# A diagrammatic axiomatisation of finite-state automata

Robin Piedeleu & Fabio Zanasi

arXiv:2009.14576

Séminaire PPS, Novembre 2020

## String diagrams for open systems



## **Compositional modelling**

• Compositional = **functorial** semantics:



## **Compositional modelling**

• Compositional = **functorial** semantics:



## **Compositional modelling**

• Compositional = **functorial** semantics:



Symmetric monoidal functor

 One step further: complete equational theory, aka axiomatisation.



## Today

- 1. Background
  - Finite-state automata
  - Regular expressions and Kleene algebra
- 2. Kleene diagrams: first attempt
  - Syntax and Semantics
  - Encoding regexes and NFA
  - Equational theory: the problem with iteration
- 3. Kleene diagrams: reprise
  - Bringing back regexes
  - Axiomatisation
  - Sketch of the completeness proof
- 4. Discussion and future work

#### Background

#### Nondeterministic finite automata

• NFA are traditionally encoded by a tuple

$$(\Sigma, Q, \delta \subseteq Q \times \Sigma \times Q, q_0 \in Q, F \subseteq Q)$$

(alphabet of basic actions, states, transition relation, initial state, accepting states)

• Example.  $(\{a, b\}, \{q_0, q_1, q_2\}, \{(q_0, a, q_1), (q_1, b, q_2), (q_2, a, q_1), more conveniently:$  $b <math>(q_2, a, q_2)\}, q_0, \{q_2\})$ 



• Recognised language: set of strings  $w = a_1 \dots a_n$  for which there exists a sequence of states  $r_0, r_2, \dots r_n$  such that  $r_0 = q_0$ ,  $(r_i, a_{i+1}, r_{i+1}) \in \delta$  and  $r_n \in F$ .





$$\begin{split} \llbracket e + f \rrbracket_R &= \llbracket e \rrbracket_R \cup \llbracket f \rrbracket_R & \llbracket 0 \rrbracket_R = \varnothing \\ \llbracket e f \rrbracket_R &= \{ vw \mid v \in \llbracket e \rrbracket_R, w \in \llbracket f \rrbracket_R \} & \llbracket 1 \rrbracket_R = \{ \varepsilon \} \\ \llbracket e^* \rrbracket_R &= \bigcup_{n \in \mathbb{N}} \llbracket e \rrbracket_R^n & \llbracket a \rrbracket_R = \{ a \} \end{split}$$

**Kleene theorem.** A language is regular if and only if it is recognised by some NFA.

## Kleene algebra

- Equational presentation of regular expressions:
  - Sum and concatenation (with their units) form an idempotent semiring
  - $e^*$  is the least-fixed point of, e.g., X = 1 + eX. But what axioms?
- **Not finitely-based**: no finite set of equations can capture all equalities in the language model [Redko, 1964]
- Finite implicational theory [Kozen, 1994]:

 $1 + ee^* \le e^*$ (star is a fixed-point) $f + ex \le x \Rightarrow e^*f \le x$ (star is the least one)

 Other axiomatisations (some infinitary): Conway, Krob, Salomaa, Kozen, Bloom, Ésik...

#### Kleene diagrams: first attempt















#### Compositionality

...means that

$$\begin{bmatrix} -c & -d & - \end{bmatrix} = \left\{ (L, K) \mid \exists M (L, M) \in \begin{bmatrix} -c & - \end{bmatrix} \right\}$$
$$(M, K) \in \begin{bmatrix} -d & - \end{bmatrix} \right\}$$

$$\begin{bmatrix} -c_1 \\ -c_2 \end{bmatrix} = \{ ((L_1, L_2), (K_1, K_2)) \\ | (L_i, K_i) \in \llbracket -c_i \end{bmatrix}, i = 1, 2 \}$$

## Sanity check: NFA

- Formal encoding from tuples definition is tedious.
- Intuition via graphical notation:



• **Theorem.** Given an NFA which recognises a language *L*, the semantics of its associated diagram, constructed as above, is

 $\big\{(K,K')\mid KL\subseteq K'\big\}$ 

### Sanity check: regexes

We can encode regexes as follows



• **Proposition**. The encoding preserves the semantics, i.e., for any expression *e*,

emantic functor
$$\llbracket \langle e \rangle \rrbracket = \{ (L, K) \mid L \llbracket e \rrbracket_R \subseteq K \}$$
Regex encodingStandard regex interpretation

Se

#### What else?

• Benefits of (de)compositionality



- Gives formal status to automata with **multiple inputs/outputs**.
- But **nc expressive**: every diagram is fully Beware! Do not necessarily coincide with initial/accepting states in the usual definition.

#### A more concrete view

A diagrammatic language to specify systems of **linear** language inequalities, i.e. for which concatenation is restricted to left-action of letters.



#### Equational theory

## Plain wires

• We have a **compact closed** category: we can bend/straighten wires at will, keeping track of only their orientation



• ... and we can eliminate isolated loops

## Copy and Sum

Cocommutative comonoid



Commutative monoid



Bimonoid







## Copy and Sum

• Idempotent



• Getting rid of trivial feedback



#### Concatenation

• Letters can be copied and deleted...



...merged and spawned



## The problem with iteration

- Recall: Kleene algebra **not finitely-based** in the standard algebraic setting. The main obstacle is iteration (represented by the star).
- Here it is a **derived** notion, made up of more primitive components:



• But the problem did not disappear.

## The problem with iteration

 Simple check: we should be able to copy/delete/merge/spawn an expression in a loop. For example,



- Incompleteness: we cannot prove this with just the current axioms.
- Even if we add it, we need to be deal with **arbitrary nestings** of loops with other operations.

### One solution

- Impose global (so infinitary) axiom schemes.
- **Definition.** A diagram is *left-to-right* if it has all inputs in its domain and all outputs in codomain.
- For any left-to-right diagram *d*, we want



• By *fiat*: similar to **matricial iteration theories** [Bloom and Ésik, 93] although, even relative to this setting, they did not produce a finitary axiomatisation for regular languages.

## Semantics of least fixed-points

- Monotone maps embed into monotone relations: f is sent to  $\{(x,y) \mid f(x) \le y\}$ .
- A relation satisfies copying and deleting,



*iff* it is the image of a monotone **map**.

- The semantics of  $e^*$  is the least fixed-point of the language map  $f = \lambda Z$ . X U eZ. This is still (the image of) a monotone map in X, *i.e.*,
  - (del) means the least fixed-point **exists** for every X;
  - (cpy) means it is **unique**.

#### Kleene diagrams: reprise

## A trick: bringing back regexes

 Extend the syntax with regular expressions on a separate wire type: \_\_\_\_\_



• Note that this is **just syntax**. Their interpretation is the free term algebra of regexes.

 $\begin{bmatrix} --- \end{bmatrix} = \{((e, e) \mid e \in \mathsf{RegExp}\} \\ \begin{bmatrix} --- \end{bmatrix} = \{(e, e^*) \mid e \in \mathsf{RegExp}\} \\ \begin{bmatrix} --- \end{bmatrix} = \{(e, e) \mid e \in \mathsf{RegExp}\} \\ \begin{bmatrix} --- \end{bmatrix} = \{(e, e) \mid e \in \mathsf{RegExp}\} \\ \begin{bmatrix} --- \end{bmatrix} = \{(e, f) \mid e, f \in \mathsf{RegExp}\} \\ \begin{bmatrix} --- \end{bmatrix} = \{(e, f) \mid e, f \in \mathsf{RegExp}\} \\ \begin{bmatrix} --- \end{bmatrix} = \{(e, f) \mid e, f \in \mathsf{RegExp}\} \\ \begin{bmatrix} --- \end{bmatrix} = \{(e, f) \mid e, f \in \mathsf{RegExp}\} \\ \begin{bmatrix} --- \end{bmatrix} = \{(e, f) \mid e, f \in \mathsf{RegExp}\} \\ \begin{bmatrix} --- \end{bmatrix} = \{(e, f) \mid e, f \in \mathsf{RegExp}\} \\ \begin{bmatrix} --- \end{bmatrix} = \{(e, f) \mid e, f \in \mathsf{RegExp}\} \\ \begin{bmatrix} --- \end{bmatrix} = \{(e, f) \mid e, f \in \mathsf{RegExp}\} \\ \begin{bmatrix} --- \end{bmatrix} = \{(e, f) \mid e, f \in \mathsf{RegExp}\} \\ \begin{bmatrix} --- \end{bmatrix} = \{(e, f) \mid e, f \in \mathsf{RegExp}\} \\ \begin{bmatrix} --- \end{bmatrix} = \{(e, f) \mid e, f \in \mathsf{RegExp}\} \\ \begin{bmatrix} --- \end{bmatrix} = \{(e, f) \mid e, f \in \mathsf{RegExp}\} \\ \begin{bmatrix} --- \end{bmatrix} = \{(e, f) \mid e, f \in \mathsf{RegExp}\} \\ \begin{bmatrix} --- \end{bmatrix} = \{(e, f) \mid e, f \in \mathsf{RegExp}\} \\ \begin{bmatrix} --- \end{bmatrix} = \{(e, f) \mid e, f \in \mathsf{RegExp}\} \\ \begin{bmatrix} --- \end{bmatrix} = \{(e, f) \mid e, f \in \mathsf{RegExp}\} \\ \begin{bmatrix} --- \end{bmatrix} = \{(e, f) \mid e, f \in \mathsf{RegExp}\} \\ \begin{bmatrix} --- \end{bmatrix} = \{(e, f) \mid e, f \in \mathsf{RegExp}\} \\ \begin{bmatrix} --- \end{bmatrix} = \{(e, f) \mid e, f \in \mathsf{RegExp}\} \\ \begin{bmatrix} --- \end{bmatrix} = \{(e, f) \mid e, f \in \mathsf{RegExp}\} \\ \begin{bmatrix} --- \end{bmatrix} = \{(e, f) \mid e, f \in \mathsf{RegExp}\} \\ \begin{bmatrix} --- \end{bmatrix} = \{(e, f) \mid e, f \in \mathsf{RegExp}\} \\ \begin{bmatrix} --- \end{bmatrix} = \{(e, f) \mid e, f \in \mathsf{RegExp}\} \\ \begin{bmatrix} --- \end{bmatrix} = \{(e, f) \mid e, f \in \mathsf{RegExp}\} \\ \begin{bmatrix} ---- \end{bmatrix} = \{(e, f) \mid e, f \in \mathsf{RegExp}\} \\ \end{bmatrix} = \{(e,$ 

## A trick: bringing back regexes

- Syntax: replace -a with general action of **any regex** (not just the letters) via
- Semantics: regex acting on languages by concatenation on the left

$$\rightarrow \rightarrow \rightarrow \parallel = \{ ((e, L), K) \mid L \llbracket e \rrbracket_R \subseteq K \text{ and } e \in \mathsf{RegExp}, L, K \subseteq \Sigma^* \}$$

Interpretation of the regex *e* (a regular language)

Free (uninterpreted) term algebra of regexes

- We recover the atomic actions as  $-a = \rightarrow \rightarrow \rightarrow \rightarrow$
- String diagrams for **generalised automata** with transitions labelled by arbitrary regexes:



## Axiomatising the action (1/2)

Capturing the behaviour of the action:

- Concatenation and empty word



#### Unfolding/compiling regexes

**Example.**  $ab(a+ab)^*$ 



## Axiomatising the action (2/2)

Back to the original problem:

- Copy and delete arbitrary regexes



- Merge and spawn arbitrary regexes



**Theorem (Completeness).** Two diagrams are equal iff they are mapped to the same monotone relation.

## Completeness proof outline

- Normal form argument: diagrammatic counterpart of constructing the minimal deterministic automaton that recognises the same language
  - An automaton is deterministic (DFA) if its transition relation is the graph of a function  $Q \times \Sigma \to Q$ .
  - Among the finite-state automata that recognise a given language, there is a unique DFA with the smallest number of states. This is our normal form.
- Obtained via **Brzozowski's algorithm**, implemented as equational reasoning:



## Completeness proof outline

- Normal form argument: diagrammatic counterpart of constructing the minimal deterministic automaton that recognises the same language
  - An automaton is deterministic (DFA) if its transition relation is the graph of a function  $Q \times \Sigma \to Q$ .
  - Among the finite-state automata that recognise a given language, there is a unique DFA with the smallest number of states. This is our normal form.
- Obtained via **Brzozowski's algorithm**, implemented as equational reasoning:



#### Determinisation, traditionally

For an NFA given by the tuple  $(\Sigma, Q, \delta, q_0, F)$ 

an equivalent (i.e. that recognises the same language) DFA is given by  $(\Sigma, 2^Q, \Delta, \{q_0\}, G)$ 

where

$$(S,a,S') \in \Delta \text{ iff } S' = \{q' \in Q \mid \exists q \in S, (q,a,q') \in \delta\}$$

and G is the set of subsets of Q that contain at least one accepting state.



## Determinisation, diagrammatically

• **Nondeterministic** transitions of automata correspond to subdiagrams of the form

$$\rightarrow \underbrace{e}_{e} \quad \text{(or } \rightarrow \underbrace{e}_{e} \quad \text{where } -e := \underbrace{e}_{e} \quad \text{)}$$

• Useless states (those that cannot reach an accepting state/ contribute to the semantics) correspond to subdiagrams of the form

• To get rid of them, just apply (not haphazardly, check the paper for details):



## Diagrammatic determinisation example



## Left-to-right diagrams again

• Now we can prove that, for any left-to-right diagram d



- Subcategory of left-to-right diagrams maps to a category of matrices over the semiring of regular languages, with matrix product as composition and direct sum as product.
- Two uses: 1) reduces the completeness proof to diagrams with one input and one output; 2) is the engine of the diagrammatic determinisation procedure.

## Diagrammatic determinisation example



## Bonus: context-free languages

- Recall that we designed a language to specify systems of **linear** language inequations.
- Remove the linearity constraint: unconstrained concatenation gives systems of polynomial language inequations.
- Diagrammatically, turn —— into  $\longrightarrow$  and  $\longrightarrow$  into  $\xrightarrow{}$  into  $\xrightarrow{}$  with

$$\begin{bmatrix} \rightarrow \\ \rightarrow \end{pmatrix} = \{ ((M,L),K) \mid ML \subseteq K \}$$

• We can specify context-free languages. For example, the language of properly matched parentheses:



## Discussion

- What's new? A finite presentation of a symmetric monoidal category (SMC) that axiomatises automata equivalence.
- In what sense is it finite? Debatable: not in the usual sense of algebraic theories, but relative to the equational theory of SMCs.
  - If we encode terms using only ; and  $\otimes$ , it is infinite.
  - But we should encode them as graphs (and equations as graph rewrites).
- Is it really new? All previous work was in a traced symmetric monoidal setting (*iteration theories* of Bloom & Ésik or *network algebra* of Stefanescu). But:
  - Still no finite axiomatisations of regular languages in these settings.
  - The trace is a global operation that cannot be finitely axiomatised relative to the theory of symmetric monoidal categories.

## Discussion

- Why does this work? Slightly mysterious, perhaps better compositionality.
  - Not the first time that a finite axiomatisation of a theory that is provably not finitely-based in the standard algebraic setting: graphical conjunctive queries vs. allegories, for example.
  - Proofs of negative results in the algebraic setting rely on showing the correspondence between terms and certain graphs. The two-dimensional syntax allows to represent all graphs/automata.

### Future work

- This category of monotone relations over languages has very rich hidden structure:
  - Cartesian **bicategory** with inclusions of relations as 2-cells;
  - adjoints (in the bicategorical sense) to copy and sum that represent the Boolean lattice structure of languages (not just the order).
- Using this more expressive language, a generalisation of **bisimulation** can be defined and is sufficient to prove that two diagrams corresponding to equivalent automata are equal.
- But **not yet a complete equational theory** for the whole extended syntax.
- This seems to correspond to **alternating automata**.

#### Questions?