

Morphotactics as Tier-Based Strictly Local Dependencies

Alëna Aksënova, Thomas Graf,
and Sedigheh Moradi

Stony Brook University

SIGMORPHON 14
Berlin, Germany
11. August 2016

Our goal

Received view

Phonology

regular

Kaplan&Kay (1994)

Recent research

subregular

Heinz (2015)

Morphology

regular

Beesley&Karttunen (2003)

?

Our goal

	Phonology	Morphology
<i>Received view</i>	regular Kaplan&Kay (1994)	regular Beesley&Karttunen (2003)
<i>Recent research</i>	subregular Heinz (2015)	?

- Show that **morphotactics is subregular**
- More precisely: **Tier-Based Strictly Local**
- **Consequences**
 - parallels to phonology
 - learnable in the limit from positive text
 - explain typological gaps

Outline

- 1 Preliminaries
- 2 SL Patterns In Morphology
- 3 Tier-Based Strictly Local
 - TSL is necessary
 - TSL is sufficient
- 4 Typological Gaps

Morphotactics

Definition (Morphotactics)

Restrictions on the linear ordering of morphemes.

- Our focus: morphotactics in underlying representations (English) OK STEM-PL * PL-STEM
- \Rightarrow allomorphy (dogs, peaches) is not considered yet

Computational nature of morphotactics

	Phonology	Morphology
<i>Received view</i>	regular	regular
<i>Recent research</i>	subregular	?

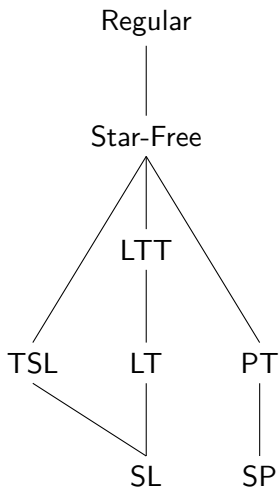
Advantages of (some) subregular languages:

- resolves learnability issues
- describes potential cognitive mechanisms
- uses less powerful generating device

Subregular Phonology and Morphology

not all languages exploit full power of
finite-state machinery
⇒ subregular hierarchy

Subregular hierarchy



Subregular Phonology and Morphology

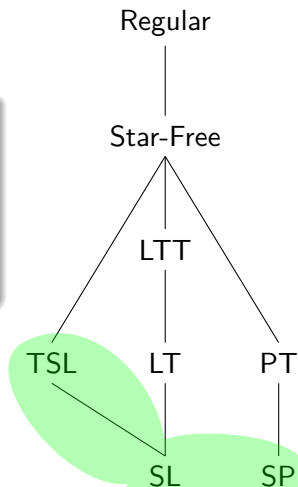
not all languages exploit full power of finite-state machinery
⇒ subregular hierarchy

Strong Subregular Hypothesis

All **phonological dependencies** are

- strictly local (SL)
- tier-based strictly local (TSL)
- strictly piecewise (SP)

Subregular hierarchy



Subregular Phonology and Morphology

not all languages exploit full power of finite-state machinery
 ⇒ subregular hierarchy

Strong Subregular Hypothesis

All **phonological dependencies** are

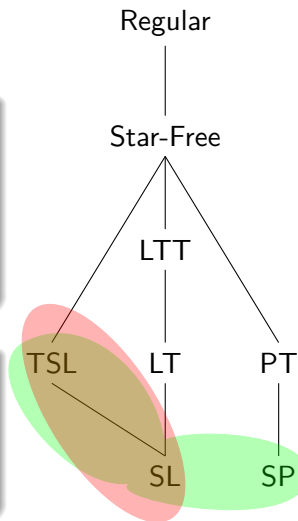
- strictly local (SL)
- tier-based strictly local (TSL)
- strictly piecewise (SP)

Subregular Morphotactics

All **morphotactic dependencies** are

- strictly local (SL)
- tier-based strictly local (TSL)

Subregular hierarchy



Strictly Local languages

- SL and TSL are generated by ***k*-gram models**.
- A *k*-gram model is a finite set of blocked *k*-grams.

Strictly Local languages

- SL and TSL are generated by ***k*-gram models**.
- A *k*-gram model is a finite set of blocked *k*-grams.

Example (Strictly Local Grammar for $(ab)^*a$)

$\Sigma = \{a, b\}$

Grammar = $\{\times b, bb, aa, b\times, \times\times\}$

Accepted strings: $\times a\times, \times aba\times, \times ababa\times, \text{etc.}$

Rejected strings: $\times ab\times, \times ba\times, \times abba\times, \text{etc.}$

Strictly Local languages

- SL and TSL are generated by ***k*-gram models**.
- A *k*-gram model is a finite set of blocked *k*-grams.

Example (Strictly Local Grammar for $(ab)^*a$)

$\Sigma = \{a, b\}$

Grammar = $\{\times b, bb, aa, b\times, \times\times\}$

Accepted strings: $\times a\times, \times aba\times, \times ababa\times, \text{etc.}$

Rejected strings: $\times ab\times, \times ba\times, \times abba\times, \text{etc.}$

Strictly Local languages

- SL and TSL are generated by ***k*-gram models**.
- A *k*-gram model is a finite set of blocked *k*-grams.

Example (Strictly Local Grammar for $(ab)^*a$)

$\Sigma = \{a, b\}$

Grammar = $\{\times b, bb, aa, b\times, \times\times\}$

Accepted strings: $\times a\times, \times aba\times, \times ababab\times, \text{etc.}$

Rejected strings: $\times ab\times, \times ba\times, \times abba\times, \text{etc.}$

Definition

Strictly *k*-Local (*k*-SL) grammar consists of a set of blocked *k*-grams over an alphabet Σ .

Tier-Based Strictly Local languages

Example (Tier-Based Strictly Local Grammar for $c^*(ac^*bc^*)^*ac^*$)

$\Sigma = \{a, b, c\}$

Grammar:

$G(a,b_tier) = \{\times b, bb, aa, b\times, \times\times\}$

Accepted strings: $\times a\times, \times accba\times, \times cacbacccba\times, \text{etc.}$

Tier-Based Strictly Local languages

Example (Tier-Based Strictly Local Grammar for $c^*(ac^*bc^*)^*ac^*$)

$\Sigma = \{a, b, c\}$

Grammar:

$G(a,b_tier) = \{\times b, bb, aa, b\times, \times\times\}$

Accepted strings: $\times a\times$, $\times accba\times$, $\times cacbaccbccba\times$, etc.

a,b_tier : $\times a\times$ $\times aba\times$ $\times ababa\times$

Tier-Based Strictly Local languages

Example (Tier-Based Strictly Local Grammar for $c^*(ac^*bc^*)^*ac^*$)

$\Sigma = \{a, b, c\}$

Grammar:

$G(a,b_tier) = \{\times b, bb, aa, b\times, \times\times\}$

Accepted strings: $\times a\times$, $\times accba\times$, $\times cacbaccbccba\times$, etc.

a,b_tier : $\times a\times$ $\times aba\times$ $\times ababa\times$

Rejected strings: $\times accccaba\times$, $\times abcccaccbcc\times$, etc.

Tier-Based Strictly Local languages

Example (Tier-Based Strictly Local Grammar for $c^*(ac^*bc^*)^*ac^*$)

$\Sigma = \{a, b, c\}$

Grammar:

$G(a,b_tier) = \{\times b, bb, aa, b\times, \times\times\}$

Accepted strings: $\times a\times$, $\times accba\times$, $\times cacbaccbccba\times$, etc.

a,b_tier : $\times a\times$ $\times aba\times$ $\times ababa\times$

Rejected strings: $\times accccaba\times$, $\times abcccaccbcc\times$, etc.

a,b_tier : $\times \color{red}aa\times$ $\times abab\color{red}b\times$

Tier-Based Strictly Local languages

Example (Tier-Based Strictly Local Grammar for $c^*(ac^*bc^*)^*ac^*$)

$\Sigma = \{a, b, c\}$

Grammar:

$G(a,b_tier) = \{\times b, bb, aa, b\times, \times\times\}$

Accepted strings: $\times a\times$, $\times accba\times$, $\times cacbaccbccba\times$, etc.

a,b_tier : $\times a\times$ $\times aba\times$ $\times ababa\times$

Rejected strings: $\times accccaba\times$, $\times abcccaccbcc\times$, etc.

a,b_tier : $\times \color{red}{aa}ba\times$ $\times abab\color{red}{b}\times$

Definition

A Tier-Based Strictly k -Local grammar is a k -SL grammar that operates over a *tier*, a specific substructure of the string.

Learnability

Learning of SL and TSL

- learning \equiv memorizing finite number of k -grams + tier induction
- learnable in the limit from positive text

Jardine & Heinz (2016)

Mappings we use

General assumption: we assume stem not to be bound in length:

- There is no limit on the length of the stem in languages.
- The stem can be result of the compounding.
whiteboard, whiteboard marker, whiteboard marker cleaning fluid,
whiteboard marker cleaning fluid purchase receipt
- Mapping of the stem to a single symbol will result in insensibility to compounds.

Mappings we use

General assumption: we assume stem not to be bound in length:

- There is no limit on the length of the stem in languages.
- The stem can be result of the compounding.
whiteboard, whiteboard marker, whiteboard marker cleaning fluid,
whiteboard marker cleaning fluid purchase receipt
- Mapping of the stem to a single symbol will result in insensibility to compounds.

- Affixes: affix-to-symbol mapping
- Stems: symbol-to-symbol mapping

Strictly Local Morphology: affixation

Example (prefix 'za-', Russian)

- exat'
'go, drive'
xxxx
- zaexat'
'call on the way'
axxxx

Bigram *xa ensures that 'za' is a prefix.

Strictly Local Morphology: affixation

Example (prefix 'za-', Russian)

- exat'
'go, drive'
xxxx
- zaexat'
'call on the way'
axxxx

Bigram *xa ensures that 'za' is a prefix.

Example (suffix '-s', English)

- dog
xxx
- dogs
xxxs

Bigram *bx ensures that 's' is a suffix.

Strictly Local Morphology: affixation [cont.]

Example (affixation, English)

● lock xxxx	● blacklist xxxxxxxxxx
● unlockable axxxx b	● unblacklistable axxxxxxxxx b

Strictly Local Morphology: affixation [cont.]

Example (affixation, English)

- lock

xxxx

- unlockable

axxxx**b**

- blacklist

xxxxxxxxx

- unblacklistable

axxxxxxxxx**b**

$$\text{SLG} = \{\times b, ba, bx, xa, a\times\}$$

This grammar necessarily generates the following forms of English, too: $\times axxxx\times$ and $\times xxxxb\times$.

Strictly Local Morphology: affixation [cont.]

Example (affixation, English)

- lock
xxxx
- unlockable
axxxx**b**
- blacklist
xxxxxxxxx
- unblacklistable
axxxxxxxxx**b**

$$\text{SLG} = \{\times b, ba, bx, xa, a\times\}$$

This grammar necessarily generates the following forms of English, too: $\times axxxx\times$ and $\times xxxxb\times$.

Indeed, this prediction is correct:

Example (affixation, English)

- unleash
axxxxx
- breakable
xxxx**b**

SL is not enough: Indonesian circumfixation

- English **un-...-able** are prefix and suffix that can co-occur.
- However, two parts of a *circumfix* cannot occur independently:

Consider the following example from Indonesian:

Example (circumfix 'ke-an', Indonesian)

- | | |
|---|--|
| <ul style="list-style-type: none"> ● tinggi
'high'
xxxxxxx | <ul style="list-style-type: none"> ● mahasiswa
'big pupil (student)'
xxxxxxxxxxx |
| <ul style="list-style-type: none"> ● ketinggian
'altitude'
axxxxxxxxb | <ul style="list-style-type: none"> ● kemahasiswaan
'student affairs'
axxxxxxxxxxxb |
| <ul style="list-style-type: none"> ● *axxxxxxxx | <ul style="list-style-type: none"> ● *xxxxxxxxb |

SL is not enough: Indonesian circumfixation [cont.]

Example (circumfix 'ke-an', Indonesian)

● tinggi

xxxxxxx

● ketinggian

axxxxxxb

● *axxxxxx

● mahasiswa

xxxxxxxxx

● kemahasiswaan

axxxxxxxxxxb

● *xxxxxxx

SLG = { \times b, ba, bx, xa, a \times }String language = \times xxxx \times , \times axxxxxb \times , \times axxxx \times , \times xxb \times ...

SL is not enough: Indonesian circumfixation [cont.]

Example (circumfix 'ke-an', Indonesian)

- | | |
|--------------------------|---------------------------------|
| ● tinggi
xxxxxx | ● mahasiswa
xxxxxxxxx |
| ● ketinggian
axxxxxxb | ● kemahasiswaan
axxxxxxxxxxb |
| ● *axxxxxx | ● *xxxxxxxxb |

SLG = { $\times b$, ba , bx , xa , $a \times$ }

String language = $\times xxx \times$, $\times axxxxxb \times$, $\times axxxx \times$, $\times xxb \times \dots$

Problem:

- SL languages can only capture local dependencies
- Circumfixes introduce non-local ones

Morphotactics is TSL

Example (circumfix 'ke-an', Indonesian)

- tinggi

xxxxxx

- ketinggian

axxxxxxb

- *axxxxxx

- mahasiswa

xxxxxxxxx

- kemahasiswaan

axxxxxxxxxxb

- *xxxxxxxxb

$$\text{TSLG}(\text{circumfix_tier}) = \{\times b, ba, a\times\}$$

Morphotactics is TSL

Example (circumfix 'ke-an', Indonesian)

- tinggi

xxxxxx

- ketinggian

axxxxxxb

- *axxxxxx

- mahasiswa

xxxxxxxxx

- kemahasiswaan

axxxxxxxxxxb

- *xxxxxxxxb

$TSLG(\text{circumfix_tier}) = \{\times b, ba, a\times\}$

Licit strings:

- $\times \text{xxxxxx} \times$

$\times \times$

- $\times \text{axxxxxxb} \times$

$\times ab \times$

Illicit strings:

- $\times \text{axxxxx} \times$

$\times a \times$

- $\times \text{bxxxxa} \times$

$\times ba \times$

Morphotactics is TSL

Example (circumfix 'ka-an', Ilocano)

In Ilocano, it is impossible to do embedded circumfixation:

- bigát
'morning'
xxxxx

- kabigátan
'the next morning'
axxxxxb

- *aaxxxxxb

Morphotactics is TSL [cont.]

Example (circumfix 'ka-an', Ilocano)

- bigát
xxxxx

- k**abigátan**
a**xxxxxb**

- ***aa**xxxx**bb**

$$\text{TSLG}(\text{circumfix_tier}) = \{\times b, ba, a\times, aa, bb\}$$

Morphotactics is TSL [cont.]

Example (circumfix 'ka-an', Ilocano)

- bigát
xxxxx

- k**abigátan**
a**xxxxxb**

- ***aa**xxxx**bb**

$TSLG(\text{circumfix_tier}) = \{\times b, ba, a\times, aa, bb\}$

Licit strings:

- \times xxxxxx \times
 $\times\times$

- \times axxxxx**b** \times
 \times ab \times

Illicit strings:

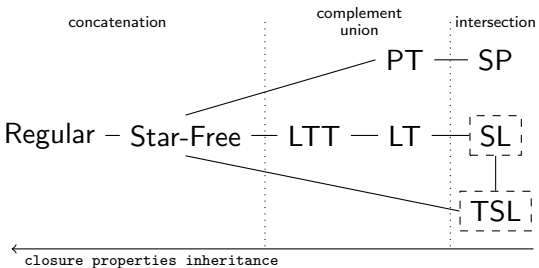
- \times aa**xxxxbb** \times
 \times **abb** \times

- \times **b**xxxxa \times
 \times **ba** \times

Interim Summary

- SL enforces local dependencies
- TSL enforces local dependencies on the determined tier
- Most of morphotactics is SL, some of it is TSL
- Learning of TSL languages is possible from positive data only
- Can morphotactics be more than TSL?

Can morphotactics be more than TSL?



Can morphotactics be more than TSL?

- **Closure under concatenation**: Frenglish contains only words whose first part is a word of French and the second a word of English.

Can morphotactics be more than TSL?

- **Closure under concatenation**: Frenglish contains only words whose first part is a word of French and the second a word of English. **X**

Can morphotactics be more than TSL?

- **Closure under concatenation**: Frenglish contains only words whose first part is a word of French and the second a word of English. **X**
- **Closure under union**: If a Mandaresian word violates rules of Mandarin Chinese, it must obey the rules of Indonesian.

Can morphotactics be more than TSL?

- **Closure under concatenation**: Frenglish contains only words whose first part is a word of French and the second a word of English. **X**
- **Closure under union**: If a Mandaresian word violates rules of Mandarin Chinese, it must obey the rules of Indonesian. **X**

Can morphotactics be more than TSL?

- **Closure under concatenation**: Frenglish contains only words whose first part is a word of French and the second a word of English. **X**
- **Closure under union**: If a Mandaresian word violates rules of Mandarin Chinese, it must obey the rules of Indonesian. **X**
- **Closure under relative complement**: Hsilgne contains all words that are ill-formed in English.

Can morphotactics be more than TSL?

- **Closure under concatenation**: Frenglish contains only words whose first part is a word of French and the second a word of English. **X**
- **Closure under union**: If a Mandaresian word violates rules of Mandarin Chinese, it must obey the rules of Indonesian. **X**
- **Closure under relative complement**: Hsilgne contains all words that are ill-formed in English. **X**

Can morphotactics be more than TSL?

- **Closure under intersection:** Russenorsk is created by combination of elements of Russian and Norwegian.

Can morphotactics be more than TSL?

- **Closure under intersection:** Russenorsk is created by combination of elements of Russian and Norwegian. ✓
(spoken in Northern Norway, 18th-19th centuries)

Can morphotactics be more than TSL?

- **Closure under intersection:** Russenorsk is created by combination of elements of Russian and Norwegian. ✓
(spoken in Northern Norway, 18th-19th centuries)

Example (Closure under intersection)

- A language allows complex nuclei and blocks codas (Supyire)
- A language forbids complex nuclei and allows codas (Russian)

Can morphotactics be more than TSL?

- **Closure under intersection:** Russenorsk is created by combination of elements of Russian and Norwegian. ✓
(spoken in Northern Norway, 18th-19th centuries)

Example (Closure under intersection)

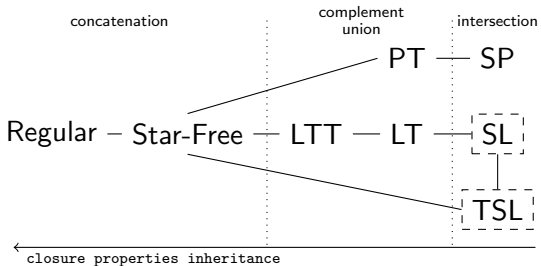
- A language allows complex nuclei and blocks codas (Supyire)
- A language forbids complex nuclei and allows codas (Russian)
- Then there will be a language that blocks complex nuclei and codas (Hawaiian, Senufo)

Can morphotactics be more than TSL?

- ✗ Closure under concatenation
- ✗ Closure under union
- ✗ Closure under relative complement
- ✓ Closure under intersection

Can morphotactics be more than TSL?

- ✗ Closure under concatenation
- ✗ Closure under union
- ✗ Closure under relative complement
- ✓ Closure under intersection



Typological gaps

Basic Logic of Argument

- All attested morphotactic patterns must be TSL.
- So if pattern A is TSL, and pattern B is TSL, but their combination A+B is not, we get a typological gap.

Some predicted gaps:

- No embedded circumfixation;
- No cases when amount of affixes A depends on the amount of affixes B;
- In general, no $a^n b^n$ pattern and its derivatives.

Typological gap I: Impossible compounding

Russian pattern – (stem-o)*-stem

Example (compounding, Russian)

● vodovoz

'water carrier'

xxxovxxx

● vodovozovoz

'carrier of water carriers'

xxxovxxxovxxx

Typological gap I: Impossible compounding

Russian pattern – (stem-o)*-stem

Example (compounding, Russian)

- | | |
|---|---|
| <ul style="list-style-type: none"> ● vodovoz
'water carrier'
xxx○xxx | <ul style="list-style-type: none"> ● vodovozovoz
'carrier of water carriers'
xxx○xxx○xxx |
|---|---|

Turkish pattern – stem-(stem⁺-o)

Example (compounding, Turkish)

- | | |
|---|--|
| <ul style="list-style-type: none"> ● bahçe kapı-sı
'garden gate'
xxxxxxxxx○ ● türk kahve-si
'Turkish coffee'
xxxxxxxxx○ | <ul style="list-style-type: none"> ● türk bahçe kapı-sı
'Turkish garden gate'
xxxxxxxxxxxxxxxxx○ ● *türk bahçe kapı-sı-sı
*xxxxxxxxxxxxxxxxx○○ |
|---|--|

Typological gap I: Impossible compounding

Russian pattern – (stem-o)*-stem

Turkish pattern – stem-(stem⁺-o)

Turkussian pattern: amount of compound markers is equal to the amount of added stems, stem-(stemⁿ-oⁿ)

Typological gap I: Impossible compounding

Russian pattern – (stem-o)*-stem

Turkish pattern – stem-(stem⁺-o)

Turkussian pattern: amount of compound markers is equal to the amount of added stems, stem-(stemⁿ-oⁿ)

- This pattern is not regular because it has infinite number of “good continuations”. (*Myhill-Nerode theorem*)
- It appears to be non-existent.

Typological gap II: Recurrent affixation

Sometimes languages allow some affixes to be iterated: **a*-stem**.

Consider example of such pattern in German:

Example (prefix 'über', German)

- morgen
'tomorrow'
xxxxxx
- übermorgen
'the day after tomorrow'
axxxxxx
- überübermorgen
'the day after the day after tomorrow'
aaxxxxxx

Typological gap II: Recurrent affixation

German pattern: **a*-stem**.

The same meaning can be expressed in another language differently, consider Ilocano (Austronesian) temporal circumfix *ka-...-an* 'next'.

Example (circumfix 'ka-an', Ilocano)

- bigát
'morning'
xxxxx

- k**abigátan**
'the next morning'
a**xxxxxb**

Typological gap II: Recurrent affixation

German pattern: **a*-stem**.

The same meaning can be expressed in another language differently, consider Ilocano (Austronesian) temporal circumfix *ka-...-an* 'next'.

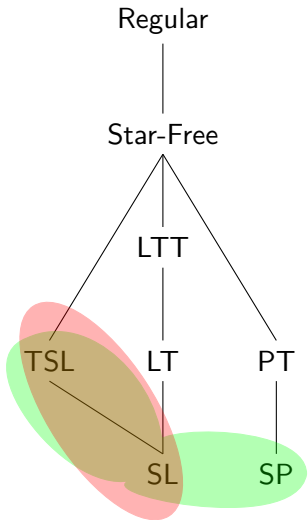
Example (circumfix 'ka-an', Ilocano)

- bigát
'morning'
xxxxx

- k**abigátan**
'the next morning'
axxxxx**b**

However, word **kakabigátanan** doesn't appear to be possible word in Ilocano: a^n -stem- b^n pattern is not regular.

Conclusion



- Morphotactics is at most Tier-Based Strictly Local
- Positive data is enough for morphological learning
- Set of typological gaps can be explained due to the subregular nature of morphology
- Same formal tools can be used for **morphology** and **phonology**

Future work

- Try to find SP patterns in morphotactics
- Look at more typologically diverse languages
- Extend to mappings from underlying to surface forms
- Work with representations of internal structure
- The elephant in the room: reduplication

Thank you!

References I



Beesley, Kenneth R. and Lauri Karttunen (2003)

Finite State Morphology.

CSLI Publications.



Chandlee, Jane (2014)

Strictly Local Phonological Processes.

PhD Thesis, University of Delaware



Chandlee, Jane, Rémi Eyraud and Jeffrey Heinz (2014)

Learning Strictly Local Subsequential Functions.

Transactions of the Association for Computational Linguistics 2, 491 – 503.



Galvez Rubino, Carl R. (1998)

Ilocano: Ilocano-English, English-Ilocano: Dictionary and Phrasebook.

Hippocrene Books Inc., U.S.



Heinz, Jeffrey (2015)

The Computational Nature of Phonological Generalizations.

Ms., University of Delaware



Heinz, Jeffrey, Chetan Rawal and Herbert G. Tanner (2011)

Tier-Based Strictly Local Constraints in Phonology.

Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics, 58 – 64.



Jardine, Adam (2015)

Computationally, Tone is Different.

Ms., University of Delaware

References II



Jardine, Adam and Jeffrey Heinz (2016)

Learning Tier-based Strictly 2-Local Languages.

Transactions of the Association for Computational Linguistics 4, 87 – 98.



Jurafsky, Daniel and James H. Martin. (2009)

Speech and language processing: an introduction to natural language processing, computational linguistics, and speech recognition.

Upper Saddle River, N.J. : Pearson Prentice Hall.



Kaplan, Ronald M. and Martin Kay (1994)

Regular Models of Phonological Rule Systems.

Computational Linguistics 20(3), 331 – 378.



Mahdi, Waruno (2012)

Distinguishing Cognate Homonyms in Indonesian.

Oceanic Linguistics 51(2), 402 – 449.



Rogers, James and Geoffrey Pullum (2007)

Aural Pattern Recognition Experiments and the Subregular Hierarchy.

Mathematics of Language 10, 1 – 16.



Sneddon, James Neil (1996)

Indonesian Comprehensive Grammar.

Routledge, London and New York.



Stump, Greg (2016)

Rule composition in an adequate theory of morphotactics.

Manuscript.