

# REGULAR TRANSFORMATIONS IN LINGUISTICS

Jeffrey Heinz



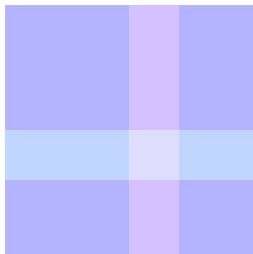
Stony Brook University

Schloss Dagstuhl  
16.05.2023

# MOTIVATION

- 1 What kind of computational resources and information are both necessary and sufficient to account for linguistic transformations?
  - 2 How can these transformations be learned from examples?
- ★ How can we characterize successful learning, including characterizing the amount and kind of data and the kinds of target concepts for such learning to occur?

Concepts



Data

# OVERVIEW OF ANSWERS

- 1 Regular transformations provide an upper bound.
- 2 Particular subclasses appear to be sufficient, and these subclasses are generally (much) less than what  $\text{FO}(<)$  provides.
- 3 These subclasses are learnable with relatively low time and data complexity in ways the full class of regular transformations is not.

# QUIZTIME! POSSIBLE LINGUISTIC PROCESSES

- **Prefixation**  
abbb  $\mapsto$  cabbb
- **Affixation Mod 2**  
abbb  $\mapsto$  cabbb  
abb  $\mapsto$  abbc
- **Bounded Spreading**  
abbb  $\mapsto$  aabb
- **Unbounded Spreading**  
abbb  $\mapsto$  aaaa  
abbbdb  $\mapsto$  aaaadb
- **Projected Unbounded Spreading**  
 $c^{10}ac^{10}bc^{10}bc^{10}bc^{10}$   
 $\mapsto c^{10}ac^{10}ac^{10}ac^{10}ac^{10}$
- **Sour Grapes**  
abbbb  $\mapsto$  aaaaa  
abccb  $\mapsto$  abccb
- **Two-sided Unbounded Spread**  
abba  $\mapsto$  aaaa  
abbb  $\mapsto$  abbb
- **Majority Rules**  
abbaa  $\mapsto$  aaaaa  
abbab  $\mapsto$  bbbbbb
- **Partial Copying**  
abcd  $\mapsto$  ababcd
- **Full Copying**  
abcd  $\mapsto$  abcdabcd
- **Triplication**  
abcd  $\mapsto$  abcdabcdabcd
- **Squaring**  
abcd  $\mapsto$  abcdabcdabcdabcd
- **Iterated Prefix Copying**  
abcd  $\mapsto$  a ab abc abcd

# QUIZTIME! POSSIBLE LINGUISTIC PROCESSES

- ✓ Prefixation  
abbb  $\mapsto$  cabbb
- ✗ Affixation Mod 2  
abbb  $\mapsto$  cabbb  
abb  $\mapsto$  abbc
- ✓ Bounded Spreading  
abbb  $\mapsto$  aabb
- ✓ Unbounded Spreading  
abbb  $\mapsto$  aaaa  
abbbdb  $\mapsto$  aaaadb
- ✓ Projected Unbounded Spreading  
 $c^{10}ac^{10}bc^{10}bc^{10}bc^{10}$   
 $\mapsto c^{10}ac^{10}ac^{10}ac^{10}ac^{10}$
- ✓✗ Sour Grapes  
abbbb  $\mapsto$  aaaaa  
abccb  $\mapsto$  abccb
- ✓ Two-sided Unbounded Spread  
abba  $\mapsto$  aaaa  
abbb  $\mapsto$  abbb
- ✗ Majority Rules  
abbaa  $\mapsto$  aaaaa  
abbab  $\mapsto$  bbbbb
- ✓ Partial Copying  
abcd  $\mapsto$  ababcd
- ✓ Full Copying  
abcd  $\mapsto$  abcdabcd
- ✓ Triplication  
abcd  $\mapsto$  abcdabcdabcd
- ✗ Squaring  
abcd  $\mapsto$  abcdabcdabcdabcd
- ✗ Iterated Prefix Copying  
abcd  $\mapsto$  a ab abc abcd

# EXAMPLES 1: PHONOLOGY

In Navajo, sibilant sounds like {s, z} are pronounced {ʃ, ʒ} if they are followed anywhere by {ʃ, ʒ} (Sapir and Hoijer 1967).

- a. /sì-ʔá/    ↦    sì-ʔá    ‘a round object lies’
- b. /sì-tí/    ↦    sì-tí    ‘he is lying’
- c. /sì-yìʃ/   ↦    ʃì-yìʃ   ‘it lies bent, curved’
- d. /sì-te:ʒ/ ↦    ʃite:ʒ   ‘we (dual) are lying’

Projected Unbounded Spreading from Right to Left

# EXAMPLES 1: PHONOLOGY

In Luganda, every toneless vowel between two high-toned vowels is pronounced with a high tone forming a high-tone ‘plateau’ (Hyman and Katamba 2010).

- a. /mu-tund-a+bi-kópo/     $\mapsto$     mutunda-bikópo    ‘cup-seller’  
L L L L H L L                     $\mapsto$     L L L L H L L
- b. /mu-tém-a+bi-sikí/     $\mapsto$     mutémá-bísíkî    ‘log-chopper’  
L H L L H L L                     $\mapsto$     L H H H H L F

Two-sided Unbounded Spread

## EXAMPLES 2: MORPHOLOGY

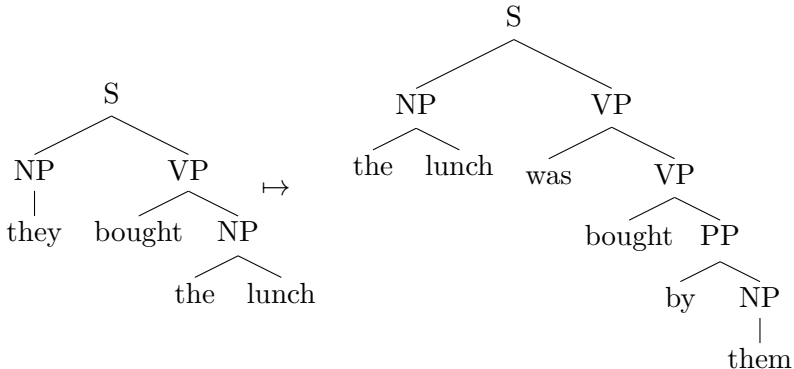
Morphological copying, “Reduplication”, (83% of languages have full copying, 74% have partial bounded copying)

- **Full copying** (Indonesian, Cohn 1989)  
buku  $\mapsto$  buku  $\sim$  buku                                ‘book’  $\mapsto$  ‘books’  
wanita  $\mapsto$  wanita  $\sim$  wanita                            ‘woman’  $\mapsto$  ‘women’
- **Initial-CVC copying** (Panganisan, Rubino 2005:11)  
baley  $\mapsto$  bal  $\sim$  baley                                        ‘town’  $\mapsto$  ‘towns’
- **Initial-CVC copying and opposite-edge placement**  
(Koryat, Riggle 2004:3)  
qanga  $\mapsto$  qanga  $\sim$  qan                                        ‘fire’  $\mapsto$  ‘fire (ABS)’
- **Root-final CVC copying and word-initial placement**  
(Madurese, Brown 2017:964)  
pa-jalan-an  $\mapsto$  lan  $\sim$  pa-jalan-an  
‘pedestrian’  $\mapsto$  ‘pedestrians’



## EXAMPLES 3: SYNTAX

The active to passive transformation can be analyzed in terms of transformations over syntactic trees.



# SUB-REGULARITY IN LANGUAGE

S Non Regular

---

P Regular

---

CNL(X) / QF(X)  
(Appropriately Subregular)

---

strings

20th century view

# SUB-REGULARITY IN LANGUAGE

S

Non Regular

---

Regular

---

P

CNL(X) / QF(X)  
(Appropriately Subregular)

---

strings

(Heinz 2018)

# SUB-REGULARITY IN LANGUAGE

Non Regular

---

S

Regular

---

P

CNL(X) / QF(X)  
(Appropriately Subregular)

---

trees strings

(Stabler 2019)

# SUB-REGULARITY IN LANGUAGE

Non Regular

---

Regular

---

S

P

CNL(X) / QF(X)  
(Appropriately Subregular)

---

trees

strings

(Graf 2022)

# WHERE WE STARTED

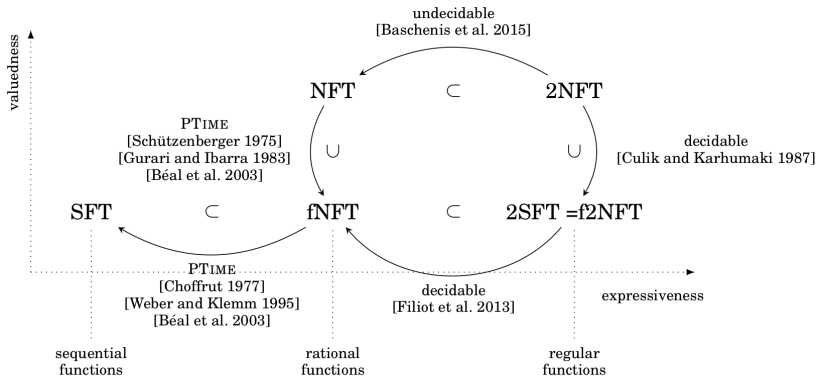


Fig. 3. A landscape of transducers of finite words.

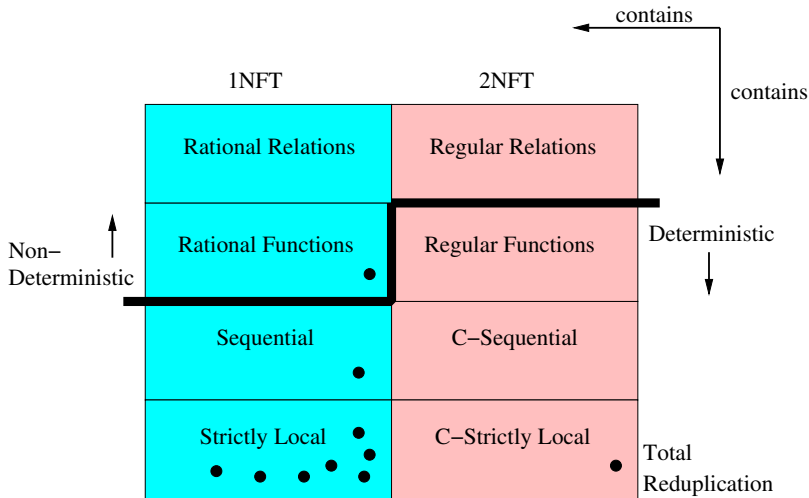
(Filiot and Reynier 2016 ACM SIGLOG)

# ROUND 1 SUMMARY OF STRING FUNCTIONS IN MORPHOLOGY AND PHONOLOGY

- 1 Many morphological and phonological operations are either left or right sequential, except for **sour grapes**, **two-sided unbounded spreading** and **full copying**
- 2 Natural Language processing and computational linguistics have focused on rational relations, leaving full copying as a “problem” (Roark and Sproat 2007).
- 3 Two reasons:
  - Those communities were unaware of 2-way transducers, Courcelle-style logical transductions, and streaming string transducers
  - Rational relations are invertible, which provides one model for both generation and analysis.

# STRING RELATIONS: ROOM AT THE BOTTOM

A more articulated view



(Filiot and Reynier 2016, Chandler 2017, Dolatian and Heinz 2020)



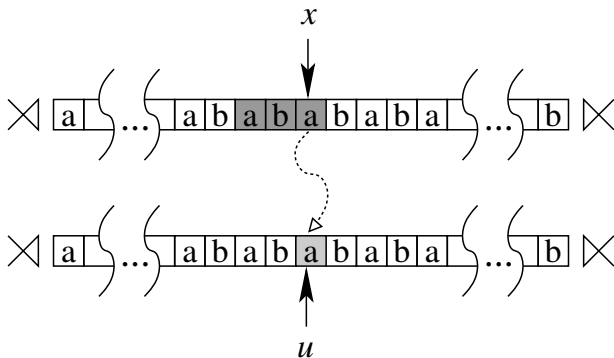
# LOCAL FUNCTIONS

Dans cet article, nous introduisons les fonctions  $p$ -locales et les fonctions  $p$ -sous locales (où  $p$  est un entier strictement positif) et nous les caractérisons par une propriété simple de leur semi-groupe syntactique : ce semigroupe doit satisfaire l'équation  $yx_1 \dots x_p = x_1 \dots x_p$ . Nous en déduisons quelques propriétés des fonctions  $p$ -locales.

(Vaysse 1986, see also Schiffler 1973)

# LOCAL FUNCTIONS AS SCANNERS WITH OUTPUTS

Vaysse defines local functions in terms of a sliding window over the input which produces corresponding outputs.



# LOCAL FUNCTIONS, NERODE/TRANSDUCERS

- 1 Independently, Chandlee et al. 2014 define these functions in terms of their Nerode equivalence classes: two inputs  $x$  and  $y$  are Nerode equivalent (that is have the same functional residuals) whenever they share the same  $k$ -sized suffix

$$\text{TAILS}_f(x) = \{(y, v) \mid f(xy) = \text{LCP}(f(x\Sigma^*))v\}$$

$$x_1 \sim_f^N x_2 \leftrightarrow \text{TAILS}_f(x_1) = \text{TAILS}_f(x_2)$$

$$\text{SUFF}_k(x_1) = \text{SUFF}_k(x_2) \Rightarrow \text{TAILS}_f(x_1) = \text{TAILS}_f(x_2)$$

- 2 Their transducer characterization is the same as Vaysse's transducer characterization.

# LOCAL FUNCTIONS, LOGICALLY

- 1 Lindell and Chandlee (2016) prove that quantifier-free interpretations of strings to strings (modeled with the successor and predecessor functions) are local functions (but not quite vice versa).
- 2 Provides a way to define “local” transformations over any model-theoretic representations, not just strings!
- 3 Subsequent work has explored this with respect to linguistically motivated graph representations of words (Jardine 2016, Strother-Garcia 2018, Jardine et al. 2021)

- 1 Many phonological and morphological processes can be described with local functions such as [Prefixation](#), [Suffixation](#), [Bounded Spreading](#) ... but not all.  
(Chandlee 2014, Chandlee 2017, Chandlee and Heinz 2018, Chandlee et al. 2018)
- 2 Chandlee et al. 2014 provide better time and complexity bounds for learning these functions than previously known (see also Jardine et al. 2014).
- 3 Current work: doing better with alternative representations of strings.

# LOCAL FUNCTIONS, ALGEBRAICALLY

- Vaysse 1986 showed that local functions have a *definite* algebraic structure.

$$Se = e$$

(recalling  $yx_1x_2 \dots x_p = x_1x_2 \dots x_p$ )

- Syntactic Monoids for Transducers are based on the Myhill equivalence relation:

$$\text{CONTEXTS}_f(x) = \{(w, y, v) \mid f(wxy) = \text{LCP}(f(wx\Sigma^*))v\}$$

$$x_1 \sim_f^M x_2 \leftrightarrow \text{CONTEXTS}_f(x_1) = \text{CONTEXTS}_f(x_2)$$

(Filiot, Gauwin, Lhote 2016)

# REVERSE DEFINITE STRUCTURES

$$eS = e$$

$$(x_1x_2 \dots x_p y = x_1x_2 \dots x_p)$$

- 1 Transducers come in two kinds: left and right.
- 2 **Unbounded Spreading** in natural languages includes both left and right Reverse Definite.

(Lambert 2022)

# TIER-PROJECTIONS (LAMBERT 2022)

- 1 A tier  $T$  is a subset of the alphabet which contain “salient” symbols.

In Navajo  $T$  equals the sibilants like  $\{s, z, \mathfrak{f}, \mathfrak{z}\}$ .

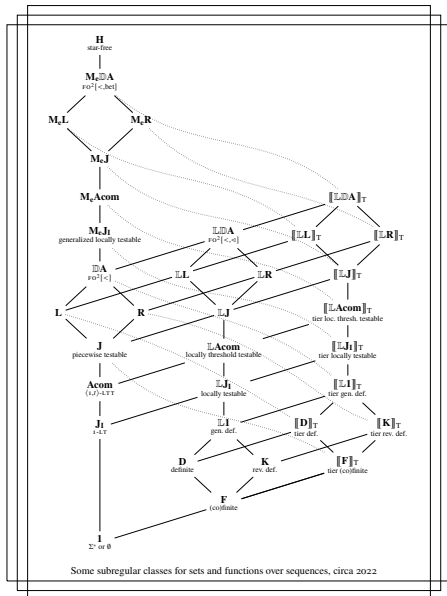
- 2 Projecting a string to a tier removes the non-salient symbols.

$$s\grave{i}-y\grave{i}\mathfrak{f} \mapsto s\mathfrak{f}$$

- 3 A function belongs to algebraic class “Tier X” if and only if its tier projection belongs to class X.
- 4 Lambert shows for each X that  $X \subsetneq \text{Tier X}$ .
  - Another flavor of **Unbounded Spreading** that is attested is Tier Definite.
  - **Projected Unbounded Spreading** is Tier Reverse Definite, with both left and right directions attested..



# LAMBERT 2022: ALGEBRAIC HIERARCHY



# ROUND 2 SUMMARY OF STRING FUNCTIONS IN MORPHOLOGY AND PHONOLOGY

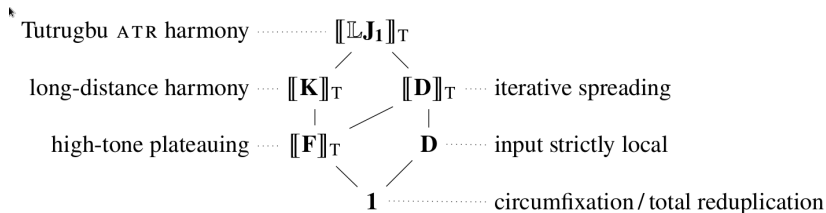


Figure 5.18: Some attested morphophonological functions.

Lambert (2022)

- Both Tutrugbu ATR Harmony ([Sour Grapes](#)) High Tone Plateauing ([Two sided Unbounded Spreading](#)) require non-deterministic or 2-way transducers.
- Their algebraic characterizations follow the methods of Carton and Dartois (2015).

# LEARNING REVISITED

- 1 Each of the classes  $X \leq \text{Tier Locally Testable}$  can be parameterized by a window size  $k$  and/or tier  $T$ .
- 2 Given  $k$  and  $T$ , the  $k$ - $X$  classes can be learned from examples with relatively low time and data complexity.
- 3 Given only  $k$ , learning the tier  $T$  is possible for languages (Jardine and Heinz 2016, Jardine and McMullin 2017) and there is some work on functions as well (Burness and McMullin 2019) that still needs to be integrated into the algebraic perspective.

# LOCAL TREE TRANSDUCTIONS

The states of the transducer correspond to local treelets of bounded depth.

- ① Jing and Heinz (2020) restricts deterministic bottom-up tree transducers.
- ② Graf (2020) takes a top-down approach.
- ③ Someone should look at QF interpretations of tree structures.

# CONCLUSION

- 1 Regular transformations provide an upper bound on morpho-phonological processes in natural languages.
- 2 Particular subclasses appear to be sufficient, and these subclasses are generally (much) less than what  $FO(<)$  provides.
- 3 These subclasses are learnable with relatively low time and data complexity in ways the full class of regular transformations is not.

# OPEN QUESTIONS

- ① Factoring transducers via composition and/or direct product
- ② Learning factored representations of transducers
- ③ Subregular classes of tree transductions
- ④ Deterministic regular relations

Thank You