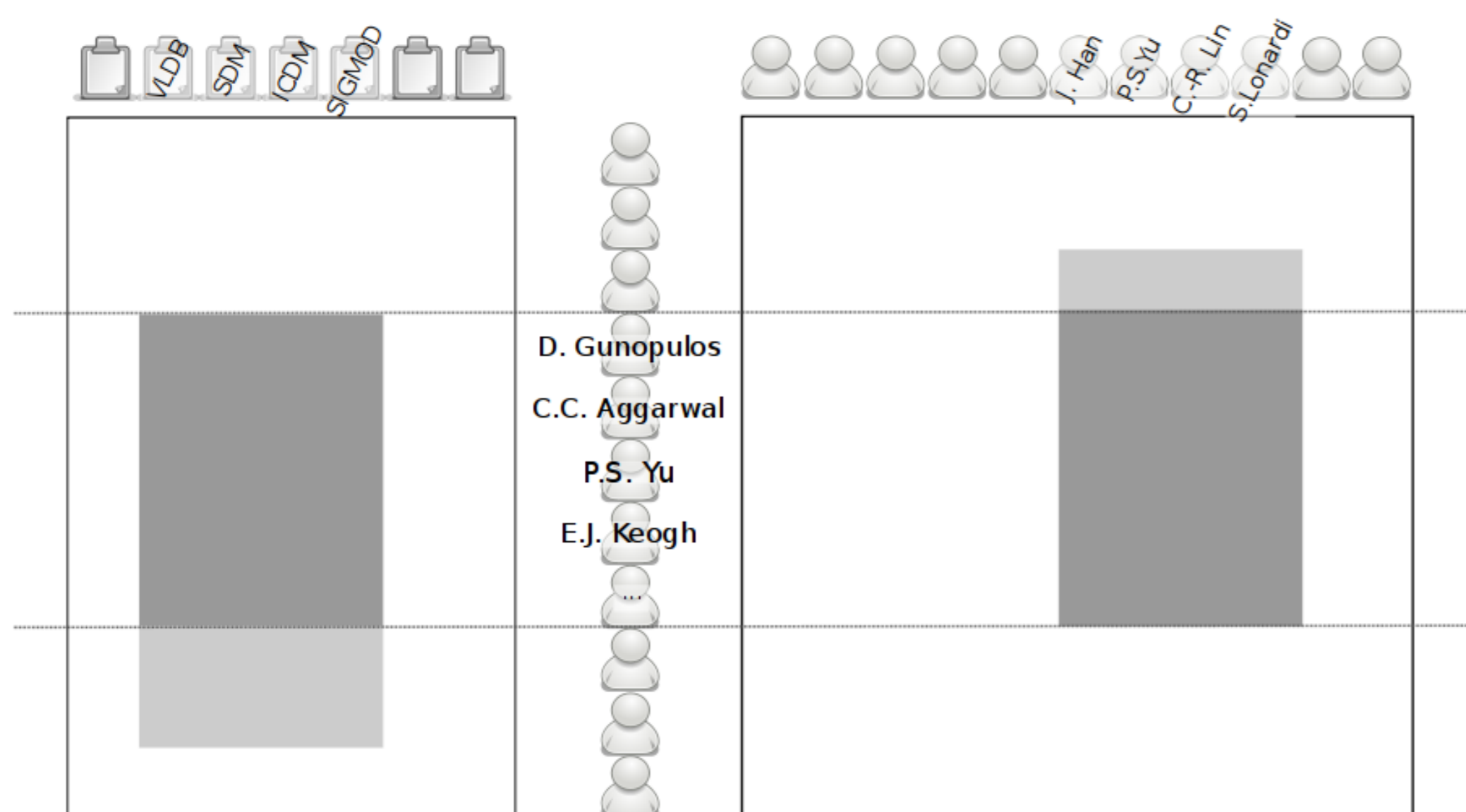
REDESCRIPTION MINING OUTSIDE THE BOOLEAN WORLD

Esther Galbrun and Pauli Miettinen

EXAMPLE

DBLP Bibliography Data

VLDB \wedge ICDM \wedge SDM \wedge SIGMOD
(J. Han \wedge P.S. Yu) \vee C.-R. Lin \vee S. Lonardi



Given two datasets with a bijection between the rows, a **redescription** is a pair of queries (q_L, q_R) over the columns characterizing approximately the same sets of rows.

In **redescription mining**, the task is to find the best redescriptions satisfying a given set of constraints.

Approaches to **boolean redescription mining** include decision trees [3], co-clusters [2], frequent itemsets [1] and greedy algorithm [1].

Jaccard coefficient:

$$J = \frac{|E_{1,1}|}{|E_{1,0}| + |E_{1,1}| + |E_{0,1}|}$$

PROBLEM DEFINITION

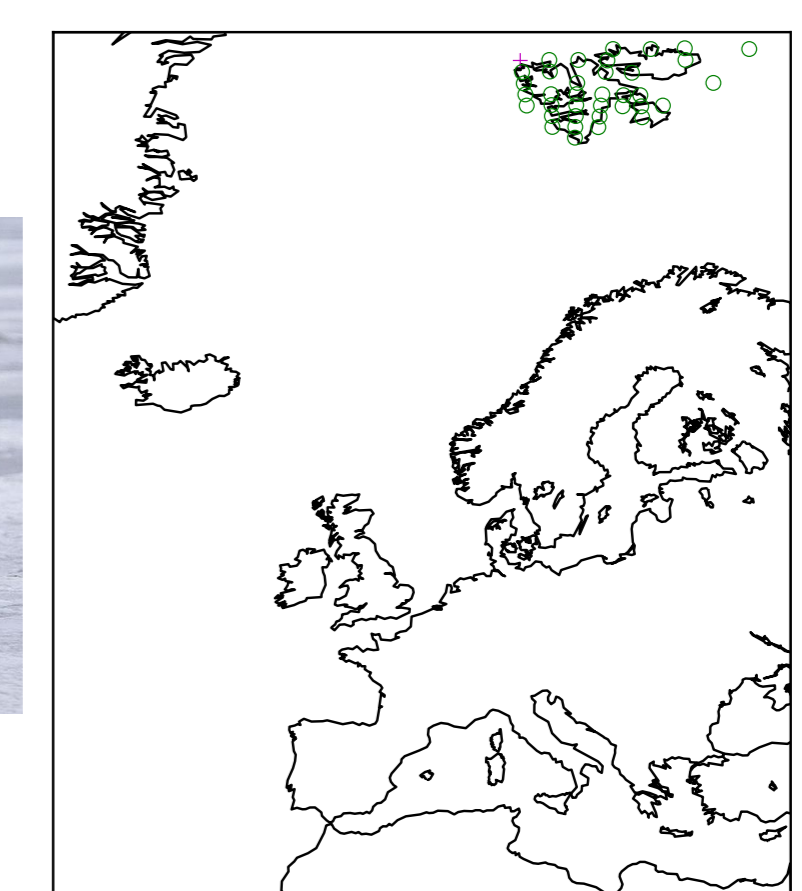
APPLICATION

Bioclimatic Niche Finding

Polar Bear

$$[-7.0727 \leq t_{\text{May}}^{\text{avg}} \leq -3.375]$$

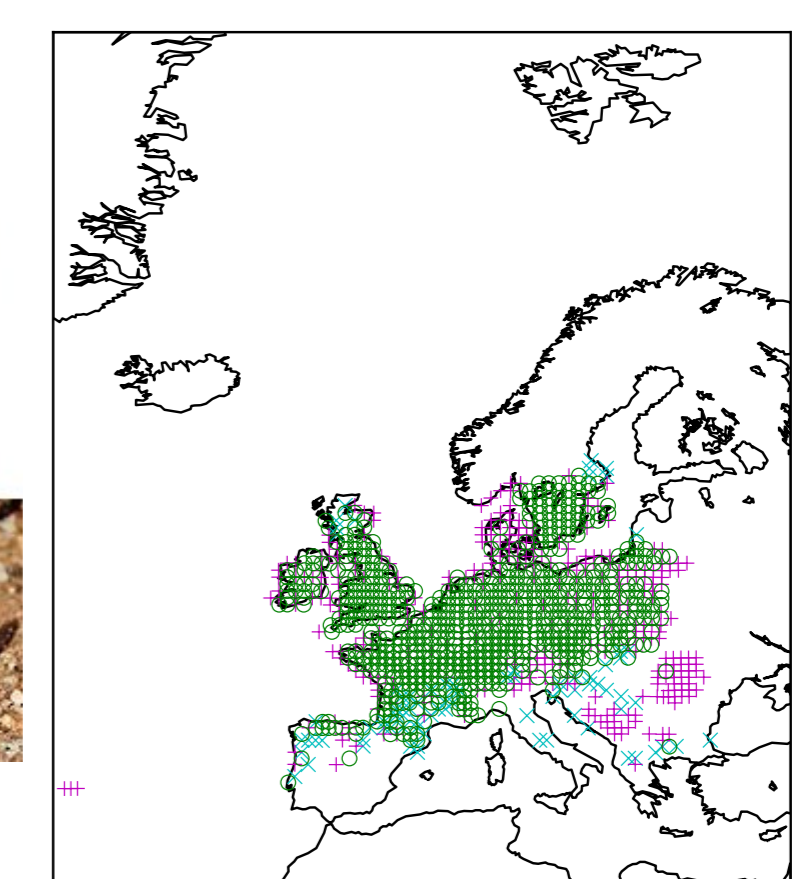
$$J = 0.973 \quad \text{supp} = 36$$



Wood Mouse \wedge Natterer's Bat
 \wedge Eurasian Pygmy Shrew

$$([3.20 \leq t_{\text{Mar}}^{\text{max}} \leq 14.50] \wedge [17.30 \leq t_{\text{Aug}}^{\text{max}} \leq 25.20]) \wedge [14.90 \leq t_{\text{Sep}}^{\text{max}} \leq 22.80] \vee [19.60 \leq t_{\text{Jul}}^{\text{avg}} \leq 19.956]$$

$$J = 0.623 \quad \text{supp} = 681$$

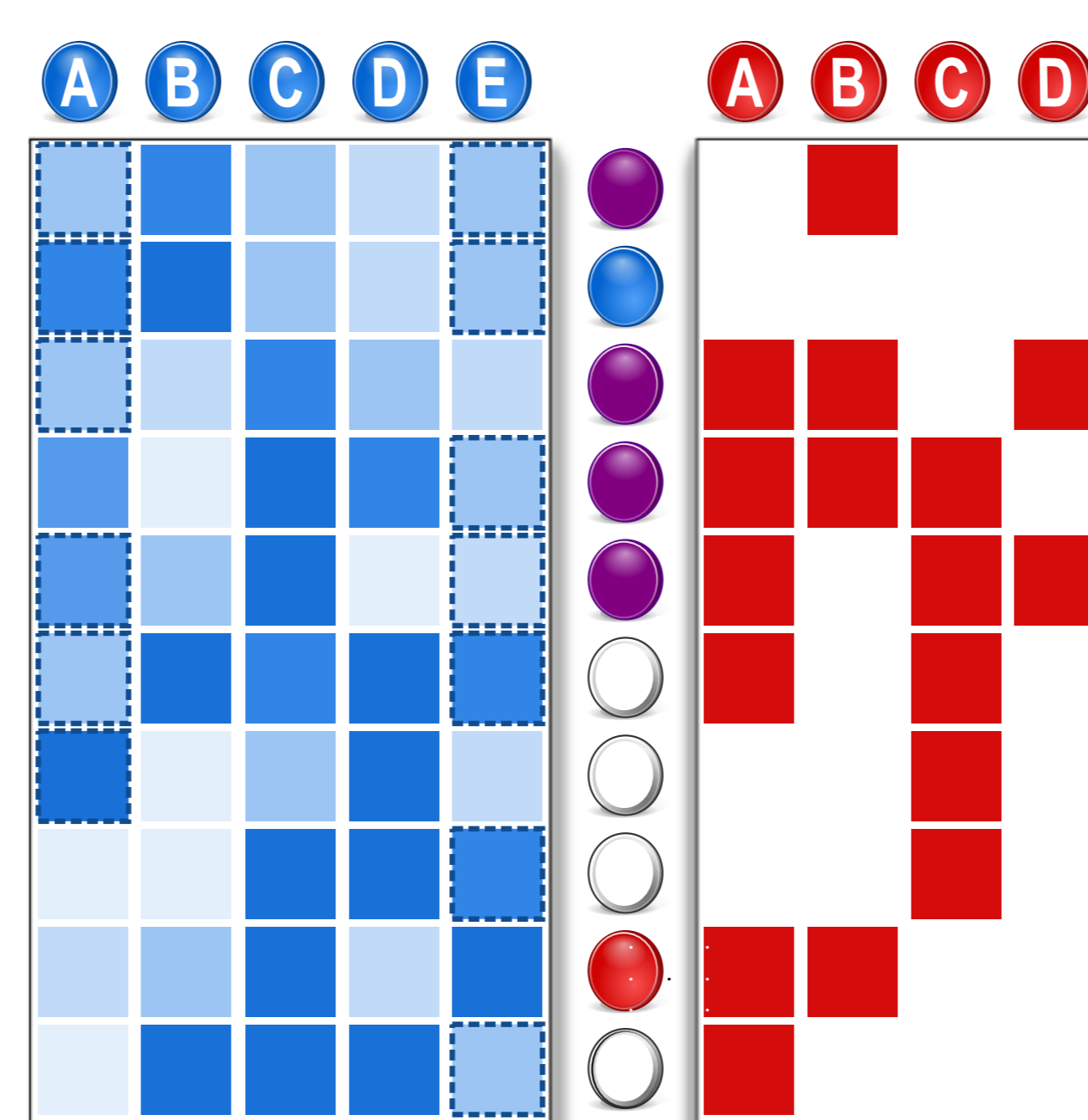


OUTSIDE THE BOOLEAN WORLD

Existing methods require **discretization** as a **pre-processing** step for real-valued data. This can lead to an explosion of the number of variables and generally requires extensive domain knowledge.

This is in contrast to our approach, where the **optimal interval** is **automatically determined** at each step during the greedy query extension.

$$[\square \leq A \leq \square] \wedge [\square \leq E \leq \square] \\ (A \wedge D) \vee B$$



FUTURE WORK

Generalizing the niche-finding problem to **traits** and identifying other domains of application (e.g. medicine) is one potential direction for future research.

Improving the algorithm and obtaining proofs of the behavior of redescription mining algorithms for real-valued data are other, more theoretical directions.

REFERENCES

- [1] Arianna Gallo et al. Finding subgroups having several descriptions: Algorithms for redescription mining. In SDM, pages 334–345, 2008.
- [2] Laxmi Parida and Naren Ramakrishnan. Redescription mining: Structure theory and algorithms. In AAAI, pages 837–844, 2005.
- [3] Naren Ramakrishnan et al. Turning cartwheels: an alternating algorithm for mining redescriptions. In KDD, pages 266–275, 2004.