

Transducers in Automata Library Mata

Author: Ing. David Chocholatý

Supervisor: doc. Mgr. Lukáš Holík Ph.D.

Brno University of Technology, Faculty of Information Technology, Brno, Czech Republic



Example of a JavaScript code vulnerable to XSS without proper string sanitization.

```
1 var text = htmlEscape(goal);
2 var action = escapeString(text);
3 element.innerHTML = '<button \
4   onclick="performAction(\
5   + action + \'\'>' + text
6   + '</button>';
1 <button onclick="performAction(
2   '&#39;);alert(1);//\'">
3   &#39;);alert(1);//\'"></button>
```

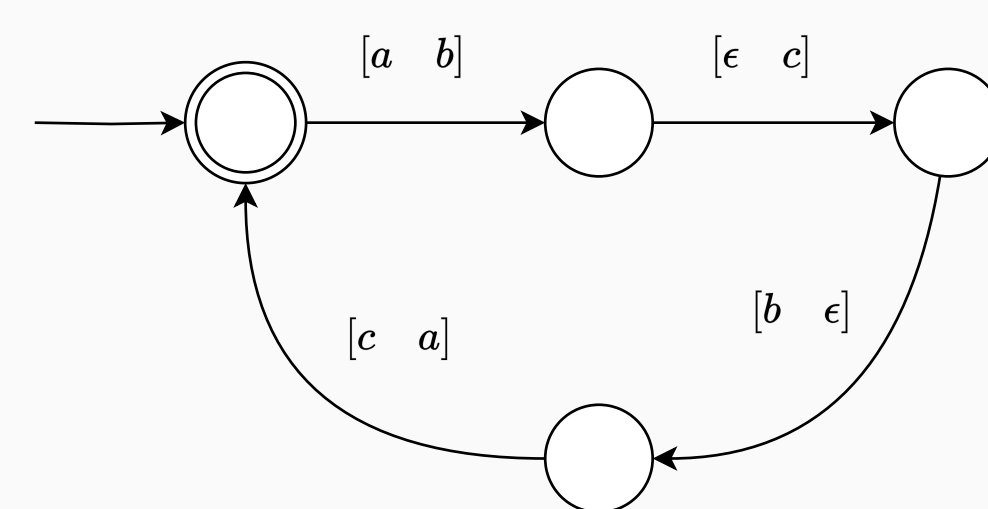
A malicious user input on a web application can lead to a XSS attack executing the malicious action alert(1).

Motivation

- Transducers are very useful in **SMT string solving** to, e.g. **verify security properties of web applications**.
- Transducers can be used to model string functions, such as `replace` and `ReplaceAll()` used to describe string constraints for **secure handling of string inputs** given by the user of web applications.
- String solving **prevents attacks** such as **cross-site scripting (XSS)** by using transducers to model *string sanitization* or *implicit browser transductions*.

Finite Transducers

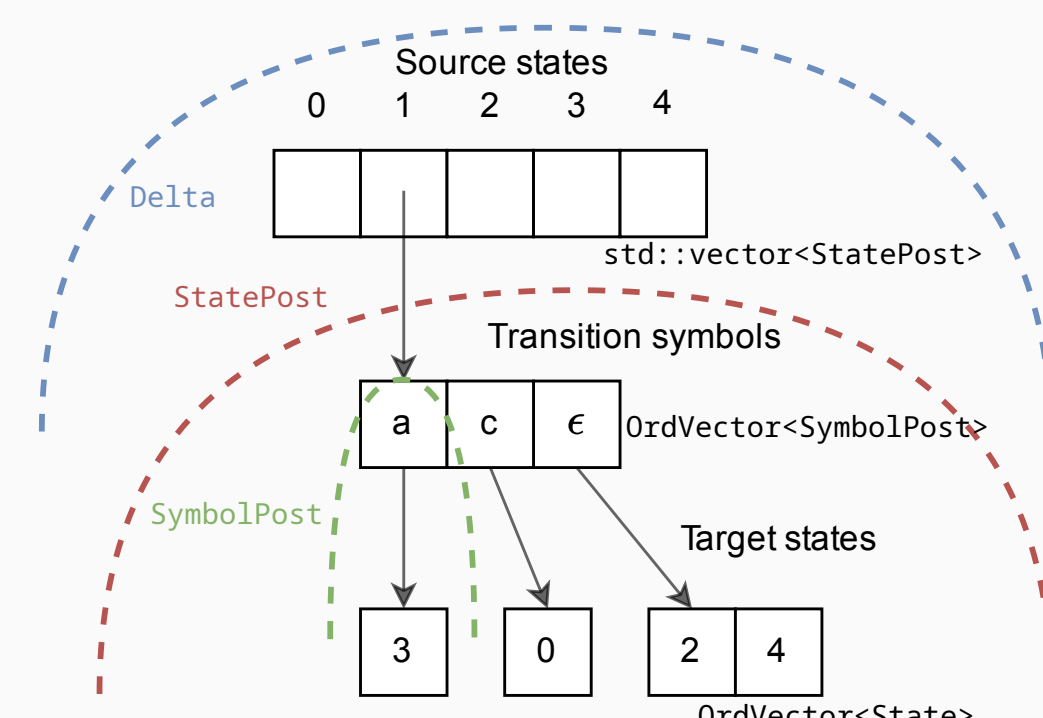
- (Non)Deterministic Finite Transducers (NFTs) are **finite state machines with n memory tapes**.
- Each run of a finite transducer represents an n -tuple of words. While a finite automaton models a regular language, finite transducers model rational relations, subsets of the Cartesian product of regular languages. That is, a finite transducer models a set of n -tuples of words.
- 2-tape NFT tapes are usually called an *input* and an *output* tape.
- 5-tuple $\mathcal{T} = (Q, \Gamma, post, I, F)$ where
 - $\Gamma = (\Sigma_\epsilon)^n$ is an n -tape alphabet of \mathcal{T} ,
 - $post : Q \times \Gamma \rightarrow 2^Q$ is a *symbol-post function*



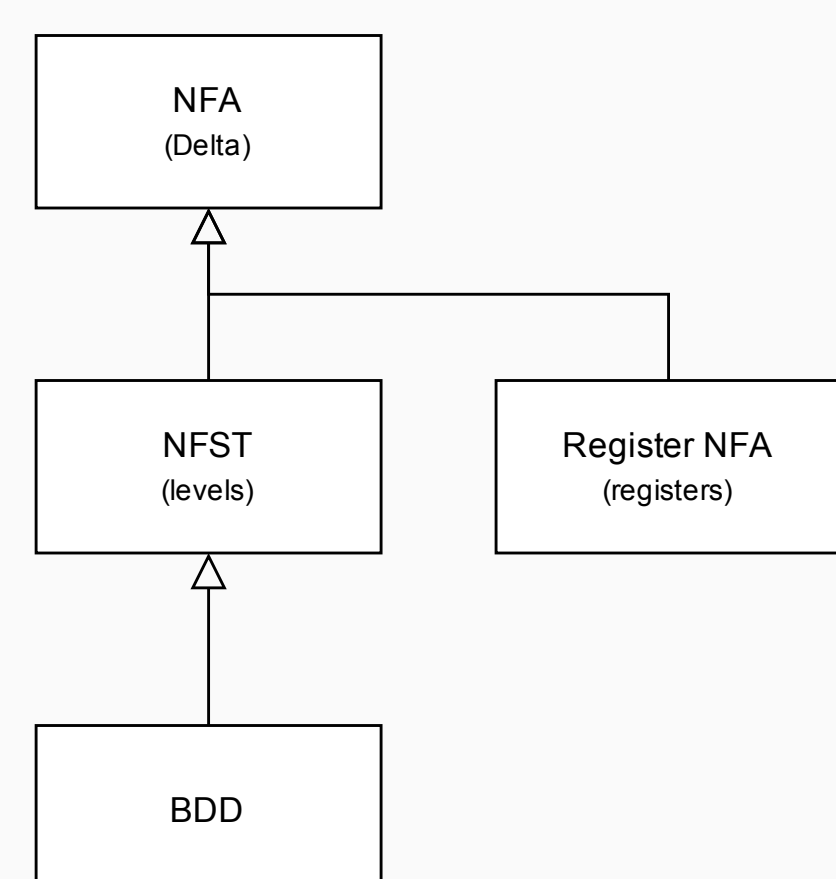
- Word *abcabc* is translated to *bcaba*.

Automata Library Mata: A Fast and Simple Finite Automata Library

- Mata is a C++ library for **manipulating finite automata**.
- Mata is designed to be **simple** to use for researchers and students, yet provide an **efficient** library for automata operations.
- The efficiency of Mata is achieved by using algorithms utilizing the features of the underlying **data structures**. E.g. the 3-level Delta representing the transition relation.
- The goal for NFTs is to **reuse the existing automata operations and data structures** in Mata with minimal modifications, namely to reuse Delta.



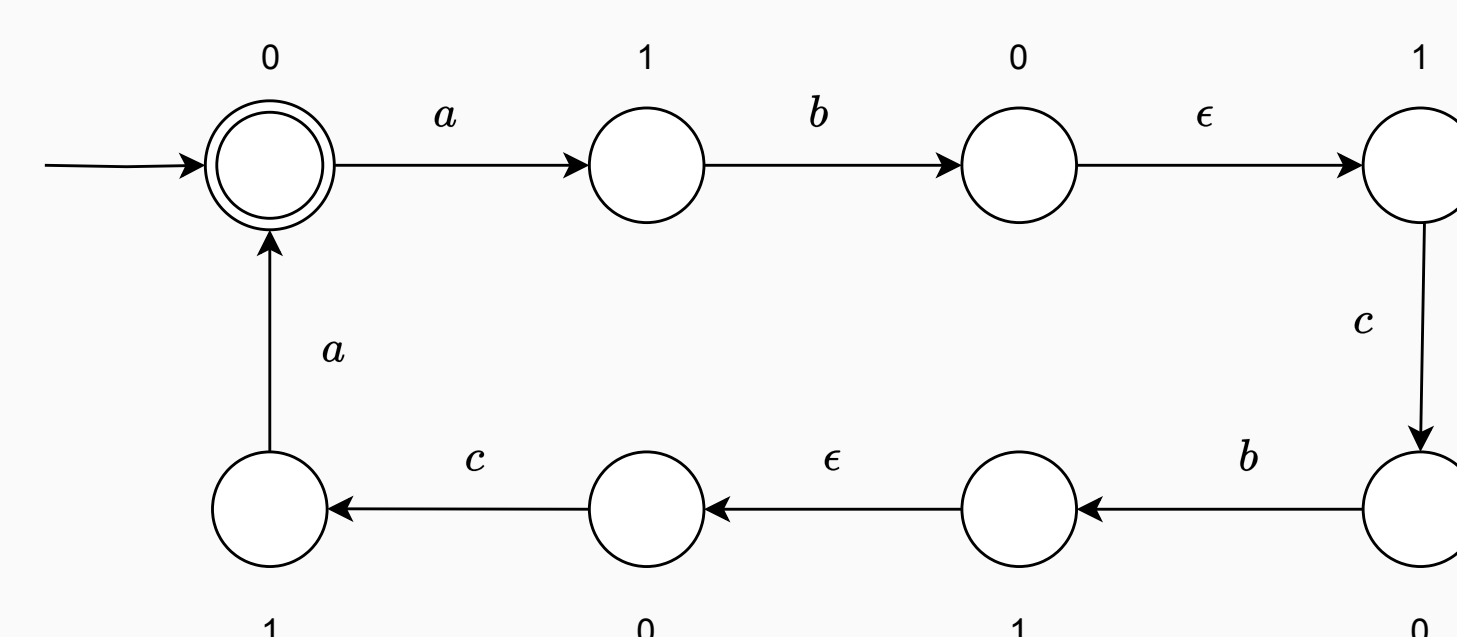
Inheritance hierarchy of automata models in Mata.



Idea

- We have implemented the representation for NFT in Mata as a **simple finite automaton with states annotated by levels where each level corresponds to a single tape in the NFT**.
- **Transition unfolding**: NFT transition is a sequence of NFA transitions.
- **Nft** inherits from the base class `Nfa` for the nondeterministic finite automata.
- This ensures most of the existing operations can be **reused and only the operations specific to NFTs need to be implemented**.
- This also gives rise to an implementation for **automata with BDD-like transitions** represented as an NFT with different interpretation for each NFT transition.

Key Contribution: Design and Implementation of Finite Transducers in Mata

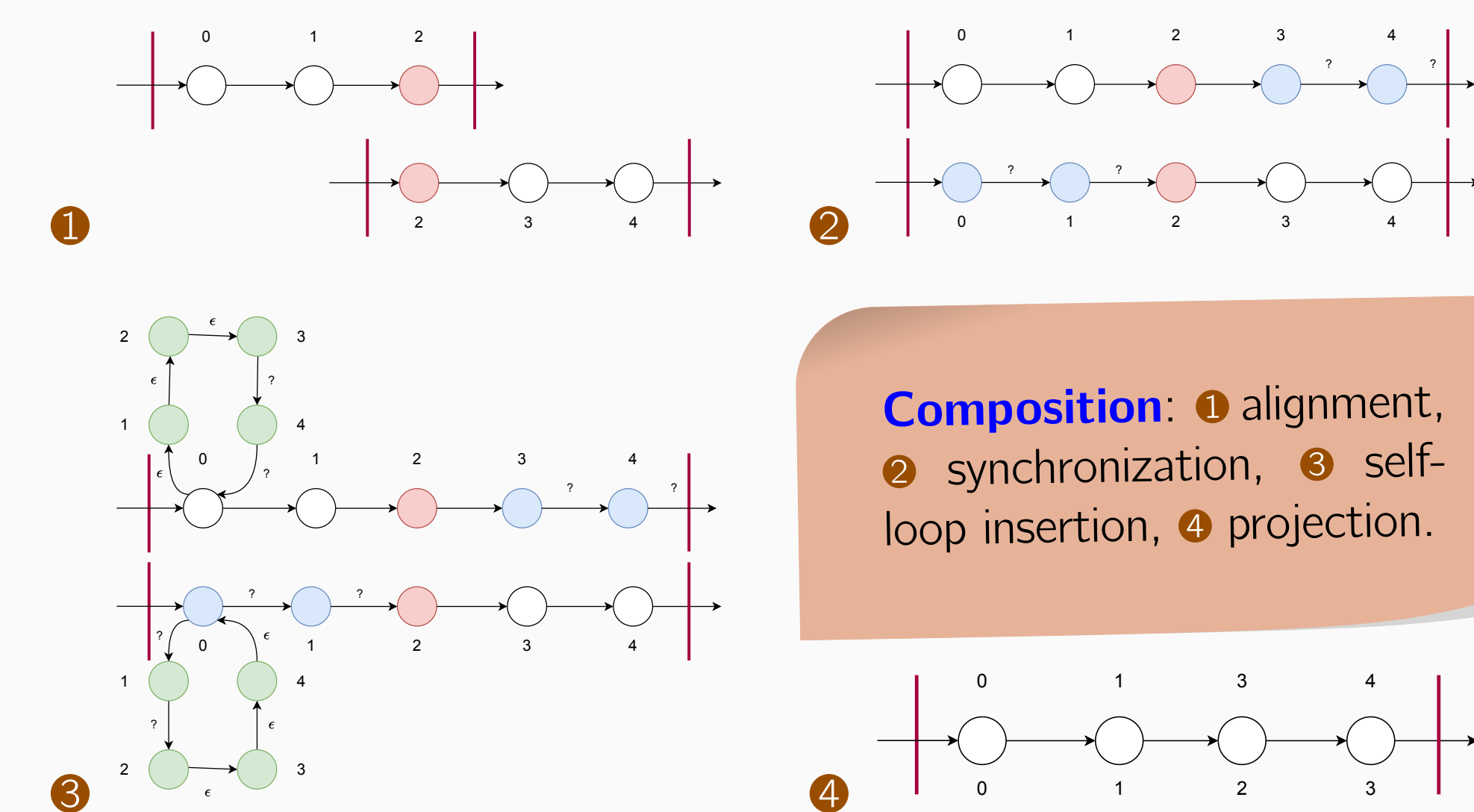


- Representation of the same NFT as earlier in Mata.
- **word unfolding**: The accepted word is an unfolded sequence of words from the NFT.

$$\bar{w} = (a_1^1, a_1^2, \dots, a_1^n) \circ (a_2^1, a_2^2, \dots, a_2^n) \circ \dots \circ (a_m^1, a_m^2, \dots, a_m^n)$$

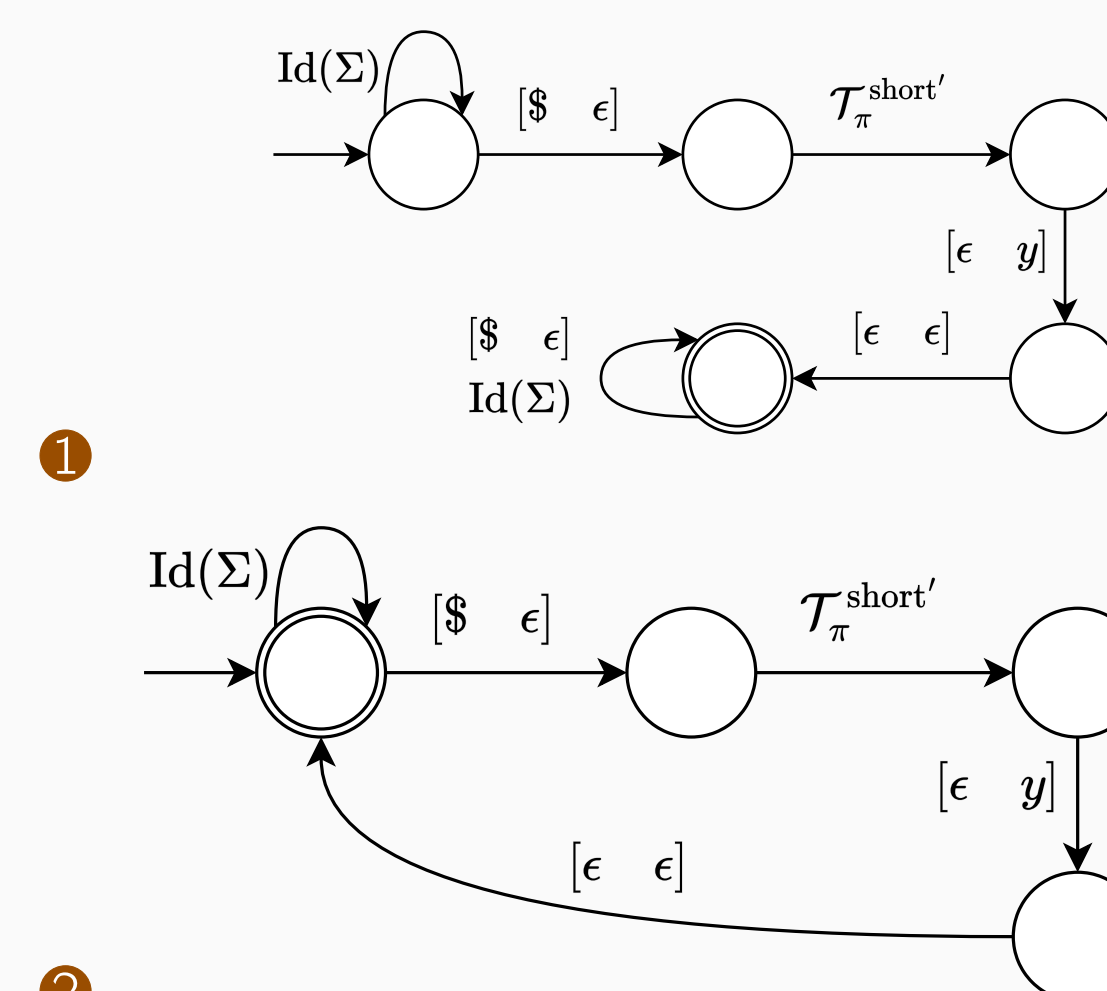
$$\bar{w} = a_1^1 \cdot a_1^2 \cdot \dots \cdot a_1^n \cdot a_2^1 \cdot a_2^2 \cdot \dots \cdot a_2^n \cdot \dots \cdot a_m^1 \cdot a_m^2 \cdot \dots \cdot a_m^n$$
- States with levels 0 represent the beginning of an NFT transition.

- NFTs support **epsilon transitions** and **don't care symbols** on the tapes as special symbols.
- These special symbols can be composed into **long jumps over several tapes**, but only up to the end of the current NFT transition.
- NFTs support operations such as: **synchronization on tapes**, **application** on a word or a language, **composition**, and various **projections**.
- Thanks to our approach, the NFT operations can be implemented as a simple combination of existing NFA operations.



Composition: ① alignment, ② synchronization, ③ self-loop insertion, ④ projection.

Support for String Solving



- NFTs can be used to model string functions such as `replace` and `ReplaceAll()`.
- Such replace operations can model **string sanitization** or **implicit browser transductions**.
- NFTs in Mata support replacement operation `replace(word, regex, replacement)` where
 - ① **single** or ② **all** occurrence(s) can be replaced, and
 - the string being replaced can be matched as a **literal** or a **regex**.

Experimental Evaluation

- Benchmarks from **SMT-LIB** (`QF_S/{webapp, transducer-plus}`) and **pattern matching** (`SymbolFromEnd`).
- Comparison with automata library `Mona`.
- NFTs in Mata have support for **nondeterminism**.

• **SymbolFromEnd**: $\# \cdot * a \cdot \{i\} \$ \rightarrow \epsilon, i \in \{1, \dots, 150\}$

method	TOs	min	max	mean	q(0.25)	median	q(0.75)	std. dev
mata	0	0.97	8070.18	2397.81	461.98	1743.10	4035.48	2233.84
mona	272	0.29	13095.06	1614.42	2.39	39.25	970.33	3485.11

method	TOs	min	max	mean	q(0.25)	median	q(0.75)	std. dev
mata	0	1.79	7900.85	2211.33	446.27	1623.44	3680.46	1983.68
mona	232	1.04	12664.70	2987.85	79.46	1304.86	4614.69	3791.32

Times are in milliseconds.

Thesis



vut.cz/en/students/final-thesis?zp_id=155757

