

ADAPTOR GRAMMARS FOR LEARNING NON-CONCATENATIVE MORPHOLOGY

Jan Botha Phil Blunsom

Department of Computer Science
University of Oxford

EMNLP 2013, Seattle





Morphology, the usual:

talk

talks

talked

talking

Morphology, the usual:

talk

talks

talked

talking

Morphology, the usual:

talk

talks

talked

talking

Gesundheitsreform (*health reform*) [German]

Morphology, the usual:

talk

talks

talked

talking

Gesundheitsreform (*health reform*) [German]

Morphology, the usual:

talk

talks

talked

talking

Gesundheitsreform (*health reform*) [German]

→ Permits concatenative view.

Morphology, the usual:

talk
talks
talked
talking

Gesundheitsreform (*health reform*) [German]

→ Permits concatenative view.

The “other” morphology:

kitAb (*book*)
kutub (*books*)

[Arabic]

Morphology, the usual:

talk
talks
talked
talking

Gesundheitsreform (*health reform*) [German]

→ Permits concatenative view.

The “other” morphology:

kitAb (*book*)
kutub (*books*)

[Arabic]

Morphology, the usual:

talk
talks
talked
talking

Gesundheitsreform (*health reform*) [German]

→ Permits concatenative view.

The “other” morphology:

kitAb	(<i>book</i>)	
kutub	(<i>books</i>)	
wakitAbi	(<i>and my book</i>)	[Arabic]

Morphology, the usual:

talk
talks
talked
talking

Gesundheitsreform (*health reform*) [German]

→ Permits concatenative view.

The “other” morphology:

kitAb	(book)	
kutub	(books)	
wakitAbi	(and my book)	[Arabic]

→ Needs non-concatenative view.

Morphology, the usual:

talk
talks
talked
talking

Gesundheitsreform (*health reform*) [German]

→ Permits concatenative view.

The “other” morphology:

kitAb	(book)	
kutub	(books)	
wakitAbi	(and my book)	[Arabic]

→ Needs non-concatenative view.

⇒ This talk deals with **both** views; focus on Arabic & Hebrew.

MOTIVATION

Rich morphology \Rightarrow novelty, sparse data

MOTIVATION

Rich morphology \Rightarrow novelty, sparse data

Rule-based & supervised methods

- the usual limitations – costly, language-dependent ...

MOTIVATION

Rich morphology \Rightarrow novelty, sparse data

Rule-based & supervised methods

- the usual limitations – costly, language-dependent ...

Unsupervised methods

- mostly limited to concatenative morphology, or:
- constrained search + dictionaries
(Darwish, 2002; Boudlal et al., 2009)
- statistics + heuristic constraints
(Rodrigues & Cavar, 2007; Daya et al., 2008)
- nonparametric Bayesian model (Fullwood & O'Donnel, 2013)

MOTIVATION

Rich morphology \Rightarrow novelty, sparse data

Rule-based & supervised methods

- the usual limitations – costly, language-dependent ...

Unsupervised methods

- mostly limited to concatenative morphology, or:
- constrained search + dictionaries
(Darwish, 2002; Boudlal et al., 2009)
- statistics + heuristic constraints
(Rodrigues & Cavar, 2007; Daya et al., 2008)
- nonparametric Bayesian model (Fullwood & O'Donnel, 2013)

Aim:

flexible model of joint segmentation and stem formation

THE PLAN

Proposal:

- encode (non-)concatenative morphology with a nice grammar formalism
- cast it as flexible nonparametric Bayesian model

THE PLAN

Proposal:

- encode (non-)concatenative morphology with a nice grammar formalism
- cast it as flexible nonparametric Bayesian model

Outcomes:

- improved segmentation of Hebrew & Arabic
- induced lexicons of roots

THE PLAN

Proposal:

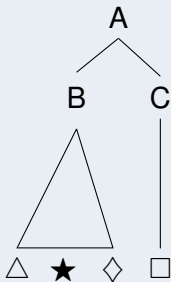
- encode (non-)concatenative morphology with a nice grammar formalism
- cast it as flexible nonparametric Bayesian model

Outcomes:

- improved segmentation of Hebrew & Arabic
- induced lexicons of roots

CONTEXT-FREE GRAMMARS

CFG

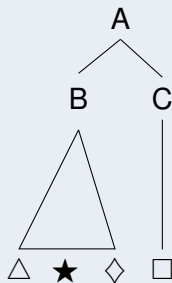


$A \rightarrow B C$

SIMPLE RANGE CONCATENATING GRAMMARS

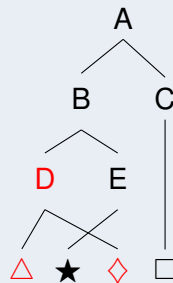
(Boullier, 2000)

CFG



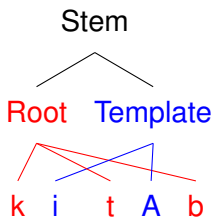
$A \rightarrow B C$

SRCG



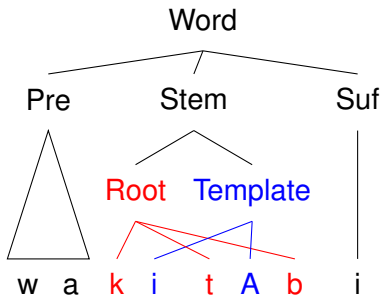
$B(\textcolor{red}{x}y\textcolor{red}{z}) \rightarrow D(\textcolor{red}{x}, \textcolor{red}{z}) E(y)$

MODELLING MORPHOLOGY WITH SRCGS



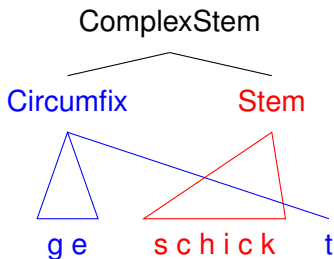
"book"

MODELLING MORPHOLOGY WITH SRCGs



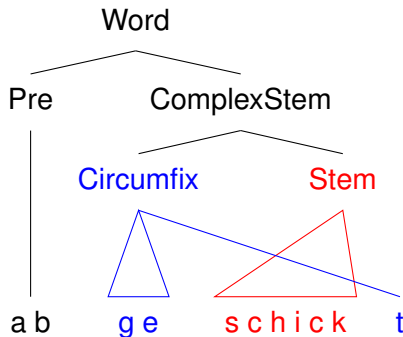
"and my book"

MODELLING MORPHOLOGY WITH SRCGs



“sent”

MODELLING MORPHOLOGY WITH SRCGS



“sent off”

PROBABILISTIC GRAMMAR

SRCGs provide the sought-after linguistic flexibility.

PROBABILISTIC GRAMMAR

SRCGs provide the sought-after linguistic flexibility.

Add **Adaptor Grammars** (Johnson et al., 2007)

- nonparametric Bayesian extension of PCFGs

PROBABILISTIC GRAMMAR

SRCGs provide the sought-after linguistic flexibility.

Add **Adaptor Grammars** (Johnson et al., 2007)

- nonparametric Bayesian extension of PCFGs
- whole subtrees reused across different inputs
e.g. stem occurring in multiple word forms

PROBABILISTIC GRAMMAR

SRCGs provide the sought-after linguistic flexibility.

Add **Adaptor Grammars** (Johnson et al., 2007)

- nonparametric Bayesian extension of PCFGs
- whole subtrees reused across different inputs
e.g. stem occurring in multiple word forms
- Expanding a node (generatively) involves:
 - a) generate its direct children via base grammar*OR*
 - b) pick a completed subtree from cache.

PROBABILISTIC GRAMMAR

SRCGs provide the sought-after linguistic flexibility.

Add **Adaptor Grammars** (Johnson et al., 2007)

- nonparametric Bayesian extension of PCFGs
- whole subtrees reused across different inputs
e.g. stem occurring in multiple word forms
- Expanding a node (generatively) involves:
 - a) generate its direct children via base grammar
 - OR*
 - b) pick a completed subtree from cache.
- choice governed by Pitman-Yor Process

PROBABILISTIC GRAMMAR

SRCGs provide the sought-after linguistic flexibility.

Add **Adaptor Grammars** (Johnson et al., 2007)

- nonparametric Bayesian extension of PCFGs
- whole subtrees reused across different inputs
e.g. stem occurring in multiple word forms
- Expanding a node (generatively) involves:
 - a) generate its direct children via base grammar*OR*
 - b) pick a completed subtree from cache.
- choice governed by Pitman-Yor Process

We apply this to SRCGs.

GRAMMAR FRAGMENT

$\text{Word}(xyz) \rightarrow \underline{\text{Pre}}(x) \underline{\text{Stem}}(y) \underline{\text{Suf}}(z)$

$\underline{\text{Stem}}(abcde) \rightarrow \underline{\text{Root}}(a, c, e) \underline{\text{Template}}(b, d)$

$\underline{\text{Root}}(f, g, h) \rightarrow \text{Char}(f) \text{Char}(g) \text{Char}(h)$

$\underline{\text{Template}}(i, j) \rightarrow \text{Char}(i) \text{Char}(j)$

GRAMMAR FRAGMENT

$\text{Word}(xyz) \rightarrow \underline{\text{Pre}}(x) \underline{\text{Stem}}(y) \underline{\text{Suf}}(z)$

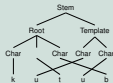
$\underline{\text{Stem}}(abcde) \rightarrow \underline{\text{Root}}(a, c, e) \underline{\text{Template}}(b, d)$

$\underline{\text{Root}}(f, g, h) \rightarrow \text{Char}(f) \text{Char}(g) \text{Char}(h)$

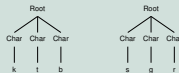
$\underline{\text{Template}}(i, j) \rightarrow \text{Char}(i) \text{Char}(j)$

CACHE

STEM



ROOT



TEMPLATE



GRAMMAR FRAGMENT

$\text{Word}(xyz) \rightarrow \underline{\text{Pre}}(x) \underline{\text{Stem}}(y) \underline{\text{Suf}}(z)$

$\underline{\text{Stem}}(abcde) \rightarrow \underline{\text{Root}}(a, c, e) \underline{\text{Template}}(b, d)$

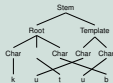
$\underline{\text{Root}}(f, g, h) \rightarrow \text{Char}(f) \text{Char}(g) \text{Char}(h)$

$\underline{\text{Template}}(i, j) \rightarrow \text{Char}(i) \text{Char}(j)$

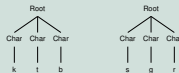
Word

CACHE

STEM



ROOT



TEMPLATE



GRAMMAR FRAGMENT

Word(xyz) \rightarrow Pre(x) Stem(y) Suf(z)

Stem($abcde$) \rightarrow Root(a, c, e) Template(b, d)

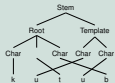
Root(f, g, h) \rightarrow Char(f) Char(g) Char(h)

Template(i, j) \rightarrow Char(i) Char(j)

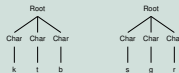
Word

CACHE

STEM



ROOT



TEMPLATE



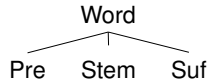
GRAMMAR FRAGMENT

$\text{Word}(xyz) \rightarrow \text{Pre}(x) \text{Stem}(y) \text{Suf}(z)$

$\text{Stem}(abcde) \rightarrow \text{Root}(a, c, e) \text{Template}(b, d)$

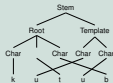
$\text{Root}(f, g, h) \rightarrow \text{Char}(f) \text{Char}(g) \text{Char}(h)$

$\text{Template}(i, j) \rightarrow \text{Char}(i) \text{Char}(j)$

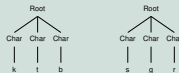


CACHE

STEM



ROOT



TEMPLATE



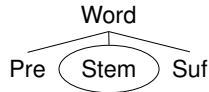
GRAMMAR FRAGMENT

$\text{Word}(xyz) \rightarrow \text{Pre}(x) \text{Stem}(y) \text{Suf}(z)$

$\text{Stem}(abcde) \rightarrow \text{Root}(a, c, e) \text{Template}(b, d)$

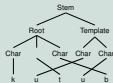
$\text{Root}(f, g, h) \rightarrow \text{Char}(f) \text{Char}(g) \text{Char}(h)$

$\text{Template}(i, j) \rightarrow \text{Char}(i) \text{Char}(j)$

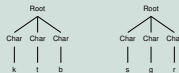


CACHE

STEM



ROOT



TEMPLATE



GRAMMAR FRAGMENT

$\text{Word}(xyz) \rightarrow \text{Pre}(x) \text{Stem}(y) \text{Suf}(z)$

$\text{Stem}(abcde) \rightarrow \text{Root}(a, c, e) \text{Template}(b, d)$

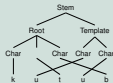
$\text{Root}(f, g, h) \rightarrow \text{Char}(f) \text{Char}(g) \text{Char}(h)$

$\text{Template}(i, j) \rightarrow \text{Char}(i) \text{Char}(j)$

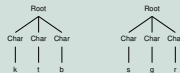
Stem

CACHE

STEM



ROOT



TEMPLATE



GRAMMAR FRAGMENT

$\text{Word}(xyz) \rightarrow \text{Pre}(x) \text{Stem}(y) \text{Suf}(z)$

$\text{Stem}(abcde) \rightarrow \text{Root}(a, c, e) \text{Template}(b, d)$

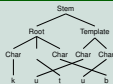
$\text{Root}(f, g, h) \rightarrow \text{Char}(f) \text{Char}(g) \text{Char}(h)$

$\text{Template}(i, j) \rightarrow \text{Char}(i) \text{Char}(j)$

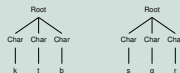
Stem

CACHE

STEM



ROOT



TEMPLATE



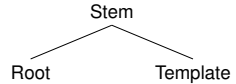
GRAMMAR FRAGMENT

$\text{Word}(xyz) \rightarrow \text{Pre}(x) \text{Stem}(y) \text{Suf}(z)$

$\text{Stem}(abcde) \rightarrow \text{Root}(a, c, e) \text{Template}(b, d)$

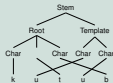
$\text{Root}(f, g, h) \rightarrow \text{Char}(f) \text{Char}(g) \text{Char}(h)$

$\text{Template}(i, j) \rightarrow \text{Char}(i) \text{Char}(j)$

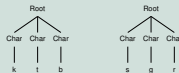


CACHE

STEM



ROOT



TEMPLATE



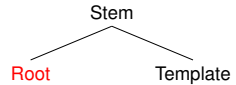
GRAMMAR FRAGMENT

$\text{Word}(xyz) \rightarrow \text{Pre}(x) \text{Stem}(y) \text{Suf}(z)$

$\text{Stem}(abcde) \rightarrow \text{Root}(a, c, e) \text{Template}(b, d)$

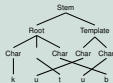
$\text{Root}(f, g, h) \rightarrow \text{Char}(f) \text{Char}(g) \text{Char}(h)$

$\text{Template}(i, j) \rightarrow \text{Char}(i) \text{Char}(j)$

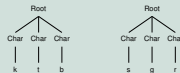


CACHE

STEM



ROOT



TEMPLATE



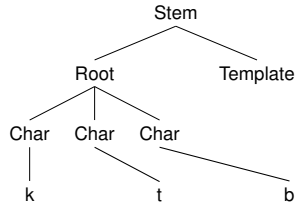
GRAMMAR FRAGMENT

$\text{Word}(xyz) \rightarrow \text{Pre}(x) \text{Stem}(y) \text{Suf}(z)$

$\text{Stem}(abcde) \rightarrow \text{Root}(a, c, e) \text{Template}(b, d)$

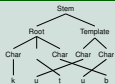
$\text{Root}(f, g, h) \rightarrow \text{Char}(f) \text{Char}(g) \text{Char}(h)$

$\text{Template}(i, j) \rightarrow \text{Char}(i) \text{Char}(j)$

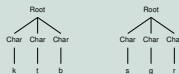


CACHE

STEM



ROOT



TEMPLATE



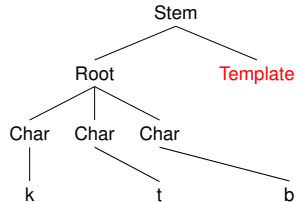
GRAMMAR FRAGMENT

$\text{Word}(xyz) \rightarrow \text{Pre}(x) \text{Stem}(y) \text{Suf}(z)$

$\text{Stem}(abcde) \rightarrow \text{Root}(a, c, e) \text{Template}(b, d)$

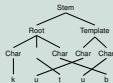
$\text{Root}(f, g, h) \rightarrow \text{Char}(f) \text{Char}(g) \text{Char}(h)$

$\text{Template}(i, j) \rightarrow \text{Char}(i) \text{Char}(j)$

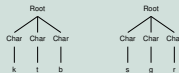


CACHE

STEM



ROOT



TEMPLATE



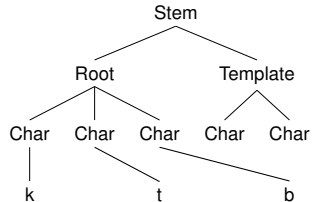
GRAMMAR FRAGMENT

$\text{Word}(xyz) \rightarrow \text{Pre}(x) \text{Stem}(y) \text{Suf}(z)$

$\text{Stem}(abcde) \rightarrow \text{Root}(a, c, e) \text{Template}(b, d)$

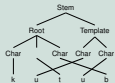
$\text{Root}(f, g, h) \rightarrow \text{Char}(f) \text{Char}(g) \text{Char}(h)$

$\text{Template}(i, j) \rightarrow \text{Char}(i) \text{Char}(j)$

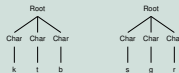


CACHE

STEM



ROOT



TEMPLATE



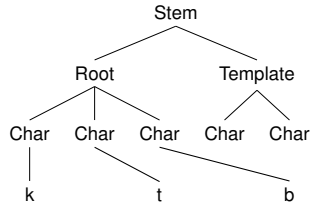
GRAMMAR FRAGMENT

$\text{Word}(xyz) \rightarrow \text{Pre}(x) \text{Stem}(y) \text{Suf}(z)$

$\text{Stem}(abcde) \rightarrow \text{Root}(a, c, e) \text{Template}(b, d)$

$\text{Root}(f, g, h) \rightarrow \text{Char}(f) \text{Char}(g) \text{Char}(h)$

$\text{Template}(i, j) \rightarrow \text{Char}(i) \text{Char}(j)$

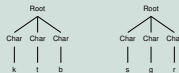


CACHE

STEM



ROOT



TEMPLATE



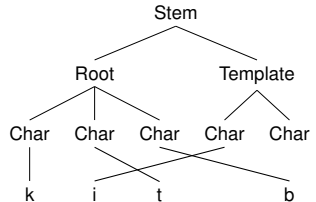
GRAMMAR FRAGMENT

$\text{Word}(xyz) \rightarrow \text{Pre}(x) \text{Stem}(y) \text{Suf}(z)$

$\text{Stem}(abcde) \rightarrow \text{Root}(a, c, e) \text{Template}(b, d)$

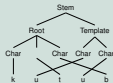
$\text{Root}(f, g, h) \rightarrow \text{Char}(f) \text{Char}(g) \text{Char}(h)$

$\text{Template}(i, j) \rightarrow \text{Char}(i) \text{Char}(j)$

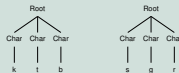


CACHE

STEM



ROOT



TEMPLATE



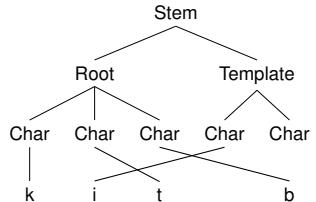
GRAMMAR FRAGMENT

$\text{Word}(xyz) \rightarrow \text{Pre}(x) \text{Stem}(y) \text{Suf}(z)$

$\text{Stem}(abcde) \rightarrow \text{Root}(a, c, e) \text{Template}(b, d)$

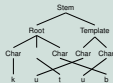
$\text{Root}(f, g, h) \rightarrow \text{Char}(f) \text{Char}(g) \text{Char}(h)$

$\text{Template}(i, j) \rightarrow \text{Char}(i) \text{Char}(j)$

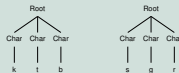


CACHE

STEM



ROOT



TEMPLATE



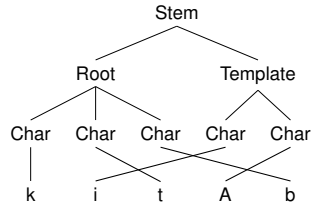
GRAMMAR FRAGMENT

$\text{Word}(xyz) \rightarrow \text{Pre}(x) \text{Stem}(y) \text{Suf}(z)$

$\text{Stem}(abcde) \rightarrow \text{Root}(a, c, e) \text{Template}(b, d)$

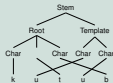
$\text{Root}(f, g, h) \rightarrow \text{Char}(f) \text{Char}(g) \text{Char}(h)$

$\text{Template}(i, j) \rightarrow \text{Char}(i) \text{Char}(j)$

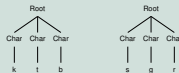


CACHE

STEM



ROOT



TEMPLATE



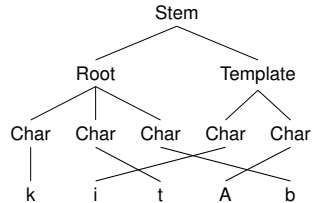
GRAMMAR FRAGMENT

$\text{Word}(xyz) \rightarrow \text{Pre}(x) \text{Stem}(y) \text{Suf}(z)$

$\text{Stem}(abcde) \rightarrow \text{Root}(a, c, e) \text{Template}(b, d)$

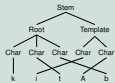
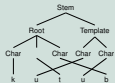
$\text{Root}(f, g, h) \rightarrow \text{Char}(f) \text{Char}(g) \text{Char}(h)$

$\text{Template}(i, j) \rightarrow \text{Char}(i) \text{Char}(j)$

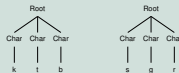


CACHE

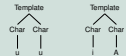
STEM



ROOT



TEMPLATE



DATASETS

Hebrew

- 5k word types from CHILDES database
- vocalised

Standard Arabic $\times 2$

- synthesised 50k word types from BAMA dictionaries
- both orthographic variants: vocalised and unvocalised

Quranic Arabic

- 18k word types from annotated Quran (Dukes et al. 2010)
- extensive diacritics

Included all parts of speech; not filtered for verbs/nouns.

METHOD

1. Run MCMC sampler over unannotated data.
2. Collect 100 posterior samples.
3. Evaluate MAP parses against references.

METHOD

1. Run MCMC sampler over unannotated data.
2. Collect 100 posterior samples.
3. Evaluate MAP parses against references.

TASK 1: SEGMENTATION

ions \rightarrow ion \cdot s

fabrications \rightarrow fabricat \cdot ion \cdot s

METHOD

1. Run MCMC sampler over unannotated data.
2. Collect 100 posterior samples.
3. Evaluate MAP parses against references.

TASK 1: SEGMENTATION

ions \rightarrow ion \cdot s

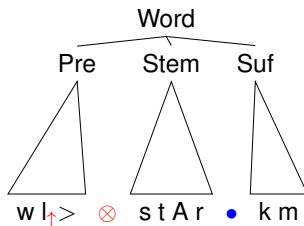
fabrications \rightarrow fabricat \cdot ion \cdot s

TASK 2: LEXICON INDUCTION

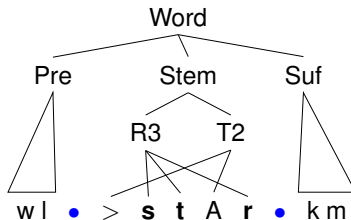
<u>Prefixes</u>	<u>Stems</u>	<u>Suffixes</u>	<u>Roots</u>
...	ion fabricat	ion s	...

EXAMPLE ANALYSES – ARABIC (UNVOCALISED)

CFG

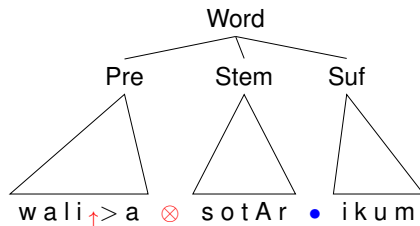


SRCGG

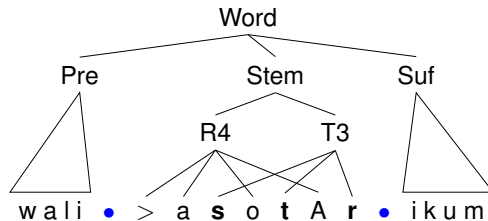


EXAMPLE ANALYSES – ARABIC (VOCALISED)

CFG

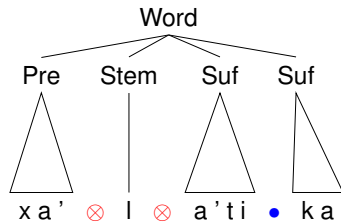


SRCGG

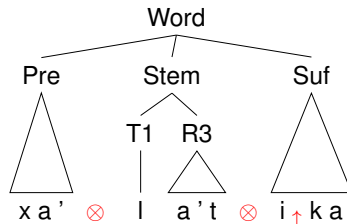


EXAMPLE ANALYSES – QURAN

CFG



SRCG



EXAMPLE ANALYSES

Top three Hebrew roots according to model

Root

1. **spr** (*tell*)
2. **lbš** (*wear*)
3. **ptx** (*open*)

EXAMPLE ANALYSES

Top three Hebrew roots according to model

Root	Correct Instances
1. spr (<i>tell</i>)	sipar•ti ye•sapr•u
2. lbš (<i>wear</i>)	li•lboš ti•lbeš•i
3. ptx (<i>open</i>)	

EXAMPLE ANALYSES

Top three Hebrew roots according to model

	Root	Correct Instances	Mistaken Instance
1.	spr (<i>tell</i>)	sipar•ti ye•sapr•u	hi⊗stap⊗ar↑t
2.	lbš (<i>wear</i>)	li•lboš ti•lbeš•i	le↑ha⊗lbiš
3.	ptx (<i>open</i>)		

SEGMENTATION RESULTS

F1-scores over word-internal morpheme boundaries

(higher better)

Morfessor
CFG

SRCG1
SRCG2
SRCG3
SRCG4

SEGMENTATION RESULTS

F1-scores over word-internal morpheme boundaries

(higher better)

Arabic (unvoc)

Morfessor	55.6
CFG	47.4
SRCG1	60.4
SRCG2	60.5
SRCG3	64.5
SRCG4	74.5

SEGMENTATION RESULTS

F1-scores over word-internal morpheme boundaries

(higher better)

	Arabic (unvoc)	Arabic
Morfessor	55.6	40.0
CFG	47.4	64.2
SRCG1	60.4	71.9
SRCG2	60.5	72.2
SRCG3	64.5	71.6
SRCG4	74.5	73.7

SEGMENTATION RESULTS

F1-scores over word-internal morpheme boundaries

(higher better)

	Arabic (unvoc)	Arabic	Hebrew
Morfessor	55.6	40.0	24.2
CFG	47.4	64.2	60.1
SRCG1	60.4	71.9	77.3
SRCG2	60.5	72.2	77.4
SRCG3	64.5	71.6	77.1
SRCG4	74.5	73.7	78.1

SEGMENTATION RESULTS

F1-scores over word-internal morpheme boundaries

(higher better)

	Arabic (unvoc)	Arabic	Hebrew	Quran
Morfessor	55.6	40.0	24.2	44.3
CFG	47.4	64.2	60.1	19.6
SRCG1	60.4	71.9	77.3	22.5
SRCG2	60.5	72.2	77.4	25.7
SRCG3	64.5	71.6	77.1	24.8
SRCG4	74.5	73.7	78.1	-

SEGMENTATION RESULTS

F1-scores over word-internal morpheme boundaries

(higher better)

	Arabic (unvoc)	Arabic	Hebrew	Quran
Morfessor	55.6	40.0	24.2	44.3
CFG	47.4	64.2	60.1	19.6
SRCG1	60.4	71.9	77.3	22.5
SRCG2	60.5	72.2	77.4	25.7
SRCG3	64.5	71.6	77.1	24.8
SRCG4	74.5	73.7	78.1	-

⇒ Modelling discontinuous substructure improves segmentation

MORPHEME LEXICON INDUCTION RESULTS

F-score in set-based evaluation against gold lexicons

(higher better)

	Prefixes	Stems	Suffixes	Triliteral Roots
Arabic (unvocalised)				P / R / F
CFG	33	44	40	—
Best SRCG	53	58	52	51 / 80 / 62

CONCLUSIONS

- Flexible modelling framework that handles concatenative *and* non-concatenative morphology
- Accounting for root-templatic stem formation improved segmentation in Hebrew and Arabic
- SRCG-variant of Adaptor Grammars potentially applicable to other SRCGs (parsing, translation)
- Further avenues for inquiry:
 - Induced roots as features in downstream tasks?
 - Success at handling other non-concatenative phenomena?
 - What does discontinuous model do on English, etc.
 - RCG with deletion $\xrightarrow{?}$ irregular patterns, weak roots?

Thank you.