Chaining Heuristic and Exact Methods for DFA Identification

1. Deterministic Finite Automata (DFA)

DFA – Set of states, transitions, alphabet -> A model that recognizes a certain regular language.



DFAs can act as surrogate models for software systems. However, the process of creating and maintaining them is costly and inefficient, and is usually omitted during software development.

2. DFA Identification



MLtut. K-Fold Cross Validation in Machine Learning. https://www.mltut.com/k-fold-cross-validation-in-machine-learning-how-does-k-fold-work/,

3. Research Question

To what extent and in what ways does **partially** learning a model heuristically and then applying exact minimization - based on the decisions already made by the heuristic - affect testperformance in DFA identification?

4. Experimental setup

Use **FlexFringe** (software for running heuristic and exact methods) for DFA identification experiment.

Take inspiration from the **STAMINA competition** setup:

- **BCR scores** for measuring test performance
- Datasets with different difficulties

Σ	Sparsity	
	100%	50%
2	0.99 (1)	0.95 (1)
5	0.97 (1)	0.78 (2)
10	0.93 (1)	0.64 (3)
20	0.91 (1)	0.63 (3)
50	0.81 (2)	0.64 (3)

Test the hybrid (heuristic, then exact) approach on:

STAMINA competition training sets with **5**fold cross validation



Independent variable:

- When to switch from heuristic to optimal **Dependent variable:**
- Resulting model (size, test-performance)

Implementations of the hybrid approach:

- Binary search on the size (DFA bound) of the partial automaton before we switch
- Binary search on the **SAT offset** after fixing the point at which we switch

In terms of state counts ₽̈́_ −1 -2



6. Conclusions and future work



5. Results



State count difference compared to a full EDSM run (Average from 5-fold cv)



Sparser datasets lead to **bigger** differences with EDSM (in terms of state count) Larger alphabets lead to **smaller** differences with EDSM (in terms of state counts)

Exact methods make for a bigger part of the hybrid approach => generally smaller models Hybrid approach identifies smaller DFAs, but that does not consistently improve accuracy

Generate more and different datasets Try stopping the greedy merging earlier (would require higher computational effort) Implement the entire DFASAT pipeline