

modeling  
morphological uncertainty  
for machine translation

**Waleed Ammar**, Victor Chahuneau, Chris Dyer, Noah A. Smith (CMU)

Jason Baldridge, Kyle Jerro, Jason Mielens (UT)

**three observations**

rich morphology is **typical**

morphology is a challenge  
(esp. into rich)

linguistic knowledge can help

# translation puzzle

## English

big man

big women

big boat

big word

big thing

big things

big cow

## Kinyarwanda

umugabo munini

? ?

ubwato bunini

ijambo rinini

ikintu kinini

ibintu binini

inka nini

## English

small man

small women

small boat

small word

small thing

small things

small cow

## Kinyarwanda

umugabo mutoya

abagore batoya

ubwato butoya

ijambo ritoya

? ?

ibintu bitoya

inka toya

# linguists

English	Kinyarwanda	English	Kinyarwanda
big man	umugabo <b>munini</b>	small man	umugabo <b>mutoya</b>
big women	? ?	small women	abagore <b>batoya</b>
big boat	ubwato <b>bunini</b>	small boat	ubwato <b>butoya</b>
big word	ijambo <b>rinini</b>	small word	ijambo <b>rito</b> ya
big thing	ikintu <b>kinini</b>	small thing	? ?
big things	ibintu <b>binini</b>	small things	ibintu <b>bitoya</b>
big cow	inka nini	small cow	inka toya
...	...	...	...

# linguists

## English

big man

big women

big boat

big word

big thing

big things

big cow

...

## Kinyarwanda

umugabo **m**unini

abagore **class2**(nini)

ubwato **b**unini

ijambo **r**inini

ikintu **k**inini

ibintu **b**inini

inka nini

...

## English

small man

small women

small boat

small word

small thing

small things

small cow

...

## Kinyarwanda

umugabo **m**utoya

abagore **b**atoya

ubwato **b**utoya

ijambo **r**itoya

ikintu **class7**(toya)

ibintu **b**itoya

inka toya

...

class 1

class 2

class 14

class 5

class 7

class 8

class 9



# linguists

## English

big man

big women

big boat

big word

big thing

big things

big cow

...

## Kinyarwanda

umugabo munini

abagore banini ✓

ubwato bunini

ijambo rinini

ikintu kinini

ibintu binini

inka nini

...

## English

small man

small women

small boat

small word

small thing

small things

small cow

...

## Kinyarwanda

umugabo mutoya

abagore batoya

ubwato butoya

ijambo ritoya

ikintu kitoya ✓

ibintu bitoya

inka toya

...

class 1

class 2

class 14

class 5

class 7

class 8

class 9

# mainstream translation models

<b>English</b>	<b>Kinyarwanda</b>	<b>English</b>	<b>Kinyarwanda</b>
big man	umugabo munini	small man	umugabo mutoya
big women	? ?	small women	abagore batoya
big boat	ubwato bunini	small boat	ubwato butoya
big word	ijambo rinini	small word	ijambo ritoya
big thing	ikintu kinini	small thing	? ?
big things	ibintu binini	small things	ibintu bitoya
big cow	inka nini	small cow	inka toya
...	...	...	...

# mainstream word alignment methods

## English

big man

big women

big boat

big word

big thing

big things

big cow

...

## Kinyarwanda

umugabo munini

? ?

ubwato bunini

ijambo rinini

ikintu kinini

ibintu binini

inka nini

...

## English

small man

small women

small boat

small word

small thing

small things

small cow

...

## Kinyarwanda

umugabo mutoya

abagore batoya

ubwato butoya

ijambo ritoya

? ?

ibintu bitoya

inka toya

...

# mainstream word alignment methods

English	Kinyarwanda	English	Kinyarwanda
big man	umugabo munini	small man	umugabo mutoya
big women	? ?	small women	abagore batoya

Problem #1: word **alignment** errors for rare inflections.

big thing	ikintu kinini	small thing	? ?
big things	ibintu binini	small things	ibintu bitoya
big cow	inka nini	small cow	inka toya
...	...	...	...

# decoding

English	Kinyarwanda	English	Kinyarwanda
big man	umugabo munini	small man	umugabo mutoya
big women	tr(women) tr(big)	small women	abagore batoya
big boat	ubwato bunini	small boat	ubwato butoya
big word	ijambo rinini	small word	ijambo ritoya
big thing	ikintu kinini	small thing	tr(thing) tr(small)
big things	ibintu binini	small things	ibintu bitoya
big cow	inka nini	small cow	inka toya
...	...	...	...

# decoding

English	Kinyarwanda	English	Kinyarwanda
big man	umugabo munini	small man	umugabo mutoya
big women	abagore munini	small women	abagore batoya
big boat	ubwato <b>banini</b>	small boat	ubwato butoya
big word	ijambo rinini	small word	ijambo ritoya
big thing	ikintu kinini	small thing	ikintu mutoya
big things	ibintu binini	small things	ibintu <b>kitoya</b>
big cow	inka nini	small cow	inka toya
...	...	...	...

# decoding

English	Kinyarwanda	English	Kinyarwanda
big man	umugabo <b>banini</b>	small man	umugabo <b>mutoya</b>
big women	abagore munini	small women	abagore batoya

Problem #2: inability to **generate** unseen inflections.

big thing	ikintu kinini	small thing	ikintu mutoya
big things	ibintu binini	small things	ibintu <b>kitoya</b>
big cow	inka nini	small cow	inka toya
...	...	...	...

# decoding

English

Kinya

real problems

Kinyarwanda

Problem #1: word **alignment** errors for rare inflections.

Problem #2: inability to **generate** unseen inflections.

Problem #3: n-gram LM fails to **predict** rare inflections.

Goal: use linguistic knowledge to overcome #1, #2 and #3.

...

...

...

...



# linguistic resources

- morphological analyzer(s) for the rich language

# road map

Problem #1: word alignment errors for rare inflections.

Problem #2: inability to generate unseen inflections.

Problem #3: n-gram LM fails to predict rare inflections.



Problem #1: word alignment errors for rare inflections.

# unsupervised word alignment with morphological features

need to compute feature expectations  
too slow, even with small amounts of data

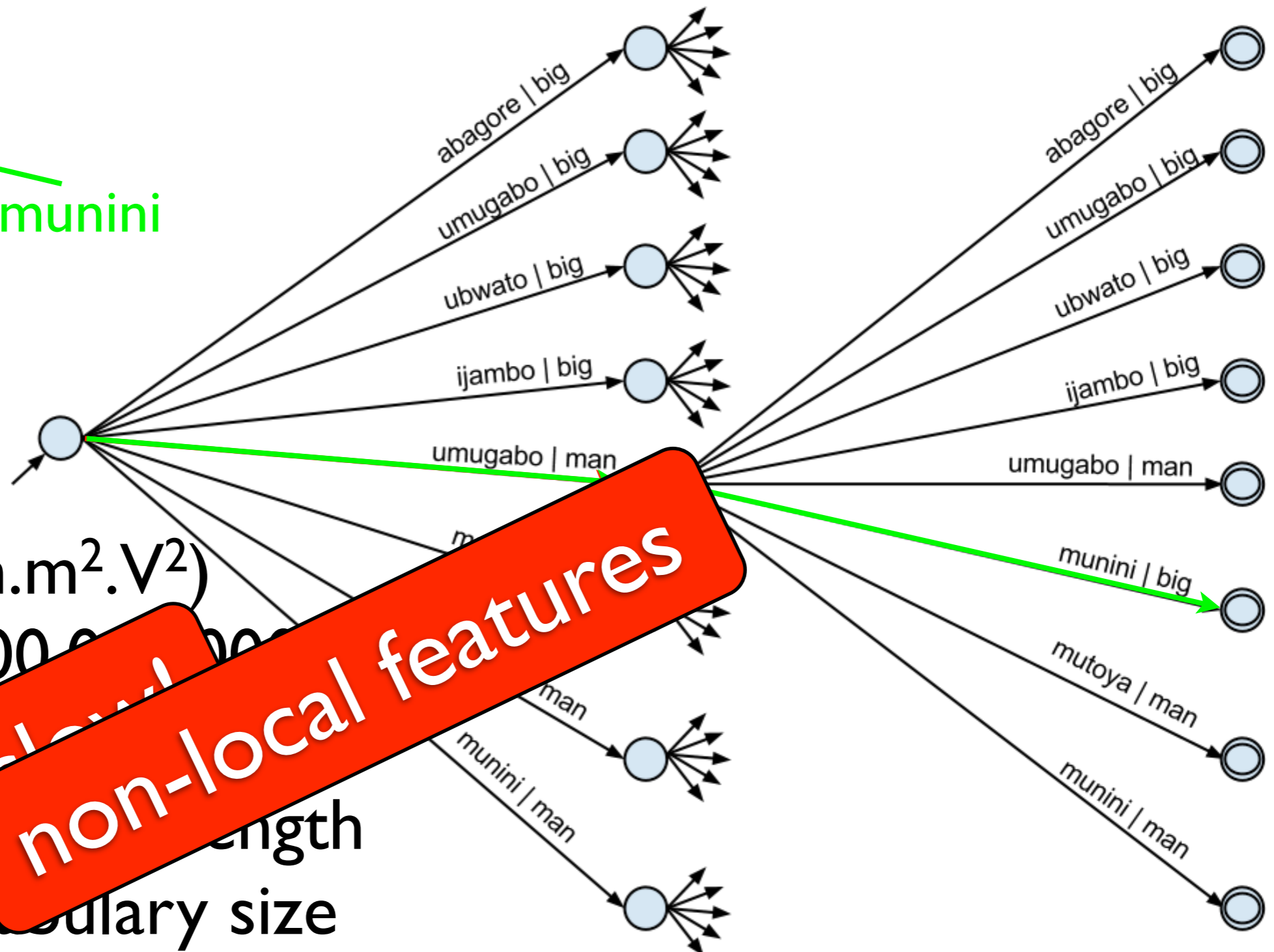
- 2011:
  - finite state transducers
  - discriminative lexicon
- 2012:
  - approximate inference
  - non-local features
  - OpenFST

Problem #1: word alignment errors for rare inflections.

## alignment FST

Eng: big man

Kin: umugabo munini



# arcs =  $O(n.m^2.V^2)$

$\sim 1,000 \times 1,000 \times 1,000$

n: source

m: target length

V: target vocabulary size

# Problem #1: word alignment errors for rare inflections.

## alignment FST

Eng: **big man**

proposal distribution

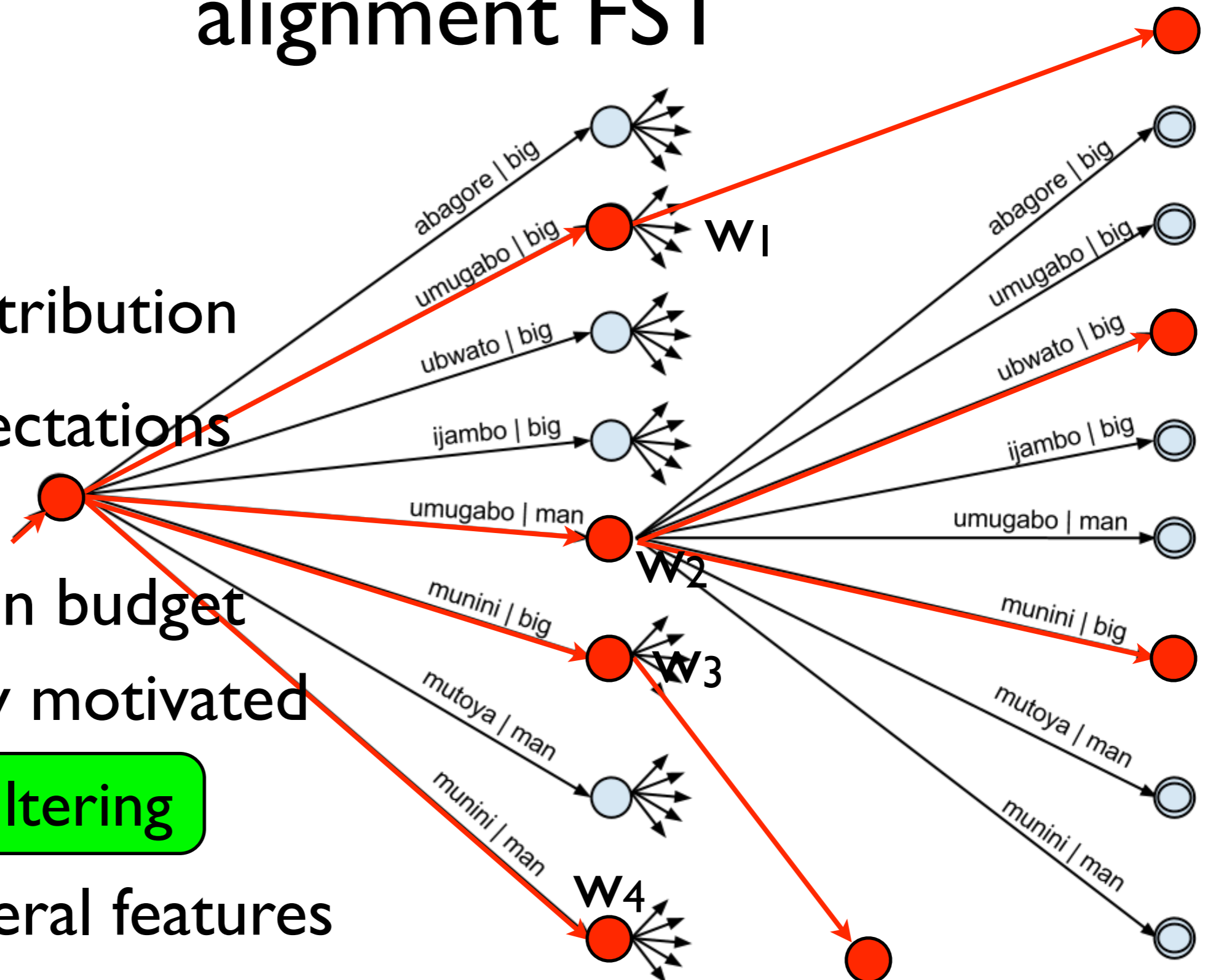
feature expectations

computation budget

theoretically motivated

**particle filtering**

enables general features



# road map

Problem #1: word alignment errors for rare inflections.

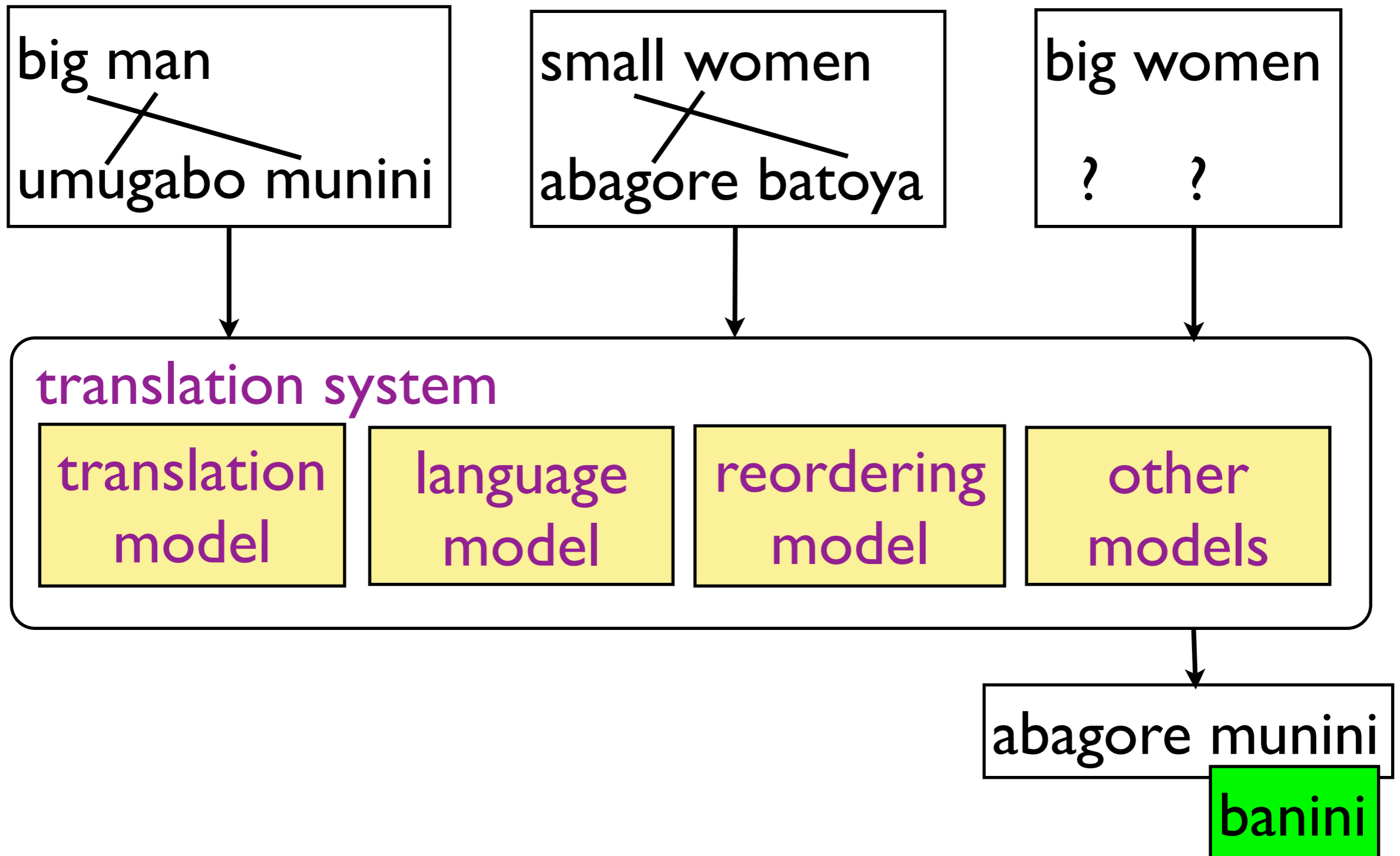
→ use morphological features + approximate inference

Problem #2: inability to generate unseen inflections.

Problem #3: n-gram LM fails to predict rare inflections.

Problem #2: inability to generate unseen inflections.

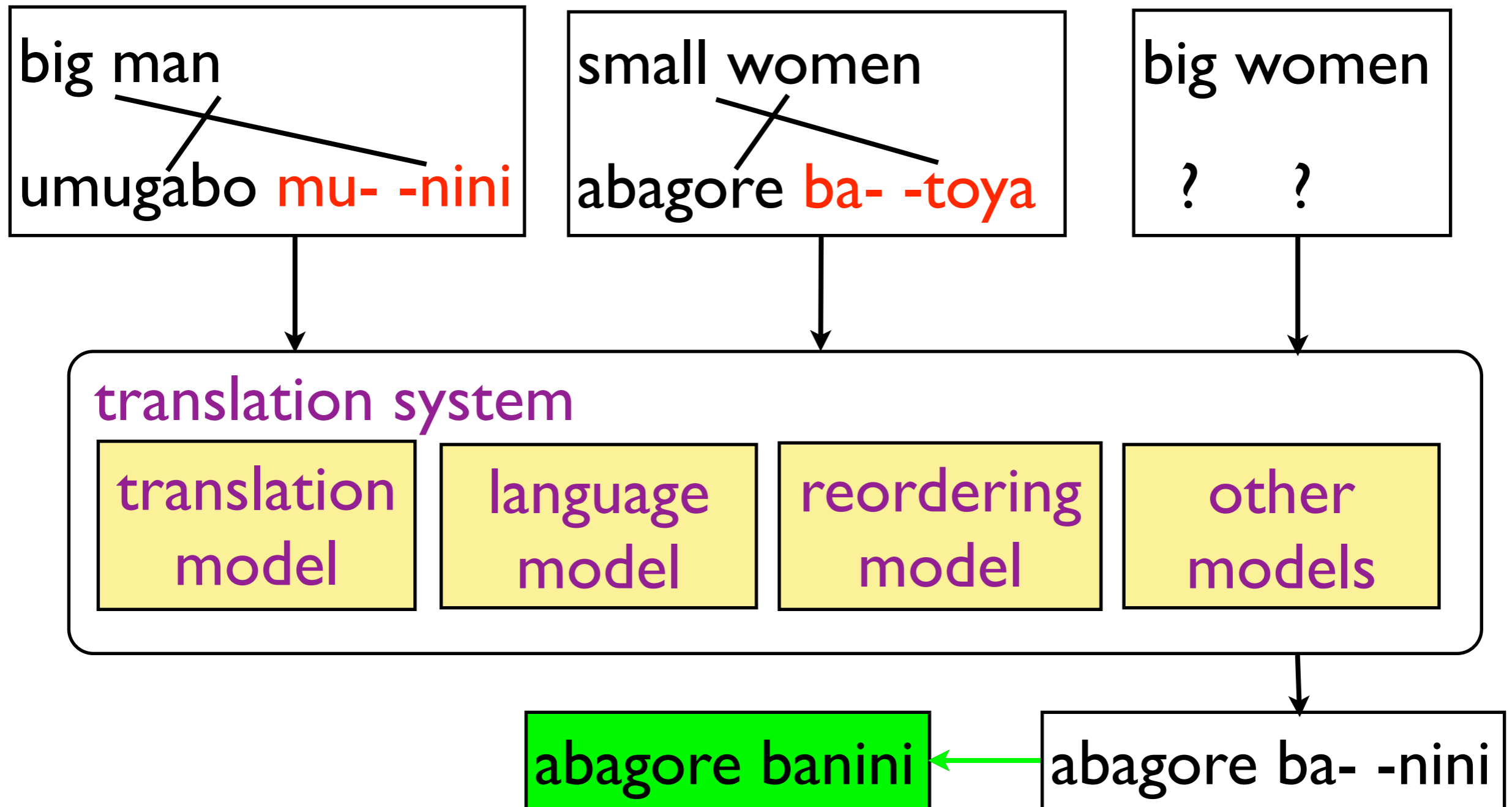
## default lexical representation





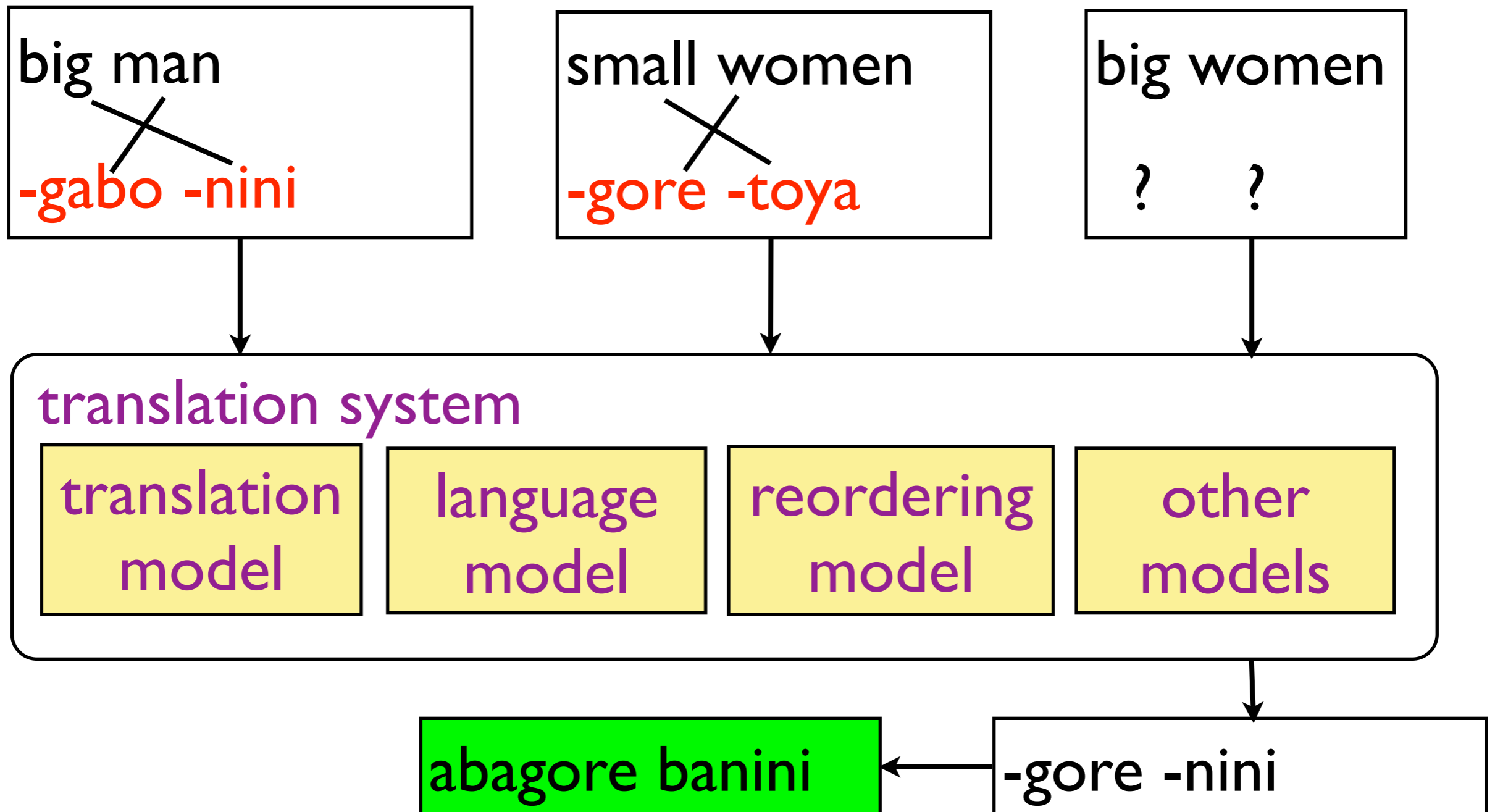
Problem #2: inability to generate unseen inflections.

## alternative lexical representation (a)



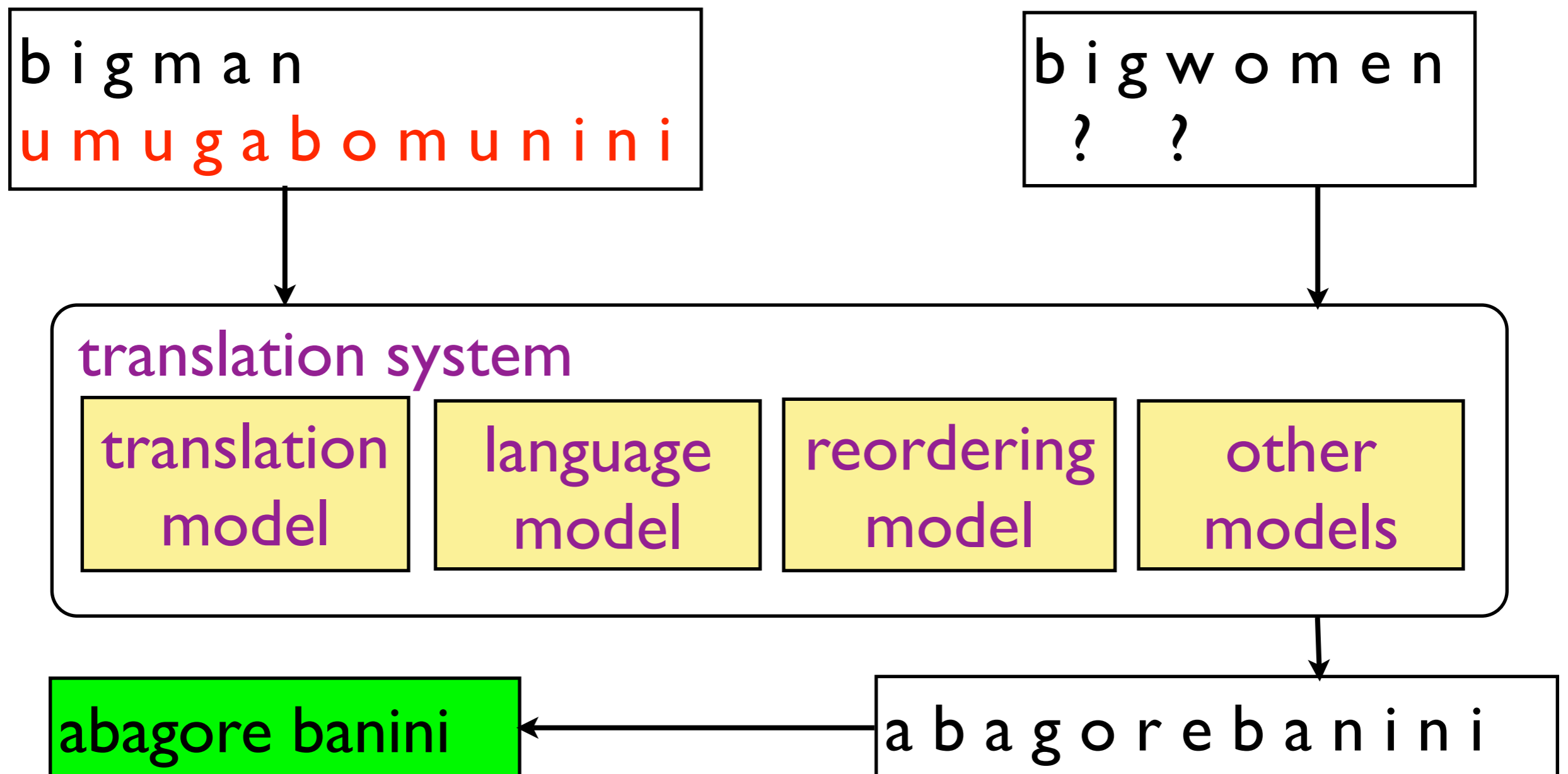
Problem #2: inability to generate unseen inflections.

## alternative lexical representation (b)



Problem #2: inability to generate unseen inflections.

## alternative lexical representation (c)



Problem #2: inability to generate unseen inflections.

## alternative lexical representations

- some proposed representations:
  - affix transformations Habash and Sadat 2006, Al-Haj and Lavie 2010
  - artificially inflected English words Oflazer and El-Kahlout 2007
  - each character is a separate token Xu et al. 2004
- problems:
  - strong assumptions
  - subset of useful representations

Problem #2: inability to generate unseen inflections.

## representation variables

English

big man

translation-friendly  
lexical representation

? ?

**possible values**

umugabo mu nini

umu gabo munini

umu gabo mu nini

u m u g a b o m u n i n i

umugabo munini

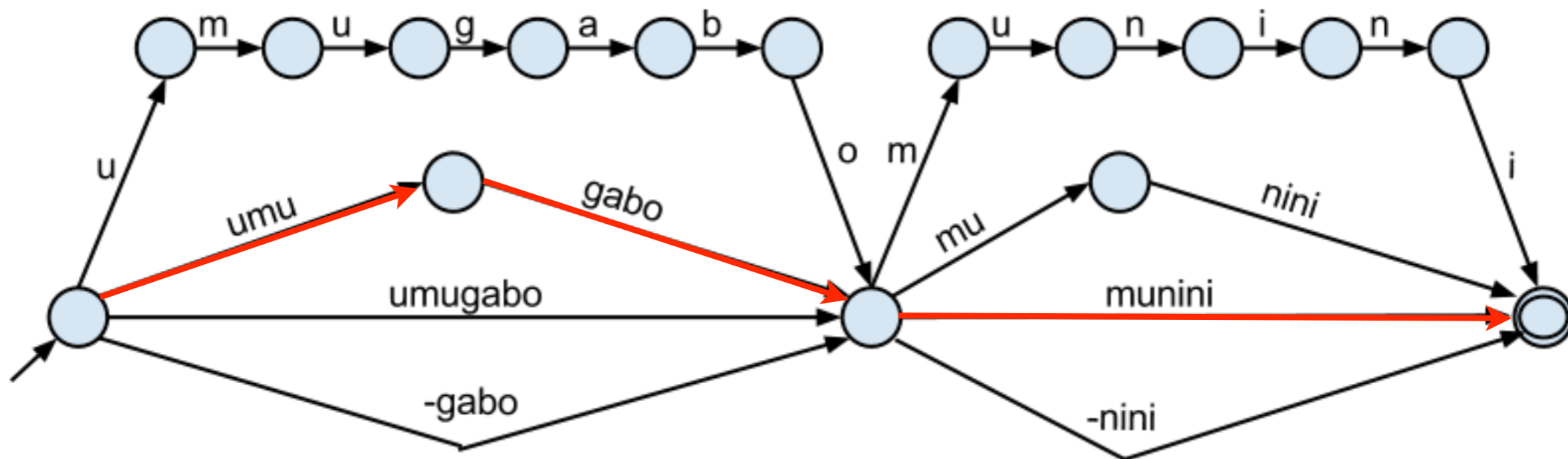
gabo nini

Kinyarwanda

umugabo munini

# Problem #2: inability to generate unseen inflections.

## representation variables



- $p(a,t|s)$
- $p(a,\mathbf{r},t | s)$  **representation**
- approximate inference

# road map

Problem #1: word alignment errors for rare inflections.

→ use morphological features + approximate inference

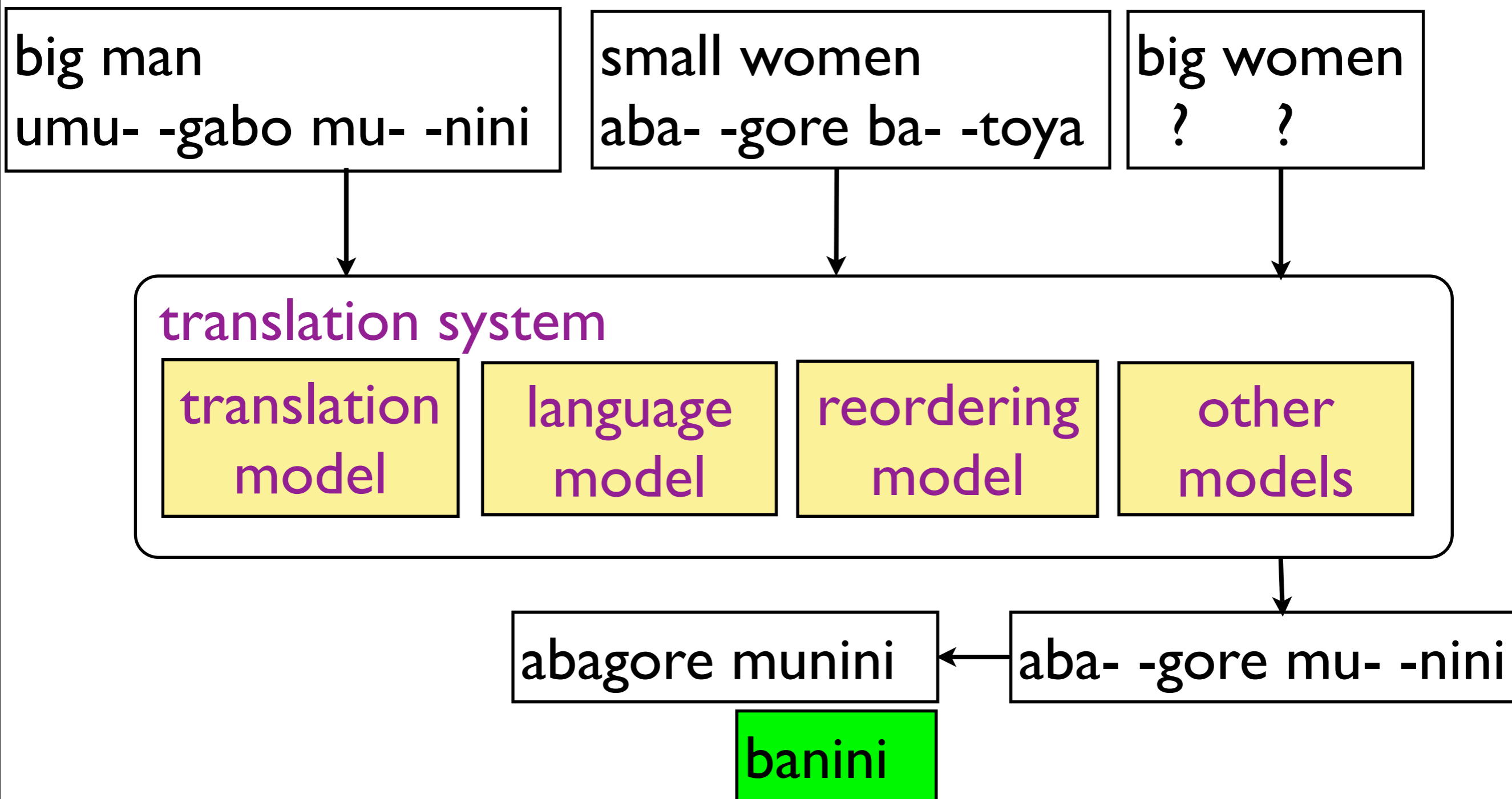
Problem #2: inability to generate unseen inflections.

→ learn translation-friendly lexical representations

Problem #3: n-gram LM fails to predict rare inflections.

Problem #3: n-gram LM fails to predict correct inflections.

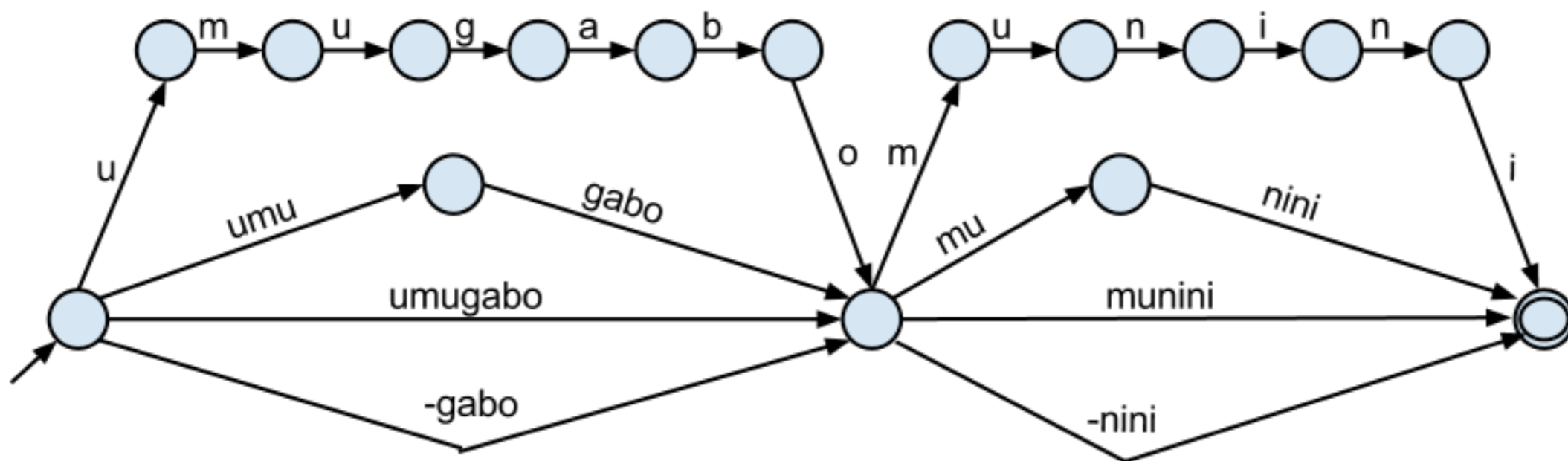
## n-gram LM favors popular sequences





Problem #3: n-gram LM fails to predict correct inflections.

## morphologically aware LM



# road map

Problem #1: word alignment errors for rare inflections.

→ use morphological features + approximate inference

Problem #2: inability to generate unseen inflections.

→ learn translation-friendly lexical representations

Problem #3: n-gram LM fails to predict rare inflections.

→ morphologically-aware language model

models that use linguistic knowledge

inside translation model

**questions?**