

Computing and classifying reduplication with 2-way finite-state transducers

Hossep Dolatian and Jeffrey Heinz
Stony Brook University

ABSTRACT

This article describes a novel approach to the computational modeling of reduplication. Reduplication is often treated as a stumbling block within finite-state treatments of morphology because they cannot adequately capture the productivity of *unbounded* copying (total reduplication) and because they cannot describe *bounded* copying (partial reduplication) without a large increase in the number of states. We provide a comprehensive typology of reduplicative processes and show that an understudied type of finite-state machine, 2-way deterministic finite-state transducers (2-way D-FSTs), captures virtually all of them. Furthermore, the 2-way D-FSTs have few states, are in practice easy to design and debug, and are linguistically motivated in terms of the transducer's origin semantics or segment alignment. Most of these processes, and their corresponding 2-way D-FSTs, are available in an online database of reduplication (RedTyp). We classify these 2-way D-FSTs according to the concatenation of known subclasses of regular relations and show that the majority fall into the Concatenated Output Strictly Local (C-OSL) class. Other cases require higher subclasses but are still definable by 2-way D-FSTs.

Keywords:
reduplication,
2-way finite state
transducer, finite
state morphology

INTRODUCTION

Reduplication is a cross-linguistically common word-formation process involving *copying*. Given a word, reduplication can copy either a bounded (1a) or unbounded (1b) number of segments. The symbol \sim marks the boundary between the two copies.

- (1) a. **Partial Reduplication** *Agta* (Moravcsik 1978, 311)
 takki \rightarrow tak \sim takki ‘leg’ \rightarrow ‘legs’
- b. **Total Reduplication** *Indonesian* (Cohn 1989, 308)
 buku \rightarrow buku \sim buku ‘book’ \rightarrow ‘books’

Reduplication is used in the majority of the world’s languages, and total reduplication is more common than partial reduplication. The World Atlas of Language Structures (WALS) database documents that 278 out of 368 (75%) languages have total and partial reduplication (Rubino 2013). 35 additional languages (10%) use only total, not partial, reduplication. The 55 (15%) remaining languages do not have productive reduplication, but this figure is debatable.¹ Therefore, developing analyzable and efficient computational models of reduplication is important.

Although reduplication is well-studied, it is a computationally challenging process (Sproat 1992). In computational linguistics, most morphological and phonological processes can be analyzed with finite-state calculus in terms of rational languages and transductions (Kaplan and Kay 1994; Beesley and Karttunen 2003). However, reduplicative processes cannot be easily modeled with the same finite-state systems. For total reduplication, this is because those finite-state systems cannot express unbounded copying in the first place (Culy 1985). As for partial reduplication, those finite-state systems are often described as unwieldy because of the state explosion caused by partial reduplication (Roark and Sproat 2007, 54). Section 2 of this article explains why reduplication is computationally challenging while reviewing previous computational approaches to this linguistic phenomenon.

¹ Most of the exceptional languages are Indo-European, but some argue that these languages still use total reduplication (Ghameshi *et al.* 2004; Stolz *et al.* 2011).

In this context, the primary contribution of this article is to show that a specific understudied type of finite-state technology *can* account for virtually all reduplicative processes. This type of transducer is known as a 2-way Finite-State Transducer or 2-way FST (Savitch 1982; Engelfriet and Hoogeboom 2001; Filiot and Reynier 2016).² In theoretical computer science, 2-way FSTs are known to be able to model unbounded copying (Engelfriet and Hoogeboom 2001). To our knowledge, we are the first to apply 2-way FSTs to computational linguistics.³

The FSTs used in most of computational linguistics are more accurately called *1-way* FSTs. They can only read the input once in one direction. 2-way FSTs are more expressive because the read head can move back and forth on the input tape. On the other hand, the write head can only move forward on the output tape. For this reason, they are less expressive than Turing machines. It is this back-and-forth movement of the read head that allows 2-way FSTs to adequately model reduplication without the difficulties faced by 1-way FSTs. This article introduces *deterministic* 2-way finite-state transducers (2-way D-FSTs) in Section 3, along with their formal definition (Section 3.1), illustrative examples of reduplication (Section 3.2), and their computational properties (Section 3.3).

The fact that the 2-way FSTs used in this article are deterministic is significant. It is well known that deterministic 1-way FSTs are less expressive than non-deterministic 1-way FSTs (Elgot and Mezei 1965; Schützenberger 1975; Choffrut 1977; Mohri 1997; Heinz and Lai 2013). Similarly, 2-way D-FSTs are less expressive than non-deterministic 2-way FSTs (Culik and Karhumäki 1986). Consequently, the empirical result that reduplication can be modeled with deterministic 2-way FSTs is in line with work which shows that various phonological and morphological processes can be described with deterministic finite-state technology (Chandlee *et al.* 2012; Gainor *et al.* 2012; Chandlee and Heinz

²This article builds off of our previous work on using 2-way FSTs for reduplication (Dolatian and Heinz 2018a,b, 2019a,b).

³2-way finite-state *acceptors* (2-way FSAs) have been used to model non-concatenative Semitic morphology (Narayanan and Hashem 1993) and to parse dependency grammars (Nelmarkka *et al.* 1984).

2012; Heinz and Lai 2013; Chandlee 2014; Luo 2017; Payne 2014, 2017).

In the later part of this article, we provide a comprehensive cross-linguistic survey of reduplicative processes based on earlier typological studies (Moravcsik 1978; Rubino 2005; Inkelas and Downing 2015a), with reduplication defined as an input-to-output function (McCarthy and Prince 1995). This survey is documented in a database we constructed, which we call The RedTyp database. It contains entries for 138 reduplicative processes from 91 languages and a 2-way D-FST for each entry. Aspects of this survey are presented in Sections 3–4, and discussed in detail in Section 6.

In Section 4, we compare 2-way D-FSTs to 1-way FSTs in terms of empirical coverage (Section 4.1), practical utility (Section 4.2), and linguistic motivation (Section 4.3). We argue that 2-way D-FSTs are linguistically motivated in that they capture the correspondence relations underlying the base and the reduplicant in a linguistically natural way. These correspondence relations are couched in terms of origin semantics (Bojańczyk 2014). We use origin semantics as a diagnostic for the strong generative capacity of reduplicative functions. (We do not claim that origin semantics matches linguistic intuitions exactly in every case, but rather that it approximately does so in many instructive cases.)

The final contribution of this article is an attempt to classify reduplicative processes according to *subclasses* of 2-way D-FSTs in Section 6. The first result we already mentioned: the full typology of reduplicative processes can be modeled with *deterministic* 2-way D-FSTs. These subclasses are defined in terms of concatenations of subclasses of 1-way FSTs. Our next result is that approximately three-quarters of the typology is expressible with the concatenation of Output Strictly Local (OSL) functions (Chandlee *et al.* 2015). The remainder of the typology is expressible with the concatenation of sequential functions, with some arguably requiring sweeping transducers or unrestricted 2-way D-FSTs.

We review these contributions and conclude in Section 7.

BACKGROUND ON COMPUTATION OF REDUPLICATION

2

Within computational linguistics, reduplication has been a challenging process to model (Culy 1985; Sproat 1992; Roark and Sproat 2007; Hulden 2009a; Chandlee 2014, 2017). Finite-state technology, as currently practiced, cannot adequately and elegantly describe many cases of productive reduplication, especially *unbounded* total reduplication. There are three kinds of issues: empirical coverage, practical utility, and matching the intensional description of reduplication. We discuss these challenges in Section 2.1. In response to these problems, some have proposed finite-state approximations for reduplication (Section 2.2) or developing more expressive systems just for total reduplication (Section 2.3). The latter approach however implies that total reduplication is ontologically different from partial reduplication, and thus should be computed differently. In Section 2.4, we discuss this implication and show that the evidence for it is inconclusive. We summarize in Section 2.5.

Why reduplication is challenging

2.1

Reduplication is challenging because segmental copying entails multiple crossing dependencies between the two copies. When the number of copied segments (and thus the number of crossing dependencies) is bound to some maximum number n , the outcome is partial reduplication. When there is no bound, the outcome is total reduplication.

Partial reduplication *can* be modeled with 1-way FSTs (Roark and Sproat 2007; Chandlee and Heinz 2012). However, as we explain in more detail in Section 4.3, these machines are understood as memorizing all finitely-many possible forms of the partial reduplicant. Consequently, the transducers suffer from an explosion of states and become unwieldy. For example, in a language with a medium-sized phonemic inventory of 22 consonants and 5 vowels, partial reduplication with a CVCV template would require at least $22 + 22 \times 5 + 22 \times 5 \times 22 = 2552$ states to memorize the first C (22 states), the first V (22×5 states), and the second C ($22 \times 5 \times 22$ states). 1-way FSTs likewise arguably

do not match the intensional description of reduplication as a *copying* process because the FSTs simply memorize all possible reduplicants in the language (Roark and Sproat 2007, 54). This is discussed in detail in Section 4.3.

On the other hand, total reduplication cannot be modeled at all with 1-way FSTs (Culy 1985). This inability is due to the fact that the output language of total reduplication is not a regular language. Rather, the copying process of total reduplication can create output languages that are identical to the non-context-free $L_{ww} = \{ww \mid w \in \Sigma^*\}$ (Hopcroft and Ullman 1969). Thus the copying function $w \mapsto ww$ (sometimes called the squaring function) is beyond the expressivity of 1-way FSTs. In fact, virtually all attested morphological processes can be described with 1-way finite-state acceptors and transducers, *except* for total reduplication (Langendoen 1981; Gazdar and Pullum 1985; Roark and Sproat 2007). In response to this problem, computational morphologists have often resorted to using either finite-state approximations (Section 2.2) or non-finite-state tools (Section 2.3).

2.2

Finite-state approximations

The literature on finite-state morphology contains many finite-state approximations to reduplication (Walther 2000; Beesley and Karttunen 2003; Cohen-Sygal and Wintner 2006; Hulden and Bischoff 2009). Roark and Sproat (2007, 57) and Cohen-Sygal and Wintner (2006, 52) provide reviews. In general, finite-state approximations are designed to lessen the burden for the developer in designing reduplication rules. They introduce new operations or tools over 1-way FSTs but they do *not* increase their expressivity. In other words, they are designed to improve on practical utility but they don't improve on empirical coverage or intensional description.

Here we briefly review the two main sets of approaches with details following. One set of approaches checks for identity between the two copies (Cohen-Sygal and Wintner 2006; Hulden and Bischoff 2009). Another set of approaches essentially 'postpones' reduplication to a run-time task (Walther 2000; Beesley and Karttunen 2000, 2003). Both try to reduce state complexity either by making a trade-off with time complexity, by implementing reduplication with a unique 1-way transducer for each morpheme in the finite lexicon, or both.

Cohen-Sygal and Wintner (2006) augment 1-way FSAs with finitely many registers (FSRA). These registers keep track of a *bounded* number of segments previously seen in the input. In order to model the total reduplication of a *given* word like *buku* → *buku~buku* (1b), the FSRA has at least as many registers as segments in the base *buku*: four. The registers check that the string *buku~buku* can be broken down into identical copies. Similarly, Hulden and Bischoff (2009) design the EQ function within the foma system (Hulden 2009b) which checks if a string is divided into two identical copies.

As for run-time procedures, these systems are designed on an input by input basis. Given some input word, they create a reduplication FST for it *on the fly*. Given an input *buku*, the compile-replace operation (Beesley and Karttunen 2000, 2003) creates an intermediate representation $\{buku\}^2$ via a 1-way FST. This intermediate representation is then interpreted as a regular expression in run-time, i.e. it is *compiled*. By compiling this regular expression, the word *bukubuku* is outputted.

Within the framework of One-Level Phonology (Bird and Ellison 1994), Walther (2000) models reduplication by representing a potentially-reduplicated morpheme like *buku* as an FSA with augmentations on the types of transition arcs: *content*, *repeat*, and *skip* arcs. These transition arcs turn a linear string *buku* into a multi-linear structure where the read head can ‘move’ around the string. This enriched representation is then intersected with a reduplication FSA that is designed to ‘move’ around this enriched representation and generate *buku~buku*.⁴ Ideally, these operations should be applied in run-time. Otherwise, if these operations are applied to the entire lexicon and stored as a single FST, they then suffer from the state explosion that they were designed to avoid.

As for total reduplication, all four of the above modifications are *approximations*. This is because they impose various restrictions which contradict the linguistic generalization that total reduplication is independent of string length. Most notably, all four approximations permit only a closed finite set of input strings to undergo total reduplication. This restriction fundamentally alters total reduplication from a process

⁴Walther (2000) does not give a formal analysis. But, we think that these augmented transition arcs are similar to 2-way FSAs and are an independently developed implementation for Precedence-Based Phonology (Raimy 2000).

which in principle applies to infinitely many words to a process which applies to only finitely many. Such approximations thus fall short of capturing how total reduplication is used as a productive process in natural language.

As for partial reduplication, all of the above four approaches have the same expressivity and are able to capture the linguistic generalization that partial reduplication is independent of string length. However, although they are designed to avoid state explosion by one means or another, they can still be said to memorize the partial reduplicant as opposed to copying it (see Section 4.3). In this way they do not intensionally capture the linguistic generalization of copying.

2.3

Extending formal power

Because of the difficulty in modeling reduplication with finite-state machinery, various augmentations and expansions of context-free grammars have been proposed to handle L_{ww} and reduplication. An early augmentation is Reduplication Context-Free Grammars (Manaster-Ramer 1986; Savitch 1989) designed to handle context-free languages and reduplication by using queues instead of stacks. A more recent augmentation is Multiple-Context Free Grammars (MCFGs) which can model L_{ww} (Seki *et al.* 1991, 1993). MCFGs have been used to model reduplication (Albro 2000, 2005). As an extension of MCFGs, Parallel MCFGs have been used to model reduplication and syntactic copying (Kobele 2006; Clark and Yoshinaka 2012, 2014; Clark 2017).⁵ Crysmann (2017) explores the use of HPSG to model total reduplication.

These technologies have had considerably less attention within mainstream computational morphology than finite-state approximations. One shortcoming of these approaches is that they model formal languages, not transformations. They accept well-formed reduplicated words ww , but they do not generate a reduplicated word ww given some input w . Thus, they do not model the squaring function $w \mapsto ww$.

⁵Kobele (2006) shows that syntactic copying can generate languages of the form a^{2^n} , i.e., exponential copying. This isn't attested in morphological copying. Note the string a^{2^n} is generated as the yield language of a tree transduction over a derivation tree.

*Computational distinctions
between total and partial reduplication*

The previous sections showed that more expressive mechanisms are needed to model reduplication. Conceptually, the use of more powerful computational formalisms implies that reduplication is ‘different’ from the rest of morpho-phonology which can be modeled using 1-way FSTs (Roark and Sproat 2007, 60). This is especially the case for total reduplication which cannot be exactly modeled with 1-way FSTs, whereas partial reduplication can. This difference has caused debates over whether both types of reduplication *should* be computed with the same formalism or not. However, this debate is inconclusive.

On one hand, Chandlee (2017) suggests that the inadequacy of 1-way FSTs for total reduplication is evidence for total reduplication being ontologically different from partial reduplication. There is some empirical support for this argument. Prosodically, total reduplication resembles more ‘syntactic’ processes like compounding more often than partial reduplication (Downing 2006). The two copies in total reduplication can be stressed separately or have separate tonal contours (Downing 2003).

On the other hand, partial and total reduplication are closely related processes. Typologically, if a language has partial reduplication, then it almost always has total reduplication too (Rubino 2013). Diachronically, both types of reduplication are typically related to each other, but not always (Hurch and Mattes 2009). And in linguistic theory, both are modeled with the same tools (Steriade 1988; Raimy 2000; Inkelas and Zoll 2005; McCarthy *et al.* 2012).

Psycholinguistic work could shed more light on the issue of total reduplication vis a vis partial reduplication. Sadly, there is little to no work on the psycholinguistic processing of reduplication. To our knowledge, existing work focuses on partial reduplication, not total reduplication (Ohala *et al.* 1986; Waksler 1999).

Learnability is another factor which could tease apart these processes. It is an open question whether both partial and total reduplication can be learned in the same way with the same mechanism. In terms of stringsets, the formal language of totally reduplicated words *ww* can be learned with distributional methods for MCFGs (Clark and Yoshinaka 2012, 2014, 2016). There is also a substantial body of work

in cognitive science and connectionism on how to learn reduplicated words (Marcus *et al.* 1999; Berent *et al.* 2014, 2016, 2017; Andan *et al.* 2018; Alhama 2017; Alhama and Zuidema 2019). Here, the task is learning words which have repeated substrings (ABB or ABA) where A and B are syllables.

In contrast, there is little to no work on learning reduplication as a function ($w \rightarrow ww$), whether in machine learning or grammatical inference. To our knowledge, the only algorithm designed specifically for learning reduplication is Nevins (2004) in the principles-and-parameters tradition. There is some recent work on using neural networks to learn copying (Gu *et al.* 2016; Prickett *et al.* 2018; Wilson 2019; Nelson *et al.* 2020). We speculate that one reason for the dearth of learning results is due to the challenges outlined above for finding natural computational models for reduplication.

2.5 *Summary and consequences*

All in all, current finite-state treatments of reduplication have issues regarding their empirical coverage (total reduplication's productivity), practical utility (state space explosion), and intensional descriptions (copying vs. remembering). The present study uses a computational formalism which does not suffer from these three problems: two-way finite-state transducers (2-way FSTs).

3 2-WAY FINITE-STATE TRANSDUCERS: DEFINITION AND APPLICATION TO REDUPLICATION

1-way FSTs read the input *once* from left to right. Most applications use non-deterministic 1-way FSTs (Roark and Sproat 2007), though deterministic 1-way FSTs are largely sufficient (Chandlee 2017). (For all inputs, deterministic FSTs have at most one path through the transducer, whereas non-deterministic ones may have more than one path.) 2-way FSTs can move back and forth on the input (Rabin and Scott 1959; Hopcroft and Ullman 1969). This ability makes them more expressive than 1-way FSTs (Savitch 1982; Engelfriet and Hoogboom 2001).

It is useful to imagine a 2-way FST as a machine operating on an input tape and writing to an output tape. The symbols on the input tape are drawn from an alphabet Σ and the symbols written to the output tape are drawn from an alphabet Γ . For an input string $w = \sigma_1 \dots \sigma_n$, the initial configuration is that the FST is in some internal state q_0 , its read head on σ_1 , and its write head at the beginning of an empty output tape. After the FST reads the symbol under the read head, three things occur:

- The internal state of the FST may change.
- The FST writes some string, possibly empty, to the output tape.
- The read head moves in one of three ways: moves to the left (-1), moves to the right ($+1$), or stays (0).

This process repeats until the read head “falls off” one of the edges of the input tape. If for some input string w , the FST falls off the right edge of the input tape when the FST is in an accepting state after writing u on the output tape, we say the FST transduces, transforms, or maps, w to u . If for some input string w , the FST falls off the left edge, falls off the right edge while in a non-accepting state, or never falls off either edge, then the FST is undefined at w . The write head cannot move back along the output tape. It can only advance as strings are written.

We formalize the definition and behavior of 2-way FSTs in Section 3.1. They are illustrated for reduplication in Section 3.2. We then describe their generative capacity and computational complexity (Section 3.3).

Preliminaries and formal definition

3.1

Given a finite alphabet Σ , the set of all possible strings of finite length built from Σ is Σ^* . The empty string is represented by λ . The length of a string w is $|w|$, so $|\lambda| = 0$. For the given strings w_1, w_2 , their concatenation is $w_1 w_2$. Below is a formalization of deterministic 2-way FSTs based on Filiot and Reynier (2016) and Shallit (2008). We adopt the convention that inputs to a 2-way D-FST are flanked with the start (\times) and end (\times) boundaries. This larger alphabet is denoted by Σ_\times .

(2) **Definition:** A 2-way D-FST is a six-tuple $(Q, \Sigma_{\times}, \Gamma, q_0, F, \delta)$ where:

- Q is a finite set of states,
- $\Sigma_{\times} = \Sigma \cup \{\times, \times\}$ is the input alphabet,
- Γ is the output alphabet,
- $q_0 \in Q$ is the initial state,
- $F \subseteq Q$ is the set of final states,
- $\delta : Q \times \Sigma \rightarrow Q \times \Gamma^* \times D$ is the transition function where the direction $D = \{-1, 0, +1\}$.

A *configuration* of a 2-way D-FST T is an element of $\Sigma_{\times}^* Q \Sigma_{\times}^* \times \Gamma^*$. The meaning of the configuration (wqx, u) is that the input to T is wx and the machine is currently in state q with the read head on the first symbol of x (or has fallen off the right edge of the input tape if $x = \lambda$) and that u is currently written on the output tape.

If the current configuration is $(wqax, u)$ and $\delta(q, a) = (r, v, 0)$ then the next configuration is $(wrax, uv)$, in which case we write $(wqax, u) \rightarrow (wrax, uv)$. If the current configuration is $(wqax, u)$ and $\delta(q, a) = (r, v, +1)$ then the next configuration is $(warx, uv)$. In this case, we write $(wqax, u) \rightarrow (warx, uv)$. If the current configuration is $(waqx, u)$ and $\delta(q, a) = (r, v, -1)$ then the next configuration is $(wrax, uv)$. We write $(waqx, u) \rightarrow (wrax, uv)$. Observe that since δ is a function, there is at most one next configuration.

The transitive closure of \rightarrow is denoted with \rightarrow^+ . Thus, if $c \rightarrow^+ c'$ then there exists a finite sequence of configurations $c_1, c_2 \dots c_n$ with $n > 1$ such that $c = c_1 \rightarrow c_2 \rightarrow \dots \rightarrow c_n = c'$.

Next we define the function f_T that a 2-way D-FST T computes. For each string $w \in \Sigma_{\times}^*$, $f_T(w) = u \in \Gamma^*$ provided there exists $q_f \in F$ such that $(q_0 \times w \times, \lambda) \rightarrow^+ (\times w \times q_f, u)$. If $f_T(w) = u$ then u is unique because the sequence of configurations is determined deterministically.

If the configurations of a 2-way D-FST T halt the computation of T on some input w , then we say T is undefined on w . If the configuration is (qax, u) and $\delta(q, a) = (r, -1, v)$ then the derivation crashes and the transduction $f_T(ax)$ is undefined. Likewise, if the configuration is (wq, u) and $q \notin F$ then the transducer crashes and the transduction f_T is undefined on input w . Another way that f_T may be undefined for some input is if the input causes the transducer to go into an infinite

loop.⁶ This occurs for input $wx \in \Sigma_{\times}^*$ whenever there exist $q \in Q$ and $u, v \in \Gamma^*$ such that $(q_0wx, \lambda) \rightarrow^+ (wqx, u) \rightarrow^+ (wqx, uv)$.

Illustration of two-way transducers for reduplication

3.2

Having established what 2-way D-FSTs are, this section illustrates how they can be used to model reduplication. We provide two examples: total reduplication and partial initial-CVC reduplication. Both examples use deterministic 2-way FSTs.

Some useful terms are ‘passes’ and ‘rewinds’. A pass (rewind) is when a 2-way D-FST moves left-to-right (right-to-left) from some position to another over the input.

Total reduplication is cross-linguistically the most common reduplicative process (Rubino 2005), and it is used in an estimated 85% of the world’s languages (Rubino 2013). We illustrate it with data from Indonesian where total reduplication marks plurality (Cohn 1989).

- | | | | |
|-----|----|------------------------------------|----------------------------------|
| (3) | a. | buku \rightarrow buku~buku | ‘book’ \rightarrow ‘books’ |
| | b. | wanita \rightarrow wanita~wanita | ‘woman’ \rightarrow ‘women’ |
| | c. | hak \rightarrow hak~hak | ‘right’ \rightarrow ‘rights’ |
| | d. | kəra \rightarrow kəra~kəra | ‘donkey’ \rightarrow ‘donkeys’ |

Figure 1 shows a 2-way D-FST that captures total reduplication. The boundary symbol \sim is a symbol in the output alphabet Γ , and is not necessary. We include it only for illustration. The 2-way D-FST in Figure 1 operates as follows:

1. **First pass:** It reads the input tape from left to right and outputs the first copy.
2. **Rewind:** When it reaches the end boundary \times , it ‘rewinds’ or goes back to the start of the input tape by moving left until the start boundary \times is reached.
3. **Second pass:** It reads the input tape once more from left to right and outputs the second copy.

⁶In practice, infinite loops are not a problem. It can be checked whether an input leads the 2-way D-FST into an infinite loop during run-time, in which case the computation can be halted.

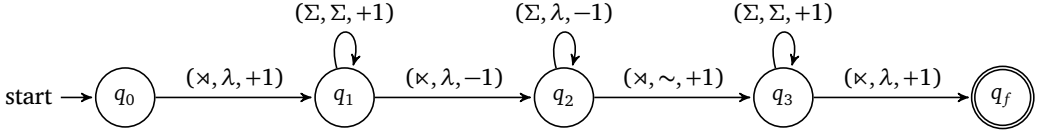


Figure 1: 2-way D-FST for total reduplication

The transition arcs are interpreted as follows. The symbol Σ is a variable representation of any alphabet symbol except for $\{x, \kappa\}$. The arrow from q_1 to itself $(\Sigma, \Sigma, +1)$ means this 2-way D-FST reads a symbol from Σ , writes that same symbol, and advances the read head one step to the right on the input tape.

Table 1 shows an example derivation for $buku \rightarrow buku \sim buku$ using the 2-way D-FST in Figure 1. The derivation shows the step-by-step configurations for the computation. The tuples in Table 1 consist of

Table 1: Derivation of $/buku/ \rightarrow [buku \sim buku]$

Outputting the first copy					
1.	$(q_0 \underline{x} buku \kappa, \lambda,$	N/A)	2.	$(x q_1 \underline{b} u \kappa, \lambda,$	$q_0 \xrightarrow[x+1]{x:\lambda} q_1)$
3.	$(x b q_1 \underline{u} \kappa, b,$	$q_1 \xrightarrow[+1]{\Sigma:\Sigma} q_1)$	4.	$(x b u q_1 \underline{\kappa} \kappa, bu,$	$q_1 \xrightarrow[+1]{\Sigma:\Sigma} q_1)$
5.	$(x buk q_1 \underline{u} \kappa, buk,$	$q_1 \xrightarrow[+1]{\Sigma:\Sigma} q_1)$	6.	$(x buku q_1 \underline{\kappa} \kappa, buku,$	$q_1 \xrightarrow[+1]{\Sigma:\Sigma} q_1)$
Going back to the start of the tape					
7.	$(x buk q_2 \underline{u} \kappa, buku,$	$q_1 \xrightarrow[-1]{\kappa:\lambda} q_2)$	8.	$(x bu q_2 \underline{\kappa} u \kappa, buku,$	$q_2 \xrightarrow[-1]{\Sigma:\lambda} q_2)$
9.	$(x b q_2 \underline{u} \kappa \kappa, buku,$	$q_2 \xrightarrow[-1]{\Sigma:\lambda} q_2)$	10.	$(x q_2 \underline{b} u \kappa \kappa, buku,$	$q_2 \xrightarrow[-1]{\Sigma:\lambda} q_2)$
11.	$(q_2 \underline{x} buku \kappa, buku,$	$q_2 \xrightarrow[-1]{\Sigma:\lambda} q_2)$			
Outputting the second copy					
12.	$(x q_3 \underline{b} u \kappa, buku \sim,$	$q_2 \xrightarrow[+1]{x:\sim} q_3)$	13.	$(x b q_3 \underline{u} \kappa \kappa, buku \sim b,$	$q_3 \xrightarrow[+1]{\Sigma:\Sigma} q_3)$
14.	$(x bu q_3 \underline{\kappa} u \kappa, buku \sim bu,$	$q_3 \xrightarrow[+1]{\Sigma:\Sigma} q_3)$	15.	$(x buk q_3 \underline{u} \kappa, buku \sim buk,$	$q_3 \xrightarrow[+1]{\Sigma:\Sigma} q_3)$
16.	$(x buku q_3 \underline{\kappa} \kappa, buku \sim buku,$	$q_3 \xrightarrow[+1]{\Sigma:\Sigma} q_3)$	17.	$(x buku x q_f, buku \sim buku,$	$q_3 \xrightarrow[+1]{\kappa:\kappa} q_f)$

three parts. The first two represent the configuration and the third part shows the transition exercised to reach this configuration from the previous one. The underlined input symbol is what the FST will read next. In the first tuple, there is no transition used (N/A). Transitions in the other tuples are given in the form shown below.

$$input\ state \xrightarrow[\text{direction}]{\text{input symbol:output string}} output\ state$$

Partial reduplication processes are also very common. A common example is initial-CVC reduplication as in Agta (Moravcsik 1978, 311).

- (4) a. takki → tak~takki ‘leg’ → ‘legs’
- b. uffu → uf~uffu ‘thigh’ → ‘thighs’

The 2-way D-FST in Figure 2 expresses partial initial-CVC reduplication. An example derivation of *takki* → *tak~takki* using our 2-way D-FST is provided in Table 2. For illustrative purposes, we assume that the function is undefined for V-initial inputs.

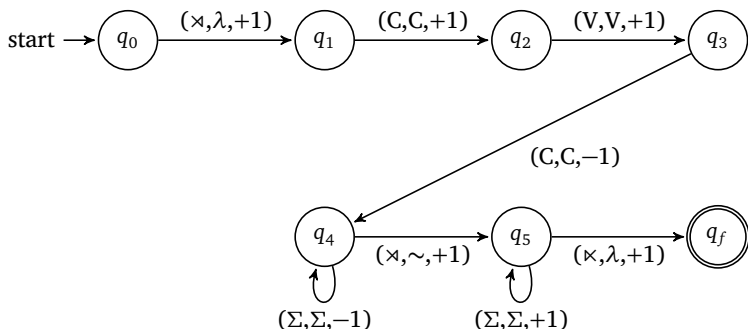


Figure 2:
2-way D-FST
for initial-CVC
reduplication

Generative capacity and computational complexity

3.3

With respect to acceptors, 1-way and 2-way finite-state acceptors are equivalent in expressive power. Both define the regular languages (Hopcroft and Ullman 1969; Shallit 2008). However, with respect to transducers, 1-way FSTs are strictly less expressive than 2-way D-FSTs (Savitch 1982; Aho *et al.* 1969; Filiot and Reynier 2016). For a 1-way

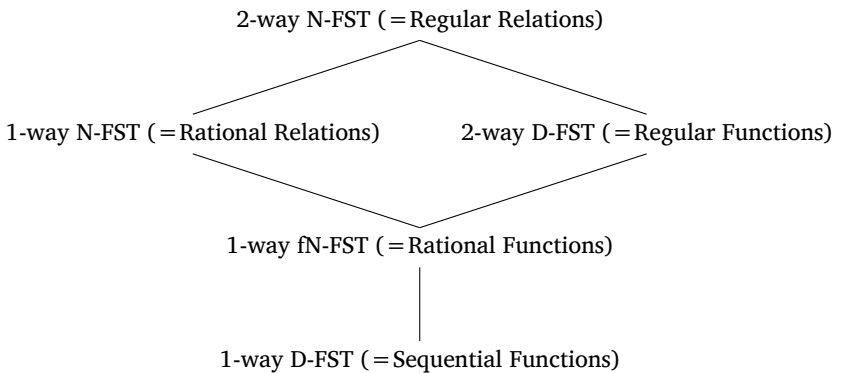
Table 2: Derivation of /takki/ → [tak~takki]

Outputting reduplicant			Outputting the base		
1.	$(q_0 \times \text{takki} \times, \lambda, \text{N/A})$		8.	$(\times q_5 \text{takki} \times, \text{tak} \sim, q_4 \xrightarrow[\text{+1}]{\times: \sim} q_5)$	
2.	$(\times q_1 \text{takki} \times, \lambda, q_0 \xrightarrow[\text{+1}]{\times: \lambda} q_1)$		9.	$(\times \text{t} q_5 \text{akki} \times, \text{tak} \sim \text{t}, q_5 \xrightarrow[\text{+1}]{\Sigma: \Sigma} q_5)$	
3.	$(\times \text{t} q_2 \text{akki} \times, \text{t}, q_1 \xrightarrow[\text{+1}]{\text{C: C}} q_2)$		10.	$(\times \text{t} a q_5 \text{kki} \times, \text{tak} \sim \text{ta}, q_5 \xrightarrow[\text{+1}]{\Sigma: \Sigma} q_5)$	
4.	$(\times \text{t} a q_3 \text{kki} \times, \text{ta}, q_2 \xrightarrow[\text{+1}]{\text{V: V}} q_3)$		11.	$(\times \text{t} a k q_5 \text{kki} \times, \text{tak} \sim \text{tak}, q_5 \xrightarrow[\text{+1}]{\Sigma: \Sigma} q_5)$	
5.	$(\times \text{t} a k q_4 \text{kki} \times, \text{tak}, q_3 \xrightarrow[\text{-1}]{\text{C: C}} q_4)$		12.	$(\times \text{t} a k k q_5 \text{i} \times, \text{tak} \sim \text{takk}, q_5 \xrightarrow[\text{+1}]{\Sigma: \Sigma} q_5)$	
Going back to the start of the tape			13.	$(\times \text{t} a k k i q_5 \times, \text{tak} \sim \text{takki}, q_5 \xrightarrow[\text{+1}]{\Sigma: \Sigma} q_5)$	
6.	$(\times q_4 \text{takki} \times, \text{tak}, q_4 \xrightarrow[\text{-1}]{\Sigma: \lambda} q_4)$		14.	$(\times \text{t} a k k i \times q_f, \text{tak} \sim \text{takki}, q_5 \xrightarrow[\text{+1}]{\times: \lambda} q_f)$	
7.	$(q_4 \times \text{takki} \times, \text{tak}, q_4 \xrightarrow[\text{-1}]{\Sigma: \lambda} q_4)$				

FST, both the input language and the output language must be regular languages. Therefore a 1-way FST cannot have its output language be the non-regular copy language $L_{ww} = \{ww | w \in \Sigma^*\}$. In contrast, the output language of a 2-way D-FST can be a non-regular language such as L_{ww} .

Figure 3 shows the hierarchy of FSTs, adapted from Filiot and Reynier (2016, p.8). Different FSTs have different generative capacity, based on whether the FST is deterministic (D-FST), non-deterministic (N-FST), 1-way, 2-way, and/or functional (f-FST).

Figure 3:
Hierarchy
of FSTs



2-way D-FSTs are equivalent in expressivity to string transductions that are defined in Monadic Second Order logic (Engelfriet and Hoogetboom 2001) and to streaming string transducers (Alur 2010).⁷ 2-way D-FSTs are less powerful than Turing machines because they cannot move back and forth on the output tape. They are closed under composition (Chytil and Jákł 1977) and some important classes are closed under inverse (Courcelle and Engelfriet 2012, 526).

Because of the difference in expressivity between 1-way and 2-way D-FSTs, it makes sense to give different names to the classes of functions that they compute. We follow Filiot and Reynier (2016) who identify the class of functions describable with 1-way deterministic FSTs as ‘sequential functions’, with 1-way functional non-deterministic FSTs as ‘rational functions’, and with 2-way deterministic FSTs as ‘regular functions’. The non-deterministic counterparts for 1-way and 2-way D-FSTs are respectively the ‘rational relations’ and ‘regular relations’.

1-way D-FSTs run in time linear to the length of the input string. As for 2-way D-FSTs, one useful metric for measuring their complexity is in terms of the number of times the 2-way D-FST passes through the input (Baschenis *et al.* 2016). In the case of the reduplication examples in Section 3.2, the 2-way D-FSTs used only two passes through the input, one for each copy. Thus, the run time for those 2-way D-FSTs is at most $2n \cdot m$ where n is the number of passes and m is the length of the input. Since n here is fixed at 2, the run time is still linear in the size of the input string. To our knowledge existing applications of regular functions have been efficient (Alur and Černý 2011; Alur *et al.* 2014).

CONTRASTING 2-WAY D-FSTs WITH 1-WAY FSTs

4

Having illustrated how 2-way D-FSTs can model reduplication, here we contrast 2-way FSTs with 1-way FSTs on three criteria: empirical coverage, practical utility, and intensional description.

⁷ A streaming-string transducer (SST) is a 1-way FST that uses finitely many registers of unbounded size. These registers allow the SST to keep track of previous information on the input tape, thus simulating 2-way D-FSTs.

4.1 *Empirical coverage of the typology and productivity*

In terms of empirical coverage, 2-way D-FSTs can effectively model *virtually* the entire typology of reduplication as described by Moravcsik (1978), Hurch (2005), Inkelas and Zoll (2005), Rubino (2005), and Samuels (2010). We review part of this typology in Section 6. This stands in stark contrast to 1-way FSTs discussed in Section 2. We say *virtually* because there are two cases in the literature which require further discussion. These are discussed in Section 6.6.2.

4.2 *Practical utility and the RedTyp database*

To showcase the empirical coverage of 2-way D-FSTs and their practical utility, we have constructed the RedTyp database (Dolatian and Heinz 2019a).⁸ It contains entries for 138 reduplicative processes from 91 languages. These were gleaned from various surveys (Rubino 2005; Inkelas and Downing 2015a). 50 of these processes were from Moravcsik (1978), an early survey which is representative of the cross-linguistically most common reduplicative patterns.

RedTyp contains 57 distinct 2-way D-FSTs that model the 138 processes. Each 2-way D-FST was designed manually, implemented in Python, and checked for correctness. On average, these 2-way D-FSTs had 8.8 states. This shows that 2-way D-FSTs are concise and convenient computational descriptions and models for reduplicative morphology. This is in contrast to 1-way FSTs which suffer from an explosion of states when modeling partial reduplication.⁹

To our knowledge, the only other database on reduplication is the Graz Database on Reduplication (Hurch 2005 ff.). However, RedTyp differs from the Graz Database because the latter does not include computational representations or implementations of its entries.

⁸It can be found on the first author's GitHub page <https://github.com/jhdeov/RedTyp>.

⁹The largest 2-way D-FST in RedTyp is for verbal reduplication in Kinande (Downing 2000) with 29 states. This pattern depends on the size of the root and the number and type of suffixes and prefixes around it. In contrast, we estimate a deterministic 1-way D-FST would require over 1,000 states for this pattern of partial reduplication.

A comparison between the two databases is provided in Dolatian and Heinz (2019a).

RedTyp offers a useful corpus of reduplicative patterns for research. For example as described in Section 6, we have used this database to identify subclasses of 2-way D-FSTs for classifying the typology of reduplication. This corpus could be used to test for other universal computational properties of reduplication. Since it contains 2-way D-FSTs, it can also be used to generate reduplicated forms. Such data sets can be used to test morphological learning algorithms.

One shortcoming is that RedTyp under-represents cases of opacity in reduplication because our main source, Moravcsik (1978), did not list opaque cases. As discussed further in Sections 6.4.3–6.5, opacity can be said to occur when phonological processes exceptionally apply either across both copies or across neither copy because of a drive to maintain identity between the two copies (McCarthy and Prince 1995). Only 5% of RedTyp displays opacity. Furthermore, RedTyp focuses on morphological copying, not syntactic copying (cf. Kobele 2006).

Linguistic motivation with origin semantics

4.3

Importantly, using 2-way D-FSTs for reduplication is linguistically motivated and matches the intensional descriptions behind the linguistic generalizations on reduplication.

2-way D-FSTs do not approximate reduplication like 1-way FSTs do. 2-way D-FSTs do not copy by remembering strings of segments (see Section 2). Instead they *actively and literally copy*.

This contrast between copying and remembering can be formalized with the notion of the *origin semantics* of a transduction (Bojańczyk 2014).¹⁰ Given a string-to-string function, the origin semantics of a function is the origin information of each symbol o_n in the output string. This is the position i_m of the read head on the input tape when the transducer had outputted o_n . To illustrate, consider a partial string-to-string function f_{ab} which maps ab to itself:

$$f(x) = \{(ab, ab)\}$$

¹⁰ For an application of origin semantics to MCFGs and potentially to machine translation, see Nederhof and Vogler (2019).

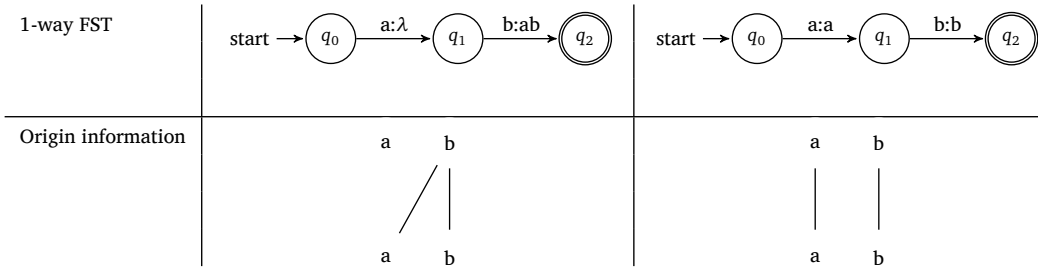


Figure 4: Pair of 1-way FSTs for the function f_{ab} and the origin information created by them for the mapping $ab \rightarrow ab$

As shown in the top row of Figure 4, this function can be modeled with at least two different 1-way FSTs which differ in *when* they output the output symbols a , b . In the bottom row of Figure 4, we visualize the origin information created by the two FSTs for the mapping (ab, ab) as graphs called *origin graphs* (Bojańczyk *et al.* 2017). The FSTs model the same function and are equivalent in their general semantics of what they output; however, they are not equivalent in their origin semantics because they use different *origin information* for their outputs.

This notion of origin semantics can be used to contrast how 1-way FSTs and 2-way FSTs model reduplication. Consider the toy example of initial-CV reduplication with a small alphabet $\Sigma = \{p,a,t\}$. This function can be modeled by either a 1-way or 2-way FST as in Figure 5. The two transducers in Figures 5 are equivalent in their general semantics because they can output the same string. For example, given the input pat , both FSTs will output $pa\sim pat$. However, the two FSTs differ in their origin semantics for the mapping $pat \rightarrow pa\sim pat$. Setting aside boundary symbols \bowtie, \bowtie, \sim , the 1-way FST associates the second pa string of the output with the vowel a of the input as in the bottom middle column of Figure 5. This is because the second pa was outputted when the 1-way FST was reading the a in the input. In contrast, the 2-way FST associates each segment in the output with an identical segment in the input as in the bottom right column of Figure 5.

The origin information created by the 2-way FST matches theoretical treatments of how the reduplicant’s segments are individually associated with identical segments in the input (Marantz 1982; Inke-

Reduplication with 2-way FSTs

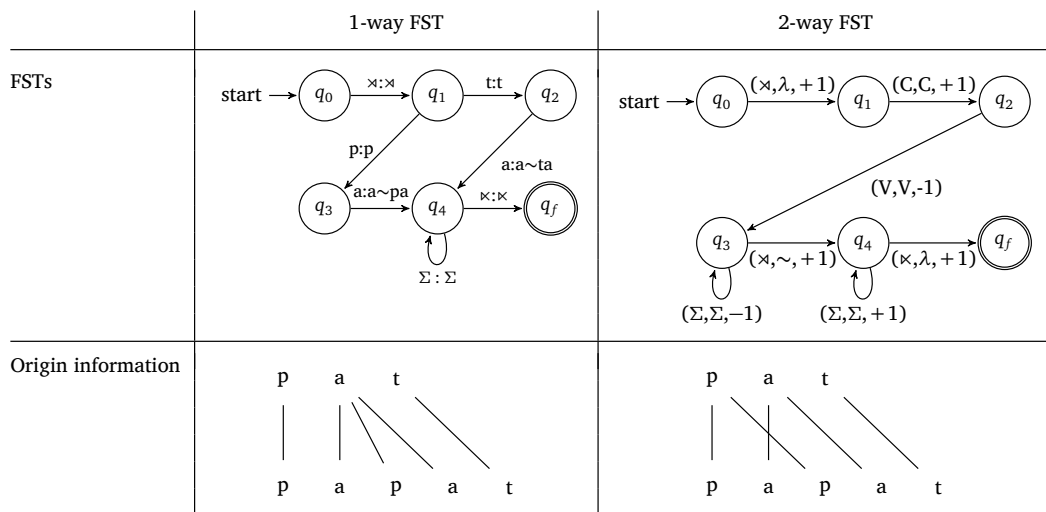


Figure 5: 1-way and 2-way FSTs for initial-CV reduplication and the origin information created by them for the mapping $pat \rightarrow pa\sim pat$

las and Zoll 2005).¹¹ In contrast, the origin information created by the 1-way FST does not match linguistic intuitions of reduplication because non-identical segments are associated. This difference in the origin semantics of the 1-way FST and 2-way FST formalizes their behavior: the 1-way FST simply *remembers* what strings of segments to output twice (Roark and Sproat 2007, 54), while the 2-way FST actively copies.

LINGUISTIC MOTIVATIONS
FOR SUBCLASSES OF TRANSDUCERS

5

Having shown the utility of 2-way D-FSTs for reduplication, the next two sections show that reduplication does not require the full power

¹¹ In Base-Reduplicant correspondence theory or BRCT (McCarthy and Prince 1995), what matters for reduplication is not the relationship or correspondence between the input and output segments, but between the two copies in the output. Origin semantics might be able to formalize the intuition behind BRCT with finite-state technology, e.g. output symbols with the same origin are in correspondence. The only computational implementation of BRCT to our knowledge (Albro 2000, 2005) uses MCFGs to do so. Note however that the empirical validity of BRCT is questionable (Inkelas and Zoll 2005; McCarthy *et al.* 2012).

of 2-way D-FSTs but falls within certain subclasses. This means that reduplication has a demarcable generative capacity or complexity. In this section, we discuss the subclasses of 1-way FSTs that have been proposed to model segmental phonology (Section 5.1), specifically the Output-Strictly Local (OSL) functions and Sequential (Seq) functions. In Section 5.2, we discuss subclasses of 2-way FSTs and design new subclasses based on the concatenation of OSL and Seq functions. We explain the intuition behind using concatenation-based subclasses for reduplication (Section 5.3). The next Section 6 goes over the typology of reduplication and shows how it fits into these subclasses.

5.1 *Computational typology of phonology and 1-way transducers*

It is known that 1-way finite-state machines can model all attested phonological processes (Johnson 1972; Kaplan and Kay 1994; Mohri 1997). However, phonological processes do not require the full power of 1-way finite-state machines (Heinz 2007; Chandlee 2014). Subclass hierarchies have been discovered for 1-way FSAs (McNaughton and Papert 1971; Rogers and Pullum 2011; Heinz and Idsardi 2013) and 1-way FSTs (Garcia *et al.* 1990; Gainor *et al.* 2012; Heinz and Lai 2013; Chandlee *et al.* 2014). Some of these subclasses have been argued to characterize different types of phonological well-formedness conditions and transformations (Heinz 2018; Chandlee and Heinz 2018; Chandlee *et al.* 2018). We give a brief and informal overview.

A common intuition in linguistic theory is that phonological processes are local or subject to adjacency constraints (Odden 1994). For example, a common phonological process is post-nasal voicing (5a) whereby voiceless stops are voiced after nasals. This process is local in the sense that the trigger for voicing (the nasal) is within a finite bound from the target of voicing (the stop). The symbols N, T, D represent nasals, voiceless stops, and voiced stops. Another local process is nasal spread whereby a vowel becomes nasalized after a nasal or nasalized vowel (5b). Nasal spread is iterative in that when a nasal triggers nasalization on a subsequent vowel, the newly nasalized vowel can then nasalize its subsequent vowel (5b-iii). The symbols V, \tilde{V} represent vowels and nasalized vowels.

- (5) a. **Post-nasal voicing**
 [+stop, -voice] → [+voice]/[+nasal] _ or T → D/N_
 i. /ata/ → [ata]
 ii. /anta/ → [anda]
- b. **Nasal spread**
 [+vowel] → [+nasal]/[+nasal] _ or V → \tilde{V} /\{N, \tilde{V} \}_
 i. /atapa/ → [atapa]
 ii. /anapa/ → [anãpa]
 iii. /anaapa/ → [anããpa]

Both processes are intuitively local. This intuition corresponds to Strict Locality in formal language theory (McNaughton and Papert 1971; Vaysse 1986; Rogers and Pullum 2011; Chandlee 2014). Formal definitions can be found in Chandlee *et al.* (2014, 2015). Informally, given an input string w , a function is Input-Strictly Local for a natural number k (k -ISL) if generating the output correspondent of some input symbol w_i relies on information about the current input symbol w_i and the $k - 1$ most recently seen input symbols. Post-nasal voicing is a 2-ISL function and it is computed by the 1-way FST in Figure 6. The symbol ? marks any other segment which isn't in an existing transition arc (Beesley and Karttunen 2003). The state labels are interpreted as keeping track of the last seen input symbol. The state labeled as N is where the post-nasal consonant is generated as voiced. The state label is interpreted as saying that a nasal was recently seen.

Output-Strictly Local functions for a number k (k -OSL) are analogously understood. A function is k -OSL if generating the output correspondent of w_i relies on information about the current input symbol w_i and the $k - 1$ most recently seen output symbols. An OSL function is L-OSL (R-OSL) if we read the input from the left (right), and write the output from the left (right). Nasal spread is a 2-L-OSL function and is computed by the 1-way FST in Figure 6. The state labels are interpreted as keeping track of the last recently generated output symbol. The state labeled as N, \tilde{V} is where a nasalized vowel is generated; the state label is interpreted as saying that the most recently outputted symbol was a nasal or nasalized vowel.

track of the last 2 symbols on the output tape and the current input symbol. The 3-OSL 1-way FST function in Figure 8 outputs up until the first VC of the input; it then stops outputting anything after that.¹³

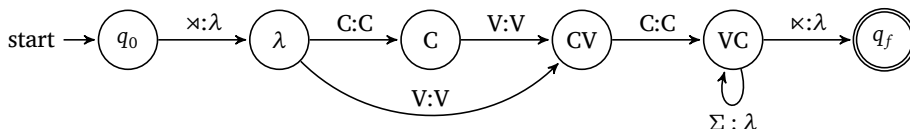


Figure 8: OSL 1-way FST for English nickname formation

A significant proportion of segmental phonology can be modeled with ISL and OSL functions (Chandlee 2014; Chandlee and Heinz 2018; Chandlee *et al.* 2018). Long-distance processes in phonology are however neither ISL or nor OSL. For example, Kikongo nasal harmony (7) requires the higher subclass of Sequential functions (Gainor *et al.* 2012). In Kikongo, alveolar stops like *d* or *l* surface as *n* if a nasal precedes them anywhere in the input. There can be any number of vowels and consonants intervening between the triggering nasal and the target alveolar (7c). This long-distance information means that the 1-way FST must keep track of whether a nasal consonant was seen *anywhere* in the input stem before it will output the alveolar.

(7) **Kikongo nasal harmony**

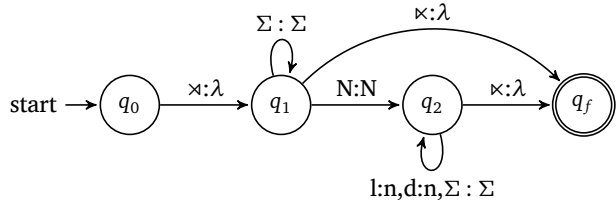
- | | | | |
|----|-------------|---------------|-----------------------|
| a. | /sakid-ila/ | → [sakid-ila] | ‘to congratulate for’ |
| b. | /mant-ila/ | → [mant-ina] | ‘to climb for’ |
| c. | /tunik-idi/ | → [tunik-ini] | ‘we ground’ |

The above pattern cannot be modeled by an ISL or OSL function but requires a Sequential 1-way FST as in Figure 9. A Seq 1-way FST is a deterministic 1-way FST that will read the input in only one direction (here left-to-right) and can use any information that it had found in the input string when processing the next input symbol.

To summarize, different types of phonological processes are computed by different subclasses of rational functions and with different subclasses of 1-way FSTs. The relatively low complexity of this subclasses has opened doors to understanding the cognitive limitations and learnability of phonological processes (Heinz 2018). In the next

¹³See Chandlee (2017) on why this function is necessarily OSL and not ISL.

Figure 9:
Sequential 1-way FST
for Kikongo nasal harmony



section, we show that extending OSL and Sequential functions into 2-way FSTs opens similar doors for the typology of reduplication.

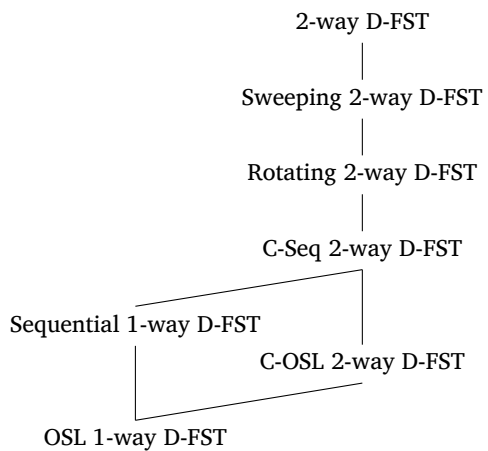
5.2 Subclasses of 2-way finite-state transducers

Unlike for 1-way FSTs, there is much less work on subclasses for 2-way FSTs. Some intuitive subclasses have been proposed in the literature. The typology of reduplication inspired us to devise additional subclasses. All of these subclasses, shown in Figure 10, restrict:

1. where the 2-way FST can rewind in the input,
2. what it can output while it is rewinding, and
3. what information can be transferred across multiple passes, i.e., if a later pass depends on an earlier pass.

At the top of the hierarchy are 2-way D-FSTs which correspond to regular functions. In regards to the first restriction, a 2-way FST is a **sweeping** transducer if the read head can change direction only at

Figure 10:
Hierarchy of subclasses
for 2-way FSTs



the ends of the input (Baschenis *et al.* 2015, 2016, 2018). A sweeping transducer is a generalization of similarly defined sweeping automata (sweeping 2-way FSAs) (Sipser 1980). For example, the 2-way FST for total reduplication in Figure 1 is a sweeping transducer. The only time the FST moves right-to-left is going from the end boundary \times to the start boundary \times . In contrast, the 2-way FST for initial-CVC partial reduplication in Figure 2 is not a sweeping transducer. It rewinds from the third input segment C to the beginning \times . However, the partial reduplication function computed by this 2-way FST can be computed by a sweeping transducer, which we show in the next section.

2-way D-FSTs are more expressive than deterministic sweeping 2-way D-FSTs. Consider the function $u_1\#u_2\#\dots\#u_n \mapsto u_n\dots u_2u_1$ where the input is a sequence of strings u_i separated by the special symbol $\#$. The output is formed by reversing these strings and deleting the $\#$'s. This function can be computed by a deterministic 2-way D-FST but not by a sweeping transducer. See Baschenis *et al.* (2016) for discussion.

In regards to the second restriction, a sweeping transducer is a **rotating** transducer if it does not output anything while it's moving right-to-left (Baschenis *et al.* 2017). The 2-way FST for total reduplication is a rotating transducer because it outputted nothing while moving right-to-left. Sweeping transducers are more expressive than rotating transducers. A sweeping transducer can compute the mirror function $w \rightarrow ww^r$, but a rotating transducer cannot.

As for the third restriction, we develop a set of concatenated-based subclasses of functions.

(8) *Subclasses of Regular Functions*

- a. **C-Seq function:** A Concatenated-Sequential function f is the concatenation of n Sequential functions s_i , e.g. $f(x) = s_1(x) \cdot s_2(x) \cdot \dots \cdot s_n(x)$. f is C-L-Seq (C-R-Seq) if the component Seq functions read the input left-to-right (right-to-left).¹⁴

¹⁴In terms of function combinatorics for regular string transformations (Alur *et al.* 2014; Dave *et al.* 2018), the class of C-Seq functions involves the use of a 'sum combinator' \otimes that concatenates the output of two or more Seq functions: $f(x) = s_1(x) \otimes s_2(x) \otimes \dots \otimes s_n(x)$ where s_i is a Seq function. This is similar to the use of product automata. See Alur *et al.* (2014) for details.

- b. **C-OSL function:** A Concatenated-OSL function f is the concatenation of n Output-Strictly Local functions o_i , e.g. $f(x) = o_1(x) \cdot o_2(x) \cdot \dots \cdot o_n(x)$. f is C-L-OSL (C-R-OSL) if the component OSL functions read the input left-to-right (right-to-left).

Rotating transducers are more expressive than C-Seq transducers, which are more expressive than C-OSL transducers. Examples witnessing these separations are drawn from the typology of reduplication in Section 6. C-Seq functions are more expressive than sequential functions (= 1-way D-FSTs) which are more expressive than OSL functions. A set of definitions is provided below for easier reference.

(9) *Subclasses of 2-way D-FSTs*

- a. **Sweeping 2-way FST:** A 2-way FST which can change direction only at the ends of the input
- b. **Rotating 2-way FST:** A sweeping 2-way FST which outputs nothing while moving right-to-left
- c. **C-Seq 2-way FST:** A rotating 2-way FST that computes a Concatenated Sequential function
- d. **C-OSL 2-way FST:** A rotating 2-way FST that computes a Concatenated Output-Strictly Local function

We have found that virtually the entire typology of reduplication can be modeled with deterministic rotating 2-way FSTs. Further, the bulk of the typology can be modeled with C-OSL functions. A few minor cases require C-Seq functions; these mostly involve infixal or internal reduplication. A smaller set of cases require the full power of rotating transducers; though these cases are not clear-cut. Before going through the typology in detail in Section 6, we illustrate the insight behind C-OSL functions and how they compute reduplicative processes.

5.3

Illustrating C-OSL

Intuitively, a C-OSL function is a function that takes as input the string x , gives x to n many separate 1-way FSTs which are OSL, and concatenates their output. To illustrate the insight behind C-OSL functions,

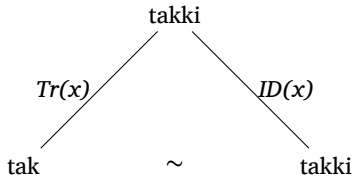


Figure 11:
Initial-CVC reduplication
as a concatenation of functions

consider initial-CVC partial reduplication from Agta again (Moravcsik 1978, 311) from (1a) repeated in (10a).

- (10) a. takki \rightarrow tak~takki ‘leg’ \rightarrow ‘legs’
b. takki \rightarrow takki~takki

As an input-to-output function, reduplication may be viewed as submitting the same input to two separate functions in parallel and concatenating their output as in Figure 11. The first function, here labeled $Tr(x)$, truncates the input to the first CVC while the second function, $ID(x)$, is the identity function. The outputs of these two functions, *tak* and *takki*, are concatenated to form the reduplicated output: *tak~takki*. In (10b), we explicitly show how initial-CVC reduplication can be seen as truncating the first copy: *takki* \rightarrow ~~*takki*~~~*takki* where truncated material is shown in strike-through.

The truncation function $Tr(x)$ is a 3-L-OSL function because it outputs a truncation of the input to just the first CVC. This is similar to English nickname formation from Section 5.1. The identity function $ID(x)$ is both 1-L-OSL and 1-R-OSL. Thus both $Tr(x)$ and $ID(x)$ are L-OSL and hence their concatenation is C-OSL. Figure 12 illustrates a 2-way D-FST for initial-CVC reduplication which is formulated as a concatenation of these two OSL functions. Contrast this model of initial-CVC reduplication (shown in Figure 12) with the non-rotating 2-way D-FST shown in Figure 2. (It is the the additional state CV_1 and its transition arcs in Figure 12 which make this D-FST rotating.)

To summarize this section, understanding reduplicative processes as C-OSL and C-Seq functions is intuitive. This analysis echoes Steriade (1988)’s treatment of partial reduplication as total reduplication followed by truncation and Inkelas and Zoll (2005)’s treatment of total reduplication as morphological doubling.

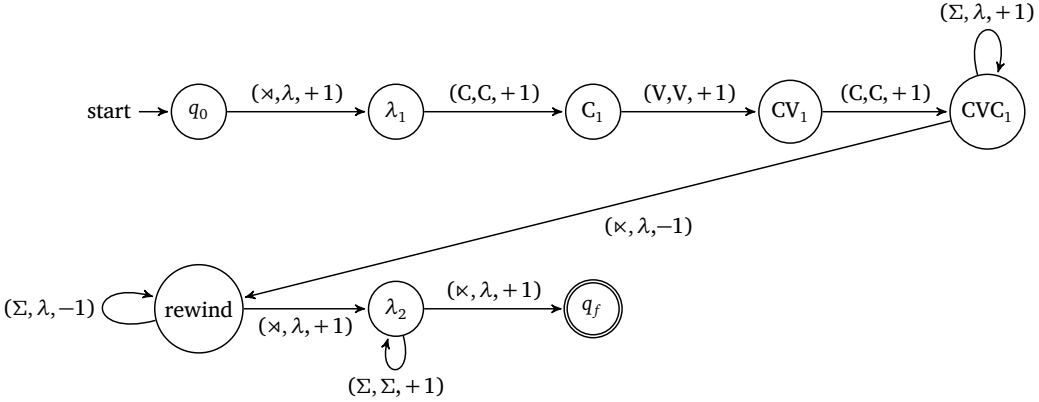


Figure 12: C-OSL 2-way D-FST for initial-CVC partial reduplication

6 COMPUTATIONAL TYPOLOGY OF REDUPLICATION

This section provides a detailed, comprehensive review of the RedTyp typology and classifies RedTyp’s reduplicative processes according to the computational classification introduced in the last section. The main finding is that most processes are C-OSL (Section 6.2) and most of the ones that are not are C-Seq (Section 6.3). There are few (and questionable) cases where reduplication needs the full power of 2-way D-FSTs (Section 6.4). We give an overview in Section 6.5.

Note that all partial reduplicative processes can be computed by 1-way FSTs. However, in order to get the linguistically-motivated origin semantics right (Section 4.3), we need the additional power of 2-way FSTs. Thus in this section, when we discuss how partial reduplicative cases fit into the subclasses of 2-way FSTs, we mean in terms of them generating the right origin semantics.

6.1 Preliminaries to the typology

Although reduplication is cross-linguistically ubiquitous, there is a wide cross-linguistic variation in a) what substring or subsequence gets repeated, b) where the copied substring or subsequence is placed

in the output, and c) whether and how phonological processes interact with copying. This section provides a brief but representative typology of reduplication compiled from various surveys (Moravcsik 1978; Rubino 2005; Inkelas and Downing 2015a).

We emphasize that our reported typology is descriptive and not theoretical. Various theoretical frameworks have been developed to account for the range of variation on reduplication (Marantz 1982; McCarthy and Prince 1995; Spaelti 1997; Raimy 2000; Inkelas and Zoll 2005; Frampton 2009; Samuels 2010; McCarthy *et al.* 2012; Saba Kirchner 2010, 2013). The reader is referred to elsewhere for theoretical overviews (Raimy 2011; Urbanczyk 2007, 2011; Inkelas and Downing 2015a,b).

We define the following *descriptive* terms which will be useful in categorizing different reduplicative processes:

(11) *Terminology for categorizing the typology:*

- **reduplicant:** the substring in the output which was created via copying
- **base:** the substring in the output which was not created via copying
- **target:** the substring in the input which will be copied or duplicated
- **anchor point:** the position in the input where the target starts or ends
- **source:** the morphological or phonological constituent in the input which contains the target

The output-based terms *base* and *reduplicant* are common in the literature on reduplication (McCarthy and Prince 1995) though their definition is problematic (Shaw 2005; Haugen 2009). Anchor points have been proposed for reduplication (Fitzpatrick 2006; Raimy 2009) and other non-concatenative processes (Yu 2007; Samuels 2010). We introduce the input-based terms *source* and *target* in order to better fully describe reduplication as an input-to-output function. This section goes through the typology of reduplication, organized in terms of how they vary in the source, target, and/or reduplicant. These variations align with what type of 2-way FST is needed to compute them.

6.2

Most reduplication is C-OSL

Most reduplicative processes are C-OSL. We go through common and some uncommon reduplicative processes and show they are C-OSL. For a function to be C-OSL, the two copies must be independent of each other, and the two passes over the input must likewise be independent of each other. Informally, some criteria for a C-OSL function are that each of the component functions:

1. reference only a finite and bounded number of the most recently generated output symbols, meaning that each of the functions,
2. do not depend on any long-distance information in the input,
3. do not use any finite lookahead or finite lookback on the input,
4. do not rely on deleted material, and
5. do not rely on any information from the other function.

6.2.1

Total and word-initial partial reduplication

Total reduplication and word-initial partial reduplication are C-L-OSL, which means they are the concatenation of two L-OSL functions.

Consider total reduplication first.

- (12) a. **Total reduplication**
Indonesian (Cohn 1989)
buku → buku~buku ‘book’ → ‘books’
- b. **Initial-CVC reduplication**
Agta (Moravcsik 1978, 311)
takki → tak~takki ‘leg’ → ‘legs’

For total reduplication, the two OSL functions are identity $ID(x)$. (Total reduplication is also C-R-OSL.)

For partial reduplication, there is limited variety in the shape of the copied material, the *reduplicant*. Some languages have a partial reduplicative process that copies the first *C* or consonant of the word (13a), first CV or consonant-vowel sequence (CV) of the word (13b), first CVC sequence (13c), or first CV(C)CV sequence (13d). In general, the copied material has to fit into some template of a particular size.

(13) *Common prefixal partial reduplicative patterns*

a. **Initial-C reduplication**

Shilh (Moravcsik 1978, 308)

gen → g~gen ‘to sleep’ → ‘to be sleeping’

b. **Initial-CV reduplication**

Sundanese (Moravcsik 1978, 319)

guyon → gu~guyon ‘to jest’ → ‘to jest repeatedly’

c. **Initial-CVC reduplication**

Panganisan (Rubino 2005, 11)

baley → bal~baley ‘town’ → ‘towns’

d. **Initial-CV(C)CV reduplication**

Dyirbal (McCarthy *et al.* 2012, 187)

a. baniju → bani~baniju ‘come’

b. balgan → balga~balgan ‘laugh’

As for the partial reduplication functions in (13), they are all C-L-OSL just like initial-CVC reduplication from Section 5.3. They involve the concatenation of a truncation $Tr(x)$ and identity function $ID(x)$. The truncation function varies in terms of how much word-initial material is faithfully outputted.

Table 3 illustrates these examples where the truncated material is shown in strike-through. The outputs of the two component OSL function are separated by ~ for illustration.

Table 3: C-OSL treatment for total and initial partial reduplication

	Total (12a)	Initial-C (13a)	Initial-CV (13b)	Initial-CVC (13c)	Initial-CV(C)CV (13d)
Input x	buku	gen	guyon	baley	balgan
Components	$ID(x) \cdot ID(x)$	$Tr(x) \cdot ID(x)$	$Tr(x) \cdot ID(x)$	$Tr(x) \cdot ID(x)$	$Tr(x) \cdot ID(x)$
Output	buku~buku	gen~gen	guyon~guyon	baley~baley	balgan~balgan
Subclass	C-L-OSL C-R-OSL	C-L-OSL	C-L-OSL	C-L-OSL	C-L-OSL

Variation in the number and placement of copies

6.2.2

The typology is larger than the above examples. Some languages create three copies of the input (triplication) instead of just two (14a).

Some reduplicative processes are suffixal; they specify that the location of the target be a word-final substring (14b) instead of word-initial substring (13d). Some reduplicative processes are wrong-sided by making the target and the reduplicant not adjacent in the output, i.e., copying the final CVC and placing it at the beginning of the output (14c vs. 13c). There are likewise cases where both the base and the reduplicant are shortened or truncated in the output, e.g., truncating both copies to CV (14d).

(14) *Variation in number and reduplicant placement*

a. **Total triplication**

Mokilese (Moravcsik 1978, 301)

roar → roar~roar~roar

‘give a shudder’ → ‘continue to shudder’

b. **Final-CVCV reduplication**

Siriono (Moravcsik 1978, 308)

erasi → erasi~rasi

‘he is sick’ → ‘he continues being sick’

c. **Initial-CVC reduplication and opposite-edge or wrong-sided placement**

Koryat (Riggle 2004, 3)

qanga → qanga~qan ‘fire’ → ‘fire (ABS)’

d. **Abbreviated reduplication (*Kager-Hamilton Problem*)**

Guarijio (Caballero 2006)¹⁵

toni → to~to ‘to boil’ → ‘to start boiling’

muhiba → mu~mu ‘to throw’ → ‘to start throwing’

All these processes are still C-OSL, however. They differ in the number and order of concatenated functions, the direction in which the input is read, and whether all or none of the functions are identity. Their computation is visualized in Table 4. Triplication is C-L-OSL and C-R-OSL; it involves concatenating three identity functions. Suffixal reduplication like final-CVCV reduplication is C-R-OSL because

¹⁵Such reduplication is often argued to be unattested and is called the *Kager-Hamilton Problem* (Idsardi and Raimy 2008). See Caballero (2006) for discussion on what prosodic and morphological factors condition this rare type of reduplication.

Table 4: C-OSL treatment for less common reduplication patterns

	Triplication (14a)	Final-CVCV (14b)	Wrong-sided (14c)	Abbreviated (14d)
Input x	roar	erasi	qanga	toni
Components	$ID(x) \cdot ID(x) \cdot ID(x)$	$ID(x) \cdot Tr(x)$	$ID(x) \cdot Tr(x)$	$Tr(x) \cdot Tr(x)$
Output	roar~roar~roar	erasi~erasi	qanga~qanga	toni~toni
Subclass	C-L-OSL C-R-OSL	C-R-OSL	C-L-OSL	C-L-OSL

it is the concatenation of an identity function and an R-OSL truncation function. The truncation function reads the input right-to-left and deletes everything to the left of the final CVCV. Wrong-sided initial-CVC reduplication is C-L-OSL. It differs from initial-CVC reduplication by ordering identity before truncation. Abbreviated reduplication is C-L-OSL. Unlike initial-CV copying, it is composed of two L-OSL truncation functions instead of just one.

Copying a morphological subconstituent

6.2.3

In the above examples, the *source* was the entire input. But unlike concatenative morphology, reduplication is often sensitive to word-internal morphological constituents, contra Bracket Erasure (Kiparsky 1982). In these cases, the semantic function of reduplication builds on the meaning of the entire input while the location of the reduplicant is word-internal (cf. Aronoff 1988). For example, some languages have reduplication target a morphological subconstituent within the input as the source, such as a root/stem (15a, 15b) or affix (15c), and whether for total reduplication (15a, 15c) or partial reduplication (15b). The source and reduplicant are usually adjacent; though there are some cases where the two copies are non-adjacent in the output, e.g., Madurese copies the root-final CVC and places it at the beginning of the output (15d).

(15) Copying from a morphological subconstituent

a. Total reduplication of the stem

KiHehe (Aronoff 1988, 8)

ku-haata → ku-haata~haata

'to ferment' → 'to start fermenting'

Table 5: C-OSL treatment for copying morphological subconstituents

	KiHehe	Bikol	Hungarian	Madurese
	(15a)	(15b)	(15c)	(15d)
Input x	$ku\{r,haata\}_r$	$na\{r,murak\}_r$	$_p\{el\}_p\text{megy}$	$pa\{r,jalan\}_r,an$
Components	$L(x) \cdot R(x)$	$L(x) \cdot R(x)$	$L(x) \cdot R(x)$	$L(x) \cdot R(x)$
Output	$ku\{r,haata\}_r \sim$ $k\ddot{u} \{r,haata\}_r$	$na\{r,murak\}_r \sim$ $n\ddot{a} \{r,murak\}_r$	$\{p,el\}_p\text{megy} \sim$ $\{p,el\}_p\text{megy}$	$p\ddot{a}\{r,jalan\}_r,an \sim$ $pa\{r,jalan\}_r,an$
Subclass	C-L-OSL C-R-OSL	C-L-OSL	C-L-OSL C-R-OSL	C-R-OSL

b. Initial CV reduplication of the stem

Bikol (Mattes 2007, 84)

$na\text{-}murak \rightarrow na\text{-}mu\sim murak$

‘to flower’ \rightarrow ‘decorating with flowers’

c. Total reduplication of an affix (prefix)

Hungarian (Inkelas and Downing 2015a, 505)

a. $el\text{-}megy \rightarrow el\sim el\text{-}megy$

‘he goes there’ \rightarrow ‘he occasionally goes there’

b. $bele\text{-}nez \rightarrow bele\sim bele\text{-}nez$

‘he looks into it’ \rightarrow ‘he occasionally looks into it’

d. Root-final CVC reduplication and word-initial placement

Madurese (Brown 2017, 964)

$pa\text{-}jalan\text{-}an \rightarrow lan\sim pa\text{-}jalan\text{-}an$

‘pedestrian’ \rightarrow ‘pedestrians’

More cases of reduplication targeting a morphological subconstituent are well-attested (Shaw 2005; Inkelas and Zoll 2005; Haugen 2009; Hyman 2009; Inkelas 2014). The above cases are C-OSL if the relevant morphological boundaries are present in the input. Their computation is visualized in Table 5. Each process uses two functions $L(x), R(x)$ which generate the two copies, reference the morphological boundaries, and they crucially output these boundaries.

For total copying in KiHehe and Hungarian, the function is C-L-OSL and C-R-OSL. For total stem copying in KiHehe, the first function $L(x)$ outputs everything up until the root right-boundary ‘}’.

ond function $R(x)$ outputs nothing until it sees the root left-boundary $\{r\}$; it outputs this and everything after it. Both functions are both L-OSL and R-OSL, thus KiHehe is C-L-OSL and C-R-OSL. Prefix copying in Hungarian is similarly defined but for the prefix boundaries $\{p\}$ and $\}p\}$. Partial stem copying in Bikol is only C-L-OSL. The function $L(x)$ outputs everything up until it outputs the string $\{rCVC\}$; it deletes everything after that. The function $R(x)$ outputs nothing up until it sees the root left-boundary $\{r\}$; it outputs this and everything after it. Non-local copying in Madurese is C-R-OSL. The function $L(x)$ reads the input right-to-left; it outputs nothing until it sees the root right-boundary $\}r\}$; it outputs this and the first CVC that it sees. After that, it deletes everything. The function $R(x)$ is the identity function.

Even though some have argued against the use of morpheme boundaries in morpho-phonological representations (Anderson 1992; Stump 2001), morphological boundaries must be part of the input for a finite-state systems like ours (e.g. Karttunen 1983; Koskeniemi 1984; Roark and Sproat 2007).¹⁶ Consider partial stem copying in Bikol: $na\{r,murak\}_r \rightarrow na\{r,mu\{r,murak\}_r$. Without the root boundaries, we could not distinguish the prefixed input $na\{r,murak\}_r$ from a hypothetical mono-morphemic input $namurak$. Without some way to encode the relevant morphological constituents in the input, we simply cannot define this reduplication function with any type of 1-way or 2-way FSTs. The use of such boundaries in finite-state morphology is standard practice.

Copying a prosodic subconstituent

6.2.4

Besides morphological subconstituents, the source can also be a metrical or prosodic subconstituent such as the stressed syllable (16a), the first syllable (16b), or the first foot (16c). The source can also

¹⁶It should be noted though that HPSG-based approaches to computational morphology (Bonami and Crysmann 2013, 2016; Crysmann and Bonami 2016) do not need morpheme boundaries as symbols in their alphabet. One reason is because they can essentially directly capture hidden morphological structure or constituency. Another strategy is to temporally apply reduplication to the subconstituent and then later add the other affixes, e.g., $haata \rightarrow haata\sim haata \rightarrow ku\sim haata\sim haata$. This a common strategy in handling morphology-semantics bracketing paradoxes (cf. Stump 1995, 2001).

be a morphophonological constituent, e.g. a prosodic stem (16d). In Chumash, the prosodic stem (underlined) consists of all the segments in the morphological stem alongside any prefixal consonants that are syllabified with the morphological stem.

- (16) *Copying from a metrical or prosodic subconstituent*
- a. **CV-reduplication of the stressed syllable**
Chamorro (Inkelas and Downing 2015a, 507)
 hu.gán.do → hu.gá~gan.do ‘play’ → ‘playing’
 - b. **Total reduplication of the initial syllable**
Hiaki (Haugen 2009)
 vu.sa → vu~vu.sa ‘awaken’
 vam.se → vam~vam.se ‘hurry’
 - c. **Total reduplication of the initial foot**
Yidiny (Marantz 1982, 453)
 (gindal)ba → gindal~gindalba ‘lizard sp.’ → ‘lizards’
 - d. **Initial-CVC reduplication of the prosodic stem**
Chumash (Downing 1998, 101)
 s + tʃeq → s-tʃeq~tʃeq ‘it is very torn’
 s + ikuk → s + ik~s-ikuk ‘he is chopping’

If the relevant prosodic boundaries are in the input, the computation is C-OSL. The computation proceeds the same as for copying a morphological constituent. Table 6 shows this with syllable boundaries $(_s)_s$, foot boundaries $(_f)_f$, stressed syllable boundaries $(_s)_s$, and prosodic stem boundaries $(_{PS})_{PS}$.¹⁷

Given an unsyllabified input, these prosodic boundaries can be generated via a 1-way FST (Hulden 2006; Yu 2017) which can be ISL because it uses finite lookahead on the input (see also Strother-Garcia 2018, 2019).

However, if the input to reduplication lacks boundaries, then reduplication is C-Seq because we need finite lookahead to know if some consonant is part of the relevant prosodic constituent. Consider

¹⁷The second function $R(x)$ in stressed syllable CV-copying must change stressed \acute{a} to unstressed a . Stressed syllable copying is also C-R-OSL if the first function $L(x)$ is R-OSL and outputs the right-boundary $)_s$. Generating the prosodic stem in Chumash requires reference to morphological boundaries too.

Table 6: C-OSL treatment for copying prosodic subconstituents

	Chamorro (16a)	Hiaki (16b)		Yidiny (16c)	Chumash (16d)
Input x	hugándo	vusa	vamse	gindalba	s + ikuk
Syllabify x	hu(_s gán) _s do	(_s vu) _s (_s sa) _s	(_s vam) _s (_s se) _s	(_f gindal) _f ba	(_{ps} s + ikuk) _{ps}
Components	$L(x) \cdot R(x)$	$L(x) \cdot R(x)$	$L(x) \cdot R(x)$	$L(x) \cdot R(x)$	$L(x) \cdot R(x)$
Output	hu(_s gá) _s ndə ~ h̥u(_s gan) _s do	(_s vu) _s (_s sa) _s ~ (_s vu) _s (_s sa) _s	(_s vam) _s (_s se) _s ~ (_s vam) _s (_s se) _s	(_f gindal) _f b̥a ~ (_f gindal) _f ba	(_{ps} sik uk̥) _{ps} ~ (_{ps} sikuk) _{ps}
Subclass	C-L-OSL C-R-OSL	C-L-OSL	C-L-OSL	C-L-OSL	C-L-OSL

initial syllable copying in Hiaki $vusa \rightarrow vu \sim vusa$. In the first function $L(x)$, the consonant s is not generated because it is part of the next syllable. We know because it precedes a vowel. In contrast for $vamse \rightarrow vam \sim vamse$, the consonant m is copied because it precedes a consonant. The use of such information from finite lookahead on the input cannot be computed by an OSL function.

This section presented cases in RedTyp which are C-OSL. They comprise the bulk of reduplicative typology. Of the 138 reduplicative processes in RedTyp (Section 4.2), 121 (87%) were C-OSL.

Some reduplication is C-Seq

6.3

This section goes through some types of reduplication which are not C-OSL but are instead C-Seq. Informally, a reduplicative function is C-Seq if its component functions do not rely on any information from the other function or its output. A component function can use finite lookback, finite lookahead, or even long-distance information in the input.

Internal reduplication and gray areas between C-OSL or C-Seq

6.3.1

One problematic area for C-OSL are internal reduplication functions which seem infixal (Broselow and McCarthy 1983; Gafos 1998; Spaelti 1997). Some of these are C-OSL, some are not. These functions are C-OSL if the truncation functions can uniquely determine what segments to delete based on only what was outputted.

In the previous sections, the reduplication's target can be thought of as a contiguous substring that is determined by scanning either the left or right edge of the source. In those cases, the target was at the edge of the source and the reduplicant was placed at the left or right edge of the base. However, there are cases of *internal* or *infixal* reduplication where the target is a substring *inside* the source such that the substring is not strictly adjacent to the source's edges (17a, 17b) and the reduplicant is placed inside the base in the output (17c). In (17c), the word-initial C is copied and placed after the first vowel.

(17) *Internal reduplication cases which are arguably not C-OSL*

a. **Leftmost-VCC* reduplication**

Mangarayi (Raimy 2000, 135)

gabuji → g-ab~abuji 'old person' → 'old persons'

b. **Rightmost-CV reduplication**

Chamorro (Inkelas and Zoll 2005, 107)

nalaŋ → nala~laŋ 'hungry' → 'very hungry'

c. **Initial-C reduplication and internal placement**

Quileute (Broselow and McCarthy 1983, 44)

t^siko → t^si~tko 'he failed sp.' → 'he failed (freq.)'

Table 7 visualizes these processes. A traditional analysis is that the reduplicant is infixal (Broselow and McCarthy 1983) where < > marks infixation, e.g., Mangarayi *gabuji* → *gab*<*ab*>*buji*. However, their treatment as C-OSL is somewhat counter-intuitive because a C-OSL function models these processes as concatenating two truncation functions: *gabuji*~*gabuji*. The first function *L(x)* outputs the first C*VC* substring and deletes everything after that. The second function *R(x)* deletes all word-initial strings of consonants C*; once it sees a vowel V, it outputs it and everything after it.

6.3.2 Internal or non-contiguous reduplication which is C-Seq

The main reason why Mangarayi and the other functions in Table 7 are C-OSL is because the target and deleted materials do not have the same shape. Knowing what to delete or generate doesn't need any finite lookahead or lookback over the input, just over the output. However, other cases of internal and non-contiguous reduplication do require such finite lookback/lookahead over the input. This makes them

Case	Mangarayi (17a)	Chamorro (17b)	Quileute (17c)
Input <i>x</i>	gabuji	nalaŋ	t ^s iko
Output	gababuji	nalalaŋ	t ^s itko
Infixed treatment	gab<ab>uji	nala<la>ŋ	t ^s i<t>ko
C-OSL treatment	gabuji~gabuji	nalaŋ~nalaŋ	t ^s iko~tiko
Subclass	C-L-OSL	C-R-OSL	C-L-OSL

Table 7:
Infixed vs. C-OSL
treatment of
internal
reduplication

C-Seq. In (18a), the penultimate syllable is reduplicated. In (18b), the word-initial CV is copied and placed before the final C. In most of these cases, the target is a contiguous substring in the input. In some cases, the target is not contiguous (18c). In (18c), the input's first CV and final C are copied and placed together at the beginning of the output.

(18) *Internal reduplication which are C-Seq*

a. **Penultimate syllable reduplication**

Samoan (Moravcsik 1978, 301,310)

a.lo.fa → a.lo~.lo.fa 'he loves' → 'they love'

ta.o.to → to.o~o.to 'he lies' → 'they lie'

b. **Initial-CV reduplication and internal placement**

Creek (Riggle 2004, 3)

fayatk + i: → fayat~fa-k + i: 'crooked' → 'crooked (pl.)'

c. **Double-sided reduplication**

Nisgha (Urbanczyk 2007, 474)¹⁸

lú:t'ux^w → lúx^w~lút'ux^w 'to value' → 'to value (pl.)'

These C-Seq processes are visualized in Table 8. Consider penultimate syllable copying in Samoan with two truncation functions *a.lo.fa* → *a.lo.fa~a.lo.fa*. If the input is read left-to-right, the first function must output everything up until the penultimate syllable: *a.lo.fa*. This is not OSL because knowledge about whether some syllable is penultimate or not requires finite lookahead on the input. If the input is instead read right-to-left, the first truncation function is still not OSL. The function would delete the last vowel and the last consonant; but once it sees the penultimate vowel, the function cannot

¹⁸We set aside issues in predicting the quality of the vowel (Shaw 2005).

Table 8:
Non-C-OSL
patterns of
internal or
non-contiguous
reduplication

Case	Samoan (16a)	Creek (18b)	Nisgha (18c)
Input x	a.lo.fa	fayatk	lú:t'ux ^w
Components	$L(x) \cdot R(x)$	$L(x) \cdot R(x)$	$L(x) \cdot R(x)$
Infixal treatment	a.lo < lo > fa	fayat < fa > k	lúx ^w ~ lút'ux ^w
C-Seq treatment	a.lo.fa ~ a.lo.fa	fayatk ~ fayatk	lút'ux ^w ~ lút'ux ^w
Subclass	C-Seq	C-Seq	C-Seq

determine if this vowel is penultimate or not based on only what it has outputted. The function would need lookback access to the input. In both left-to-right and right-to-left cases, the truncation functions are Seq. The other processes in Table 8 are not C-OSL for similar reasons.

Although penultimate syllable copying is not C-OSL, Samoan has penultimate stress (Zuraw *et al.* 2014). Stressed syllable copying is C-OSL. This is an argument for reanalyzing Samoan as instead copying the stressed syllable. This relates to Nelson's (2003, 117) hypothesis that any references to the penultimate position in reduplication must be prosodic. Note that for Creek and Nisgha, the component functions are Seq; however they do not generate the right origin semantics because of unbounded word-internal deletion. The first function for Nisgha deletes everything except for the first CV and last C $lut'ux^w$. If read left-to-right, in order to generate the final C, we need to move to the next symbol and check if it is \times ; thus the final C is generated as an output correspondent for the end-boundary \times . If read left-to-right, in order to generate the initial CV, we move on to the preceding symbol and check if it is \times ; thus the first CV are generated as output correspondents to the start-boundary \times . This is because Seq functions are deterministic. In contrast, the right origin semantics would be generated if the component functions were non-deterministic or if the function was computed by a full 2-way FST.¹⁹

¹⁹To exactly capture the right origin semantics, it is possible that a subclass of *Streaming-string transducers* (1-way FSTs with registers) (Alur and Černý 2011; Alur and Deshmukh 2011) are a suitable alternative for modeling infixal reduplication. Discovering subclasses of SSTs and their relations to subclasses of 2-way D-FSTs is a worthwhile open question.

The above cases with infixation showed that capturing the right origin semantics might require classes which are more expressive than C-OSL and C-Seq. In this section, we go through more cases. Some are ambiguously C-Seq depending on the analysis; others must use rotating or even unrestricted 2-way D-FSTs in order to capture the right origin semantics.

Reduplication with syllable-count

6.4.1

Reduplication that is sensitive to syllables may involve iteration (19a) or minimality requirements on what is reduplicated (19b). Both examples are from Mandarin for different reduplicative processes. In (19a), reduplication is iterative because each syllable undergoes total reduplication: an input of the form A.B has A.A.B.B as the output. In (19b), a word undergoes total reduplication if it is monosyllabic, otherwise the morpheme *meei* is added.

(19) *Reduplication and syllable number*a. **Iterative reduplication of syllables**

Mandarin (Moravcsik 1978, 314)

huang.jang → huang~huang-jang~jang

‘flustered’ → ‘flustered (vivid form)’

b. **Minimality in reduplication**

Mandarin (Moravcsik 1978, 305-6).

jang → jang~jang ‘sheet’ → ‘every sheet’

jia.luen → meei-jia.luen ‘gallon’ → ‘every gallon’

Iterative copying (19a) is C-OSL if the number of iterations (= number of syllables) is bounded. The individual syllables must also be uniquely identifiable in the input. For Mandarin, the function is made up of four concatenated OSL truncation functions. The first two functions output everything up until the medial syllable boundary ‘.’ while the latter two delete everything up until the ‘.’ boundary.

(20) *Mandarin iterative reduplication as C-OSL:*

$$L1(x) \cdot L2(x) \cdot R1(x) \cdot R2(x)$$

huang.jang~huang.jang~huang.jang~huang.jang.

We are unaware of any examples showing reduplicative processes which iterate over inputs that have at least three syllables. So while Mandarin provides examples of $A \rightarrow A\sim A$ and $A.B \rightarrow A.A\sim B.B$, we have no examples of $*A.B.C \rightarrow A.A\sim B.B\sim C.C$. We likewise have not seen cases of trisyllabic iterative copying in other languages. This is computationally significant. If Mandarin allowed iterative copying over trisyllabic words, then generating the right origin semantics would need as many passes over the input as there are syllables in the input. The function would either need the full power of a 2-way D-FST in order to generate the right origin semantics, otherwise we could use a 1-way FST that has the linguistically-unmotivated origin semantics.

As for minimality requirements (19b), this cannot be computed by a C-Seq transducer with the right origin semantics. In order to reduplicate a monosyllabic input: $jang \rightarrow jang\sim jang$, we use two concatenated identity functions. But to block reduplication in a bisyllabic input $jia.luen \rightarrow meei-jia.luen$, we need to check that the input does not contain any medial syllable boundaries. The first function would need to use to finite lookahead before choosing to output the first segment j or the prefix $meei$. As with the infixation cases in Section 6.3.2, a C-Seq 2-way FST can do so but it then generates the wrong origin semantics because it associates the output segment j with the input syllable boundary ‘.’. Generating the right origin semantics needs a rotating 2-way D-FST that involves three passes. The first pass reads the input and checks if it is monosyllabic or not. If yes, the second and third passes apply the identity function: $jang\sim jang$. If no, the second pass outputs the prefix and the base $meei-jia.luen$; there is no third pass.

6.4.2

Phonological changes to the reduplicant

The previous section illustrated how C-Seq 2-way FSTs are distinct from rotating 2-way FSTs. In the latter, a pass can transfer information (e.g., *is the input monosyllabic*) to a later pass. Similar information transfer is required in certain cases where phonological processes interact with reduplication. We first go over cases where we arguably do not need such information transfer.

Reduplicative patterns do not only involve copying. In addition to copying segments, a reduplicative process may involve a host of other

L-OSL, their composition is L-OSL, and the concatenation with $ID(x)$ is C-L-OSL. Complex onset reduction in Tagalog and echo reduplication in Turkish are likewise C-L-OSL and consist of the concatenation of a composed L-OSL function with some other L-OSL function.

However, this does not mean that all hypothetical cases of reduplicant modifications are C-OSL. Such processes can be C-Seq or higher if the composition of a truncation or identity function with the modification function is Seq or higher. For example, if complex onset reduction in the reduplicant deleted the first consonant: $mag\text{-}trabaho \rightarrow mag\text{-}ra\sim trabaha$, this process would be C-Seq and not C-OSL. In the first copy, the truncation function would generate $mag\{,trabaho\}_r$, while the modification function would generate $mag\{,tra\}$. Deleting only the root-initial consonant t if it precedes a consonant is not OSL because we need finite lookahead on the input. Interestingly, this type of cluster reduction is argued to be unattested in reduplication (Zukoff 2017, 25). This may either be an accidental gap or evidence that reduplication modification must be C-OSL. To our knowledge, there is no typological survey of attested reduplicant modifications to settle this.

6.4.3 Phonological changes to or across both copies

Phonological changes may likewise affect both copies or apply across the boundary between the copies (22). Some involve a phonological process which is productive in the language (22a), others involve a phonological process which is not found anywhere else in the language outside of reduplication (22b). The former set of cases are often called

Table 9: C-OSL treatment for phonological changes to the reduplicant

	Papago (21a)	Tagalog (21b)	Turkish (21c)
Input x	bana	$mag\{,trabaho\}_r$	kitap
Components	$M(Tr(x)) \cdot ID(x)$	$M(Tr(x)) \cdot R(x)$	$ID(x) \cdot M(ID(x))$
Innermost function	ba $\#$ a	$mag\{,tra\ \#aho\}_r$	kitap
Composition	baa	$mag\{,t\ \# a$	mitap
Concatenation	baa \sim bana	$mag\text{-}ta \sim mag\{,trabaho\}_r$	kitap \sim mitap
Subclass	C-L-OSL	C-L-OSL	C-L-OSL

normal application of phonological rules, while the latter are *juncture effects* which are morpheme-specific phonological processes.

(22) *Phonological modifications across the two copies*

a. **Normal application of nasal substitution**

Balangao (McCarthy and Prince 1995, 85)

- i. /maN + tagtag/ → [ma + nagtag] ‘running’
- ii. /maN-RED + tagtag/ → [ma + nagta~tagtag], ‘running’
*[ma + nagta~nagtag] everywhere’

b. **Phonology across the boundary (*juncture effects*)**

Dakota (Inkelas and Zoll 2005, 101)

- i. /skokpá → o-skókpa~kpa ‘to be scooped out’
- ii. /čap/ → čap~čap-a ‘trot’
- iii. /žat/ → žag~žat-a ‘curved’

In (22a), the prefix *maN-* can trigger reduplication of the root/stem. Nasal substitution combines the prefix’s nasal with an adjacent voiceless consonant into a single nasal that has the place of articulation of the consonant. Nasal substitution applies only to the segment next to the prefix, regardless of whether that consonant is part of the reduplicant or not. In (22b), the final syllable of the root is copied and placed at the left edge of the input. If there are two coronals across the reduplicative boundary (b), then the first coronal becomes dorsal. The final /a/ is epenthesized.

We likewise find phonological processes or rules interacting *differently* in the context of reduplication. For example in Madurese, there is a phonological process of nasal spread in which nasality is spread from nasals onto sequences of glides and vowels (23a). Reduplication copies the final CVC and places it at the beginning of the output (23b). If a vowel in the base is nasalized by a nasal, its nasality will transfer to the reduplicant as well. Because the reduplicant does not contain any nasals to trigger nasal spread, nasal spread in the reduplicant is treated as an over-application of the phonological process of nasal spread.

(23) **Over-application of nasal spread**

Madurese (McCarthy and Prince 1995, 30; Cohn 1993, 358)

- a. /neyat/ → [nēỹāt] ‘intention’
- b. [ỹāt~nēỹāt] ‘intentions’

Traditionally, these cases can be thought as a composition of a morphological rule of reduplication (C-OSL) and a phonological rule (that is independently ISL, OSL, or Sequential) (Raimy 2000; Inkelas and Zoll 2005). If the morphological rule precedes the phonological rule, then we have normal application. If the morphological rule outputs reduplicant boundaries and precedes the phonological rule, then we have juncture effects. And, if the phonological rule precedes the morphological rule then we have over-application. Table 10 visualizes these three types of interactions as rule or function composition.

Table 10: Order of compositions for different reduplication-phonology interactions

	Normal application	Juncture effect	Over-application
Language	Balangao (22a)	Dakota (22b)	Madurese (23)
Input x	$\text{maN}\{\text{,tagtag}\}_r$	žat	neyat
Order of composition	1. Copy 2. Modify	1. Copy 2. Modify	1. Modify 2. Copy
Components	$M(L(x)) \cdot R(x)$	$M(L(x)) \cdot R(x)$	$L(M(x)) \cdot R(M(x))$
Innermost functions	$\text{maN}\{\text{,tagtag}\}_r \sim \text{maŋ}\{\text{,tagtag}\}_r$	žat \sim žata	něyāt
Outer function	$\text{maN}\{\text{,ŋagtag}\}_r \sim \{\text{,tagtag}\}_r$	žag \sim žata	něyāt \sim něyāt

Computationally, we can treat all these cases as composition of a C-OSL/C-Seq function for reduplication with an OSL/Seq function for phonology in either order. We conjecture C-OSL/C-Seq functions are not closed under composition, but we do not prove it. This means that composition may create a rotating 2-way FST.

Whether a case of normal application, juncture effect, or over-application is C-OSL, C-Seq, or higher depends on the complexity of the individual functions. In fact, the above three examples can be done with a C-Seq transducer. To illustrate, consider over-application in Madurese. Nasal spread is an L-OSL function $M(x)$. Reduplication is the concatenation of an R-OSL truncation function $L(x)$ and an L/R-OSL identity function $R(x)$. To generate overapplication, reduplication is instead the concatenation of two *modified* functions $L(M(x)) \cdot R(M(x))$. The first function is the composition of truncation over nasal spread. The composition of these two OSL rules of different directions is L-Seq because nasalization relies on deleted information

from the input.²¹ The second function is the composition of identity over nasal spread; this is L-OSL. Together, Madurese overapplication is C-L-Seq.

It is an open question if there are cases of normal application, juncture effects, and over-application which *cannot* be treated with a C-Seq formalization but require an unrestricted rotating 2-way FST. Solving this requires an in-depth knowledge of both the morphology and phonology of any such example (Inkelas and Zoll 2005).

Under-application and Back-copying In contrast to the over-application of phonological processes in reduplication, we likewise find cases of *under-application*. For example in Akan, velar consonants become palatalized before nonlow front vowels: /k,g/ → [tɕ, dʒ]/ _ /i,e/ as in (24a). Akan likewise has a process of initial-CV reduplication where the reduplicant V is a pre-specified non-low front vowel /ɪ/ (24b).²² However if the reduplicant C is a velar, it will not be palatalized before the reduplicant's non-low front vowel /ɪ/ (24c). Thus the rule under-applies. The velar will only palatalize if both copies of the velar in the reduplicant and base are preceded by a non-low front vowel (24d).

(24) **Under-application of palatalization in reduplication**

<i>Akan</i>	(McCarthy and Prince 1995, 83-93)
	(Schachter and Fromkin 1968, 89))
a. /ke/ → [tɕe]	‘divide’
b. /si/ → [si~siʔ]	‘stand’
c. /kaʔ/ → [kɪ~kaʔ], *[tɕɪ~kaʔ]	‘bite’
d. /ge/ → [dʒɪ~dʒe]	‘receive’

Cases of apparent under-application or over-application in reduplication are termed *opacity effects* (cf. the transparency of normal application (22a)). They are often understood as being caused by a need to maintain identity between the two copies that reduplication cre-

²¹ As with infixal reduplication (Section 6.3.2), the C-Seq transducer needs finite look-ahead into the end boundary \times and this makes it not have the exact origin semantics that we want.

²² The reduplicant V in Akan gets its front/back features from vowel harmony. For illustration, we represent it simply as /ɪ/.

ated (Wilbur 1973; McCarthy and Prince 1995). A more drastic version of identity is back-copying whereby the reduplicant undergoes some phonological rule, and then the effects of this rule are transferred onto the base. It is reported that in Malay, nasality spreads from a nasal consonant onto a sequence of vowels. Nasality can spread over glides and /h/. Plurality is marked by total reduplication. If nasal spread applies across the two copies, nasality will transfer onto both copies.

(25) **Back-copying of nasal spread**

Malay (McCarthy and Prince 1995, 85)

/hamə/ → [hām̃ə~hām̃ə], *[ham̃ə~hām̃ə] ‘germ’ → ‘germs’

These opacity effects are controversial both theoretically and empirically (Inkelas and Zoll 2005; Samuels 2010; Kiparsky 2010; McCarthy *et al.* 2012). Many cases of under-application have been re-analyzed as either unproductive (McCarthy *et al.* 2012) or due to morpheme-specific rules (Inkelas and Zoll 2005).²³ In fact, Akan palatalization (24) is the classical case study on under-application but it is likely a synchronically unproductive and fossilized rule (Silverman 2002; Adomako 2018). Empirically, there have been little if any convincing cases of back-copying (Bruening 1997) and some are arguably due to morphological factors outside of reduplication (McLaughlin 2005). The Malay data itself has not been successively reproduced (Kiparsky 2010).

Because under-application and back-copying have weak empirical backing, there is a limited attested typology of these processes. It is thus unclear whether we can make any computational generalizations about them. But putting aside these empirical problems, Akan under-application can be modeled with a C-Seq function which uses finite lookahead on the input. It is the concatenation of a modified truncation function and a modified identity function. The first function truncates the input $C_1V_2\Sigma^*$ to C_1i and applies palatalization if V_2 is /i,e/. The second function applies palatalization to the input.

Malay back-copying is not C-Seq. This is because nasalization requires unbounded lookahead on the input. The function requires an

²³An exception is Tonkawa (Gouskova 2007) which is arguably a bona fide case of under-application.

unrestricted rotating transducer with three passes over the input.²⁴ In the first pass, we output nothing but we check if the input ends in a $N(V+G)^*$ sequence where G stands for glides and $/h/$: $hamə$. If yes, the second pass applies nasal spread starting from the first segment: $hamə \sim \tilde{h}am\tilde{ə}$. The third pass does the same: $hamə \sim \tilde{h}am\tilde{ə} \sim \tilde{\tilde{h}}am\tilde{\tilde{ə}}$.

Overview of the typology summary

6.5

To summarize, we cataloged a wide variety of attested reduplicative patterns. All of it can be computed with deterministic 2-way FSTs. Most common and uncommon types of reduplication can be computed with the subclass of C-OSL functions, including total reduplication (12a), common partial reduplication patterns (13), triplication (14a), suffixal reduplication (14b), non-local reduplication (14c), and abbreviated reduplication (14d). Subconstituent reduplication is likewise C-OSL if the relevant morphological (15) or prosodic boundaries (16) are present in the input. In fact, of the 138 reduplicative processes in RedTyp (§5.2), 121 (87%) are C-OSL.²⁵

We analyzed the typology in terms of generating the right origin semantics. To do so, some less common types of reduplication are C-Seq or higher. This is largely because of the need for finite lookahead on the input. Some but not all types of infixal or non-contiguous reduplication are C-OSL (Section 6.3.1) and some are C-Seq (Section 6.3.2). In the latter case, generating the right origin semantics can require full 2-way FSTs because of the need for finite lookahead. Some cases like iterative reduplication (19a) are C-OSL if the input is at most bisyllabic; otherwise generating the right origin semantics needs an unrestricted 2-way D-FST. Minimality requirements (19b) likewise require

²⁴ Malay back-copying can likewise be treated as the composition of a C-OSL function for triplication $hamə \sim hamə \sim hamə$, followed by an OSL function for nasal spread $ham\tilde{ə} \sim \tilde{h}am\tilde{ə} \sim \tilde{\tilde{h}}am\tilde{\tilde{ə}}$, followed by an OSL function that deletes everything before the first \sim boundary $\tilde{h}am\tilde{ə} \sim \tilde{\tilde{h}}am\tilde{\tilde{ə}}$. This analysis is inspired by Reiss and Simpson (2009).

²⁵ Although the cross-linguistic typology on reduplication is overwhelmingly C-OSL, our numbers from RedTyp do not mean that we estimate that 13% of the cross-linguistic typology of reduplicative processes is not C-OSL. RedTyp likely under-represents cases of opacity. Such cases can be non-C-OSL.

full 2-way D-FSTs. When reduplication interacts with phonological processes, the computation can range anywhere from C-OSL to full 2-way D-FSTs depending on the individual phonological process and the order of function composition. We suspect the finite lookahead in these cases may be resolved with more sophisticated representations and logical transductions (Dolatian 2020).

This concludes the section on how various subclasses of 2-way D-FSTs map to certain divisions in the reduplicative typology. The above cases are representative of common and uncommon reduplicative processes. There are other subtle variations for reduplication in natural language, such as cases of allomorphy (Spaelti 1997), or multiple reduplicants (Urbanczyk 1999, 2001; Fitzpatrick and Nevins 2004; Fitzpatrick 2006), among others. We will not discuss these cases because a full typology is beyond the scope of this paper. However, virtually all attested reduplicative processes can be modeled with 2-way FSTs. Fitting the *entire* attested typology into the right subclasses is a fruitful research direction. The next section looks at cases where 2-way D-FSTs arguably over-generate or under-generate the typology, even with these well-defined subclasses.

6.6 *Issues in over- and under-generation*

Here, we address the questions whether and how 2-way D-FSTs under- and over-generate reduplicative processes.

6.6.1 Over-generation with 2-way D-FSTs

One way to interpret the contribution we have made is that we are advocating the following hypothesis:

- (26) **2-way Hypothesis:** Reduplication is anything that can be computed with 2-way D-FSTs.

This is not, in fact, a position we advocate. We think this hypothesis is false because it overgenerates in ways we consider linguistically bizarre. For example, 2-way D-FSTs can map words to their reverse ($w \mapsto w^r$) and to a copy of itself and its mirror image ($w \mapsto ww^r$). None

of these transformations are attested morphologically. Some overgeneration can be avoided by hypothesizing stronger computational properties; that is, focusing on *subclasses* of 2-way D-FSTs which cannot generate the above unattested patterns.

This leads to another hypothesis.

- (27) **C-OSL Hypothesis:** Reduplication is anything that can be computed with C-OSL 2-way D-FSTs.

The C-OSL hypothesis is well supported because it covers the *bulk* of reduplicative typology as shown in Section 6. Rarer reduplicative patterns require more powerful subclasses of 2-way D-FSTs.

Even this hypothesis can be said to suffer from overgeneration. For example, while this excludes the reversal and mirror image processes above, it permits total reduplication of a word up to some large natural number n ($w \mapsto w^n$), or partial reduplication up to some natural number n of segments.

Nonetheless, not all issues in overgeneration can be reduced to computation or computability. Some are certainly due to external factors.²⁶ To illustrate, total reduplication in most *spoken* languages creates at most two copies. The creation of three copies (= triplication) is relatively rare in spoken languages, e.g. Thao (Blust 2001). In sign languages, we find the reverse situation: creating two copies is rare but triplication is common, e.g. ASL (Wilbur 2005). The difference between sign and spoken reduplication is more likely due to modality and not to the computation.²⁷

Under-generation with 2-way D-FSTs

6.6.2

Non-computational factors can also help us understand apparent cases of under-generation. There are two cases we discuss here: abstract morphemic copying and reduplication with haplogy. Both of these

²⁶ Like Potts and Pullum (2002, 375), “we are extremely sceptical of the idea that formalisms exist that correspond exactly to what linguists wish to say.”

²⁷ A similar point can be made for the role of pivot or anchor points in reduplication. Cross-linguistically, most reduplicative processes target specific positions in the word which are perceptually or psycho-linguistically more salient (Samuels 2010; Raimy 2009; Idsardi and Raimy 2008), e.g. the first syllable and not the third syllable. The choice of these pivots is likely functional, not computational.

can be explained as involving interactions between reduplication and other linguistic modules (the lexicon) or processes (filters).

Undergeneration of abstract morphemic copying Abstract morphemic copying is when the input to the copying mechanism is not a string of phonological segments but a more abstract morphological entity, i.e. a morpheme or morpho-syntactic feature bundle (Inkelas and Zoll 2005). This is in contrast to examples reviewed earlier, where the source of reduplication was a string of segments which may contain morpheme boundaries. Such a case occurs in Sye in Table 11.

In Sye, a stem may have multiple suppletive allomorphs used in different morphological contexts. For example, the abstract morpheme $\sqrt{\text{FALL}}$ in Table 11a has two allomorphs *amol* and *omol*, such that *amol* is used after future morpheme and certain other tense morphemes while *omol* is used elsewhere. As for reduplication, total reduplication is used to mark intensification (Table 11b). When total reduplication applies in a context that requires using one of the allomorphs, we have an allomorph mismatch between the two copies (Table 11c).

Table 11:
Abstract
morphemic
copying in Sye

Morphemes	a. $\sqrt{\text{FALL}}$	b. $\sqrt{\text{FALL + RED}}$	c. $\text{FUT-}\sqrt{\text{FALL-RED}}$
Output	omol	omol~omol	cw- <i>amol</i> ~ <i>omol</i> *cw- <i>omol</i> ~ <i>omol</i> *cw- <i>amol</i> ~ <i>amol</i>
Gloss	'fall'	'fall all over'	'they will fall all over'

Inkelas and Zoll (2005) analyze Sye as involving morphological copying. The copies are not in phonological correspondence because they are different allomorphs of the same morpheme. What was copied was an abstract morpheme $\sqrt{\text{FALL}}$. Its two copies were later spelled-out as two different allomorphs. Inkelas and Zoll's (2005) analysis for Sye is controversial (Frampton 2009); but there are a few other languages which show that the reduplicant is copying an abstract morphological entity (Inkelas and Downing 2015a,b; Hyman *et al.* 2009).

Cases of morphological copying for *suppletive* roots can be modeled with a 2-way D-FST that copies an abstract pre-spelled-out morphological entity, e.g. a root morpheme $\sqrt{\text{FALL}}$ or a root index (Harley 2014) which can be represented as a finite string of symbols. This is followed by a 1-way FST that models spell-out such that it is equipped

with knowledge over what all the *finite* pairs of morphemes and their suppletive allomorphs are. We make the safe assumption that the number of morphemes in a language that show suppletion is finite.²⁸

Undergeneration of reduplication with haplology Another case of potential under-generation is when reduplication is affected by anti-homophony constraints or haplology, i.e. when reduplication is blocked because it would create a sequence of identical syllables or feet that is dis-preferred by speakers (Yip 1995; Nevins 2012).

For example in Kanuri (Moravcsik 1978, 313), total reduplication is used to form glossonyms (28b). However reduplication is blocked if it creates sequence of identical syllables/feet (28b).

(28) Reduplication and haplology in Kanuri

- | | | |
|----|--------------------|----------------------------------|
| a. | kanəmbu | ‘Kanembu tribe’ |
| | kanəmbu~kanəmbu | ‘language of the Kanembu tribe’ |
| b. | karekare | ‘Karekare tribe’ |
| | *karekare~karekare | ‘language of the Karekare tribe’ |

A 2-way D-FST can encode this requirement that the input must not itself be a sequence of identical syllables or feet. However that would require the 2-way D-FST to know what all the finitely possible sequences of syllables and feet are in the language.

Two alternatives to this solution are possible. One is using a copy-and-filter mechanism (Golston 1995). The 2-way D-FST would handle the copying. The output of the copying process would be fed to a phonological system which would filter out any homophonous sequences of syllables or feet. The other alternative is to argue that the Kanuri input stems which contain a repeated sequence of symbols /karekare/ are underlyingly already reduplicated via lexical reduplication /kare + RED/. Such arguments have been brought up for superficially similar haplology effects in Manam (Buckley 1997). From this approach, there is then no haplology problem for 2-way D-FSTs.

In sum, although virtually the entire typology of reduplication can be modeled with 2-way D-FSTs, there are complications if one

²⁸ An FST which handles spell-out would resemble the lexical transducer used in the *xfst* finite-state package (Beesley and Karttunen 2003).

wishes to model the full interface of reduplicative morphology with other systems. However, as a reviewer point outs, it may not be desirable, from a linguistic perspective at least, to model the interfaces in this way. It is disputable whether complications occurring at the interface of morphology with syntax or phonology should be addressed within an FST that is intended to account for the computational complexity of the morphology itself. On the other side of the coin, some conceptual problems arise with over-generation of 2-way D-FSTs because of the power that they require to handle copying in the first place.

7

CONCLUSION

The present study has taken a step in formalizing the wide typology of reduplicative processes in formal-language theoretic terms. We showed that 2-way D-FSTs, which are an understudied type of finite-state transducer, can easily model reduplication because they can reread their input multiple times in multiple directions. Computationally, this means that *recognizing* whether strings belong to the copy language $\{ww \mid w \in \Sigma^*\}$ (so for any $w \in \Sigma^*$ determining whether there is a $v \in \Sigma^*$ such that $w = \nu v$) is a harder problem than the one that takes any $w \in \Sigma^*$ as input and returns ww as output (*copying*). Reduplication studied as recognition is computationally more complex than reduplication studied as copying.

In addition to modeling reduplicative morphology as copying, 2-way D-FSTs do not suffer from state explosion nor do they assume finite bounds on the input, unlike 1-way FSTs. In terms of strong generative capacity, 2-way FSTs actively copy segments instead of memorizing segments. A diagnostic for copying vs. remembering is the origin semantics of the function. This article also presented the RedTyp database, which provides concrete examples of 2-way DFSTs modeling a range of cross-linguistic reduplicative morphemes.

Furthermore, we showed that the typology of reduplication can be modeled with subclasses of 2-way FSTs that are essentially defined as concatenations of simple subclasses of 1-way FSTs. Thus, our work

showed the role of computational subclasses in carving out the generative capacity of morphological processes, whether reduplicative or not. To give more context, most morphological processes can be computed by 1-way finite state automata and transducers (Koskenniemi 1983; Beesley and Karttunen 2003). In fact, substantially less expressive subregular classes are capable of computing most of these morphological processes (Aksënova *et al.* 2016; Chandlee 2017). So far, these subclasses have been identified based on considerations of locality (ISL, OSL) and determinism (Seq, sequentiality). At first, reduplication looks like an outlier in that it requires the more expressive generative capacity of 2-way transducers. However, even within this larger class of 2-way FSTs, we argued that reduplication only needs certain subclasses which are also based on the same considerations of locality (C-OSL) and sequentiality (C-Seq). These subclasses reinforce the role of locality and determinism as general constraints in linguistic processes (cf. Heinz 2018).

Having showcased the utility of 2-way D-FSTs for modeling reduplication, we conclude with three avenues of future research.

First, we have approached reduplication from the perspective of morphological generation. Given an input *buku*, a 2-way D-FST can generate the output *buku~buku* easily. On the other hand, it is an open question as to how to do morphological *analysis* with 2-way FSTs to get the inverse relation of *buku~buku* \rightarrow *buku*. As a class, deterministic 2-way FSTs are not invertible. We are currently developing algorithms for inverting the subclasses (C-OSL, C-Seq) that we have set up.

A second area of research is the integration of 2-way FSTs into natural language processing. This obviously has many aspects. A first step may be the integration of 2-way FSTs into existing platforms such as *xfst* (Beesley and Karttunen 2003), *open-fst* (Allauzen *et al.* 2007), *foma* (Hulden 2009b), and *pynini* (Gorman 2016).²⁹

A third promising area of research is developing learning models based on the computational models that we proposed here. One approach builds on Chandlee *et al.*'s 2015 learning results of OSL func-

²⁹In fact, the team behind Thrax (Tai *et al.* 2011) have recently been exploring the use of multi-pushdown transducers (MPDT) to generate reduplication (Richard Sproat, p.c.). An open question is comparing the generative capacity of MPDTs and 2-way FSTs.

tions (Dolatian and Heinz 2018a). Another approach probes the learnability of reduplicative patterns with neural networks (Nelson *et al.* 2020).

REFERENCES

- Kwasi ADOMAKO (2018), Velar palatalization in Akan: A reconsideration, *Journal of West African Languages*, 45(2).
- Alfred V. AHO, John E. HOPCROFT, and Jeffrey D. ULLMAN (1969), A general theory of translation, *Mathematical Systems Theory*, 3(3):193–221.
- Alëna AKSËNOVA, Thomas GRAF, and Sedigheh MORADI (2016), Morphotactics as tier-based strictly local dependencies, in *Proceedings of the 14th SIGMORPHON Workshop on Computational Research in Phonetics, Phonology, and Morphology*, pp. 121–130.
- Daniel M. ALBRO (2000), Taking Primitive Optimality Theory beyond the finite state, in Jason EISNER, Lauri KARTTUNEN, and Alain THÉRIAULT, editors, *Finite-State Phonology: Proceedings of the 5th Workshop of SIGPHON*, pp. 57–67, Luxembourg, <http://aclanthology.coli.uni-saarland.de/pdf/W/W00/W00-1806.pdf>.
- Daniel M. ALBRO (2005), *Studies in Computational Optimality Theory, with Special Reference to the Phonological System of Malagasy*, Ph.D. thesis, University of California, Los Angeles.
- Raquel G. ALHAMA (2017), *Computational Modelling of Artificial Language Learning: Retention, Recognition & Recurrence*, Ph.D. thesis, Universiteit van Amsterdam.
- Raquel G. ALHAMA and Willem ZUIDEMA (2019), A review of computational models of basic rule learning: The neural-symbolic debate and beyond, *Psychonomic Bulletin & Review*, 26(4):1–21.
- Cyril ALLAUZEN, Michael RILEY, Johan SCHALKWYK, Wojciech SKUT, and Mehryar MOHRI (2007), OpenFst: A general and efficient weighted finite-state transducer library, in Jan HOLUB and Jan ŽĎÁREK, editors, *Implementation and Application of Automata*, pp. 11–23, Springer, Berlin, Heidelberg.
- Rajeev ALUR (2010), Expressiveness of streaming string transducers, in *Proceedings of the 30th Annual Conference on Foundations of Software Technology and Theoretical Computer Science*, volume 8, p. 1–12, doi:10.4230/LIPICs.FSTTCS.2010.1.

Rajeev ALUR and Jyotirmoy V. DESHMUKH (2011), Nondeterministic streaming string transducers, in Luca ACETO, Monika HENZINGER, and Jiří SGALL, editors, *Automata, Languages and Programming*, pp. 1–20, Springer, Berlin, Heidelberg, ISBN 978-3-642-22012-8.

Rajeev ALUR, Adam FREILICH, and Mukund RAGHOTHAMAN (2014), Regular combinators for string transformations, in *Proceedings of the Joint Meeting of the Twenty-Third EACSL Annual Conference on Computer Science Logic (CSL) and the Twenty-Ninth Annual ACM/IEEE Symposium on Logic in Computer Science (LICS), Vienna, Austria, CSL-LICS '14*, pp. 9:1–9:10, Association for Computing Machinery, New York, NY, USA, ISBN 978-1-4503-2886-9, doi:10.1145/2603088.2603151.

Rajeev ALUR and Pavol ČERNÝ (2011), Streaming transducers for algorithmic verification of single-pass list-processing programs, in *Proceedings of the 38th Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages, Austin, Texas, USA, POPL '11*, pp. 599–610, Association for Computing Machinery, New York, NY, USA, ISBN 978-1-4503-0490-0, doi:10.1145/1926385.1926454.

Qatherine ANDAN, Outi BAT-EL, Diane BRENTARI, and Iris BERENT (2018), ANCHORING is amodal: Evidence from a signed language, *Cognition*, 180:279–283.

Stephen R. ANDERSON (1992), *A-morphous morphology*, volume 62 of *Cambridge Studies in Linguistics*, Cambridge University Press, Cambridge.

Mark ARONOFF (1988), Head operations and strata in reduplication: A linear treatment, in Geert BOOIJ and Jaap VAN MARLE, editors, *Yearbook of Morphology*, volume 1, pp. 1–15, Foris, Dordrecht.

Félix BASCHENIS, Olivier GAUWIN, Anca MUSCHOLL, and Gabriele PUPPIS (2015), One-way definability of sweeping transducers, in *35th IARCS Annual Conference on Foundations of Software Technology and Theoretical Computer Science (FSTTCS'15)*, Bangalore, India, <https://hal.archives-ouvertes.fr/hal-01219509>.

Félix BASCHENIS, Olivier GAUWIN, Anca MUSCHOLL, and Gabriele PUPPIS (2016), Minimizing resources of sweeping and streaming string transducers, in Ioannis CHATZIGIANNAKIS, Michael MITZENMACHER, Yuval RABANI, and Davide SANGIORGI, editors, *43rd International Colloquium on Automata, Languages, and Programming (ICALP 2016)*, volume 55 of *Leibniz International Proceedings in Informatics (LIPIcs)*, pp. 114:1–114:14, Schloss Dagstuhl-Leibniz-Zentrum für Informatik, Dagstuhl, Germany, ISBN 978-3-95977-013-2, ISSN 1868-8969, doi:10.4230/LIPIcs.ICALP.2016.114.

Félix BASCHENIS, Olivier GAUWIN, Anca MUSCHOLL, and Gabriele PUPPIS (2017), Untwisting two-way transducers in elementary time, in *32nd Annual ACM/IEEE Symposium on Logic in Computer Science, LICS 2017, Reykjavik*,

Iceland, June 20-23, 2017, pp. 1–12, ISBN 978-1-5090-3018-7,
doi:10.1109/LICS.2017.8005138.

Félix BASCHENIS, Olivier GAUWIN, Anca MUSCHOLL, and Gabriele PUPPIS (2018), One-way definability of two-way word transducers, *Logical Methods in Computer Science*, 14.

Kenneth BEESLEY and Lauri KARTTUNEN (2000), Finite-state non-concatenative morphotactics, in *Proceedings of the 38th Annual Meeting on Association for Computational Linguistics, Hong Kong, ACL '00*, pp. 191–198, Association for Computational Linguistics, Stroudsburg, PA, USA, doi:10.3115/1075218.1075243.

Kenneth BEESLEY and Lauri KARTTUNEN (2003), *Finite-State Morphology: Xerox Tools and Techniques*, CSLI Publications, Stanford, CA.

Iris BERENT, Outi BAT-EL, Diane BRENTARI, Amanda DUPUIS, and Vered VAKNIN-NUSBAUM (2016), The double identity of linguistic doubling, *Proceedings of the National Academy of Sciences*, 113(48):13702–13707.

Iris BERENT, Outi BAT-EL, and Vered VAKNIN-NUSBAUM (2017), The double identity of doubling: Evidence for the phonology-morphology split, *Cognition*, 161:117–128.

Iris BERENT, Amanda DUPUIS, and Diane BRENTARI (2014), Phonological reduplication in sign language: Rules rule, *Frontiers in Psychology*, 5:560.

Steven BIRD and T. Mark ELLISON (1994), One-level phonology: Autosegmental representations and rules as finite automata, *Computational Linguistics*, 20(1):55–90.

Robert A. BLUST (2001), Thao triplication, *Oceanic Linguistics*, 40(2):324–335.

Mikołaj BOJAŃCZYK (2014), Transducers with origin information, in Javier ESPARZA, Pierre FRAIGNIAUD, Thore HUSFELDT, and Elias KOUTSOPIAS, editors, *Automata, Languages, and Programming*, pp. 26–37, Springer, Berlin, Heidelberg.

Mikołaj BOJAŃCZYK, Laure DAVIAUD, Bruno GUILLON, and Vincent PENELLE (2017), Which classes of origin graphs are generated by transducers, in Ioannis CHATZIGIANNAKIS, Piotr INDYK, Fabian KUHN, and Anca MUSCHOLL, editors, *44th International Colloquium on Automata, Languages, and Programming (ICALP 2017)*, volume 80 of *Leibniz International Proceedings in Informatics (LIPIcs)*, pp. 114:1–114:13, Schloss Dagstuhl–Leibniz-Zentrum für Informatik, Dagstuhl, Germany, ISBN 978-3-95977-041-5, ISSN 1868-8969, doi:10.4230/LIPIcs.ICALP.2017.114.

Olivier BONAMI and Berthold CRYSMANN (2013), Morphotactics in an information-based model of realisational morphology, in Stefan MÜLLER, editor, *Proceedings of HPSG 2013*, pp. 27–47, CSLI Publications, Stanford, CA.

- Olivier BONAMI and Berthold CRYSMANN (2016), The role of morphology in constraint-based lexicalist grammars, in Andrew HIPPISEY and Gregory T. STUMP, editors, *The Cambridge Handbook of Morphology*, p. 609–656, Cambridge University Press, Cambridge.
- Ellen BROSELOW and John MCCARTHY (1983), A theory of internal reduplication, *The Linguistic Review*, 3(1):25–88.
- Jason BROWN (2017), Non-adjacent reduplication requires spellout in parallel, *Natural Language & Linguistic Theory*, 35(4):1–23.
- Benjamin BRUENING (1997), Abkhaz mabkhaz: M-reduplication in Abkhaz, weightless syllables, and base-reduplicant correspondence, in Benjamin BRUENING, Yoonjung KANG, and Martha MCGINNIS, editors, *PF: Papers at the Interface*, volume 30, MIT Working Papers in Linguistics, Cambridge, MA.
- Eugene BUCKLEY (1997), Integrity and correspondence in Manam double reduplication, in *Proceedings of NELS*, volume 28, pp. 59–67.
- Gabriela CABALLERO (2006), “Templatic backcopying” in Guarijio abbreviated reduplication, *Morphology*, 16(2):273–289.
- Jane CHANDLEE (2014), *Strictly local phonological processes*, Ph.D. thesis, University of Delaware, Newark, DE.
- Jane CHANDLEE (2017), Computational locality in morphological maps, *Morphology*, 27(4):1–43.
- Jane CHANDLEE, Angeliki ATHANASOPOULOU, and Jeffrey HEINZ (2012), Evidence for classifying metathesis patterns as subsequential, in Jaehoon CHOI, E. Alan HOGUE, Jeffrey PUNSKE, Deniz TAT, Jessamyn SCHERTZ, and Alex TRUEMAN, editors, *The Proceedings of the 29th West Coast Conference on Formal Linguistics*, pp. 303–309, Cascillida Press, Somerville, MA.
- Jane CHANDLEE, Rémi EYRAUD, and Jeffrey HEINZ (2014), Learning strictly local subsequential functions, *Transactions of the Association for Computational Linguistics*, 2:491–503, <http://aclweb.org/anthology/Q14-1038>.
- Jane CHANDLEE, Rémi EYRAUD, and Jeffrey HEINZ (2015), Output strictly local functions, in *14th Meeting on the Mathematics of Language*, pp. 112–125.
- Jane CHANDLEE and Jeffrey HEINZ (2012), Bounded copying is subsequential: Implications for metathesis and reduplication, in *Proceedings of the 12th Meeting of the ACL Special Interest Group on Computational Morphology and Phonology*, SIGMORPHON '12, pp. 42–51, Association for Computational Linguistics, Montreal, Canada.
- Jane CHANDLEE and Jeffrey HEINZ (2018), Strict locality and phonological maps, *Linguistic Inquiry*, 49(1):23–60.
- Jane CHANDLEE, Jeffrey HEINZ, and Adam JARDINE (2018), Input strictly local opaque maps, *Phonology*, 35(2):171–205.

Christian CHOFFRUT (1977), Une caractérisation des fonctions séquentielles et des fonctions sous-séquentielles en tant que relations rationnelles, *Theoretical Computer Science*, 5(3):325–337, ISSN 0304-3975, doi:[https://doi.org/10.1016/0304-3975\(77\)90049-4](https://doi.org/10.1016/0304-3975(77)90049-4).

Michal P. CHYTIK and Vojtěch JÁKL (1977), Serial composition of 2-way finite-state transducers and simple programs on strings, in Arto SALOMAA and Magnus STEINBY, editors, *Automata, Languages and Programming*, pp. 135–147, Springer, Berlin, Heidelberg, ISBN 978-3-540-37305-6.

Alexander CLARK (2017), Computational learning of syntax, *Annual Review of Linguistics*, 3:107–123.

Alexander CLARK and Ryo YOSHINAKA (2012), Beyond semilinearity: Distributional learning of parallel multiple context-free grammars, in *International Conference on Grammatical Inference*, pp. 84–96.

Alexander CLARK and Ryo YOSHINAKA (2014), Distributional learning of parallel multiple context-free grammars, *Machine Learning*, 96(1-2):5–31.

Alexander CLARK and Ryo YOSHINAKA (2016), Distributional learning of context-free and multiple context-free grammars, in Jeffrey HEINZ and José M. SEMPERE, editors, *Topics in Grammatical Inference*, pp. 143–172, Springer, Berlin, Heidelberg.

Yael COHEN-SYGAL and Shuly WINTNER (2006), Finite-state registered automata for non-concatenative morphology, *Computational Linguistics*, 32(1):49–82.

Abigail C. COHN (1989), Stress in Indonesian and bracketing paradoxes, *Natural language & linguistic theory*, 7(2):167–216.

Abigail C. COHN (1993), The status of nasalized continuants, in Marie K. HUFFMAN and Rena A. KRAKOW, editors, *Nasals, Nasalization, and the Velum*, volume 5 of *Phonetics and Phonology*, pp. 329–367, Academic Press, Inc., San Diego, CA.

Bruno COURCELLE and Joost ENGELFRIET (2012), *Graph Structure and Monadic Second-Order Logic, a Language Theoretic Approach*, Cambridge University Press, Cambridge.

Berthold CRYSMANN (2017), Reduplication in a computational HPSG of Hausa, *Morphology*, 27(4):527–561.

Berthold CRYSMANN and Olivier BONAMI (2016), Variable morphotactics in information-based morphology, *Journal of Linguistics*, 52(2):311–374.

Karel CULIK and Juhani KARHUMÄKI (1986), The equivalence of finite valued transducers (on HDTOL languages) is decidable, *Theoretical Computer Science*, 47:71–84, ISSN 0304-3975, doi:[https://doi.org/10.1016/0304-3975\(86\)90134-9](https://doi.org/10.1016/0304-3975(86)90134-9).

- Christopher CULY (1985), The complexity of the vocabulary of Bambara, *Linguistics and Philosophy*, 8:345–351.
- Vrunda DAVE, Paul GASTIN, and Shankara Narayanan KRISHNA (2018), Regular transducer expressions for regular transformations, in *Proceedings of the 33rd Annual ACM/IEEE Symposium on Logic in Computer Science, LICS '18*, pp. 315–324, Association for Computing Machinery, New York, NY, USA, ISBN 978-1-4503-5583-4, doi:10.1145/3209108.3209182.
- Hossep DOLATIAN (2020), *Computational Locality of Cyclic Phonology in Armenian*, Ph.D. thesis, Stony Brook University.
- Hossep DOLATIAN and Jeffrey HEINZ (2018a), Learning reduplication with 2-way finite-state transducers, in Olgierd UNOLD, Witold DYRKA, , and Wojciech WIECZOREK, editors, *Proceedings of Machine Learning Research: International Conference on Grammatical Inference*, volume 93 of *Proceedings of Machine Learning Research*, pp. 67–80, Wrocław, Poland.
- Hossep DOLATIAN and Jeffrey HEINZ (2018b), Modeling reduplication with 2-way finite-state transducers, in *Proceedings of the 15th SIGMORPHON Workshop on Computational Research in Phonetics, Phonology, and Morphology*, Association for Computational Linguistics, Brussels, Belgium.
- Hossep DOLATIAN and Jeffrey HEINZ (2019a), RedTyp: a database of reduplication with computational models, in *Proceedings of the Society for Computation in Linguistics*, volume 2, article 3.
- Hossep DOLATIAN and Jeffrey HEINZ (2019b), Reduplication with finite-state technology, in *Proceedings of the 53rd Annual Meeting of the Chicago Linguistics Society*, Chicago Linguistics Society, Chicago.
- Laura J. DOWNING (1998), Prosodic misalignment and reduplication, in Geert BOOIJ and Jaap VAN MARLE, editors, *Yearbook of Morphology 1997*, pp. 83–120, Kluwer Academic Publishers, Dordrecht.
- Laura J. DOWNING (2000), Morphological and prosodic constraints on Kinande verbal reduplication, *Phonology*, 17(01):1–38.
- Laura J. DOWNING (2003), Compounding and tonal non-transfer in Bantu languages, *Phonology*, 20(1):1–42.
- Laura J. DOWNING (2006), *Canonical Forms in Prosodic Morphology*, number 12 in Oxford studies in Theoretical Linguistics, Oxford University Press, Oxford.
- Calvin C. ELGOT and Jorge E. MEZEI (1965), On relations defined by generalized finite automata, *IBM Journal of Research and development*, 9(1):47–68.
- Joost ENGELFRIET and Hendrik Jan HOOGEBOOM (2001), MSO definable string transductions and two-way finite-state transducers, *Transactions of the Association for Computational Linguistics*, 2(2):216–254, ISSN 1529-3785, doi:10.1145/371316.371512.

Emmanuel FILIOT and Pierre-Alain REYNIER (2016), Transducers, logic and algebra for functions of finite words, *ACM SIGLOG News*, 3(3):4–19, ISSN 2372-3491, doi:10.1145/2984450.2984453.

Justin FITZPATRICK (2006), Sources of Multiple Reduplication in Salish and Beyond, *Studies in Salishan* 7, pp. 211–240.

Justin FITZPATRICK and Andrew NEVINS (2004), Linearizing nested and overlapping precedence in multiple reduplication, in *University of Pennsylvania Working Papers in Linguistics*, pp. 75–88.

Jennifer FITZPATRICK-COLE (1994), *The Prosodic Domain Hierarchy in Reduplication*, Ph.D. thesis, Stanford University, Stanford, CA.

John FRAMPTON (2009), *Distributed Reduplication*, MIT Press, Cambridge.

Diamandis GAFOS (1998), A-templatic reduplication, *Linguistic Inquiry*, 29(3):515–527.

Brian GAINOR, Regine LAI, and Jeffrey HEINZ (2012), Computational characterizations of vowel harmony patterns and pathologies, in Jaehoon CHOI, E. Alan HOGUE, Jeffrey PUNSKE, Deniz TAT, Jessamyn SCHERTZ, and Alex TRUEMAN, editors, *The Proceedings of the 29th West Coast Conference on Formal Linguistics*, pp. 63–71, Cascillida Press, Somerville, MA.

Pedro GARCIA, Enrique VIDAL, and José ONCINA (1990), Learning locally testable languages in the strict sense, in *Proceedings of the Workshop on Algorithmic Learning Theory*, pp. 325–338.

Gerald GAZDAR and Geoffrey K PULLUM (1985), Computationally relevant properties of natural languages and their grammars, *New generation computing*, 3:273–306.

Jila GHOMESHI, Ray JACKENDOFF, Nicole ROSEN, and Kevin RUSSELL (2004), Contrastive focus reduplication in English (the salad-salad paper), *Natural Language & Linguistic Theory*, 22(2):307–357.

Chris GOLSTON (1995), Syntax outranks phonology: Evidence from Ancient Greek, *Phonology*, 12(3):343–368.

Kyle GORMAN (2016), Pynini: A Python library for weighted finite-state grammar compilation, in *Proceedings of the SIGFSM Workshop on Statistical NLP and Weighted Automata*, pp. 75–80, Association for Computational Linguistics, Berlin, Germany, doi:10.18653/v1/W16-2409, <http://www.aclweb.org/anthology/W16-2409>.

Maria GOUSKOVA (2007), The reduplicative template in Tonkawa, *Phonology*, 24(3):367–396.

Jiatao GU, Zhengdong LU, Hang LI, and Victor O.K. LI (2016), Incorporating copying mechanism in sequence-to-sequence learning, in *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 1631–1640, Association for Computational Linguistics, Berlin, Germany.

- Heidi HARLEY (2014), On the identity of roots, *Theoretical linguistics*, 40(3/4):225–276.
- Jason D. HAUGEN (2009), What is the base for reduplication?, *Linguistic Inquiry*, 40(3):505–514.
- Jeffrey HEINZ (2007), *The Inductive Learning of Phonotactic Patterns*, Ph.D. thesis, University of California, Los Angeles.
- Jeffrey HEINZ (2018), The computational nature of phonological generalizations, in Larry HYMAN and Frans PLANK, editors, *Phonological Typology*, Phonetics and Phonology, chapter 5, pp. 126–195, Mouton de Gruyter, Berlin.
- Jeffrey HEINZ and William IDSARDI (2013), What complexity differences reveal about domains in language, *Topics in Cognitive Science*, 5(1):111–131.
- Jeffrey HEINZ and Regine LAI (2013), Vowel harmony and subsequentiality, in Andras KORNAI and Marco KUHLMANN, editors, *Proceedings of the 13th Meeting on the Mathematics of Language (MoL 13)*, pp. 52–63, Association for Computational Linguistics, Sofia, Bulgaria, <http://www.aclweb.org/anthology/W13-3006>.
- John E. HOPCROFT and Jeffrey D. ULLMAN (1969), *Formal Languages and their Relation to Automata*, Addison-Wesley Longman Publishing Co., Inc., Boston, MA.
- Mans HULDEN (2006), Finite-state syllabification, in Anssi YLI-JYRÄ, Lauri KARTTUNEN, and Juhani KARHUMÄKI, editors, *Finite-State Methods and Natural Language Processing. FSMNLP 2005. Lecture Notes in Computer Science*, volume 4002, Springer, Berlin/Heidelberg.
- Mans HULDEN (2009a), *Finite-State Machine Construction Methods and Algorithms for Phonology and Morphology*, Ph.D. thesis, University of Arizona, Tucson, AZ.
- Mans HULDEN (2009b), Foma: a finite-state compiler and library, in *Proceedings of the Demonstrations Session at EACL 2009*, pp. 29–32, Association for Computational Linguistics, Athens, Greece, <http://www.aclweb.org/anthology/E09-2008>.
- Mans HULDEN and Shannon T. BISCHOFF (2009), A simple formalism for capturing reduplication in finite-state morphology, in Jakub PISKORSKI, Bruce WATSON, and Anssi YLI-JYRÄ, editors, *Proceedings of the 2009 conference on Finite-State Methods and Natural Language Processing: Post-proceedings of the 7th International Workshop FSMNLP 2008*, pp. 207–214, IOS Press, Amsterdam, ISBN 978-1-58603-975-2, <http://dl.acm.org/citation.cfm?id=1564035.1564059>.
- Bernhard HURCH, editor (2005), *Studies on Reduplication*, number 28 in *Empirical Approaches to Language Typology*, Walter de Gruyter, Berlin.
- Bernhard HURCH (2005 ff.), Graz database on reduplication, last accessed 10-26-2017 from <http://reduplication.uni-graz.at/redup/>.

- Bernhard HURCH and Veronika MATTES (2009), Introduction: diachrony and productivity of reduplication, *Morphology*, 19(2):107–112.
- Larry M. HYMAN (2009), The natural history of verb-stem reduplication in Bantu, *Morphology*, 19(2):177–206.
- Larry M. HYMAN, Sharon INKELAS, and Galen SIBANDA (2009), Morphosyntactic correspondence in Bantu reduplication, in Kristin HANSON and Sharon INKELAS, editors, *The Nature of the Word: Studies in Honor of Paul Kiparsky*, Current Studies in Linguistics, pp. 273–309, The MIT Press, Cambridge, MA.
- William IDSARDI and Eric RAIMY (2008), Reduplicative economy, in Bert VAUX and Andrew NEVINS, editors, *Rules, constraints, and phonological phenomena*, chapter 5, pp. 149–184, Oxford University Press, Oxford.
- Sharon INKELAS (2014), *The Interplay of Morphology and Phonology*, Oxford University Press, Oxford.
- Sharon INKELAS and Laura J. DOWNING (2015a), What is reduplication? Typology and analysis part 1/2: The typology of reduplication, *Language and Linguistics Compass*, 9(12):502–515.
- Sharon INKELAS and Laura J. DOWNING (2015b), What is Reduplication? Typology and analysis Part 2/2: The analysis of reduplication, *Language and Linguistics Compass*, 9(12):516–528.
- Sharon INKELAS and Cheryl ZOLL (2005), *Reduplication: Doubling in Morphology*, Cambridge University Press, Cambridge.
- Adam JARDINE (2016), Computationally, tone is different, *Phonology*, 33(2):247–283.
- C. Douglas JOHNSON (1972), *Formal Aspects of Phonological Description*, Mouton, The Hague.
- Ronald M. KAPLAN and Martin KAY (1994), Regular models of phonological rule systems, *Computational linguistics*, 20(3):331–378.
- Lauri KARTTUNEN (1983), KIMMO: A general morphological processor, in *Texas Linguistic Forum*, volume 22, pp. 163–186.
- Paul KIPARSKY (1982), Lexical morphology and phonology, in I.-S. YANG, editor, *Linguistics in the morning calm: Selected papers from SICOL-1981*, pp. 3–91, Hansin, Seoul.
- Paul KIPARSKY (2010), Reduplication in stratal OT, in Linda UYECHE and Lian Hee WEE, editors, *Reality Exploration and Discovery: Pattern Interaction in Language & Life*, pp. 125–142, CSLI Press, Stanford.
- Gregory Michael KOBELE (2006), *Generating Copies: An Investigation into Structural Identity in Language and Grammar*, Ph.D. thesis, University of California, Los Angeles.

- Kimmo KOSKENNIEMI (1983), *Two-level morphology: A general computational model for word-form recognition and production*, Ph.D. thesis, University of Helsinki.
- Kimmo KOSKENNIEMI (1984), A general computational model for word-form recognition and production, in *Proceedings of the 10th international conference on Computational Linguistics*, pp. 178–181, Association for Computational Linguistics.
- D. Terence LANGENDOEN (1981), The generative capacity of word-formation components, *Linguistic Inquiry*, 12(2):320–322.
- Jeffrey LIDZ (2001), Echo reduplication in Kannada and the theory of word-formation, *Linguistic review*, 18(4):375–394.
- Huan LUO (2017), Long-distance consonant agreement and subsequentiality, *Glossa: A journal of General Linguistics*, 2(1):1–25, doi:<http://doi.org/10.5334/gjgl.42>.
- Alexis MANASTER-RAMER (1986), Copying in natural languages, context-freeness, and queue grammars, in *Proceedings of the 24th annual meeting on Association for Computational Linguistics*, pp. 85–89, Association for Computational Linguistics.
- Alec MARANTZ (1982), Re reduplication, *Linguistic Inquiry*, 13(3):435–482.
- Gary F. MARCUS, Sugumaran VIJAYAN, S. Bandi RAO, and Peter M. VISHTON (1999), Rule learning by seven-month-old infants, *Science*, 283(5398):77–80.
- Veronika MATTES (2007), *Reduplication in Bikol*, Ph.D. thesis, University of Graz, Graz, Austria.
- John J. MCCARTHY, Wendell KIMPER, and Kevin MULLIN (2012), Reduplication in Harmonic Serialism, *Morphology*, 22(2):173–232.
- John J. MCCARTHY and Alan PRINCE (1994), The emergence of the unmarked: Optimality in prosodic morphology, in Mercé GONZÁLEZ, editor, *Proceedings of the North East Linguistic Society 24*, p. 333–79, Graduate Linguistic Student Association, University of Massachusetts, Amherst, MA.
- John J. MCCARTHY and Alan PRINCE (1995), Faithfulness and reduplicative identity, in Jill N. BECKMAN, Laura Walsh DICKEY, and Suzanne URBANCZYK, editors, *Papers in Optimality Theory*, Graduate Linguistic Student Association, University of Massachusetts, Amherst, MA.
- Fiona MCLAUGHLIN (2005), Reduplication and consonant mutation in the Northern Atlantic languages, in Hurch (2005), pp. 111–134.
- Robert MCNAUGHTON and Seymour A. PAPERT (1971), *Counter-free automata*, MIT Press, Cambridge, MA.
- Mehryar MOHRI (1997), Finite-state transducers in language and speech processing, *Computational Linguistics*, 23(2):269–311.

Edith MORAVCSIK (1978), Reduplicative constructions, in Joseph GREENBERG, editor, *Universals of Human Language*, volume 1, pp. 297–334, Stanford University Press, Stanford, California.

Ajit NARAYANAN and Lama HASHEM (1993), On abstract finite-state morphology, in *Proceedings of the Sixth Conference on European Chapter of the Association for Computational Linguistics, Utrecht, The Netherlands, EACL '93*, pp. 297–304, Association for Computational Linguistics, Stroudsburg, PA, USA, ISBN 90-5434-014-2, doi:10.3115/976744.976779.

Mark-Jan NEDERHOF and Heiko VOGLER (2019), Regular transductions with MCFG input syntax, in *Proceedings of the 14th International Conference on Finite-State Methods and Natural Language Processing*, pp. 56–64, Association for Computational Linguistics, Dresden, Germany, <https://www.aclweb.org/anthology/W19-3109>.

Esa NELIMARKKA, Harri JÄPPINEN, and Aarno LEHTOLA (1984), Two-way finite automata and dependency grammar: A parsing method for inflectional free word order languages, in *Proceedings of the 10th international conference on Computational linguistics*, pp. 389–392, Association for Computational Linguistics.

Max NELSON, Hossep DOLATIAN, Jonathan RAWSKI, and Brandon PRICKETT (2020), Probing RNN encoder-decoder generalization of subregular functions using reduplication, in *Proceedings of the Society for Computation in Linguistics*, volume 3.

Nicole Alice NELSON (2003), *Asymmetric Anchoring*, Ph.D. thesis, Rutgers University, New Brunswick, NJ.

Andrew NEVINS (2004), What UG can and can't do to help the reduplication learner, in Aniko CSLRMAZ, Andrea GUALMINI, and Andrew NEVINS, editors, *MIT Working Papers in Linguistics 48*, pp. 113–126, MIT Department of Linguistics and Philosophy, Cambridge, MA.

Andrew NEVINS (2012), Haplological dissimilation at distinct stages of exponence, in Jochen TROMMER, editor, *The Morphology and Phonology of Exponence*, pp. 84–116, Oxford University Press, Oxford.

Andrew NEVINS and Bert VAUX (2003), Metalinguistic, shmetalinguistic: The phonology of shmreduplication, in *Proceedings from the Annual Meeting of the Chicago Linguistic Society*, volume 39, pp. 702–721, Chicago Linguistic Society, Chicago.

David ODDEN (1994), Adjacency parameters in phonology, *Language*, 70(2):289–330.

John J. OHALA, Joseph Paul STEMBERGER, and Marshall LEWIS (1986), Reduplication in Ewe: Morphological accommodation to phonological errors, *Phonology*, 3:151–160.

Amanda PAYNE (2014), Dissimilation as a subsequential process, in Jyoti IYER and Leland KUSMER, editors, *NELS 44: Proceedings of the 44th Meeting of the North East Linguistic Society*, volume 2, pp. 79–90, Graduate Linguistic Student Association, University of Massachusetts, Amherst, MA.

Amanda PAYNE (2017), All dissimilation is computationally subsequential, *Language: Phonological Analysis*, 93(4):e353–e371, doi:doi:10.1353/lan.2017.0076.

Christopher POTTS and Geoffrey K. PULLUM (2002), Model theory and the content of OT constraints, *Phonology*, 19(3):361–393.

Brandon PRICKETT, Aaron TRAYLOR, and Joe PATER (2018), Seq2Seq models with dropout can learn generalizable reduplication, in *Proceedings of the Fifteenth Workshop on Computational Research in Phonetics, Phonology, and Morphology*, pp. 93–100.

Michael O. RABIN and Dana SCOTT (1959), Finite automata and their decision problems, *IBM Journal of Research and Development*, 3(2):114–125.

Eric RAIMY (2000), *The Phonology and Morphology of Reduplication*, Mouton de Gruyter, Berlin.

Eric RAIMY (2009), Deriving reduplicative templates in a modular fashion, in Eric RAIMY and Charles E. CAIRNS, editors, *Contemporary views on architecture and representations in phonology*, number 48 in Current Studies in Linguistics, pp. 383–404, MIT Press, Cambridge, MA.

Eric RAIMY (2011), Reduplication, in Marc VAN OOSTENDORP, Colin EWEN, Elizabeth HUME, and Keren RICE, editors, *The Blackwell Companion to Phonology*, volume 4, pp. 2383–2413, Wiley-Blackwell, Malden, MA.

Charles REISS and Marc SIMPSON (2009), Reduplication as projection, unpublished manuscript, Concordia University, Montréal.

Jason RIGGLE (2004), Nonlocal reduplication, in Kier MOULTON and Matthew WOLF, editors, *Proceedings of the 34th meeting of the North Eastern Linguistics Society*, Graduate Linguistic Student Association, University of Massachusetts, Amherst, MA.

Brian ROARK and Richard SPROAT (2007), *Computational Approaches to Morphology and Syntax*, Oxford University Press, Oxford.

James ROGERS and Geoffrey PULLUM (2011), Aural pattern recognition experiments and the subregular hierarchy, *Journal of Logic, Language and Information*, 20:329–342.

Carl RUBINO (2005), Reduplication: Form, function and distribution, in Hurch (2005), pp. 11–29.

Carl RUBINO (2013), *Reduplication*, Max Planck Institute for Evolutionary Anthropology, Leipzig, <http://wals.info/chapter/27>.

- Jesse SABA KIRCHNER (2010), *Minimal reduplication*, Ph.D. thesis, University of California, Santa Cruz.
- Jesse SABA KIRCHNER (2013), Minimal reduplication and reduplicative exponence, *Morphology*, 23(2):227–243.
- Bridget SAMUELS (2010), The topology of infixation and reduplication, *The Linguistic Review*, 27(2):131–176.
- Walter J. SAVITCH (1982), *Abstract machines and grammars*, Little Brown and Company, Boston.
- Walter J. SAVITCH (1989), A formal model for context-free languages augmented with reduplication, *Computational Linguistics*, 15(4):250–261.
- Paul SCHACHTER and Victoria FROMKIN (1968), A phonology of Akan: Akuapem, Asante, Fante, in *UCLA Working Papers in Phonetics 9*, University of California, Los Angeles, Los Angeles.
- Marcel-Paul SCHÜTZENBERGER (1975), Sur certaines opérations de fermeture dans les langages rationnels, in *Symposia Mathematica*, volume 15, pp. 245–253.
- Hiroyuki SEKI, Takashi MATSUMURA, Mamoru FUJII, and Tadao KASAMI (1991), On multiple context-free grammars, *Theoretical Computer Science*, 88(2):191–229.
- Hiroyuki SEKI, Ryuichi NAKANISHI, Yuichi KAJI, Sachiko ANDO, and Tadao KASAMI (1993), Parallel multiple context-free grammars, finite-state translation systems, and polynomial-time recognizable subclasses of lexical-functional grammars, in *Proceedings of the 31st annual meeting on Association for Computational Linguistics*, pp. 130–139, Association for Computational Linguistics.
- Jeffrey SHALLIT (2008), *A Second Course in Formal Languages and Automata Theory*, Cambridge University Press, New York, NY, USA, 1 edition, ISBN 0521865727, 9780521865722.
- Patricia A. SHAW (2005), Non-adjacency in reduplication, in Hurch (2005), pp. 161–210.
- Daniel SILVERMAN (2002), Dynamic versus static phonotactic conditions in prosodic morphology, *Linguistics*, 40(1):29–60.
- Michael SIPSER (1980), Lower bounds on the size of sweeping automata, *Journal of Computer and System Sciences*, 21(2):195–202.
- Philip SPAELTI (1997), *Dimensions of Variation in Multi-Pattern Reduplication*, Ph.D. thesis, University of California, Santa Cruz.
- Richard William SPROAT (1992), *Morphology and Computation*, MIT press, Cambridge, MA.
- Donca STERIADE (1988), Reduplication and syllable transfer in Sanskrit and elsewhere, *Phonology*, 5(1):73–155.

Thomas STOLZ, Cornelia STROH, and Aina URDZE (2011), *Total Reduplication: The Areal Linguistics of a Potential Universal*, volume 8, Walter de Gruyter, Berlin.

Kristina STROTHER-GARCIA (2018), Imdlawn Tashlhiyt Berber syllabification is quantifier-free, in *Proceedings of the Society for Computation in Linguistics*, volume 1, pp. 145–153, doi:10.7275/R5J67F4D.

Kristina STROTHER-GARCIA (2019), *Using model theory in phonology: a novel characterization of syllable structure and syllabification*, Ph.D. thesis, University of Delaware.

Gregory STUMP (1995), Two types of mismatch between morphology and semantics, in Eric SCHILLER, Elisa STEINBERG, and Barbara NEED, editors, *Autolexical Theory: Ideas and Methods*, number 85 in Trends in Linguistics: Studies and Monographs, pp. 291–318, Mouton De Gruyter, Berlin.

Gregory STUMP (2001), *Inflectional morphology: A theory of paradigm structure*, number 93 in Cambridge Studies in Linguistics, Cambridge University Press, Cambridge.

Terry TAI, Wojciech SKUT, and Richard SPROAT (2011), Thrax: An open source grammar compiler built on OpenFst, in *IEEE Automatic Speech Recognition and Understanding Workshop*, volume 12.

Suzanne URBANCZYK (1999), Double reduplications in parallel, in René KAGER, Harry VAN DER HULST, and Wim ZONNEVELD, editors, *The prosody-morphology interface*, pp. 390–428, Cambridge University Press, Cambridge.

Suzanne URBANCZYK (2001), *Patterns of reduplication in Lushootseed*, Garland, New York.

Suzanne URBANCZYK (2007), Themes in phonology, in Paul DE LACY, editor, *The Cambridge Handbook of Phonology*, pp. 473–493.

Suzanne URBANCZYK (2011), Reduplication, in Mark ARONOFF, editor, *Oxford Bibliography*,
<http://oxfordindex.oup.com/view/10.1093/obo/9780199772810-0036>.

Odile VAYSSE (1986), Addition molle et fonctions p-locales, in *Semigroup Forum*, volume 34, pp. 157–175, Springer.

Rachelle WAKSLER (1999), Cross-linguistic evidence for morphological representation in the mental lexicon, *Brain and Language*, 68(1-2):68–74.

Markus WALTHER (2000), Finite-state reduplication in one-level prosodic morphology, in *Proceedings of the 1st North American chapter of the Association for Computational Linguistics conference*, NAACL 2000, pp. 296–302, Association for Computational Linguistics, Seattle, Washington,
<http://dl.acm.org/citation.cfm?id=974305.974344>.

Ronnie B. WILBUR (1973), *The Phonology of Reduplication*, Ph.D. thesis, University of Indiana, Bloomington, Indiana.

Ronnie B WILBUR (2005), A reanalysis of reduplication in American Sign Language, in Hurch (2005), pp. 595–623.

Colin WILSON (2019), Re (current) reduplication: Interpretable neural network models of morphological copying, *Proceedings of the Society for Computation in Linguistics*, 2(1):379–380.

Moira YIP (1995), Repetition and its avoidance: The case of Javanese, in Keiichiro SUZUKI and Dirk ELZINGA, editors, *Proceedings of the South Western Optimality Theory workshop 1995. Arizona Phonology Conference Volume 5*, pp. 238–262, University of Arizona, Tucson, AZ.

Alan C.L. YU (2007), *A Natural History of Infixation*, number 15 in Oxford Studies in Theoretical Linguistics, Oxford University Press, Oxford.

Kristine YU (2017), Advantages of constituency: Computational perspectives on Samoan word prosody, in *International Conference on Formal Grammar 2017*, p. 105–124, Springer, Berlin.

Sam ZUKOFF (2017), *Indo-European Reduplication: Synchrony, Diachrony, and Theory*, Ph.D. thesis, Massachusetts Institute of Technology.

Kie ZURAW, M. Yu KRISTINE, and Robyn ORFITELLI (2014), The word-level prosody of Samoan, *Phonology*, 31(2):271–327.

Hossep Dolatian

© 0000-0001-5044-8434

hossep.dolatian@stonybrook.edu

Jeffrey Heinz

© 0000-0002-5954-3195

jeffrey.heinz@stonybrook.edu

Department of Linguistics

Institute of Advanced Computational Science

Stony Brook University

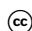
Stony Brook, NY, US

<https://you.stonybrook.edu/deovlet/>

Hossep Dolatian and Jeffrey Heinz (2020), *Computing and classifying reduplication with 2-way finite-state transducers*, *Journal of Language Modelling*, 8(1):179–250

doi <https://dx.doi.org/10.15398/jlm.v8i1.-1>

This work is licensed under the *Creative Commons Attribution 4.0 Public License*.

©  <http://creativecommons.org/licenses/by/4.0/>