

Towards a Coalgebraic Chomsky Hierarchy

(Short paper)

Sergey Goncharov¹ and Stefan Milius¹ and Alexandra Silva²

¹ Friedrich-Alexander-Universität Erlangen-Nürnberg

² Radboud University Nijmegen

The Chomsky hierarchy plays a prominent role in the foundations of theoretical computer science relating classes of formal languages of primary importance as well as the grammars and machine models that represent those language classes. It is well-known that regular languages can be captured by coalgebras. They form precisely the image of all *finite* coalgebras for the set functor $X \mapsto 2 \times X^A$ under the unique homomorphisms into the final coalgebra, which is carried by all formal languages in this case. There have also been attempts to model context-free grammars (in Greibach normal form) coalgebraically [2, 5] and hence, to capture context-free languages. But what about other language classes and other machine models in the Chomsky Hierarchy?

Here we bring together several ideas in order to treat machine models that are automata extended by extra storage; in particular, we show how to model machines with a stack or a Turing tape coalgebraically, and so we obtain the first coalgebraic account of recursively enumerable languages. Our central notion is that of a \mathbb{T} -automaton, i.e. a finite coalgebra $m : X \rightarrow B \times (TX)^A$, where A is a set of inputs, $\mathbb{T} = (T, \eta, \mu)$ is a finitary monad, which we think as modelling extra storage, and B is an (Eilenberg-Moore) \mathbb{T} -algebra $a^m : TB \rightarrow B$ of outputs. The semantics of a \mathbb{T} -automaton can then be defined via the *generalized power-set construction* [4]: notice first that the functor $LX = B \times X^A$ (being essentially given by product) lifts to the category of \mathbb{T} -algebras; then one can extend the coalgebra structure m to the free T -algebra on X to obtain an L -coalgebra $m^\sharp : TX \rightarrow B \times (TX)^A$. Observe that the final L -coalgebra is carried by the set B^{A^*} of formal power-series on B . Finally, one defines the semantic map as $\llbracket - \rrbracket_m = h \cdot \eta_X : X \rightarrow B^{A^*}$, where $h : TX \rightarrow B^{A^*}$ is the unique L -coalgebra homomorphism from (TX, m^\sharp) to the final L -coalgebra.

In order to model machines with pushdown storage we consider the *stack monad*: for a finite set Γ of stack symbols this is the submonad \mathbb{T} of the store monad $(- \times \Gamma^*)^{\Gamma^*}$ for which elements $\langle r, t \rangle$ of TX satisfy the following restriction: there exists a k (depending on r, t) such that for every $w \in \Gamma^k$ and $u \in \Gamma^*$ we have $r(wu) = r$ and $t(wu) = t(w)u$. Goncharov [1, Proposition 5] has shown that the stack monad can be presented by operations and equations by the *stack theory* w.r.t. $\Gamma = \{\gamma_1, \dots, \gamma_n\}$: this theory consists of the operations $\text{pop} : n + 1 \rightarrow 1$ and $\text{push}_i : 1 \rightarrow 1$ for $i = 1, \dots, n$. These operations are subject to the following axioms:

$$\begin{aligned} \text{push}_i(\text{pop}(x_1, \dots, x_n, y)) &= x_i, & \text{pop}(\text{push}_1(x), \dots, \text{push}_n(x), x) &= x, \\ \text{pop}(x_1, \dots, x_n, \text{pop}(y_1, \dots, y_n, z)) &= \text{pop}(x_1, \dots, x_n, z). \end{aligned}$$

A *stack \mathbb{T} -automaton* is a \mathbb{T} -automaton for the stack monad for which $B \subseteq 2^{\Gamma^*}$ consists of precisely those predicates p for which there exists a k such that $p(wu) = p(w)$ for all w with $|w| \geq k$, and a^m is given by evaluation: $a^m(r, t)(s) = r(s)(t(s))$.

Theorem 1. *Let m be a stack \mathbb{T} -automaton. Then for any $x_0 \in X$ and any $\gamma_0 \in \Gamma$, $\{w \in A^* \mid \llbracket x_0 \rrbracket_m(w)(\gamma_0)\}$ is a real-time deterministic context-free language. Conversely, any real-time deterministic context-free language can be obtained in this way.*

A similar theorem allows to capture all context-free languages by \mathbb{T} -automata. But here one needs non-deterministic machines and so one takes for \mathbb{T} the tensor product of the stack monad with the finite power-set monad \mathcal{P}_f .

In order to capture languages beyond context-free ones a treatment of silent transitions is needed (besides a monad \mathbb{T} modelling more powerful storage like e.g. a Turing tape). For this we first consider the *CPS* (“*continuation passing style*”)-transform of a given \mathbb{T} -automaton m , which is obtained as follows. First, consider the continuation monad \mathbb{T}_B given by $T_B X = B^{B^X}$. Then observe that having the \mathbb{T} -algebra structure $a^m : TB \rightarrow B$ is equivalent to having a monad morphism $\kappa : \mathbb{T} \rightarrow \mathbb{T}_B$ (see Kelly [3, Prop. 22.4]). We then extend the given \mathbb{T} -automaton m to $m_* = (\text{id}_B \times \kappa^A) \cdot m$ with $a^{m_*} = \lambda t. t(\text{id}) : T_B B \rightarrow B$. One can show that $\llbracket - \rrbracket_{m_*} = \llbracket - \rrbracket_m$. Now assume that B is ω -additive, i.e. (besides being a \mathbb{T} -algebra) it also is an algebra for the countably supported multiset monad. In other words, B is a commutative monoid with infinite summation. It is not difficult to see that the ω -additive structure extends pointwise to T_B , i.e., for each set X , $T_B X$ is an ω -additive monoid, too. This then allows us to define the elimination of silent transitions from a \mathbb{T} -automaton; indeed, suppose we have a \mathbb{T} -automaton with input set $A_\tau = A \cup \{\tau\}$, then one can define a \mathbb{T}_B -automaton by first forming m_* and then using infinite summation to obtain $m_\tau : X \rightarrow B \times (T_B X)^A$ (without τ -transitions)³. We define the (τ -free) semantics of m as $\llbracket - \rrbracket_m^\tau = \llbracket - \rrbracket_{m_\tau}$.

Now we are ready to give a coalgebraic description of recursively enumerable languages. For this we consider the *tape monad*: for a finite set Γ of tape symbols this is the submonad \mathbb{T} of the store monad $(- \times \mathbb{Z} \times \Gamma^{\mathbb{Z}})^{\mathbb{Z} \times \Gamma^{\mathbb{Z}}}$ for which the elements $\langle r, z, t \rangle : \mathbb{Z} \times \Gamma^{\mathbb{Z}} \rightarrow (X \times \mathbb{Z} \times \Gamma^{\mathbb{Z}})$ of TX satisfy the following restriction: there exists $k \geq 0$ such that for any $i, j \in \mathbb{Z}$ and any $\sigma, \sigma' : \mathbb{Z} \rightarrow \Gamma$ if $\sigma =_{i \pm k} \sigma'$ ⁴ then

$$\begin{aligned} t(i, \sigma) &=_{i \pm k} t(i, \sigma'), & r(i, \sigma) &= r(i, \sigma'), & |z(i, \sigma) - i| &\leq k, & t(i, \sigma) &=_{i \pm k} \sigma, \\ z(i, \sigma) &= z(i, \sigma'), & t(i, \sigma_{+j}) &= t(i + j, \sigma_{+j}), & r(i, \sigma_{+j}) &= r(i + j, \sigma), \\ z(i, \sigma_{+j}) &= z(i + j, \sigma) - j. \end{aligned}$$

where $\sigma_{+j}(i) = \sigma(i + j)$. A *tape automaton* is a \mathbb{T} -automaton for the tape monad with $B \subseteq 2^{\mathbb{Z} \times \Gamma^{\mathbb{Z}}}$ consisting of all those $p \in 2^{\mathbb{Z} \times \Gamma^{\mathbb{Z}}}$ for each of which there exists a k such that $p(i, \sigma) = p(i, \sigma')$ and $p(i, \sigma_{+j}) = p(i + j, \sigma)$ if $\sigma =_{i \pm k} \sigma'$ and with $a^m : TB \rightarrow B$ given by evaluation similarly as for stack automata. It can be shown that tape automata over A_τ are equivalent to deterministic 2-tape Turing machines with input alphabet A , where the first tape is a special read-only and right-only tape holding the input word at the beginning of a computation. Thus, we obtain that tape automata represent precisely the recursively enumerable-languages.

Theorem 2. *For every tape automaton m over A_τ , Γ with $|\Gamma| \geq 2$ containing a special blank symbol \square , and every state $x \in X$ the following language is recursively enumerable: $\{w \in A^* \mid \llbracket x \rrbracket_m^\tau(w)(0, \sigma_\square) = 1\}$, where σ_\square is the constant function on \square . Conversely, every recursively enumerable language can be represented in this way.*

³ Unfortunately, space limitations prevent us from giving the technical details here.

⁴ We write $\sigma =_{i \pm k} \sigma'$ ($\sigma =_{i \pm k} \sigma'$) if $\sigma(j) = \sigma'(j)$ for all j such that $|i - j| \leq k$ ($|i - j| > k$).

References

- [1] Goncharov, S.: Trace semantics via generic observations. In: Heckel, R., Milius, S. (eds.) CALCO 2013. LNCS, vol. 8089 (2013)
- [2] Hasuo, I., Jacobs, B.: Context-free languages via coalgebraic trace semantics. In: et al., J.F. (ed.) Proc. Algebra and Coalgebra in Computer Science: First International Conference (CALCO'05). Lecture Notes in Comput. Sci., vol. 3629, pp. 213–231. Springer (2005)
- [3] Kelly, G.M.: A unified treatment of transfinite constructions for free algebras, free monoids, colimits, associated sheaves, and so on. Bull. Austral. Math. Soc. 22, 1–83 (1980)
- [4] Silva, A., Bonchi, F., Bonsangue, M.M., Rutten, J.J.M.M.: Generalizing determinization from automata to coalgebras. Log. Methods Comput. Sci 9(1:9) (2013)
- [5] Winter, J., Bonsangue, M.M., Rutten, J.J.M.M.: Context-free languages, coalgebraically. In: Proc. Algebra and Coalgebra in Computer Science (CALCO'11). Lecture Notes Comput. Sci., vol. 6859, pp. 359–376. Springer (2011)