PhD Proposal: Computation of distances between transducer models and application to approximate learning of word functions and relations.

Supervisors: C Aiswarya and Benjamin Monmege

Text-to-text transductions are a fundamental aspect of computation, appearing in applications such as code-to-code translation, text preprocessing, and syntactic transformation. Their formal mathematical models, known as **finite-state transducers (FSTs)**, extend finite automata by associating input strings with output strings. Key theoretical questions in this domain include **functional equivalence** (determining whether two transducers define the same transformation) and **type checking** (ensuring that a transduction preserves membership in a formal language). These problems are particularly relevant in software verification, where one must confirm that a code transformation maintains syntactic and semantic properties.

However, exact equivalence checking is computationally hard or even undecidable in many settings, motivating the need for **weakened or approximate equivalence**. In this thesis, we propose to study thoroughly a relaxation of functional equivalence called *k*-equivalence, where two transducers are said equivalent if their outputs can be converted into one another with at most *k* edits, inspired by the famous "edit distance". Preliminary results have been obtained [1] on this notion for word-to-word transducers. Our goal is to investigate the theoretical properties of this notion on more general forms of transducers like nested-word-to-word (the so-called *visibly pushdown transducers* [2]), transducers working on multiply nested words to model multi-pushdown systems, or transducers working on data words [8]. Our goal is also to develop efficient learning algorithms for transducers under *k*-equivalence using the Angluin L* framework [3], already in the case of word-to-word transducers, but also in the generalizations like for data words [7]. We hypothesize that allowing small errors in the learning process could lead to more compact and generalizable models.

Several studies have explored equivalence checking and learning for finite-state models. Angluin's seminal work on active learning [3] provides the foundation for learning finite automata from queries. More recently, advances in approximate learning have been explored in probabilistic automata and error-tolerant string matching [6]. The complexity of exact equivalence checking for several classes of transducers has been investigated: see, [4] for a survey, and, e.g., [5] where weighted extensions of such transducers are also considered. Some other models of transducers with registers called streaming string transducers have also been explored and shown decidable [9] (that have the particularity to be deterministic, thus the name "streaming"). However approximate equivalence remains largely unexplored. The approximate equivalence question is also strongly related (and thus could lead to advances) to the search for minimal models for a given language or word-to-word function: while finite-state automata are minimizable, the search for such minimal object is left open for streaming string transducers, where there is a tradeoff between the number of states and the number of registers.

By bridging these areas, this thesis aims to develop a theoretically grounded yet practically viable framework for learning approximate transductions.

Methodology

Our approach will consist of both theoretical and experimental components:

• **Formalization of** *k***-equivalence:** We will generalize the *k*-equivalence notion introduced in [1] for word-to-word transducers, using standard edit distance metrics (insertions, deletions, substitutions), to various other models of transducers, in particular visibly pushdown transducers, and study its computational properties. In the realm of nested words (modelling

the pushdown capabilities), edits must be performed with care since it is not possible to delete a push on the stack without also deleting the associated pop.

- Algorithm Design: We will develop efficient algorithms for checking *k*-equivalence on these more general models of transducers, and also study the relationship of the approximate equivalence problem with the minimization problem, in particular in streaming string transducers.
- Learning: We will extend Angluin's classical algorithm for learning finite-state machines to work with approximate equivalence, allowing the learner to tolerate small errors in teacher responses. Extensions will then be considered, like data words.
- **Empirical Evaluation:** We will implement and benchmark our algorithms on real-world datasets, such as programming language translation tasks and natural language processing benchmarks.

We plan on publishing the results in the best international conferences and journals of theoretical computer science, but also, depending on the empirical results we will obtain, in more applied venues to demonstrate the applicability of those theoretical results to several domains like programming languages or naturel language processing.

This proposal is ambitious and vast, and the technical work required will be challenging. However the relative risk, regarding the difficult underlying questions that this PhD proposal contains, is tempered by the vast open space of questions, some of which have low-hanging fruits to balance the high-risk portions of the proposal. There are no single bottlenecks that can block progress globally.

References

- 1. C. Aiswarya, Amaldev Manuel, Saina Sunny (2024). *Edit Distance of Finite State Transducers*. ICALP 2024: 125:1-125:20
- Emmanuel Filiot, Jean-François Raskin, Pierre-Alain Reynier, Frédéric Servais, Jean-Marc Talbot (2018). *Visibly pushdown transducers*. Journal of Computer and System Sciences, 97, pp.147-181.
- 3. Angluin, D. (1987). *Learning regular sets from queries and counterexamples*. Information and Computation, 75(2), 87–106.
- 4. Filiot, E., & Reynier, P. A. (2016). *Transducers, logic and algebra for functions of finite words*. Logical Methods in Computer Science, 12(3).
- 5. Mohri, M. (2003). *Finite-state transducers in language and speech processing*. Computational Linguistics, 29(1), 23-27.
- 6. Balle, B., & Mohri, M. (2015). *Learning weighted automata*. Journal of Machine Learning Research, 16, 1625-1655.
- 7. B. Bollig, P. Habermehl, M. Leucker, and B. Monmege (2014). *A robust class of data languages and an application to learning*. Logical Methods in Computer Science, 10(4:19).
- 8. Léo Exibard, Emmanuel Filiot, and Pierre-Alain Reynier (2019). *Synthesis of Data Word Transducers*. CONCUR 2019: LIPIcs, Volume 140.
- 9. R. Alur and P. Černý (2011). *Streaming transducers for algorithmic verification of singlepass list-processing programs*. POPL 2011, p. 599–610, ACM.