

## A Formalisation of Finite Automata in Isabelle / HOL

Thomas Tuerk

29th February 2012

## Computer Aided Verification of Automata (CAVA)

- CAVA is a DFG funded project
- goals:
  - formalise algorithmic parts of automata theory
  - formalise applications in model checking
  - high level of abstraction
  - generation of reference implementations
- members
  - Javier Esparza / René Neumann
  - Tobias Nipkow / Thomas Tuerk / (Peter Lammich)
  - Jan-Georg Smaus / Alexander Schimpf
- in this talk: finite automata on finite words

## Outline

- 1 Motivation
- 2 Finite Automata
- 3 Presburger Arithmetic
- 4 Hopcroft's Algorithm
- 5 Conclusion

## Nondeterministic Finite Automata

- **nondeterministic finite automata** (NFAs) are quintuples  $(Q, \Sigma, \Delta, \mathcal{I}, \mathcal{F})$ , where
  - $Q :: \text{state set}$  set of states
  - $\Sigma :: \text{label set}$  set of labels
  - $\Delta :: (\text{state} \times \text{label} \times \text{state}) \text{ set}$  transition relation
  - $\mathcal{I} :: \text{state set}$  set of initial states
  - $\mathcal{F} :: \text{state set}$  set of accepting states
- additionally **well-formedness** conditions
  - finite  $Q$
  - finite  $\Sigma$
  - $\Delta \subseteq Q \times \Sigma \times Q$
  - $\mathcal{I} \subseteq Q$
  - $\mathcal{F} \subseteq Q$

## Deterministic Finite Automata

- **deterministic finite automata** (DFAs) are NFAs with
  - $|\mathcal{I}| = 1$
  - $\forall q \in \mathcal{Q}, a \in \Sigma. \exists! q' \in \mathcal{Q}. (q, a, q') \in \Delta$
- Notice: in this formalisation all DFAs are **complete**
- $\Delta$  can be represented as a function  
 $\delta :: \text{states} \times \text{labels} \rightarrow \text{states}$

## What is there?

- Operations
  - acceptance / language of automaton
  - Boolean operations
  - determinisation
  - minimisation
  - ...
- Concepts
  - isomorphy
  - initially connected automata
  - efficient constructions
  - reversal / left and right languages
  - ...
- Infrastructure
  - construction from lists / destruction
  - interface in SML and OCaml
  - connection with Graphviz
  - ...

## Isabelle/HOL Representation

- high-level datastructure for NFAs
  - sets and relations
  - record structure
  - well-formedness predicates
  - locales for NFAs and DFAs
- low-level, executable representation
  - *Isabelle Collection Framework* (ICF) by Peter Lammich
  - set and finite map implementations
  - e. g. red-black trees, arrays, lists

## Example: Product Automata

### Product Automata

$$\begin{aligned}
 \text{product\_NFA } \mathcal{A}_1 \mathcal{A}_2 = (& \\
 & \mathcal{Q} = \mathcal{Q}(\mathcal{A}_1) \times \mathcal{Q}(\mathcal{A}_2), \\
 & \Sigma = \Sigma(\mathcal{A}_1) \cap \Sigma(\mathcal{A}_2), \\
 & \Delta = \{((q_1, q_2), a, (q'_1, q'_2)) \mid (q_1, a, q'_1) \in \Delta(\mathcal{A}_1) \wedge \\
 & \hspace{10em} (q_2, a, q'_2) \in \Delta(\mathcal{A}_2)\}, \\
 & \mathcal{I} = \mathcal{I}(\mathcal{A}_1) \times \mathcal{I}(\mathcal{A}_2), \\
 & \mathcal{F} = \mathcal{F}(\mathcal{A}_1) \times \mathcal{F}(\mathcal{A}_2) \\
 & \left. \right)
 \end{aligned}$$

- clear definition / simple proofs
- **Problem: type of states changes, composition tricky**  
 $\text{product\_NFA} :: ('q_1, 'a) \text{NFA} \Rightarrow ('q_2, 'a) \text{NFA} \Rightarrow ('q_1 \times 'q_2, 'a) \text{NFA}$
- **Problem: too many states**

## Renaming States / Isomorphic Automata

- solving the state-type problem requires renaming
- renaming leads to equivalent automaton, provided only equivalent states are identified
- two corner cases
  - injective functions lead to **isomorphic** automaton
  - identifying all equivalent states useful for minimisation
- isomorphy is a key concept of the library
- implementations operate modulo isomorphy
- state set is finite
- type-class **NFA\_states** guarantees infinite universe
- therefore, safe renaming into **NFA\_states** possible

## Unreachable states

- **unreachable** states are states that cannot be reached from an initial state
- removing unreachable states leads to an equivalent automaton
- such an automaton is called **initially connected**
- not considering unreachable states is essential for performance

## Reconsider Product Automata

- *product\_NFA*  $\mathcal{A}_1 \mathcal{A}_2$ 
  - original definition
  - intuitive definition / straightforward proofs
- *efficient\_product\_NFA*  $\mathcal{A}_1 \mathcal{A}_2 :=$   
*NFA\_remove\_unreachable\_states*(*product\_NFA*  $\mathcal{A}_1 \mathcal{A}_2$ )
  - consider only reachable states
  - essential for efficiency
- *NFA\_product*  $\mathcal{A}_1 \mathcal{A}_2 :=$   
*NFA\_normalise\_states*(*efficient\_product\_NFA*  $\mathcal{A}_1 \mathcal{A}_2$ )
  - allows composing
- locales provide implementation interface for *NFA\_product* and *efficient\_product\_NFA* modulo isomorphy

## Efficient construction

```

NFA_construct_reachable ( $I, \Sigma, FP, \Delta$ ) =
  let  $rm$  be an injective function on  $I$ ;
  initialise  $\mathcal{A}$  with  $(\emptyset, \Sigma, \emptyset, \text{image } rm \ I, \emptyset)$  and  $ws$  with  $I$ ;
  while  $ws \neq \emptyset$  do
    choose  $q \in ws$ ;
    if  $rm \ q \notin Q(\mathcal{A})$  then
       $N := \{q' \mid \exists a \in \Sigma. (q, a, q') \in \Delta\}$ ;
      extend  $rm$  to be still injective and defined on  $N$ ;
       $\mathcal{A} := ( \{rm \ q\} \cup Q(\mathcal{A}), \Sigma,$ 
         $\Delta(\mathcal{A}) \cup \{(rm \ q, a, rm \ q') \mid a \in \Sigma \wedge (q, a, q') \in \Delta\},$ 
         $\mathcal{I}(\mathcal{A}), \text{ if } FP \ q \text{ then } \{rm \ q\} \cup \mathcal{F}(\mathcal{A}) \text{ else } \mathcal{F}(\mathcal{A}) );$ 
       $ws := N \cup ws;$ 
    end
  end
  end
  return  $\mathcal{A}$ 

```

- this constructs only reachable states and renames them
- it provides a simple interface

well-formed  $\mathcal{A} \wedge (\forall q \in Q(\mathcal{A}). FP \ q \Leftrightarrow q \in \mathcal{F}(\mathcal{A})) \implies$   
*NFA\_construct\_reachable*( $I(\mathcal{A}), \Sigma(\mathcal{A}), FP, \Delta(\mathcal{A})$ ) is isomorphic to  
*NFA\_remove\_unreachable\_states*( $\mathcal{A}$ )

### Summary

- I implemented a automaton library in Isabelle/HOL
- Boolean operations, determinisation, minimisation, ...
- abstract definitions and lemmata
- efficient implementations and simple interface available

### Future Work

- automata with  $\epsilon$ -transitions
- connection to regular expressions

# Questions?

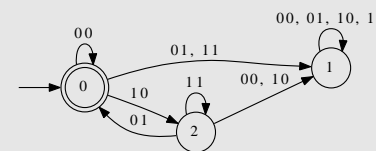
- **automata for diophantine (in)equations**
- Boolean combinations
- **existential quantification on automata**
  - projection
  - determinisation
  - (minimisation)
  - **right quotient**

- there is a Presburger Arithmetic Library by Stefan Berghofer and Markus Reiter
- available in the *Archive of Formal Proofs*
- presented at TPHOLs 2009
- they implement an automata based decision procedure
- they formalised automata on bitstrings
  - efficiently executable
  - tailored towards the application
  - little abstraction

### Case Study

Implement automata using my library.

### Automaton for $2x - y = 0$



### Example

$x = 3 = 110000\dots$   
 $y = 6 = 011000\dots$

- about 500 lines in Presburger Library
- mixed with low level definitions
- about 350 (+70) lines on my abstract level
- straightforward mapping

- existential quantification combines
  - projection (remove head of label)
  - right quotient
  - determinisation
  - (minimisation)
- implementation right quotient straightforward (about 40 lines)
- with interfaces about 200 lines
- correctness proof tricky (about 270 lines)
- about 200 lines for right quantification in original work

- straightforward
- but lengthy
- about 100 lines in my and the original work

<i>Example</i>	<i>Berghofer/Reiter</i>	<i>Tuerk</i>
stamp	0.01 s	0.01 s
example	0.10 s	0.01 s
example2	0.13 s	0.03 s
harrison1	0.03 s	0.01 s
harrison2	0.15 s	0.02 s
stamp-false	0.01 s	0.01 s
example2-false	0.46 s	0.04 s

speed of generated OCaml code on an Intel Core i7 2720QM

- straightforward to use my automata library
- about 2 days of work involved
- higher level of abstraction
- more efficient execution
- about 1200 lines in my adaption
- about 4500 lines in original library

# Questions?

## Motivation

- minimisation is an important operation
- Hopcroft's Algorithm has the best known behaviour in practise
- implemented using Peter Lammich's refinement framework
- good case study for refinement
- here: focus on refinement

## General Idea

- Hopcroft's algorithm minimises initially connected DFAs
- it computes the Myhill-Nerode equivalence relation

$$\{\{q' \mid q' \text{ is equivalent to } q\} \mid q \in Q\}$$

- this equivalence relation is used to rename states

**General Idea**

```

initialise P with part_F;
while there are C in P and (a, C_s) in Sigma x P with splittable_A(C, (a, C_s)) do
  choose such C and (a, C_s);
  update P by removing C and adding the two results of split_A(C, (a, C_s));
end
return P;

split_A(C, (a, C_s)) := ({q | q in C and delta(q, a) in C_s}, {q | q in C and delta(q, a) not in C_s})
splittable_A(C, (a, C_s)) := let (C_t, C_f) = split_A(C, (a, C_s)) in C_t not empty and C_f not empty
part_F := {F, Q - F} - {empty}
  
```

## Abstract Algorithm

Idea: maintain explicit list of splitters

**Abstract Algorithm**

```

Hopcroft_step_abstract(A, a, C_s, P, L) =
  spec (P', L'). P' = Split_A(P, (a, C_s)) ^ splitter_PA(P, (a, C_s), L, L');

Hopcroft_abstract(A) =
  if (Q = empty) then return empty else
  if (F = empty) then return {Q} else
  while_Hopcroft_abstract_invar (lambda(P, L). L not empty) (lambda(P, L). {
    (a, C_s) <- spec x. x in L;
    (P', L') <- Hopcroft_step_abstract(A, a, C_s, P, L);
    return (P', L');
  }) (part_F, {(a, F) | a in Sigma})
  
```

The proof of correctness needs about 2600 lines.

## Set Implementation

Idea: implement the step with a foreach loop

**Set Implementation**

```

Hopcroft_step_set(A, a, C_s, P, L) =
  P' <- spec P'. P' subset P ^ (forall C in P. splittable_A(C, (a, C_s)) => C in P');
  (P', L') <- foreach_Hopcroft_set_invar P' (lambda(C, L'). {
    let (C_t, C_f) = split_A(C, (a, C_s));
    if (C_t = empty v C_f = empty) then return (P', L') else {
      (C_1, C_2) <- spec x. x in {(C_t, C_t), (C_t, C_f)};
      let P' = (P' - {C}) union {C_1, C_2};
      let L' = (L' - {(a, C) | a in Sigma}) union {(a, C_1) | a in Sigma} union {(a, C_2) | (a, C) in L'};
      return (P', L');
    }
  }) (P, L - {(a, C_s)});
  return (P', L');
  
```

Refinement needs about 400 lines.

This is the level of abstraction usually used in literature.

## Precomputing Predecessors

### Precomputing Predecessors

```

Hopcroft_step_pre(A, a, Cs, P, L) =
  let pre = {q | δ(q, a) ∈ Cs};
  P' ← spec P'. P' ⊆ P ∧ (∀C ∈ P'. C ∩ pre ≠ ∅) ∧
    (∀C ∈ P. C ∩ pre ≠ ∅ ∧ splittableA(C, (a, Cs) ⇒ C ∈ P');
  (P', L') ← foreachHopcroft_set_invar P' (λC (P', L')). {
    let (Ct, Cf) = ({q | q ∈ C ∧ q ∈ pre}, {q | q ∈ C ∧ q ∉ pre});
    if (Cf = ∅) then return (P', L') else {
      let (Cmin, Cmax) = if (|Cf| < |Ct|) then (Cf, Ct) else (Ct, Cf);
      let P' = (P' - {C}) ∪ {Cmin, Cmax};
      let L' = (L' - {(a, C) | a ∈ Σ}) ∪
        {(a, Cmin | a ∈ Σ} ∪ {(a, Cmax | (a, C) ∈ L'};
      return (P', L');
    }
  } (P, L - {(a, Cs)});
return (P', L');

```

Refinement needs about 50 lines.

## Data Refinement Partitions

- representing partitions as sets of sets is inefficient
- use pair of finite maps instead
  - im* - finite map from index to classes
  - sm* - finite map from states to index of class it belongs to
  - n* - maximal index

### Abstraction Function and Invariant

```

partition_map_invar(im, sm, n) :=
  dom(im) = {i | 0 ≤ i < n} ∧ (∀0 ≤ i < n. im(i) ≠ ∅) ∧
  (∀q. sm(q) = i ⇔ (0 ≤ i < n ∧ q ∈ im(i)))
partition_map_α(im, sm, n) := {im(i) | 0 ≤ i < n}

```

- when splitting classes, the old index is reused for  $C_{max}$
- update of  $L$  becomes much simpler

## Data Refinement Partitions II

### Data Refinement Partitions

```

Hopcroft_step_map(A, a, is, (im, sm, n), L) =
  let pre = {q | δ(q, a) ∈ im(is)};
  let I = {sm(q) | q ∈ pre};
  ((im', sm', n'), L') ← foreachHopcroft_map_invar I (λi ((im', sm', n'), L')). {
    let (Ct, Cf) = ({q | q ∈ im(i) ∧ q ∈ pre}, {q | q ∈ im(i) ∧ q ∉ pre});
    if (Cf = ∅) then return ((im', sm', n'), L') else {
      let (Cmin, Cmax) = if (|Cf| < |Ct|) then (Cf, Ct) else (Ct, Cf);
      let (im', sm', n') = ( im'(i ↦ Cmax, n ↦ Cmin),
        λq. if q ∈ Cmin then n else sm'(q), n' + 1);
      let L' = {(a, n) | a ∈ Σ} ∪ L';
      return ((im', sm', n'), L');
    }
  } ((im, sm, n), L - {(a, is)});
return ((im', sm', n'), L');

```

Refinement needs about 1100 lines.

## Data Refinement Classes

- represent classes by interval of natural numbers
- use pair of finite maps for mapping between classes and natural numbers
  - pm* - finite map from states to natural numbers
  - pim* - inverse map

### Abstraction Function and Invariant

```

class_map_invar Q (pm, pim) :=
  dom(pm) = Q ∧ dom(pim) = {i | i < |Q|} ∧
  (∀q i. pm(q) = i ⇔ pim(i) = q)
class_map_α(pm, pim)(l, u) := {pim(i) | l ≤ i ≤ u}

```

## Data Refinement Classes II

### Data Refinement Classes

```

Hopcroft_step_map2( $\mathcal{A}, a, i_s, (im, sm, n), L, pm, pim$ ) =
  let  $pre = \{q \mid \delta(q, a) \in im(i_s)\};$ 
  let  $(iM, (pm, pim)) =$ 
    Hopcroft_step_map2_compute_iM( $\mathcal{A}, (pm, pim), pre, (im, sm, n)$ );
   $((im', sm', n'), L') \leftarrow$  foreachHopcroft_map2_invar  $\{(i, s) \mid iM(i) = s\}$ 
   $(\lambda(i, s) ((im', sm', n'), L')). \{$ 
    let  $(l, u) = im(i);$ 
    let  $(C_t, C_f) = ((l, s - 1), (s, u));$ 
    if  $((u + 1) - s = 0)$  then return  $((im', sm', n'), L')$  else {
      let  $(C_{min}, C_{max}) =$  if  $((u + 1) - s < s - l)$  then  $(C_f, C_t)$  else  $(C_t, C_f);$ 
      let  $(im', sm', n') =$ 
         $(im'(i \mapsto C_{max}, n \mapsto C_{min}),$ 
           $\lambda q. \text{if } q \in \text{class\_map\_}\alpha(pm, pim)(C_{min}) \text{ then } n \text{ else } sm'(q), n' + 1);$ 
      let  $L' = \{(a, n) \mid a \in \Sigma\} \cup L';$ 
      return  $((im', sm', n'), L');$ 
    }
  }  $((im, sm, n), L - \{(a, i_s)\});$ 
return  $((im', sm', n'), L'), (pm, pim));$ 

```

## Collection Framework

- the next refinement step introduces executable datastructures
- Peter Lammich's **Collection Framework** used
- red-black-trees and arrays chosen as data-structures
- 800 lines needed
- further 380 lines for connection with automata implementation (computing *pre*)

## Data Refinement Classes III

### Data Refinement Classes

```

Hopcroft_step_map2_compute_iM( $\mathcal{A}, (pm, pim), pre, (im, sm, n)$ ) =
  foreach  $pre (\lambda q (iM, (pm, pim))). \{$ 
    let  $i = im(q);$ 
    let  $s =$  if defined  $iM(q)$  then  $iM(q)$  else  $\text{fst}(im(i));$ 
    let  $iq = pm(q);$ 
    if  $iq = s$  then
      return  $(iM(i \mapsto s + 1), (pm, pim));$ 
    else
      let  $qs = pim(s);$ 
      let  $pm = pm(q \mapsto s)(qs \mapsto iq);$ 
      let  $pim = pim(s \mapsto q)(iq \mapsto qs);$ 
      return  $(iM(i \mapsto s + 1), (pm, pim));$ 
    end
  }  $(\emptyset, (pm, pim));$ 

```

Refinement needs about 1600 lines.

## Experimental Results

No. DFAs	No. states	No. labels	Baclet/Pagetti OCaml	Lammich/Tuerk OCaml	Leiß PolyML	Leiß PolyML
10000	50	2	0.17 s	7.87 s	1.91 s	5.38 s
10000	50	5	0.28 s	15.03 s	3.67 s	19.34 s
10000	100	2	0.30 s	16.21 s	4.06 s	16.41 s
10000	100	5	0.53 s	30.79 s	7.86 s	63.21 s
10000	250	2	0.76 s	46.53 s	11.48 s	83.62 s
1000	1000	2	0.44 s	20.95 s	5.67 s	134.21 s
1000	2500	2	1.44 s	58.24 s	17.83 s	905.82 s

Experimental Results (measured on an Intel Core I7 2720QM)



Motivation ○○	Finite Automata ○○○○○○○○○○	Presburger Arithmetic ○○○○○○○	Hopcroft's Algorithm ○○○○○○○○○○●	Conclusion ○
Conclusion				

- refinement worked well for Hopcroft's algorithm
- there is a non-trivial, efficient implementation
- experimental result show decent performance

## Questions?

Motivation ○○	Finite Automata ○○○○○○○○○○	Presburger Arithmetic ○○○○○○○	Hopcroft's Algorithm ○○○○○○○○○○○○	Conclusion ●
Conclusion				

- I implemented a library for finite automata in Isabelle/HOL
- this library provides abstract concepts for reasoning about automata
- there are efficient implementations of non-trivial algorithms
- the library can easily be used by other projects
- generated code in SML and OCaml exists