# Weighted Tree Automata and Transducers for Syntactic Natural Language Processing

Jonathan May
Thesis Defense
April 20, 2010

1

# How do we view natural language?
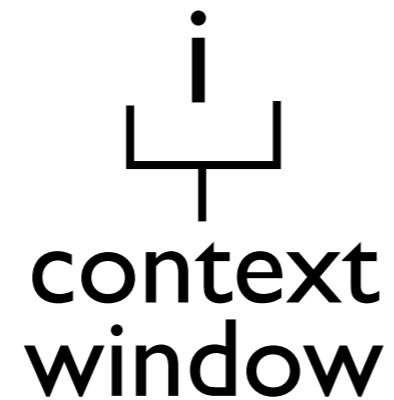
## As a string?

Monday, April 19, 2010

# How do we view natural language?

## As a string?

i

# How do we view natural language?

## As a string?

i

context window

4

# How do we view natural language?

## As a string?

i gave
context
window

5

# How do we view natural language?

## As a string?

i  gave  my

context
window

6

# How do we view natural language?
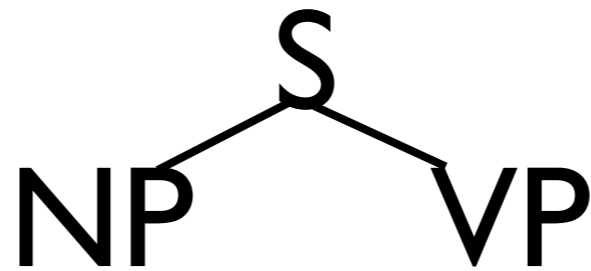
## As a string?

i gave my son

context window

Monday, April 19, 2010

# How do we view natural language?

## As a string?

i gave my son ?

context window

8

# How do we view natural language?

## As a string?

i gave **my son** ? a

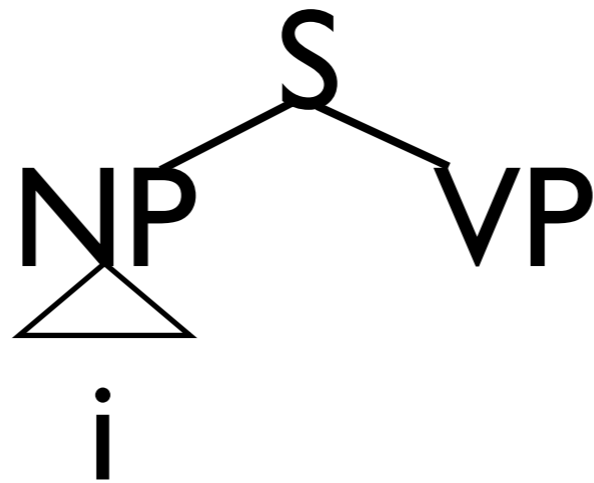context window

# How do we view natural language?

## As a string?

a  baseball bat

i gave  my son  ?

context
window

# How do we view natural language?

## As a string?

a baseball bat

i gave my son ?

context
window

is

# How do we view natural language?

## As a string?

i gave my son ? a baseball bat

context window

is three years old

12

# How do we view natural language?

## As a string?

a  baseball bat

i  gave  my  son  ?

|
context
window

is  three years old

Language is more hierarchical than this!

13

# How do we view natural language?

## Or as a tree?

S

# How do we view natural language?

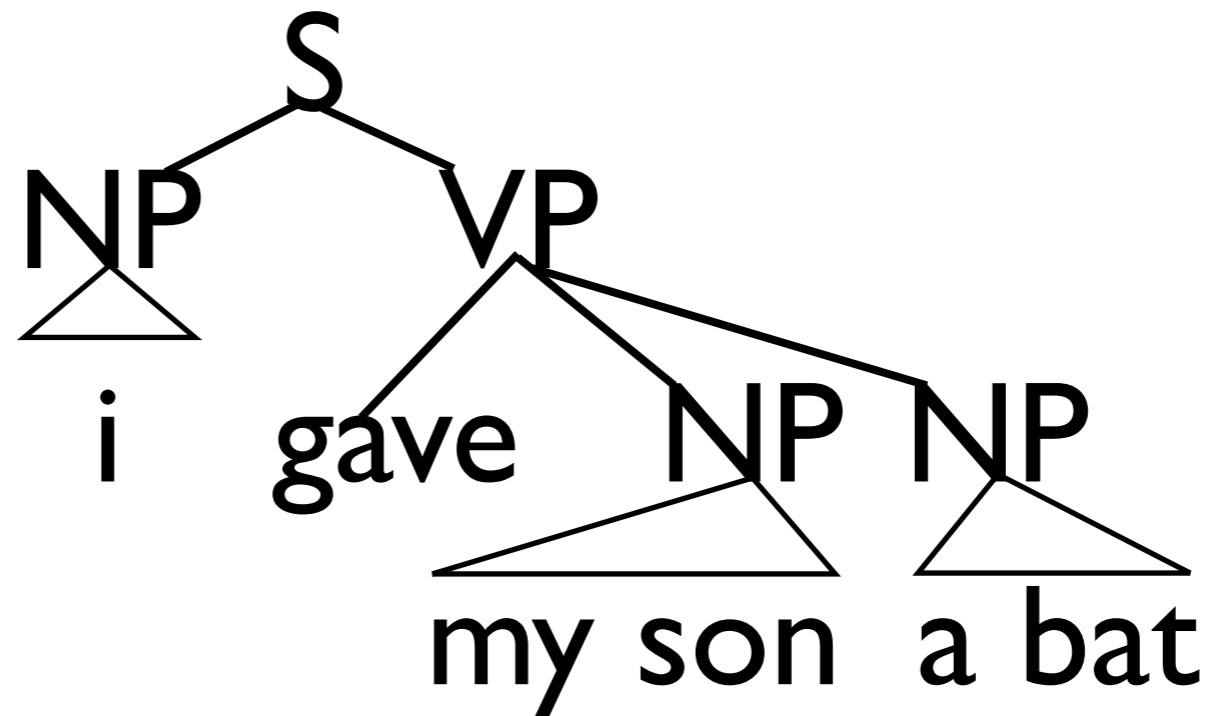## Or as a tree?

S
NP VP

# How do we view natural language?

## Or as a tree?

S
NP VP
i

# How do we view natural language?

## Or as a tree?

```
           S
      NP        VP
      /\
     i    gave    NP NP
```

17

# How do we view natural language?

## Or as a tree?

# How do we view natural language?

## Or as a tree?

# How do we view natural language?

## Or as a tree?

Trees provide syntactic context!

```
            S
        NP      VP
        △
        i    gave   NP  NP
                     △   △
                  my son  a bat
```

20

# String World vs Tree World

```
string
```

Monday, April 19, 2010

# String World vs Tree World

| |
|---|
| string |

| |
|---|
| great formalisms |

# String World vs Tree World

| |
|---|
| string |
| great formalisms |
| useful algorithms |

# String World vs Tree World

| |
|---|
| string |
| great formalisms |
| useful algorithms |
| toolkits |

# String World vs Tree World

| |
|---|
| string |
| great formalisms |
| useful algorithms |
| toolkits |
| rapid progress |

# String World vs Tree World

| |
|---|
| string |
| great formalisms |
| useful algorithms |
| toolkits |
| rapid progress |
| limited expressiveness |

21

# String World vs Tree World

| string | tree |
|---|---|
| great formalisms | |
| useful algorithms | |
| toolkits | |
| rapid progress | |
| limited expressiveness | |

21

# String World vs Tree World

| string | tree |
|---|---|
| great formalisms | great formalisms |
| useful algorithms | |
| toolkits | |
| rapid progress | |
| limited expressiveness | |

21

# String World vs Tree World

| string | tree |
|---|---|
| great formalisms | great formalisms |
| useful algorithms | |
| toolkits | |
| rapid progress | |
| limited expressiveness | powerful expressiveness |

21

# String World vs Tree World

| string | tree |
|---|---|
| great formalisms | great formalisms |
| useful algorithms | few algorithms |
| toolkits | |
| rapid progress | |
| limited expressiveness | powerful expressiveness |

21

# String World vs Tree World

| string | tree |
|---|---|
| great formalisms | great formalisms |
| useful algorithms | few algorithms |
| toolkits | no toolkits |
| rapid progress | |
| limited expressiveness | powerful expressiveness |

# String World vs Tree World

| string | tree |
|---|---|
| great formalisms | great formalisms |
| useful algorithms | few algorithms |
| toolkits | no toolkits |
| rapid progress | slow progress |
| limited expressiveness | powerful expressiveness |

21

# Contributions

| string | tree |
|---|---|
| great formalisms | great formalisms |
| useful algorithms | few algorithms |
| toolkits | no toolkits |
| rapid progress | slow progress |
| limited expressiveness | powerful expressiveness |

22

# Contributions

| string | tree |
|---|---|
| great formalisms | great formalisms |
| useful algorithms | **new algorithms!** |
| toolkits | no toolkits |
| rapid progress | slow progress |
| limited expressiveness | powerful expressiveness |

# Contributions

| string | tree |
|---|---|
| great formalisms | great formalisms |
| useful algorithms | **new algorithms!** |
| toolkits | **new toolkit!** |
| rapid progress | slow progress |
| limited expressiveness | powerful expressiveness |

22

# Contributions

| string | tree |
|---|---|
| great formalisms | great formalisms |
| useful algorithms | **new algorithms!** |
| toolkits | **new toolkit!** |
| rapid progress | **rapid progress!** |
| limited expressiveness | powerful expressiveness |

22

# Weighted finite-state string machines

### Acceptor

the/.3   blue/.8   hairy/.1
red/.2   ε/.8
green/.2   dwarf/.2   elf/.9

### Transducer

blue:ε/.8
the:el/.3   elf:duende/.8
man:hombre/1
elf:duende/.8
man:hombre/1
ε:azúl/.3
ε:triste/.2
the:la/.7

23

# Weighted finite-state string machines

## Acceptor



the/.3
blue/.8
hairy/.1
red/.2
ε/.8
dwarf/.2
green/.2
elf/.9

→

the blue dwarf/.048
green hairy elf/.0144
the red hairy hairy elf/.000432

...

## Transducer



blue:ε/.8
the:el/.3
elf:duende/.8
man:hombre/1
elf:duende/.8
man:hombre/1
ε:azúl/.3
ε:triste/.2
the:la/.7

23

# Weighted finite-state string machines

## Acceptor



the blue dwarf/.048
green hairy elf/.0144
the red hairy hairy elf/.000432

...

## Transducer



the blue elf : el duende azúl/.0576
the blue man : el duende triste/.048

...

23

# Using WFSTs for NLP

Given a string and a transducer, calculate the highest weighted transformation of the string by the transducer

24

# Using WFSTs for NLP

Given a string and a transducer, calculate the highest weighted transformation of the string by the transducer

the blue dwarf

24

# Using WFSTs for NLP

Given a string and a transducer, calculate the highest weighted transformation of the string by the transducer

the blue dwarf

Machine Translation

24

# Using WFSTs for NLP

Given a string and a transducer, calculate the highest weighted transformation of the string by the transducer

the blue dwarf → **Machine Translation** →
el enano azúl / .3
el enano triste / .1
el duende azúl / .05
el azúl duende / .01

...

24

# MT as weighted transducers

| Imagine an English sentence | → | Re-order the words | → | Translate into Spanish |



the green ball = .098

a green green horse = .0036

horse green the ❌

# MT as weighted transducers

Imagine an English sentence → Re-order the words → Translate into Spanish

the
green/.7
the/.01
ball/.2
green
horse/.75
green/.6
green/.04

ball/.05
ball/.3
the
the/.7
green/.03
start
a/.2
green
horse/.02
a
horse/.4

horse green the ✗

25

# MT as weighted transducers

Imagine an English sentence → Re-order the words → Translate into Spanish



the: ε/0.1

ball: ball the/1
green: green the/1
the: the the/1

→ the

ball: ε/0.2

ball: ball ball/1
green: green ball/1
the: the ball/1

→ ball

green: ε/0.7

ball: ball green/1
green: green green/1
the: the green/1

→ green

start

the:the/0.9
ball:ball/0.8
green:green/0.3

the green ball →
- the green ball    .216
- the ball green    .63
- green the ball    .03

26

# MT as weighted transducers



Imagine an English sentence → Re-order the words → Translate into Spanish

ball:pelota/.8
ball:bola/.2

the:el/.35
the:la/.35
the:los/.15
the:las/.15

green:verde/.7
green:campo/.3

horse:caballo/.9
horse:bayo/.1

the ball green

la pelota verde  .196

el bola campo  .021

27

# MT as weighted transducers



Generative story: we corrupt good English into (possibly bad) Spanish

28

# MT as weighted transducers



Decoding story: given some good Spanish, determine the best good English that could produce it

29

# Secret weapons

- WFST toolkits do this calculation for us:

  - AT&T FSM[1] / Google OpenFst[2]

  - USC/ISI Carmel[3]

- Generic operations for manipulation, combination, inference, training

| WFST toolkit operations |
| --- |
| k-best |
| em training |
| determinization |
| composition |
| pipeline inference |
| on-the-fly inference |

1: Mohri, Pereira, Riley, '98    2: Allauzen et al., '07    3: Graehl, '97

# Widely applicable!

the blue elf → Machine Translation (Kumar & Byrne '03) → el enano azúl / .3
el enano triste / .1
el duende azúl / .05
...

# Widely applicable!

the blue elf → **Machine Translation** (Kumar & Byrne '03) → el enano azúl / .3
el enano triste / .1
el duende azúl / .05
...

 → **Decipherment** (Ravi & Knight '09) → "i love killing people" / .2
"eggs, milk, flour, tin foil" / .002
...

# Widely applicable!

the blue elf → **Machine Translation** (Kumar & Byrne '03) →

el enano azúl / .3
el enano triste / .1
el duende azúl / .05
...

 → **Decipherment** (Ravi & Knight '09) →

"i love killing people" / .2
"eggs, milk, flour, tin foil" / .002
...

 → **Speech Recognition** (Pereira et al. '94) →

"this is the bbc" / .5
"this is the bee I see" / .1
"the sister beyoncé" / .1
...

31

# Widely applicable!

the blue elf → **Machine Translation** (Kumar & Byrne '03) →

el enano azúl / .3
el enano triste / .1
el duende azúl / .05
...

**Decipherment** (Ravi & Knight '09) →

"i love killing people" / .2
"eggs, milk, flour, tin foil" / .002
...

**Speech Recognition** (Pereira et al. '94) →

"this is the bbc" / .5
"this is the bee I see" / .1
"the sister beyoncé" / .1
...

**Poetry Generation** (Greene & Knight '10) →

"there was a computer from apple that wore a red rose in its lapel" / .1
...

31

# NLP work using WFSTs

Translation
(Kumar & Byrne '03)

Decipherment
(Ravi & Knight '09)

Speech
Recognition
(Pereira et al. '94)

Poetry
Generation
(Greene & Knight '10)

OCR
(Kolak et al. '03)

Morphology
(Karttunen et al. '92)

POS Tagging
(Church '88)

Spelling
Correction
(Boyd '09)

Transliteration
(Knight & Graehl '98)

Also see summary: book chapter of *Handbook of Weighted Automata* (Knight & May '08)

# Limitations of strings

- Can't do arbitrary long-distance reordering

- Can't maintain arbitrary long-distance dependencies

- Can't naturally integrate syntax information



33

# Limitations of strings

- Can't do arbitrary long-distance reordering

- Can't maintain arbitrary long-distance dependencies

- Can't naturally integrate syntax information



33

# Limitations of strings

- Can't do arbitrary long-distance reordering

- Can't maintain arbitrary long-distance dependencies

- Can't naturally integrate syntax information

```
          S
         / \
        VP  NP
       /△\  /△\
      /
 (1st-person-
   singular)      33
```

# Limitations of strings

- Can't do arbitrary long-distance reordering

- Can't maintain arbitrary long-distance dependencies

- Can't naturally integrate syntax information

```
            S
           / \
          VP  NP
         /\   /\
        /  \ /  \
```

*(1st-person-singular)*     33     *(1st-person-singular)*

# Limitations of strings

- Can't do arbitrary long-distance reordering

- Can't maintain arbitrary long-distance dependencies

- Can't naturally integrate syntax information

## But that's what we want!

33

# Limitations of strings

- Can't do arbitrary long-distance reordering

- Can't maintain arbitrary long-distance dependencies

- Can't naturally integrate syntax information

## But that's what we want!

Parsing
(Collins '97)

33

# Limitations of strings

- Can't do arbitrary long-distance reordering

- Can't maintain arbitrary long-distance dependencies

- Can't naturally integrate syntax information

## But that's what we want!

Parsing          Question Answering
(Collins '97)    (Echihabi & Marcu '03)

33

# Limitations of strings

- Can't do arbitrary long-distance reordering

- Can't maintain arbitrary long-distance dependencies

- Can't naturally integrate syntax information

## But that's what we want!

Parsing
(Collins '97)

Question Answering
(Echihabi & Marcu '03)

Language Modeling
(Charniak '01)

33

# Limitations of strings

- Can't do arbitrary long-distance reordering

- Can't maintain arbitrary long-distance dependencies

- Can't naturally integrate syntax information

## But that's what we want!

Parsing
(Collins '97)

Question Answering
(Echihabi & Marcu '03)

Language Modeling
(Charniak '01)

Summarization
(Knight & Marcu '03)

33

# Limitations of strings

- Can't do arbitrary long-distance reordering

- Can't maintain arbitrary long-distance dependencies

- Can't naturally integrate syntax information

## But that's what we want!

Parsing
(Collins '97)

Question Answering
(Echihabi & Marcu '03)

Machine Translation
(Yamada & Knight '01)
(Galley et al. '04)

Language Modeling
(Charniak '01)

Summarization
(Knight & Marcu '03)

(Mi et al. '08)
(Zhang et al. '08)

33

# Lots of work with tree models, but NO tree toolkit!

Parsing
(Collins '97)

Question Answering
(Echihabi & Marcu '03)

Machine Translation
(Yamada & Knight '01)
(Galley et al. '04)
(Mi et al. '08)
(Zhang et al. '08)

Language Modeling
(Charniak '01)

Summarization
(Knight & Marcu '03)

33

# Weighted finite-state tree machines

## Grammar



## Transducer

# Weighted finite-state tree machines

# Weighted finite-state tree machines

# Weighted regular tree grammars



Tree      Weight

(Berstel & Reutenauer, 1982)

35

# Weighted regular tree grammars



Tree     Weight

(Berstel & Reutenauer, 1982)

# Weighted regular tree grammars



| Tree | Weight |
|------|--------|
| NP (1 2 3) | .2 |

(Berstel & Reutenauer, 1982)

# Weighted regular tree grammars

0 →(.2) NP
         / | \
        1  2  3

1 →(.8) the

2 →(.4) blue

2 →(.2) red

2 →(.3) JJ
         / \
        2   2

3 →(.9) elf

**Tree        Weight**

NP
/ | \
1  2  3
|
1 →(.8) the                    .2

(Berstel & Reutenauer, 1982)

# Weighted regular tree grammars



| Tree | Weight |
|------|--------|
| NP<br>the  2  3 | .16 |

(Berstel & Reutenauer, 1982)

# Weighted regular tree grammars



Tree      Weight

.16

(Berstel & Reutenauer, 1982)

# Weighted regular tree grammars



Tree     Weight

.064

(Berstel & Reutenauer, 1982)

41

# Weighted regular tree grammars



Tree     Weight

.064

(Berstel & Reutenauer, 1982)

# Weighted regular tree grammars



Tree      Weight

the blue elf      .0576

(Berstel & Reutenauer, 1982)

43

# Weighted tree transducers



NP → NP   .3

④ ⑤ ⑥ ⑦

⑤ the → el   .5

⑥ elf → duende   1

⑦ blue → azúl   .4

⑦ blue → triste   .2

| Tree | Weight |
|------|--------|
| ④ NP<br>the blue elf | 1 |

(Kuich, 1998)

44

# Weighted tree transducers



Tree        Weight

(Kuich, 1998)

45

# Weighted tree transducers



Tree     Weight

(Kuich, 1998)

# Weighted tree transducers



Tree    Weight

(Kuich, 1998)

# Weighted tree transducers



Tree    Weight

(Kuich, 1998)

# Weighted tree transducers



Tree      Weight

.3

(Kuich, 1998)

49

# Weighted tree transducers



Tree      Weight

(Kuich, 1998)

# Weighted tree transducers



Tree          Weight

.15

(Kuich, 1998)

# Weighted tree transducers



Tree      Weight

(Kuich, 1998)

# Weighted tree transducers



Tree      Weight

.15

(Kuich, 1998)

53

# Weighted tree transducers



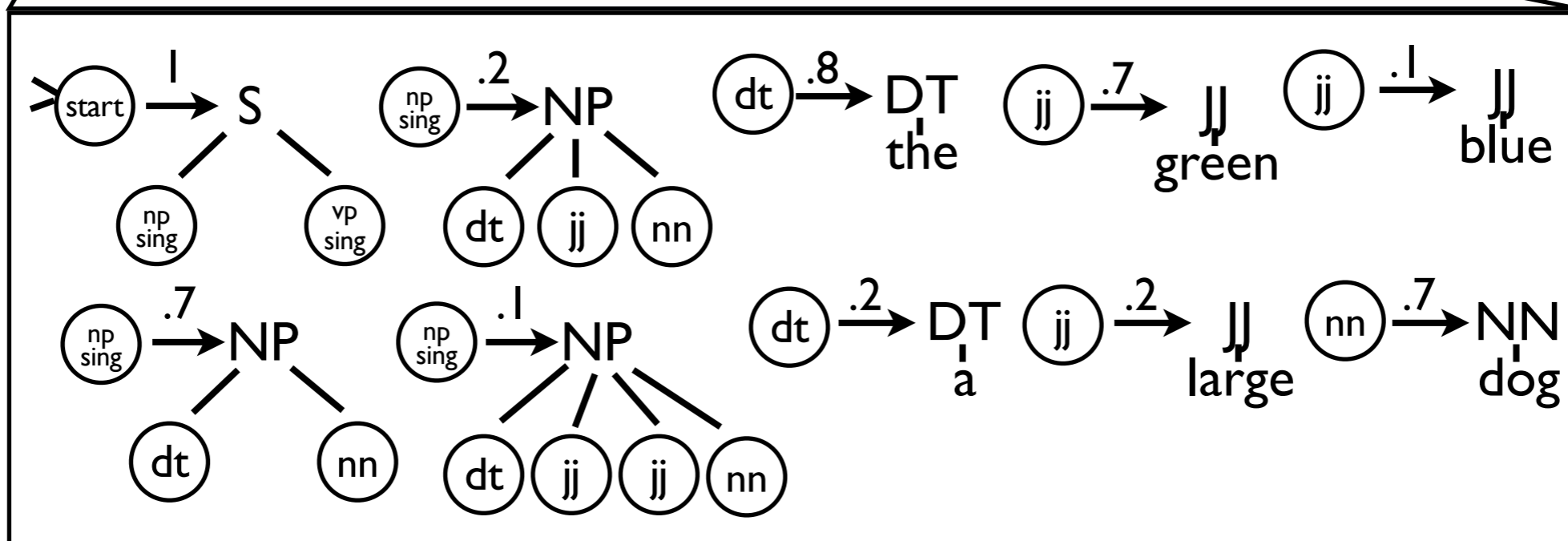Tree       Weight

.15

(Kuich, 1998)

54

# Weighted tree transducers



Tree — Weight

NP
el duende azúl      .06

(Kuich, 1998)

55

# Weighted tree-string transducers



$$NP(\bullet, \blacklozenge, \pentagon) \xrightarrow{.3} \bigcirc \ \pentagon \ \blacklozenge$$

Tree    Weight

the $\xrightarrow{.5}$ el

elf $\xrightarrow{1}$ duende

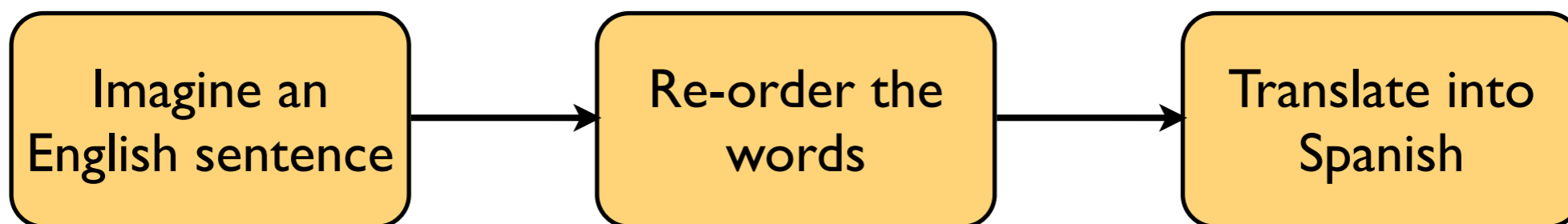blue $\xrightarrow{.4}$ azúl
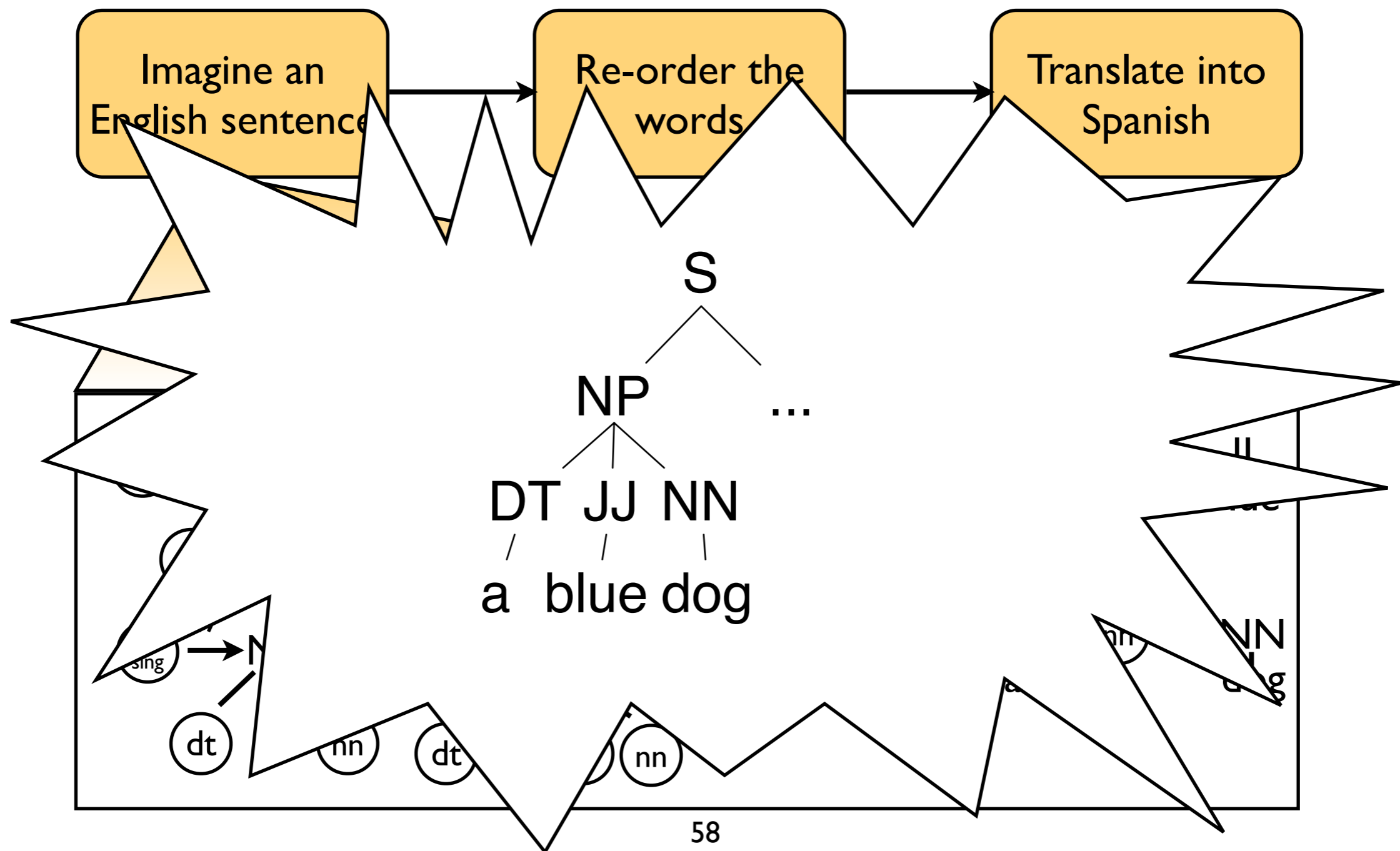
blue $\xrightarrow{.2}$ triste

NP(the, blue, elf)    1

(Kuich, 1998)

56

# Weighted tree-string transducers



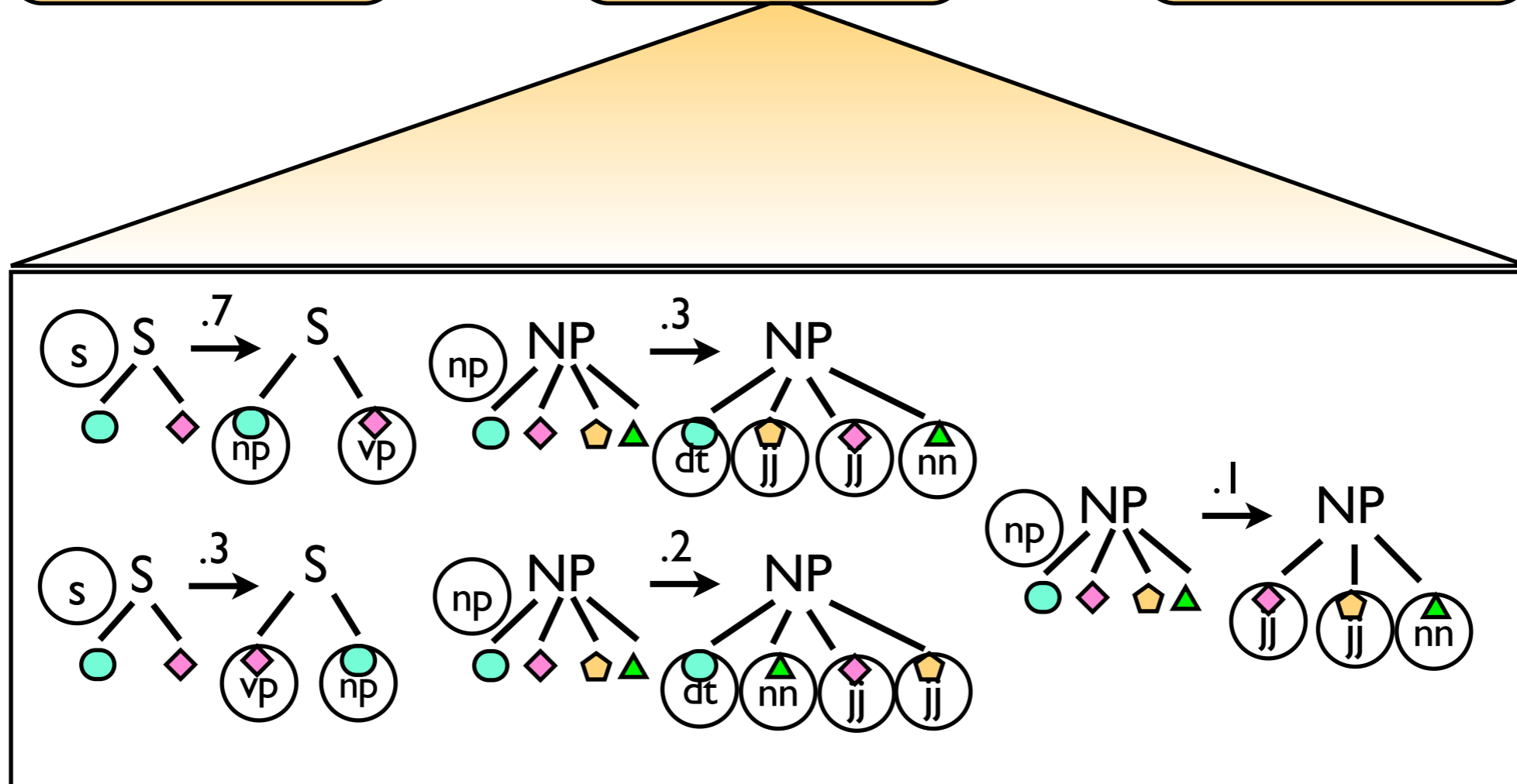| String | Weight |
|---|---|
| el duende azúl | .06 |

(Kuich, 1998)

57

# MT as weighted tree transducers
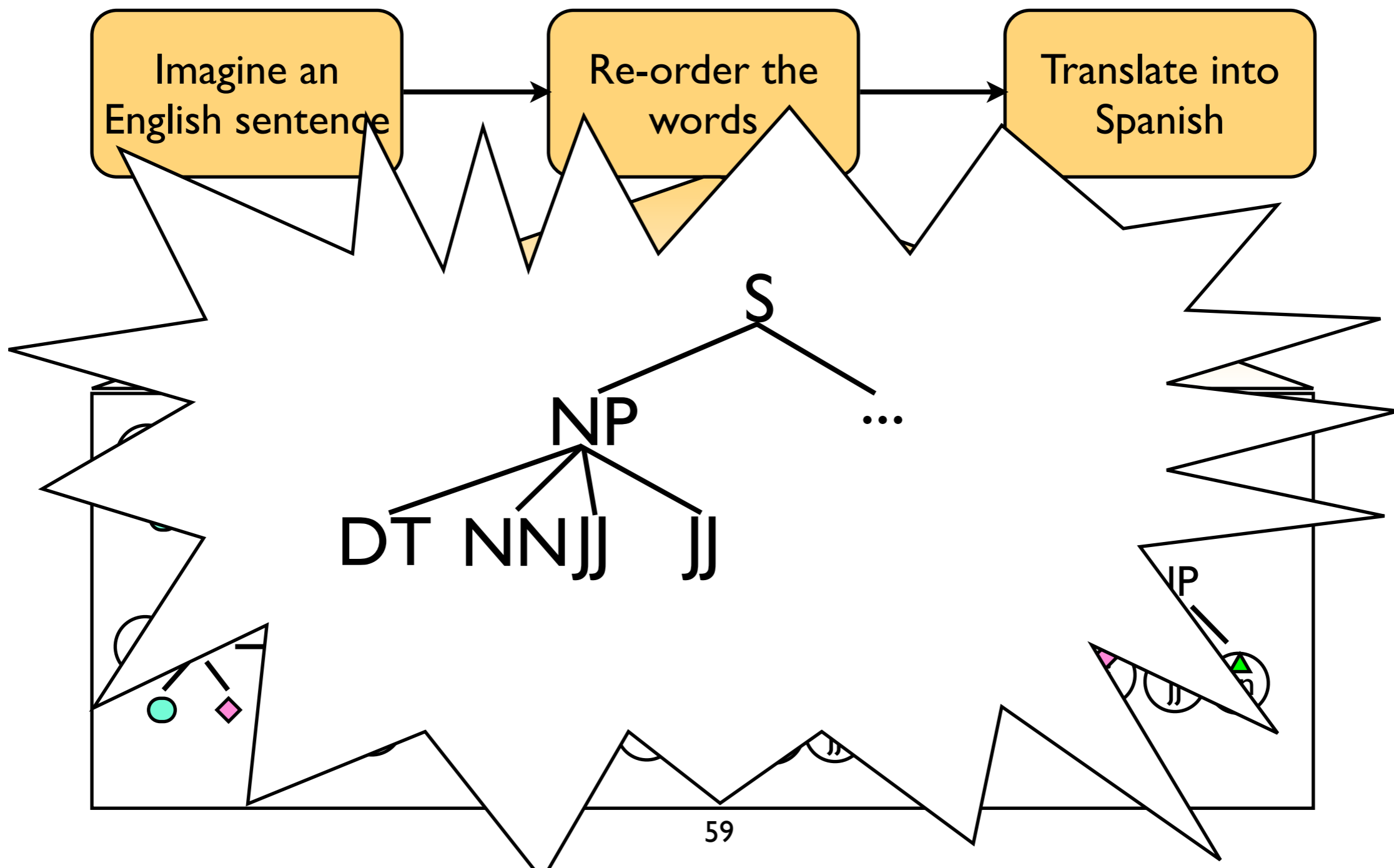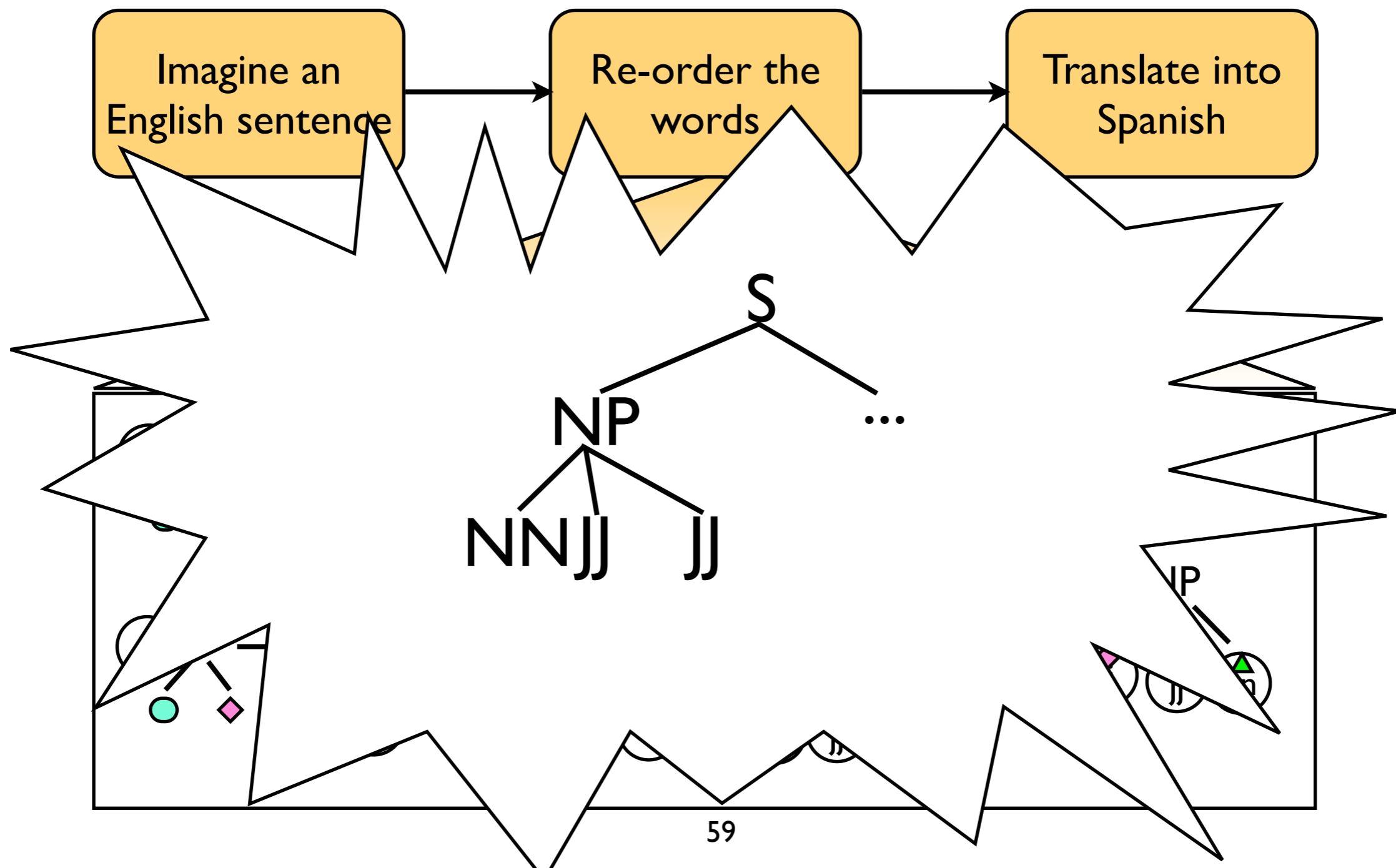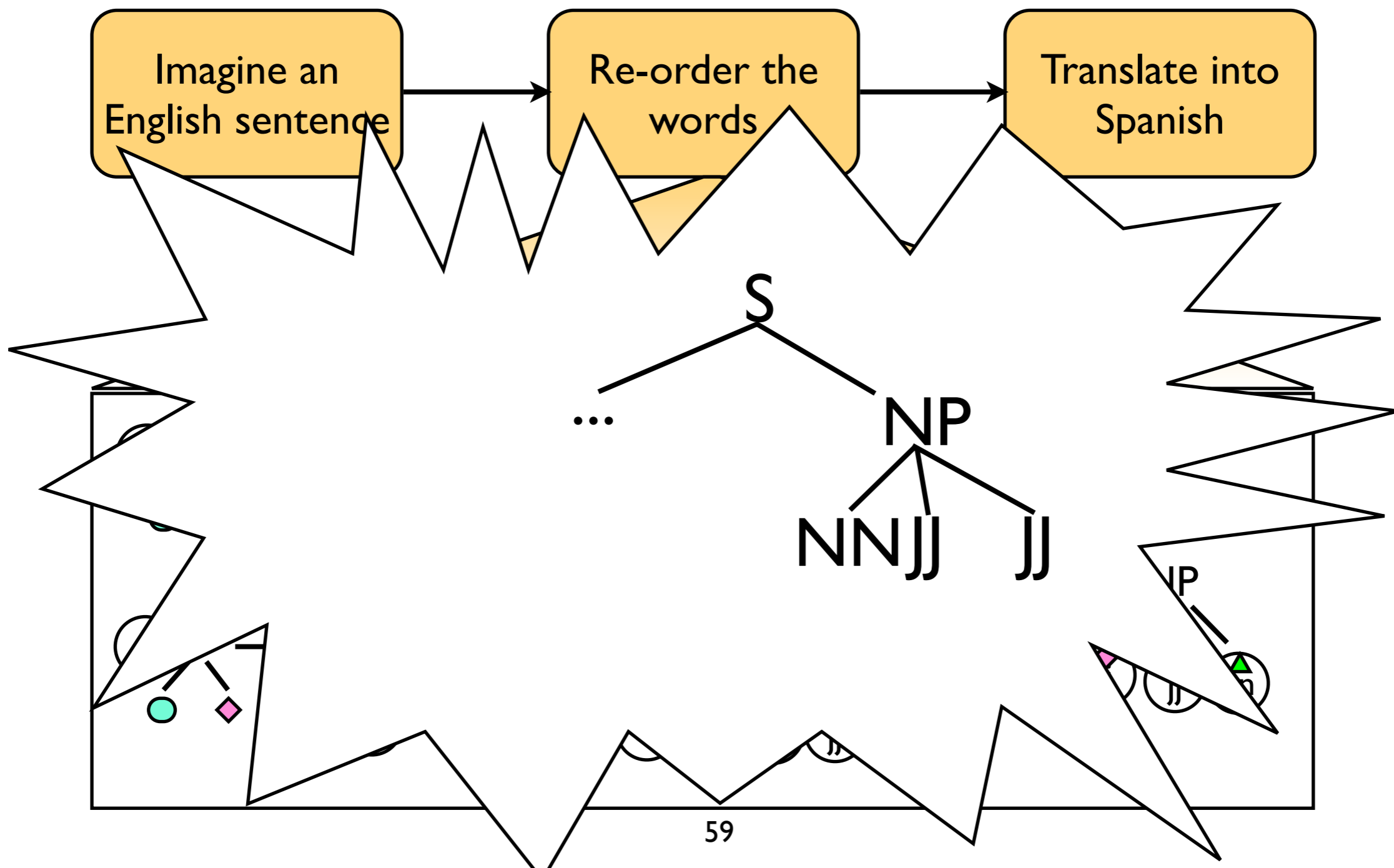
# MT as weighted tree transducers

Imagine an English sentence

Re-order the words

Translate into Spanish

```
          S
         / \
       NP   ...
      / | \
    DT JJ NN
    |   |   |
    a  blue dog
```
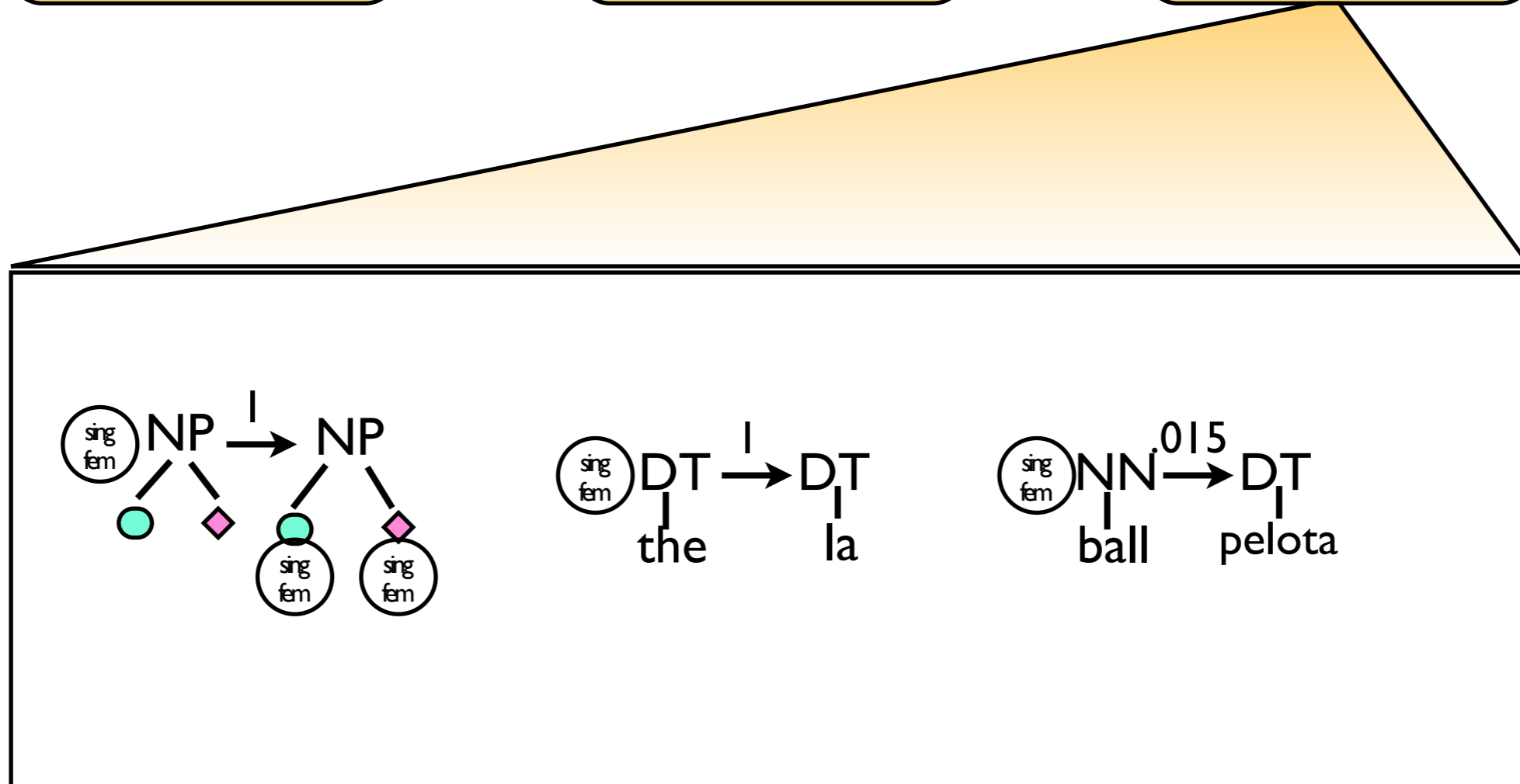
dt          nn    dt          nn

# MT as weighted tree transducers
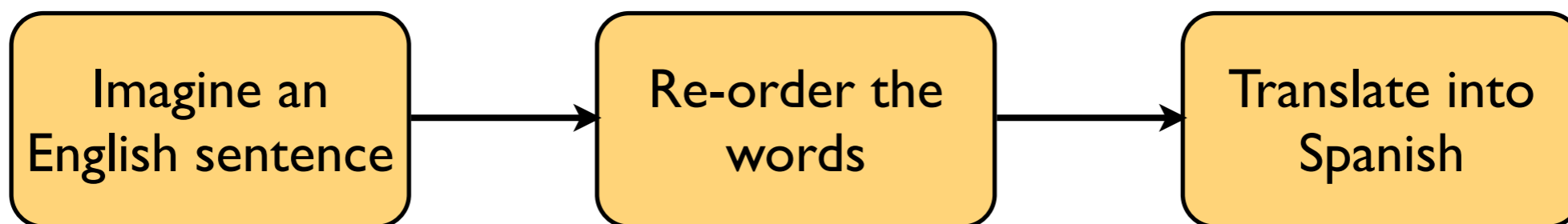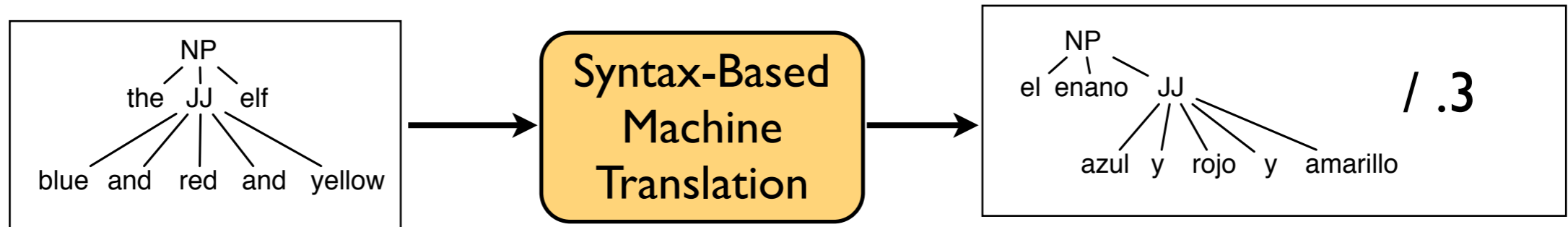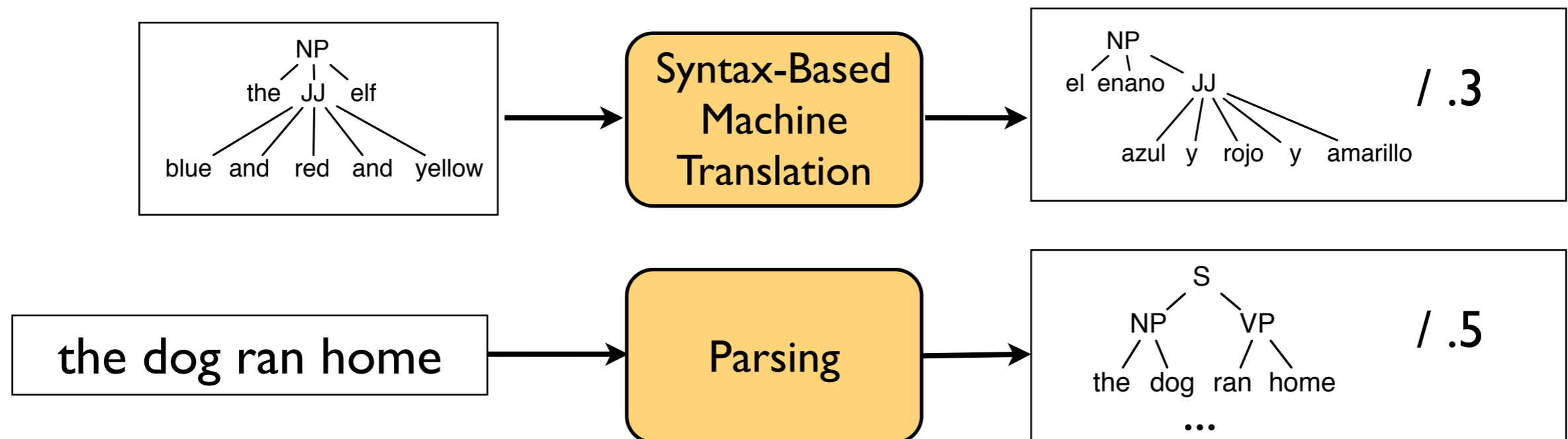
# MT as weighted tree transducers

Imagine an English sentence → Re-order the words → Translate into Spanish

S
├─ NP
│  ├─ DT  JJ  JJ  NN
└─ ...

# MT as weighted tree transducers

Imagine an English sentence → Re-order the words → Translate into Spanish

S
NP ...
DT NN JJ JJ

IP

JJ

59

# MT as weighted tree transducers

Imagine an English sentence → Re-order the words → Translate into Spanish

S
NP  ...
NN JJ  JJ

IP

59

# MT as weighted tree transducers

Imagine an English sentence → Re-order the words → Translate into Spanish

S
... NP
NN JJ   JJ

IP

# MT as weighted tree transducers

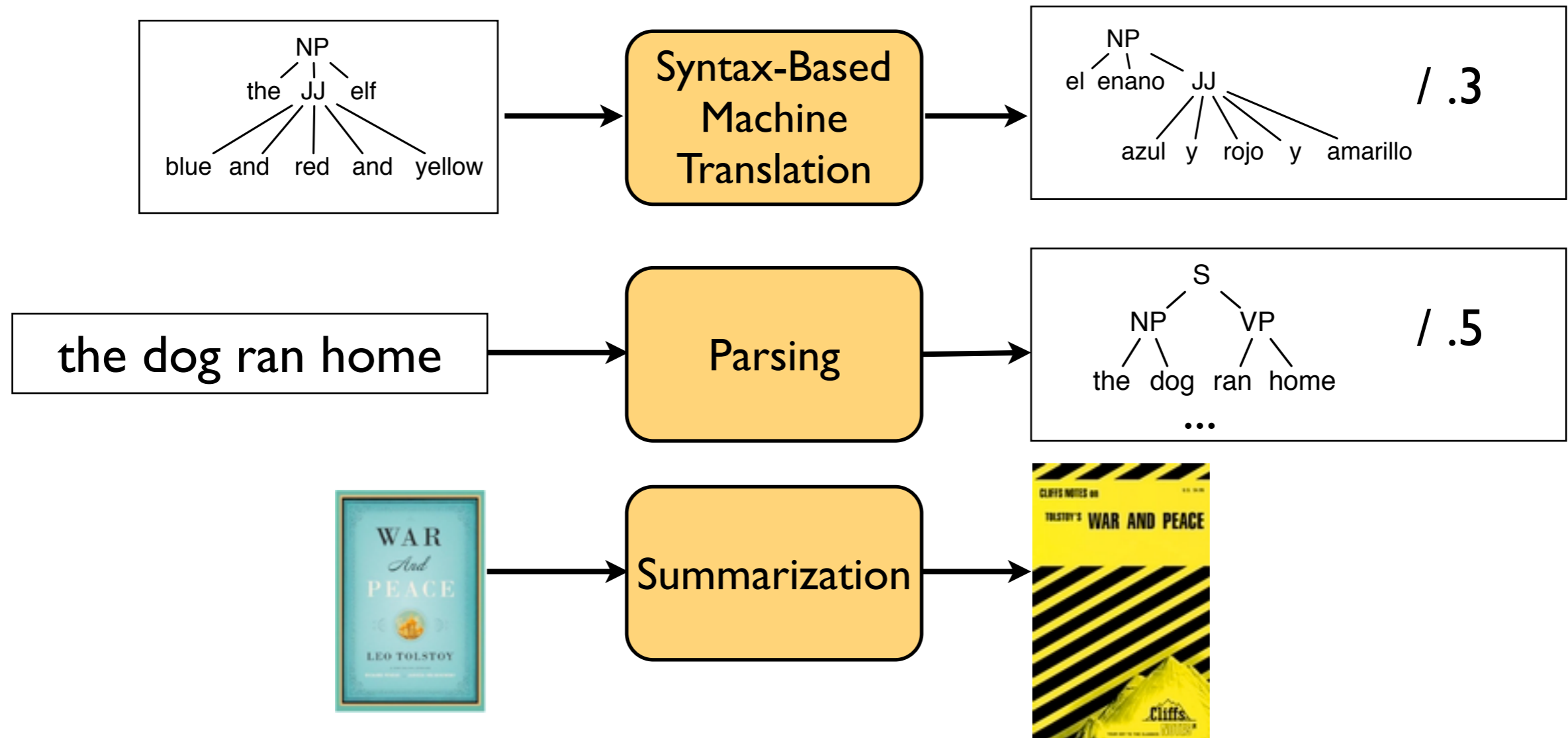| Imagine an English sentence | → | Re-order the words | → | Translate into Spanish |
|---|---|---|---|---|



60

# Great, so now we can solve harder problems!

# Great, so now we can solve harder problems!

# Great, so now we can solve harder problems!

# Great, so now we can solve harder problems!



NP
the JJ elf
blue and red and yellow

→ Syntax-Based Machine Translation →

NP
el enano JJ
azul y rojo y amarillo

/ .3

the dog ran home → Parsing →

S
NP VP
the dog ran home
...

/ .5

WAR And PEACE LEO TOLSTOY → Summarization → CLIFFS NOTES on TOLSTOY'S WAR AND PEACE

# Not so fast!

# String world has many more available operations than tree world!

| Operation | String | Tree |
|---|---|---|
| k-best | yes | alg[1] |
| em training | yes | alg[2] |
| determinization | yes | no |
| composition | yes | proof of concept[3] |
| pipeline inference | yes | proof of concept[4] |
| on-the-fly inference | yes | no |

1: Huang & Chiang, 2005
2: Graehl & Knight, 2004

3: Maletti, 2006
4: Fülöp, Maletti, Vogler, 2010

62

# Algorithmic contribution I: weighted determinization

| Operation | String | Tree |
|---|---|---|
| k-best | yes | alg |
| em training | yes | alg |
| determinization | yes | alg |
| composition | yes | proof of concept[3] |
| pipeline inference | yes | proof of concept[4] |
| on-the-fly inference | yes | no |

Algorithmic I

# Algorithmic contribution II: efficient inference

| Operation | String | Tree |
|---|---|---|
| k-best | yes | alg |
| em training | yes | alg |
| determinization | yes | alg |
| composition | yes | alg |
| pipeline inference | yes | alg |
| on-the-fly inference | yes | alg |

Algorithmic I

Algorithmic II

64

# Practical contribution I:
# weighted tree transducer toolkit

Practical I

| Operation | String | Tree |
|---|---|---|
| k-best | yes | yes |
| em training | yes | yes |
| determinization | yes | yes |
| composition | yes | yes |
| pipeline inference | yes | yes |
| on-the-fly inference | yes | yes |

Algorithmic I

Algorithmic II

65

# Practical contribution II: syntactic re-alignment

Practical I

Practical II

Algorithmic I

Algorithmic II

| Operation | String | Tree |
|---|---|---|
| k-best | yes | yes |
| em training | yes | yes |
| determinization | yes | yes |
| composition | yes | yes |
| pipeline inference | yes | yes |
| on-the-fly inference | yes | yes |

66

# Determinization of weighted tree automata

(May & Knight, HLT-NAACL '06)
(Büchse, May, Vogler, FSMNLP '09)



Elevated Mohri algorithm ('97) to tree automata
Demonstrated empirical gains in parsing and MT

67

BEFORE

$t \xrightarrow{.2} D$

D — q

D — r

$t \xrightarrow{.3} D$

D — q

D — s

$q \xrightarrow{.3} A$

$r \xrightarrow{.2} B$

$s \xrightarrow{.6} B$

$s \xrightarrow{.4} C$

AFTER

68

BEFORE

$$t \xrightarrow{.2} D$$
$$q \quad r$$

$$t \xrightarrow{.3} D$$
$$q \quad s$$

$$q \xrightarrow{.3} A \qquad r \xrightarrow{.2} B \qquad s \xrightarrow{.6} B \qquad s \xrightarrow{.4} C$$

AFTER

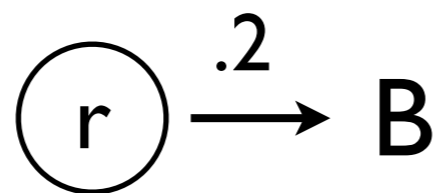Non-deterministic rules
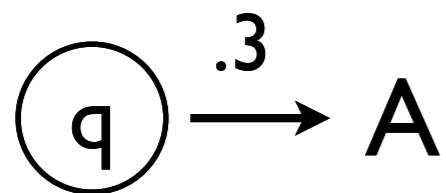(treating grammar as bottom-up acceptor)

69

BEFORE

$t \xrightarrow{.2} D$ (q, r)

$t \xrightarrow{.3} D$ (q, s)

$q \xrightarrow{.3} A$  $r \xrightarrow{.2} B$  $s \xrightarrow{.6} B$  $s \xrightarrow{.4} C$
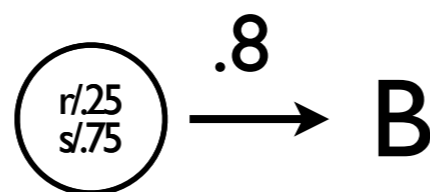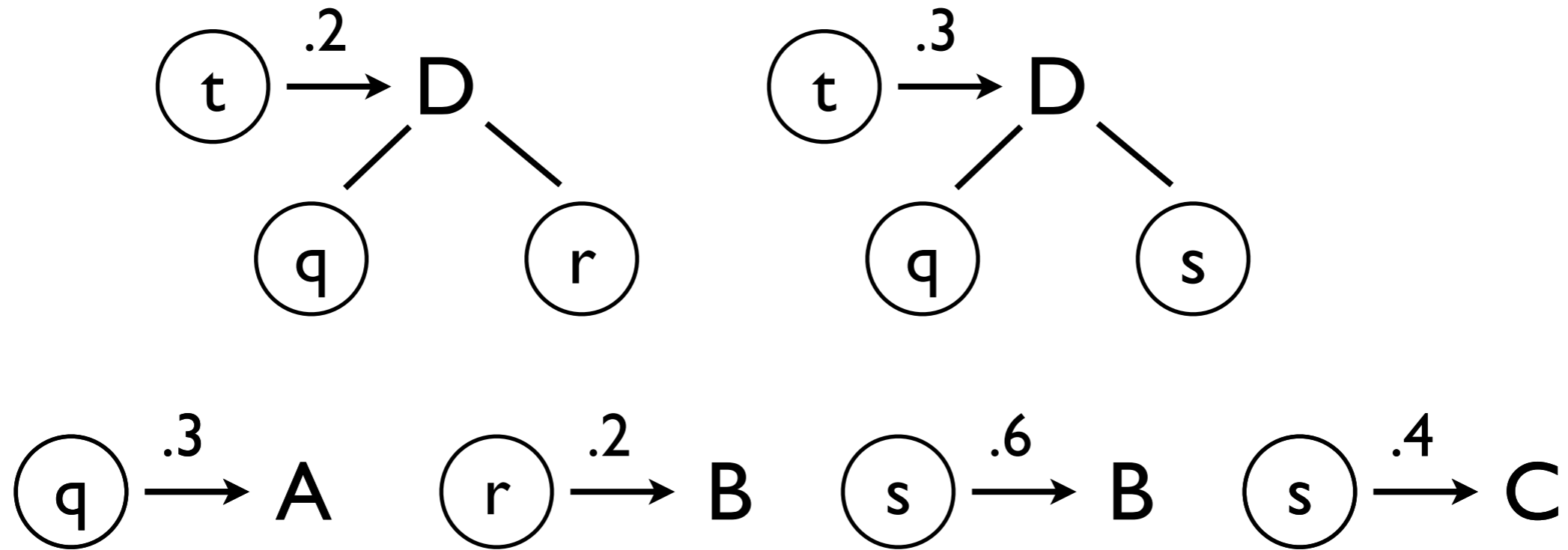
AFTER

Merge terminal rules with same right sides

$s \xrightarrow{.8} B$

70

**BEFORE**

$t \xrightarrow{.2} D$ with children $q$, $r$

$t \xrightarrow{.3} D$ with children $q$, $s$

$q \xrightarrow{.3} A$
$r \xrightarrow{.2} B$
$s \xrightarrow{.6} B$
$s \xrightarrow{.4} C$

**AFTER**

Merge terminal rules
with same right sides

$\text{r/.25, s/.75} \xrightarrow{.8} B$

71

BEFORE

$$t \xrightarrow{.2} D$$
$$q \quad r$$

$$t \xrightarrow{.3} D$$
$$q \quad s$$

$$q \xrightarrow{.3} A \qquad r \xrightarrow{.2} B \qquad s \xrightarrow{.6} B \qquad s \xrightarrow{.4} C$$

## Process the other terminal rules

AFTER

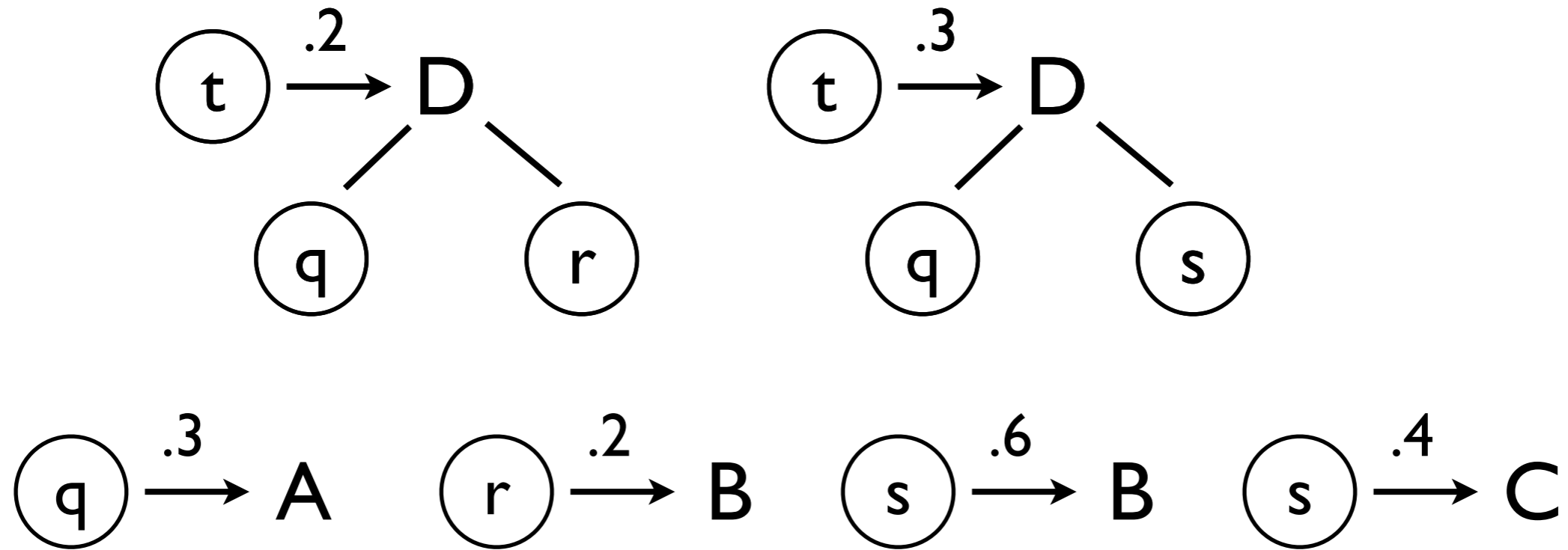$$q \xrightarrow{.3} A \qquad \text{r/.25 s/.75} \xrightarrow{.8} B \qquad s \xrightarrow{.4} C$$

73

**BEFORE**

$$t \xrightarrow{.2} D$$
$$D \text{ — } q \quad r$$

$$t \xrightarrow{.3} D$$
$$D \text{ — } q \quad s$$

$$q \xrightarrow{.3} A \qquad r \xrightarrow{.2} B \qquad s \xrightarrow{.6} B \qquad s \xrightarrow{.4} C$$

## Process the other terminal rules

**AFTER**

$$q \xrightarrow{.3} A \qquad \overset{r/.25}{\underset{s/.75}{\bigcirc}} \xrightarrow{.8} B \qquad s \xrightarrow{.4} C$$
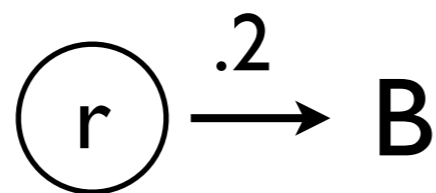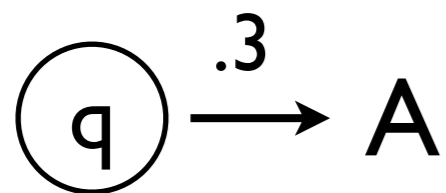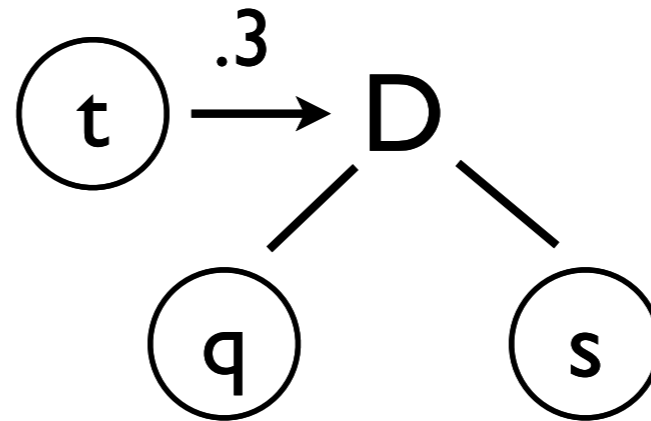
74

Process the other terminal rules

75

BEFORE

$$t \xrightarrow{.2} D \quad (q \quad r)$$

$$t \xrightarrow{.3} D \quad (q \quad s)$$

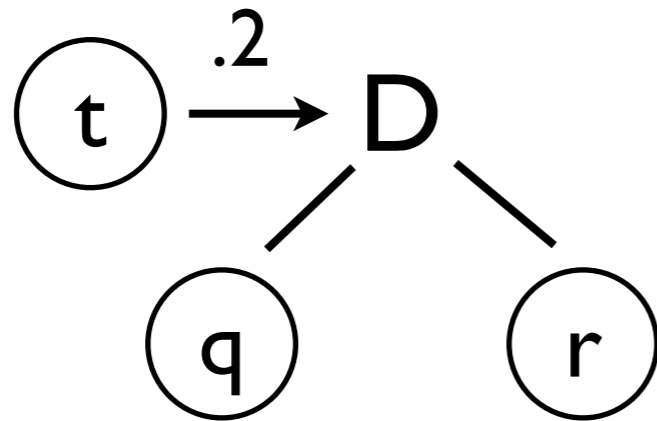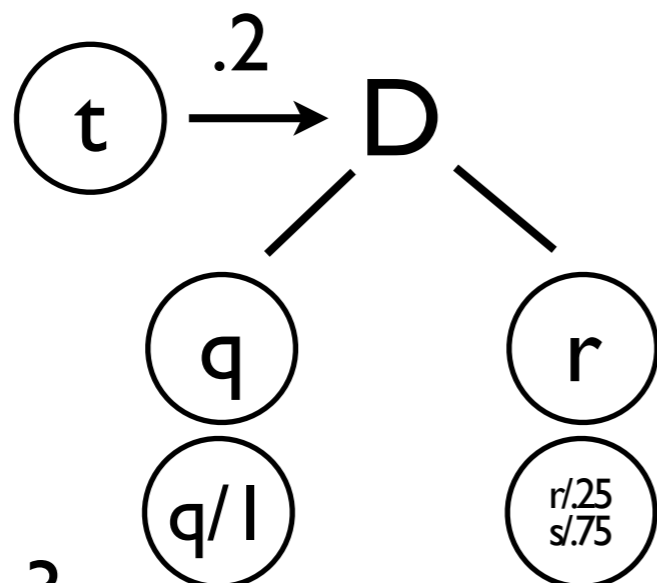$$q \xrightarrow{.3} A \qquad r \xrightarrow{.2} B \qquad s \xrightarrow{.6} B \qquad s \xrightarrow{.4} C$$

Process the other terminal rules

AFTER

$$q/1 \xrightarrow{.3} A \qquad \binom{r/.25}{s/.75} \xrightarrow{.8} B \qquad s/1 \xrightarrow{.4} C$$

76

BEFORE

t $\xrightarrow{.2}$ D
q    r

t $\xrightarrow{.3}$ D
q    s

q $\xrightarrow{.3}$ A    r $\xrightarrow{.2}$ B    s $\xrightarrow{.6}$ B    s $\xrightarrow{.4}$ C
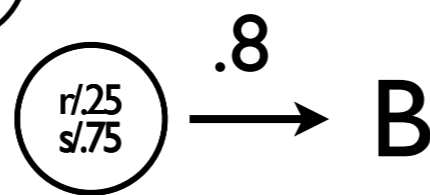
AFTER

t $\xrightarrow{.2}$ D
q    r
q/1    r/.25 s/.75
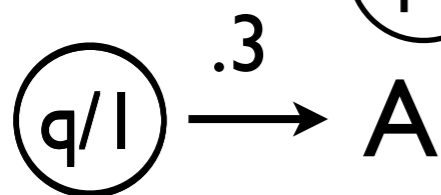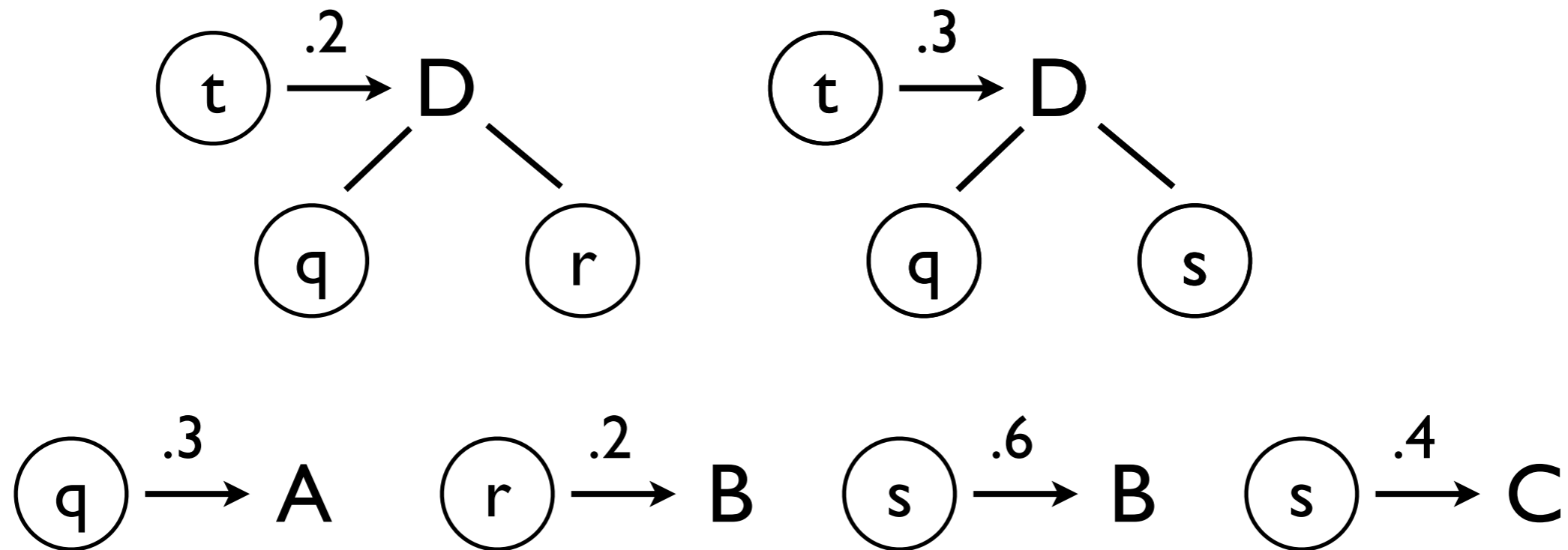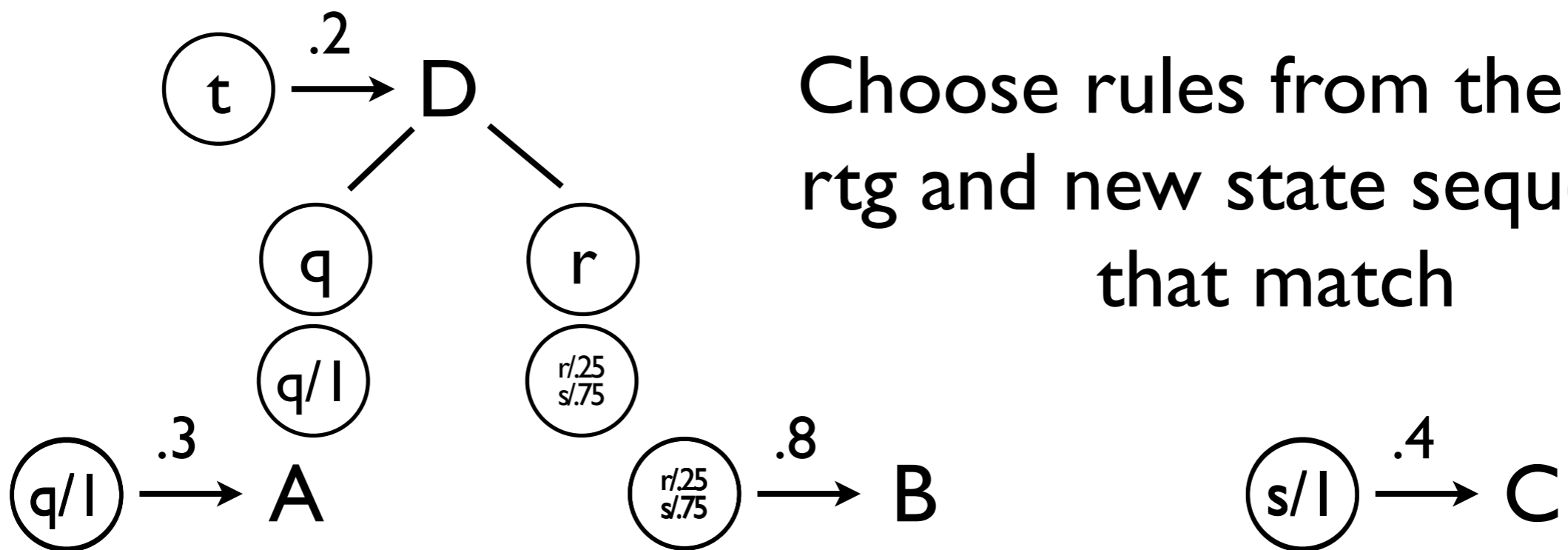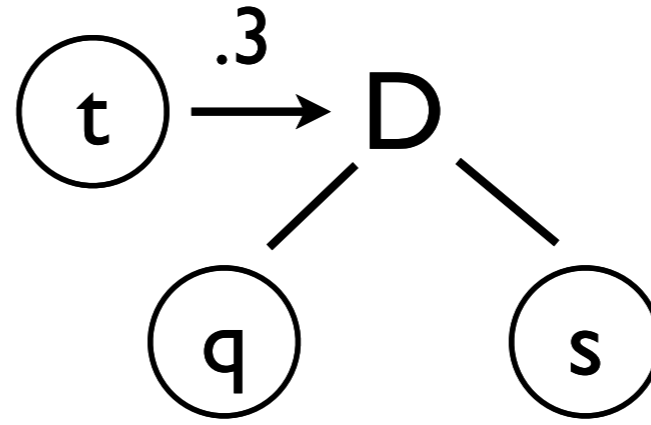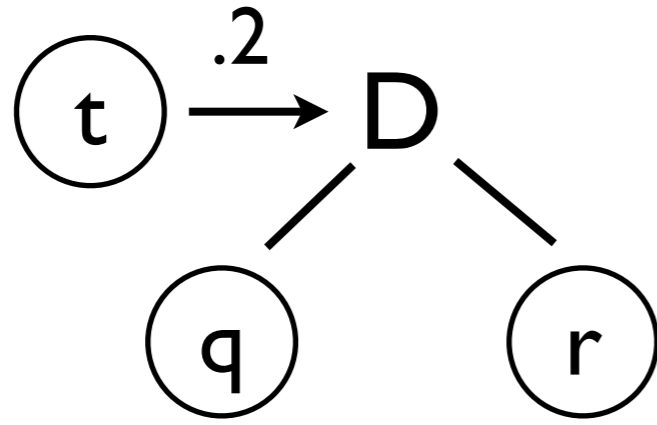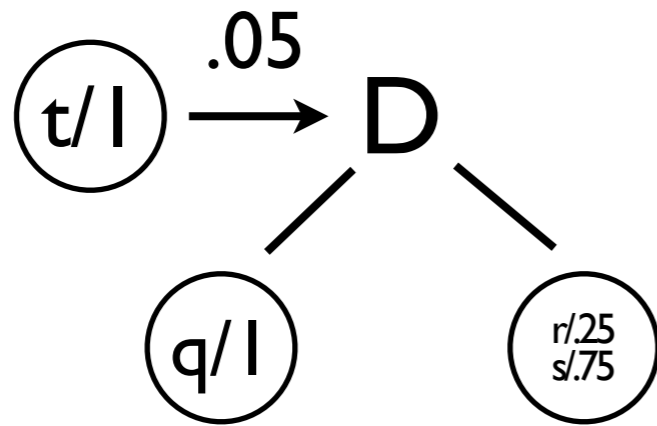
Choose rules from the input rtg and new state sequences that match

q/1 $\xrightarrow{.3}$ A    r/.25 s/.75 $\xrightarrow{.8}$ B    s/1 $\xrightarrow{.4}$ C

77

BEFORE

$$t \xrightarrow{.2} D \quad (q, r)$$

$$t \xrightarrow{.3} D \quad (q, s)$$

$$q \xrightarrow{.3} A \qquad r \xrightarrow{.2} B \qquad s \xrightarrow{.6} B \qquad s \xrightarrow{.4} C$$

AFTER

$$t \xrightarrow{.2} D \quad (q, r)$$

q/1

r/.25 s/.75

$$q/1 \xrightarrow{.3} A \qquad r/.25 \; s/.75 \xrightarrow{.8} B \qquad s/1 \xrightarrow{.4} C$$

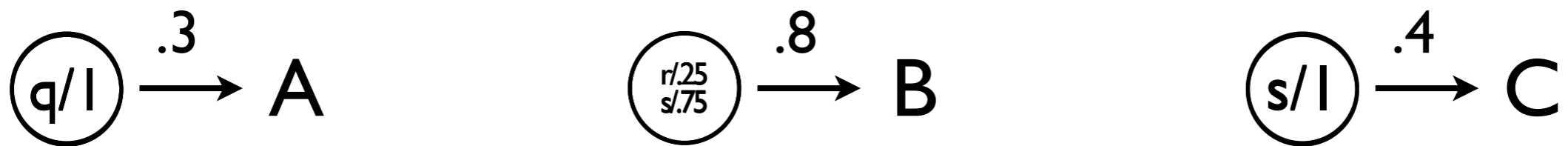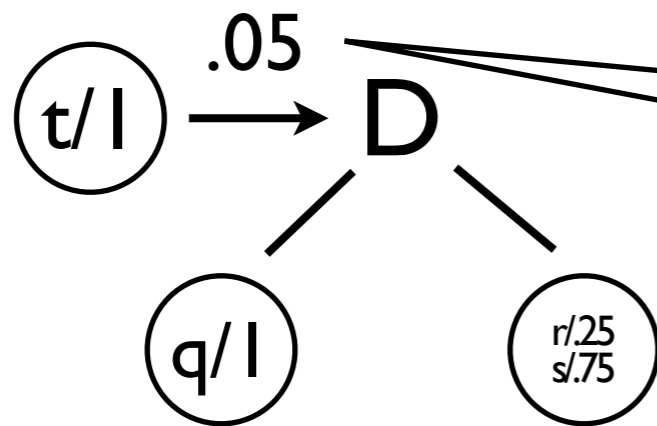Choose rules from the input rtg and new state sequences that match

78

Form new rules from these components

B**EFORE**

$t \xrightarrow{.2} D$ with children $q$, $r$

$t \xrightarrow{.3} D$ with children $q$, $s$

$q \xrightarrow{.3} A$ 　 $r \xrightarrow{.2} B$ 　 $s \xrightarrow{.6} B$ 　 $s \xrightarrow{.4} C$

A**FTER**

$t/1 \xrightarrow{.05} D$ with children $q/1$, $r/.25 \; s/.75$

rule weight of .2 times residual of .25

$q/1 \xrightarrow{.3} A$ 　 $r/.25 \; s/.75 \xrightarrow{.8} B$ 　 $s/1 \xrightarrow{.4} C$

80

rule weight of .3 times residual of .75

83

BEFORE

$t \xrightarrow{.2} D$ (q, r)

$t \xrightarrow{.3} D$ (q, s)

$q \xrightarrow{.3} A$   $r \xrightarrow{.2} B$   $s \xrightarrow{.6} B$   $s \xrightarrow{.4} C$

AFTER

$t/1 \xrightarrow{.025} D$ (q/1, r/.25 s/.75)

These rules are identical except for their weight, so we'll sum them

$q/1 \xrightarrow{.3} A$   $r/.25 \; s/.75 \xrightarrow{.8} B$   $s/1 \xrightarrow{.4} C$

84

These rules are identical except for their weight, so we'll sum them

85

# Empirical experiments

## Machine translation (Galley et al. '04, '06)



Determinization removes duplicates and re-ranks n-best lists

| Method | Bleu |
|---|---|
| Undeterminized | 21.87 |
| Top-500 "crunching" | 23.33 |
| Determinized | 24.17 |

90

# Empirical experiments
## DOP parsing (Bod '92)



Determinization removes duplicates and re-ranks n-best lists

| Method | Precision | Recall | F |
|---|---|---|---|
| Undeterminized | 80.23 | 80.18 | 80.20 |
| Top-500 "crunching" | 80.48 | 80.29 | 80.39 |
| Determinized | 81.09 | 79.72 | 80.40 |

# Efficient inference through cascades of weighted tree transducers
## (May, Knight, Vogler, Submitted)

- First presentation of algorithms for inference through weighted extended tree transducer cascades

- On-the-fly approach significantly outperforms "classic" approach

92

# Inference through *string* transducers

Given a string and a transducer, calculate the highest weighted transformation of the string by the transducer

93

# Inference through *string* transducers

Given a string and a transducer, calculate the highest weighted transformation of the string by the transducer

the blue dwarf

93

# Inference through *string* transducers

Given a string and a transducer, calculate the highest weighted transformation of the string by the transducer

the blue dwarf

Machine Translation

93

# Inference through *string* transducers

Given a string and a transducer, calculate the highest weighted transformation of the string by the transducer

the blue dwarf → **Machine Translation** →

el enano azúl / .3
el enano triste / .1
el duende azúl / .05
el azúl duende / .01

...

93

# Inference through string *cascades*

Given a string and a cascade, calculate the highest weighted transformation of the string by the cascade

(Pereira & Riley, 1997)

94

# Inference through string *cascades*

Given a string and a cascade, calculate the highest weighted transformation of the string by the cascade

A B

(Pereira & Riley, 1997)

94

# Inference through string *cascades*

Given a string and a cascade, calculate the highest weighted transformation of the string by the cascade

A:A/.9
A:B/.1
A:A/.6
B:A/.8

A B + d e

B:B/.3
B:B/.2    A:B/.4    B:A/.7

(Pereira & Riley, 1997)

94

# Inference through string *cascades*

Given a string and a cascade, calculate the highest weighted transformation of the string by the cascade



(Pereira & Riley, 1997)

94

# Inference through string *cascades*

Given a string and a cascade, calculate the highest weighted transformation of the string by the cascade



1-BEST( A B + ... + ... ) = ?

(Pereira & Riley, 1997)

94

# Pipeline approach



1-BEST( A B + d e + f ) = ?

Embed the string

(Pereira & Riley, 1997)

95

# Pipeline approach



I-BEST( > [a] [b] ((c)) + > ((d)) ((e)) + > ((f)) ) = ?

A:A/1   B:B/1

A:B/.1   A:A/.9   A:A/.6   A:C/.6   A:D/.4
B:A/.8
B:B/.3
B:B/.2   A:B/.4   B:A/.7   B:C/.7   B:D/.3

## Embed the string

(Pereira & Riley, 1997)

96

# Pipeline approach



1-BEST( a b c + d e + f ) = ?

A:A/1   B:B/1

A:B/.1   A:A/.9   A:A/.6   A:C/.6   A:D/.4
B:A/.8
B:B/.3
B:B/.2   A:B/.4   B:A/.7   B:C/.7   B:D/.3

Compose the cascade

(Pereira & Riley, 1997)

97

# Pipeline approach



1-BEST(

A:B/.1
B:B/.2
B:A/.8
B:B/.3
A:A/.9
B:A/.7

ad  bd  cd  be  ce

+

A:C/.6  A:D/.4
B:C/.7  B:D/.3

f

) = ?

Compose the cascade

(Pereira & Riley, 1997)

98

# Pipeline approach

1-BEST(

A:C/.07
A:D/.03

B:C/.14
B:D/.06

bdf

B:C/.48
B:D/.32

cdf

adf

B:C/.21
B:D/.09

bef

cef

A:C/.54
A:D/.36

B:C/.42
B:D/.28

) = ?

Compose the cascade

(Pereira & Riley, 1997)

99

# Pipeline approach



1-BEST(                    ) = ?

Graph with nodes and edges:
- adf → bdf: C/.07 D/.03
- bdf → cdf: C/.14 D/.06
- bdf → cef: C/.48 D/.32
- adf → bef: C/.54 D/.36
- bef → cdf: C/.21 D/.09
- bef → cef: C/.42 D/.28

Project the range

(Pereira & Riley, 1997)

100

# Pipeline approach



1-BEST( ) = ?

Find the 1-best path of the result

(Dijkstra, 1959)

101

# Pipeline approach

1-BEST(  ) = ?

Find the 1-best path of the result

(Dijkstra, 1959)

# Pipeline approach



1-BEST( ) = ?

C/.07
D/.03

C/.14
D/.06

bdf

C/.48
D/.32

cdf

adf

C/.21
D/.09

bef

cef

**C/.54**
D/.36

**C/.42**
D/.28

Find the 1-best path of the result

(Dijkstra, 1959)

103

# Problems with pipeline

- Extra work done to create unused arcs

- Building done without input of all cascade members



104

# On-the-fly approach



1-BEST(  ) = ?

- Build arcs in result graph as needed
- All members of cascade "vote" simultaneously
- Less total construction cost

(Mohri, Pereira, Riley, 1999)

105

# On-the-fly approach



1-BEST(    adf  ) = ?

- Build arcs in result graph as needed
- All members of cascade "vote" simultaneously
- Less total construction cost

(Mohri, Pereira, Riley, 1999)

106

# On-the-fly approach



- Build arcs in result graph as needed
- All members of cascade "vote" simultaneously
- Less total construction cost

(Mohri, Pereira, Riley, 1999)

107

# On-the-fly approach



1-BEST(  ) = ?

- Build arcs in result graph as needed

- All members of cascade "vote" simultaneously

- Less total construction cost

(Mohri, Pereira, Riley, 1999)

108

# On-the-fly approach



- Build arcs in result graph as needed

- All members of cascade "vote" simultaneously

- Less total construction cost

(Mohri, Pereira, Riley, 1999)

109

# On-the-fly approach



1-BEST(  ) = ?

- Build arcs in result graph as needed
- All members of cascade "vote" simultaneously
- Less total construction cost

(Mohri, Pereira, Riley, 1999)

110

# On-the-fly approach

A:A/1  B:B/1
a → b → c

A:B/.1  A:A/.9  A:A/.6
B:A/.8
d  e
B:B/.3
B:B/.2  A:B/.4  B:A/.7

A:C/.6  A:D/.4
f
B:C/.7  B:D/.3

C/.07
D/.03
bdf

1-BEST(  adf  ) = ?

C/.54
D/.36

bef

C/.42

cef

• Build arcs in result graph as needed

• All members of cascade "vote" simultaneously

• Less total construction cost

(Mohri, Pereira, Riley, 1999)

111

# On-the-fly approach



1-BEST( ) = ?

- Build arcs in result graph as needed

- All members of cascade "vote" simultaneously

- Less total construction cost
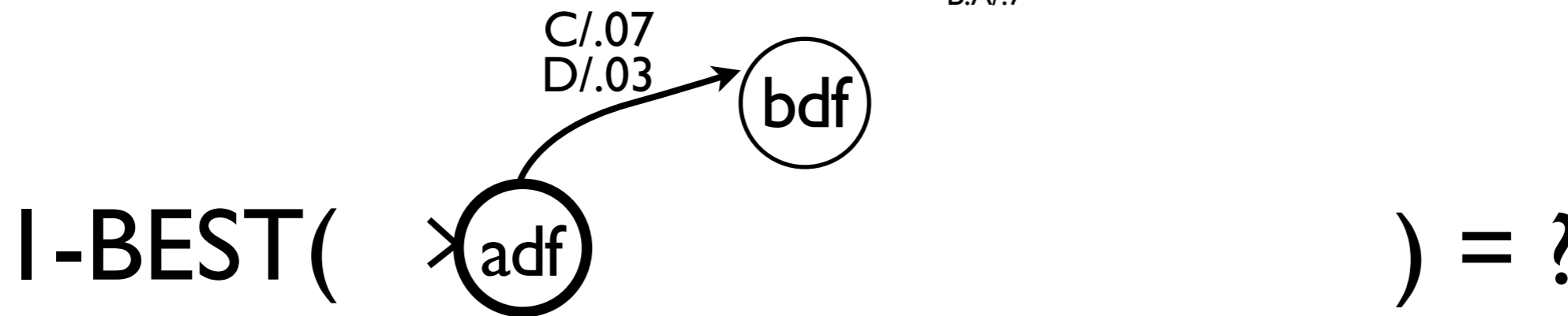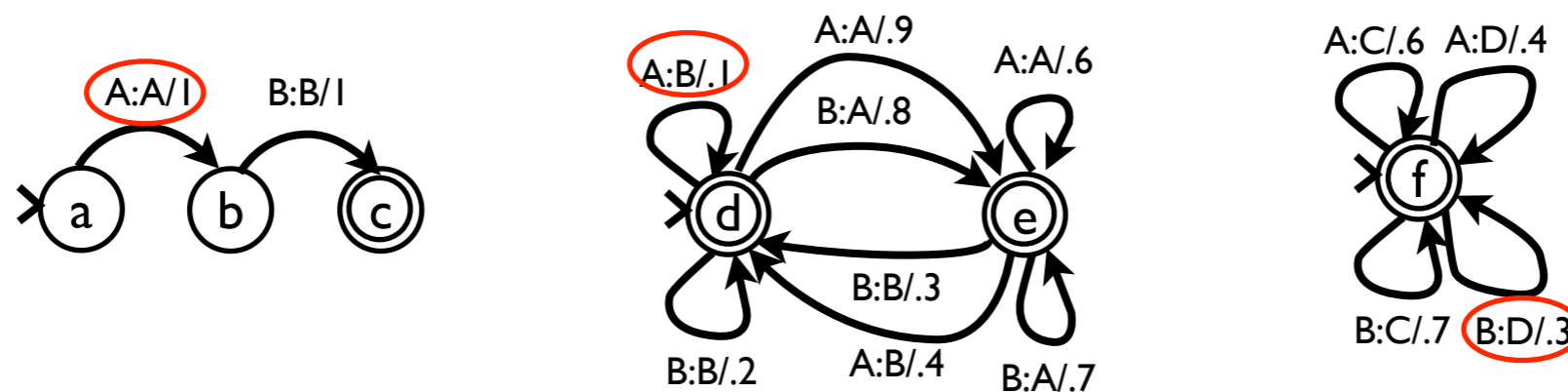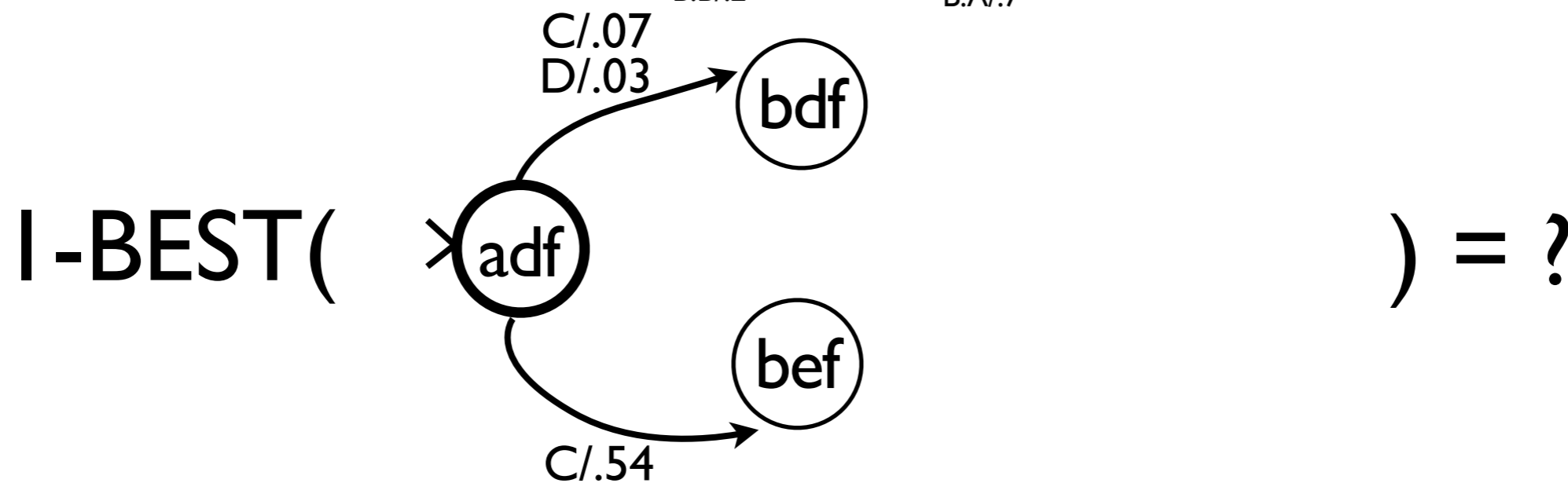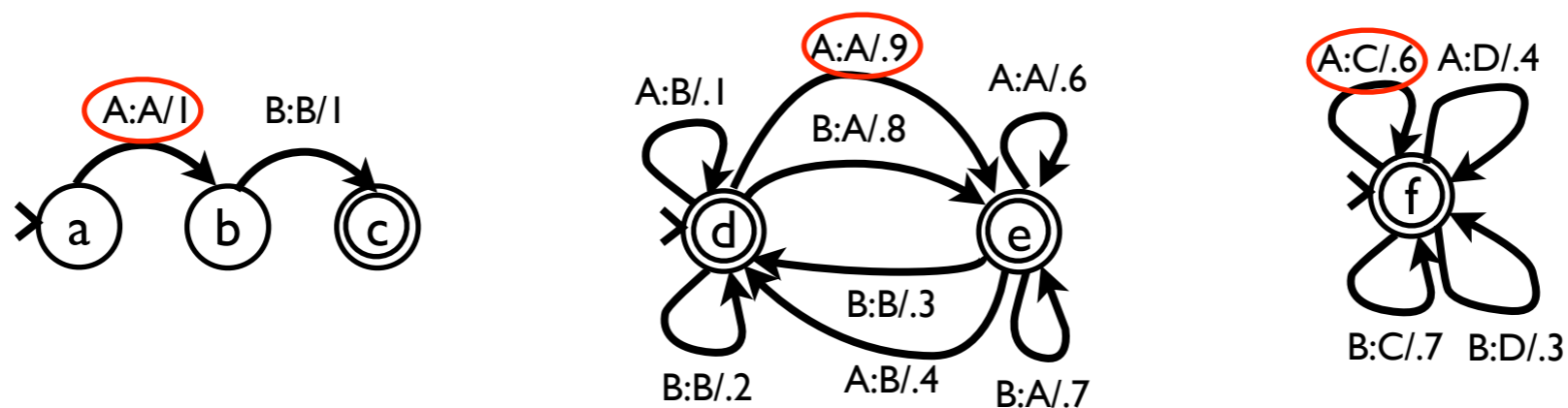
(Mohri, Pereira, Riley, 1999)

112

# On-the-fly approach



- Build arcs in result graph as needed

- All members of cascade "vote" simultaneously

- Less total construction cost

(Mohri, Pereira, Riley, 1999)

113

# Inference through *tree* cascades?

- In general, tree transducers are *not closed* under composition

- However, some classes are closed, and by adding additional steps to the process, we can conduct inference

- We provide pipeline and on-the-fly algorithms for applicable classes of weighted tree transducers

114

# Inference through *tree* cascades

Given a tree and a cascade, calculate the highest weighted transformation of the tree by the cascade

115

# Inference through *tree* cascades

Given a tree and a cascade, calculate the highest weighted transformation of the tree by the cascade



115

# Inference through *tree* cascades

Given a tree and a cascade, calculate the highest weighted transformation of the tree by the cascade



115

# Inference through *tree* cascades

Given a tree and a cascade, calculate the highest weighted transformation of the tree by the cascade



115

# Inference through *tree* cascades

Given a tree and a cascade, calculate the highest weighted transformation of the tree by the cascade



115

# Pipeline approach

# Pipeline approach



Embed the tree

117

# Pipeline approach



1-BEST( ... ) = ?

Compose adjacent transducers

118

# Pipeline approach

# Pipeline approach



New step!

Embed the grammar

# Pipeline approach



I-BEST( ... ) = ?

Embed the grammar

Identity transducer has more composition cases

# Pipeline approach



I-BEST(        ) = ?

Compose adjacent transducers

122

# Pipeline approach

1-BEST(



) = ?

Project the range

123

# Pipeline approach



.016

Find 1-best path of the result

(Knuth '77)

124

# On-the-fly approach



1-BEST( ) = ?

125

# On-the-fly approach



1-BEST(  ) = ?

126

# On-the-fly approach

# On-the-fly approach

# On-the-fly approach

# On-the-fly vs. pipeline

| Language model | Re-order English words | Insert function words | Translate to Japanese |

← 

- We recovered 1-best English tree through this cascade

- We calculated time to complete for several language models and both pipeline and on-the-fly methods

- On-the-fly was much faster and in some cases the only method that worked in the memory allotted

(Yamada & Knight, 2001)

130

# On-the-fly vs. pipeline

| language model | method | time/sentence |
|---|---|---|
| weak | pipeline<br>on-the-fly | 28s<br>17s |
| strong & large | pipeline<br>on-the-fly | >60s*<br>24s |
| strong & small | pipeline<br>on-the-fly | 2.5s<br>.06s |

\* Ran out of memory before completing

131

# Extension for tree-string transducers

What if the cascade ends in a tree-string transducer, and we want to pass a string through the cascade?

tree-tree     tree-tree     tree-**string**

```
        VB
      /  |  \
   PRP  VB   NN
    |    |    |
    i  hate snakes
```

← | Re-order English words | ← | Insert function words | ← | Translate to Japanese | ← ヘビ が 大嫌い だ

132

# Extension for tree-string transducers

What if the cascade ends in a tree-string transducer, and we want to pass a string through the cascade?

tree-tree      tree-tree      tree-**string**

VB

PRP  VB  NN

i  hate  snakes

| Re-order English words | ← | Insert function words | ← | Translate to Japanese | ← ヘビ が 大嫌い だ |

cfg parsing

132

# Extension for tree-string transducers

What if the cascade ends in a tree-string transducer, and we want to pass a string through the cascade?

tree-tree    tree-tree    tree-**string**

VB
PRP   VB   NN
i    hate  snakes

Re-order English words ← Insert function words ← Translate to Japanese ← ヘビ が 大嫌い だ

on-the-fly or pipeline

cfg parsing

132

# A weighted tree automata and transducer toolkit
## (May & Knight, CIAA '06)

- Operations for inference, manipulation, and training of tree transducers and automata

- Very easy to experiment quickly, without coding

- http://www.isi.edu/licensed-sw/tiburon

133

# Tiburon example 1: syntax MT cascade

Simplified English trees to Japanese strings

```
        TOP
         |
        VB
      / | \
   PRP  VB  NN
    |   |    |
    i  hate snakes
```

↔   ヘビ が 大嫌い だ

e-tree ↔ rotate ↔ insert ↔ translate ↔ j-string

(Yamada & Knight, 2001)

134

# Tiburon example I: syntax MT cascade

## 1) Rotate children



e-tree ↔ **rotate** ↔ insert ↔ translate ↔ j-string

(Yamada & Knight, 2001)

135

# Tiburon example I: syntax MT cascade

## 2) Insert function words



e-tree ↔ | rotate | ↔ | insert | ↔ | translate | ↔ j-string

(Yamada & Knight, 2001)

136

# Tiburon example 1: syntax MT cascade

## 3) Translate leaves

$$\overset{vb}{\bigcirc}\ hate\ \overset{.25}{\longrightarrow}\ 大嫌い$$

e-tree ↔ [ rotate ] ↔ [ insert ] ↔ [ translate ] ↔ j-string

(Yamada & Knight, 2001)

137

# Tiburon example 1: syntax MT cascade

- Task: Decode candidate sentence, get top 5 answers

- Algorithms used: inference through cascade, k-best, determinization

Candidate:　彼ら は 偽善 が 大嫌い だ

Correct answer:

TOP(VB(NN("hypocrisy") VB("is") JJ(JJ("abhorrent") TO(TO("to") PRP("them")))))

138

# Tiburon example 1: syntax MT cascade

Let's try it!

% tiburon  -k 5  -m tropical  -e euc-jp rot ins trans ej.1.f

139

# Tiburon example 1: syntax MT cascade

Let's try it!

% tiburon   -k 5  -m tropical  -e euc-jp rot ins trans  ej.1.f

program

139

# Tiburon example 1: syntax MT cascade

Let's try it!

% tiburon  <span style="color:red">-k 5</span>  -m tropical  -e euc-jp rot ins trans ej.1.f

↑

| 5 best |
|--------|

139

# Tiburon example 1: syntax MT cascade

Let's try it!

% tiburon  -k 5  -m tropical  -e euc-jp rot ins trans ej.1.f

semiring

139

# Tiburon example 1: syntax MT cascade

Let's try it!

% tiburon  -k 5  -m tropical  **-e euc-jp** rot ins trans  ej.1.f

character set

139

# Tiburon example 1: syntax MT cascade

Let's try it!

% tiburon  -k 5  -m tropical  -e euc-jp  rot ins trans  ej.1.f

cascade

139

# Tiburon example 1: syntax MT cascade

Let's try it!

% tiburon  -k 5  -m tropical  -e euc-jp rot ins trans <span style="color:red">ej.1.f</span>

input

139

# Tiburon example 1: syntax MT cascade

## First try is not so good!

% tiburon -k 5 -m tropical -e euc-jp rot ins trans ej.1.f
TOP(VB(PRP("him") VB("abominate") IN(IN("above") NN(JJ("abhorrent") NN("fanatic")))))  # 18.368
TOP(VB(PRP("them") VB("abominate") IN(IN("above") NN(JJ("abhorrent") NN("fanatic")))))  # 18.368
TOP(VB(PRP("him") VB("abominate") IN(IN("above") NN(JJ("abhorrent") NN("hypocrisy")))))  # 18.368
TOP(VB(PRP("them") VB("abominate") IN(IN("above") NN(JJ("abhorrent") NN("hypocrisy")))))  # 18.368
TOP(VB(PRP("him") VB("abominate") IN(IN("above") NN(JJ("abhorrent") NN("clouds")))))  # 18.368

140

# Tiburon example 1: syntax MT cascade

## Add in a simple PCFG-based language model

```
      IN
     /  \
   IN    NN
   /\    /\
  /__\  /__\
```

141

# Tiburon example 1: syntax MT cascade

## Add in a simple PCFG-based language model



141

# Tiburon example I: syntax MT cascade

## Add in a simple PCFG-based language model



% tiburon -k 5 -m tropical -e euc-jp pcfg-lm rot ins trans ej.1.f
TOP(VB(PRP("i") VB("abominate") JJ(JJ("abhorrent") TO(TO("to") PRP("i")))) # 33.024
TOP(VB(PRP("i") VB("is") JJ(JJ("abhorrent") TO(TO("to") PRP("i")))) # 33.718
TOP(VB(PRP("him") VB("abominate") JJ(JJ("abhorrent") TO(TO("to") PRP("i")))) # 33.718
TOP(VB(PRP("i") VB("abominate") JJ(JJ("abhorrent") TO(TO("to") PRP("him")))) # 33.718
TOP(VB(PRP("them") VB("abominate") JJ(JJ("abhorrent") TO(TO("to") PRP("i")))) # 33.718

141

# Tiburon example 1: syntax MT cascade

## Try a grandparent language model

```
        VB
       /  \
      IN    △
     /  \
   IN    NN
   △     △
```

142

# Tiburon example 1: syntax MT cascade

## Try a grandparent language model

# Tiburon example 1: syntax MT cascade

## Try a grandparent language model



% tiburon -k 5 -m tropical -e euc-jp gp-lm rot ins trans ej.1.f
TOP(VB(PRP("i") VB("abominate") JJ(JJ("abhorrent") TO(TO("to") PRP("them"))))) # 26.603
TOP(VB(PRP("i") VB("is") JJ(JJ("abhorrent") TO(TO("to") PRP("them"))))) # 27.297
TOP(VB(NN("hypocrisy") VB("abominate") JJ(JJ("abhorrent") TO(TO("to") PRP("them"))))) # 28.033
TOP(VB(PRP("i") VB("abominate") JJ(JJ("abhorrent") TO(TO("to") PRP("them"))))) # 28.071
TOP(VB(NN("hypocrisy") VB("is") JJ(JJ("abhorrent") TO(TO("to") PRP("them"))))) # 28.726

142

# Tiburon example 1: syntax MT cascade

## Try a grandparent language model



% tiburon -k 5 -m tropical -e euc-jp gp-lm rot ins trans ej.1.f

TOP(VB(PRP("i") VB("abominate") JJ(JJ("abhorrent") TO(TO("to") PRP("them"))))) # 26.603

TOP(VB(PRP("i") VB("is") JJ(JJ("abhorrent") TO(TO("to") PRP("them"))))) # 27.297

TOP(VB(NN("hypocrisy") VB("abominate") JJ(JJ("abhorrent") TO(TO("to") PRP("them"))))) # 28.033

TOP(VB(PRP("i") VB("abominate") JJ(JJ("abhorrent") TO(TO("to") PRP("them"))))) # 28.071

TOP(VB(NN("hypocrisy") VB("is") JJ(JJ("abhorrent") TO(TO("to") PRP("them"))))) # 28.726

### Correct sentence is 5th

142

# Tiburon example I: syntax MT cascade

## Try a grandparent language model

VB

IN

IN    NN

0.667

vb in → IN

in in    in nn

Duplicates

% tiburon -k 5 -m tropical -e euc-jp gp-lm rot ins trans ej.1.f
TOP(VB(PRP("i") VB("abominate") JJ(JJ("abhorrent") TO(TO("to") PRP("them"))))) # 26.603
TOP(VB(PRP("i") VB("is") JJ(JJ("abhorrent") TO(TO("to") PRP("them"))))) # 27.297
TOP(VB(NN("hypocrisy") VB("abominate") JJ(JJ("abhorrent") TO(TO("to") PRP("them"))))) # 28.033
TOP(VB(PRP("i") VB("abominate") JJ(JJ("abhorrent") TO(TO("to") PRP("them"))))) # 28.071
TOP(VB(NN("hypocrisy") VB("is") JJ(JJ("abhorrent") TO(TO("to") PRP("them"))))) # 28.726

Correct sentence is 5th

142

# Tiburon example 1: syntax MT cascade

- Combine duplicate derivations in entire search space using *weighted determinization*

143

# Tiburon example 1: syntax MT cascade

- Combine duplicate derivations in entire search space using *weighted determinization*

% tiburon -d 5 -k 5 -m tropical -e euc-jp gp-lm rot ins trans ej.1.f
TOP(VB(PRP("i") VB("abominate") JJ(JJ("abhorrent") TO(TO("to") PRP("them")))) # 26.329
TOP(VB(PRP("i") VB("is") JJ(JJ("abhorrent") TO(TO("to") PRP("them")))) # 27.023
TOP(VB(NN("hypocrisy") VB("abominate") JJ(JJ("abhorrent") TO(TO("to") PRP("them")))) # 27.759
TOP(VB(NN("hypocrisy") VB("is") JJ(JJ("abhorrent") TO(TO("to") PRP("them")))) # 28.452
TOP(VB(NN(DT("a") NN("clouds")) VB("abominate") JJ(JJ("abhorrent") TO(TO("to") PRP("them")))) # 31.250

143

# Tiburon example 1: syntax MT cascade

- Combine duplicate derivations in entire search space using *weighted determinization*

Now we're 4th

% tiburon -d 5 -k 5 -m tropical -e euc-jp gp-lm rot ins trans ej.1.f
TOP(VB(PRP("i") VB("abominate") JJ(JJ("abhorrent") TO(TO("to") PRP("them"))))) # 26.329
TOP(VB(PRP("i") VB("is") JJ(JJ("abhorrent") TO(TO("to") PRP("them"))))) # 27.023
TOP(VB(NN("hypocrisy") VB("abominate") JJ(JJ("abhorrent") TO(TO("to") PRP("them"))))) # 27.759
TOP(VB(NN("hypocrisy") VB("is") JJ(JJ("abhorrent") TO(TO("to") PRP("them"))))) # 28.452
TOP(VB(NN(DT("a") NN("clouds")) VB("abominate") JJ(JJ("abhorrent") TO(TO("to") PRP("them"))))) # 31.250

143

# Tiburon example 2: training a syntax LM

- The LMs we used before had no hidden states

- Let's introduce hidden states and learn weights with EM



(Petrov & Klein, '07)

144

# Tiburon example 2: training a syntax LM

- The LMs we used before had no hidden states

- Let's introduce hidden states and learn weights with EM



(Petrov & Klein, '07)

144

# Tiburon example 2: training a syntax LM

- The LMs we used before had no hidden states

- Let's introduce hidden states and learn weights with EM



(Petrov & Klein, '07)

144

# Tiburon example 2: training a syntax LM

% tiburon  -t 50  --randomize  trees  rtg.4split > 4split-lm

145

# Tiburon example 2: training a syntax LM

% tiburon  **-t 50**  --randomize  trees  rtg.4split > 4split-lm

50 iterations

145

# Tiburon example 2: training a syntax LM

% tiburon  -t 50  **--randomize**  trees  rtg.4split > 4split-lm

random initial weights avoids saddles

145

# Tiburon example 2: training a syntax LM

% tiburon  -t 50  --randomize  trees  rtg.4split > 4split-lm

training data

145

# Tiburon example 2: training a syntax LM

% tiburon  -t 50  --randomize  trees  rtg.4split  > 4split-lm

4-way split

145

# Tiburon example 2: training a syntax LM

% tiburon -t 50 --randomize trees rtg.4split > 4split-lm

Cross entropy with normalized initial weights is 1.868; corpus prob is e^-269.025

Cross entropy after 1 iterations is 1.190; corpus prob is e^-171.383

Cross entropy after 2 iterations is 1.138; corpus prob is e^-163.866

Cross entropy after 3 iterations is 1.036; corpus prob is e^-149.229

...

Cross entropy after 47 iterations is 0.581; corpus prob is e^-83.665

Cross entropy after 48 iterations is 0.581; corpus prob is e^-83.665

Cross entropy after 49 iterations is 0.581; corpus prob is e^-83.665

146

# Tiburon example 2: training a syntax LM

% tiburon -t 50 --randomize trees rtg.4split > 4split-lm
Cross entropy with normalized initial weights is 1.868; corpus prob is e^-269.025
Cross entropy after 1 iterations is 1.190; corpus prob is e^-171.383
Cross entropy after 2 iterations is 1.138; corpus prob is e^-163.866
Cross entropy after 3 iterations is 1.036; corpus prob is e^-149.229
...
Cross entropy after 47 iterations is 0.581; corpus prob is e^-83.665
Cross entropy after 48 iterations is 0.581; corpus prob is e^-83.665
Cross entropy after 49 iterations is 0.581; corpus prob is e^-83.665

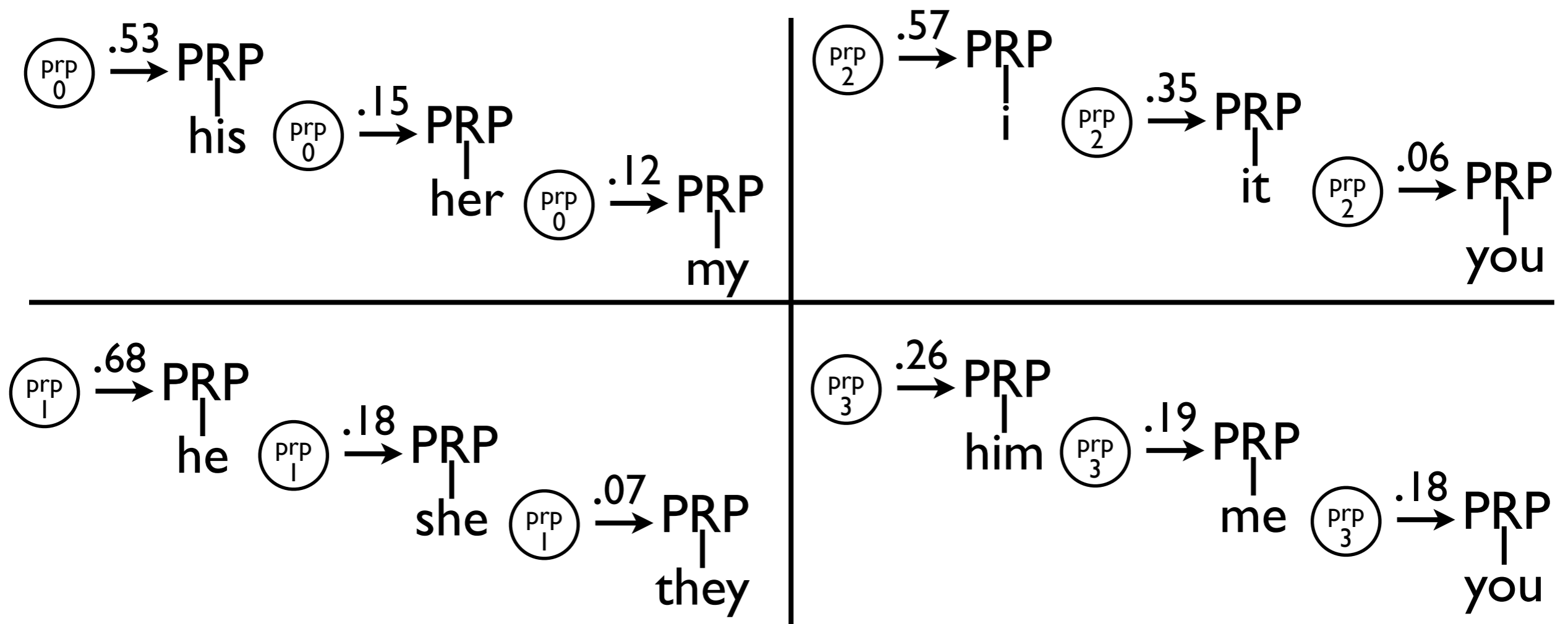## Compare with GP-PCFG

% tiburon -t 3 --randomize trees rtg.gp.pcfg > lm
Cross entropy with normalized initial weights is 0.827; corpus prob is e^-119.022
Cross entropy after 1 iterations is 0.732; corpus prob is e^-105.448
Cross entropy after 2 iterations is 0.732; corpus prob is e^-105.448

146

# Tiburon example 2: training a syntax LM

## We can subjectively see state specialization

$\text{prp}_0 \xrightarrow{.53} \text{PRP}$
his

$\text{prp}_0 \xrightarrow{.15} \text{PRP}$
her

$\text{prp}_0 \xrightarrow{.12} \text{PRP}$
my

$\text{prp}_2 \xrightarrow{.57} \text{PRP}$
i

$\text{prp}_2 \xrightarrow{.35} \text{PRP}$
it

$\text{prp}_2 \xrightarrow{.06} \text{PRP}$
you

$\text{prp}_1 \xrightarrow{.68} \text{PRP}$
he

$\text{prp}_1 \xrightarrow{.18} \text{PRP}$
she

$\text{prp}_1 \xrightarrow{.07} \text{PRP}$
they

$\text{prp}_3 \xrightarrow{.26} \text{PRP}$
him

$\text{prp}_3 \xrightarrow{.19} \text{PRP}$
me

$\text{prp}_3 \xrightarrow{.18} \text{PRP}$
you

147

# Tiburon example 2: training a syntax LM

Tied for first!

% tiburon -k 5 -m tropical -e euc-jp 4split-lm rot ins trans ej.1.f
TOP(VB(NN("hypocrisy") VB("is") JJ(JJ("abhorrent") TO(TO("to") PRP("them")))))  # 29.556
TOP(VB(NN("fanatic") VB("is") JJ(JJ("abhorrent") TO(TO("to") PRP("them")))))  # 29.556
TOP(VB(NN("clouds") VB("is") JJ(JJ("abhorrent") TO(TO("to") PRP("them")))))  # 29.556
TOP(VB(NN("fanatic") VB("is") JJ(JJ("abhorrent") TO(TO("to") PRP("them")))))  # 29.717
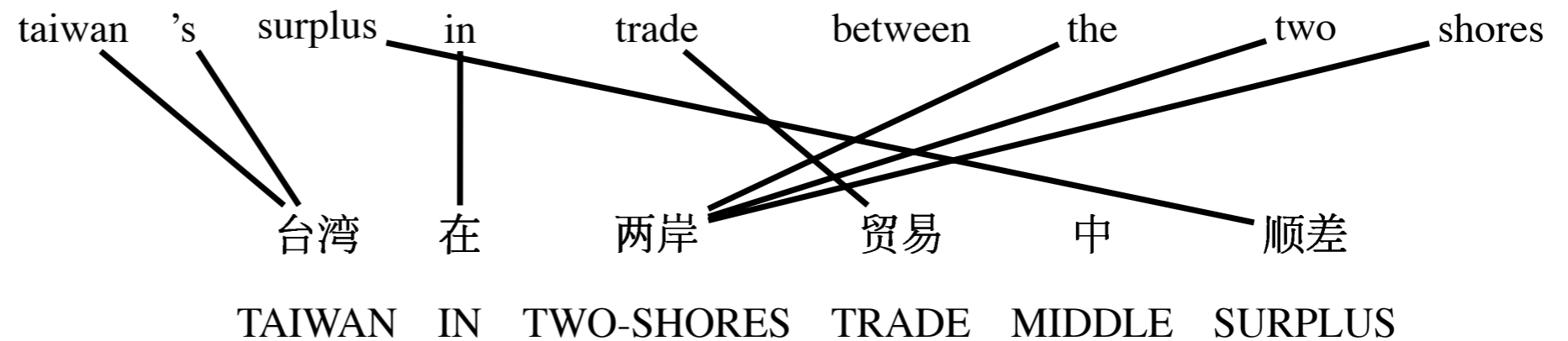TOP(VB(NN("hypocrisy") VB("is") JJ(JJ("abhorrent") TO(TO("to") PRP("them")))))  # 29.717

148

# Using tree transducers to improve machine translation

## (May & Knight, EMNLP '07)

- We will now shift focus to improving state-of-the-art syntax MT results

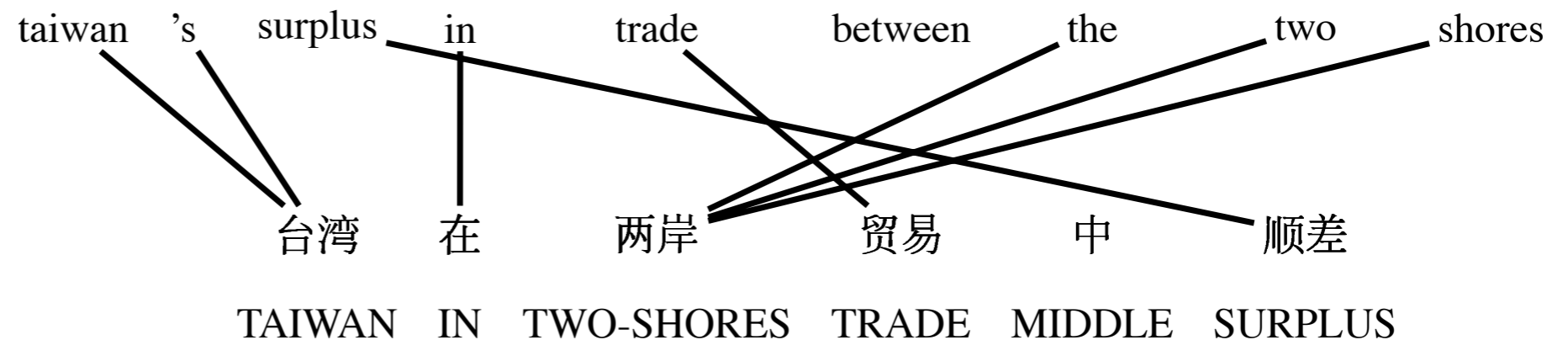- At core, we're using the power of training tree transducers to achieve gains

149

# Extracting syntactic rules
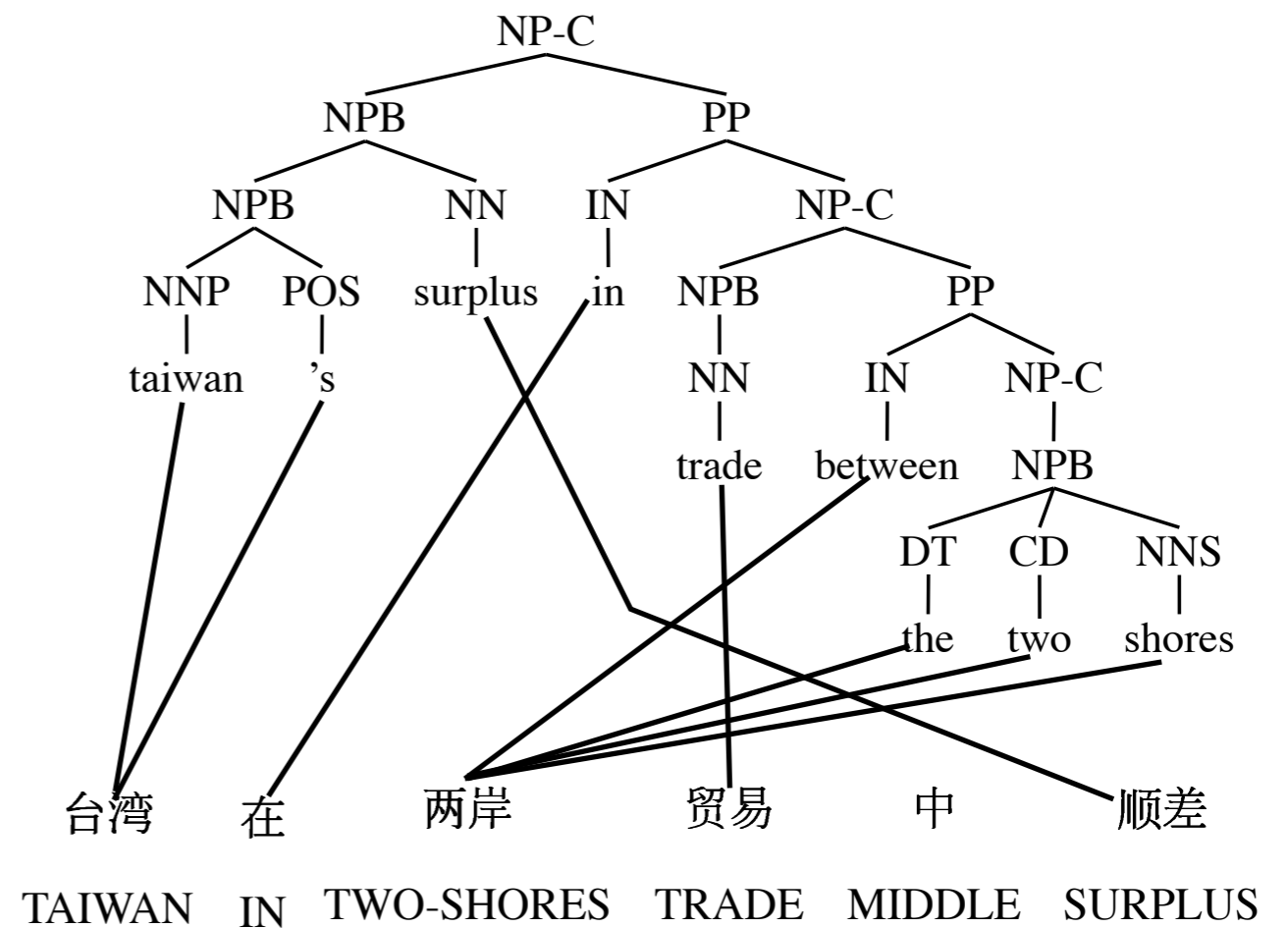
## 1) Obtain alignments



taiwan    's    surplus    in    trade    between    the    two    shores

台湾    在    两岸    贸易    中    顺差

TAIWAN    IN    TWO-SHORES    TRADE    MIDDLE    SURPLUS

# Extracting syntactic rules

1) Obtain alignments



taiwan  's  surplus  in  trade  between  the  two  shores

台湾  在  两岸  贸易  中  顺差

TAIWAN  IN  TWO-SHORES  TRADE  MIDDLE  SURPLUS

(Galley et al. '04, '06)

150

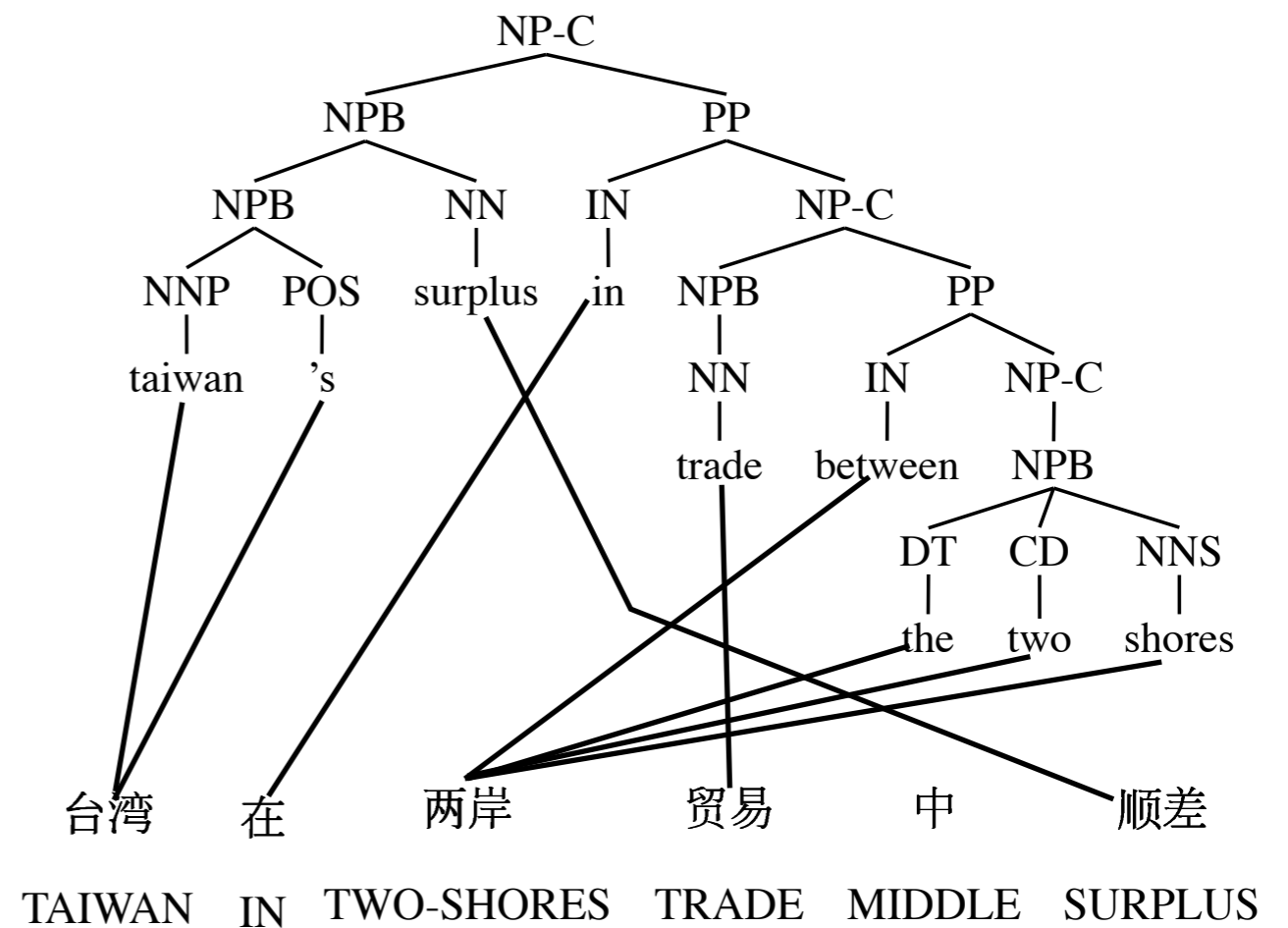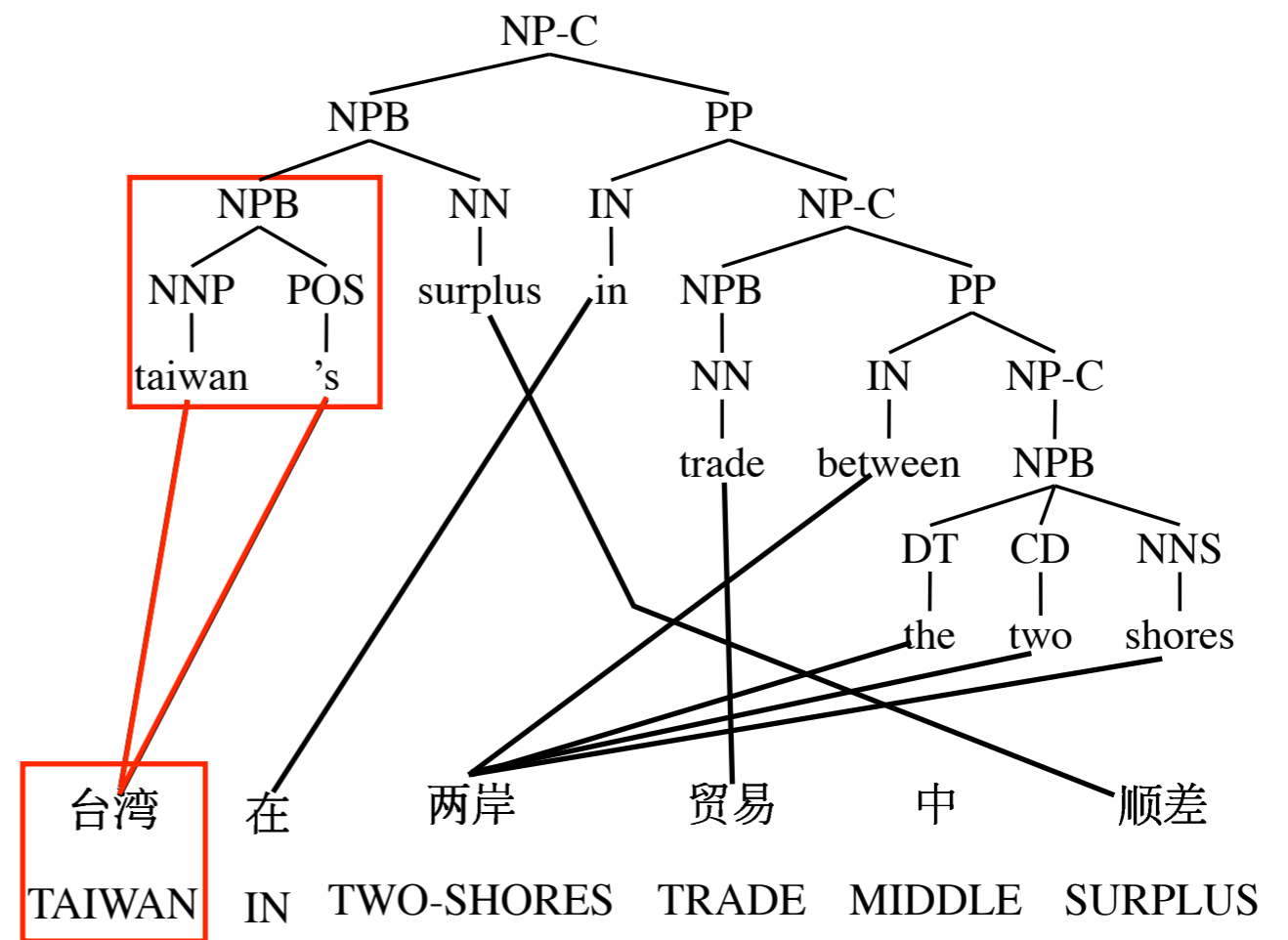# Extracting syntactic rules

## 2) Add parse tree



(Galley et al. '04, '06)

151

# Extracting syntactic rules

## 3) Extract rules


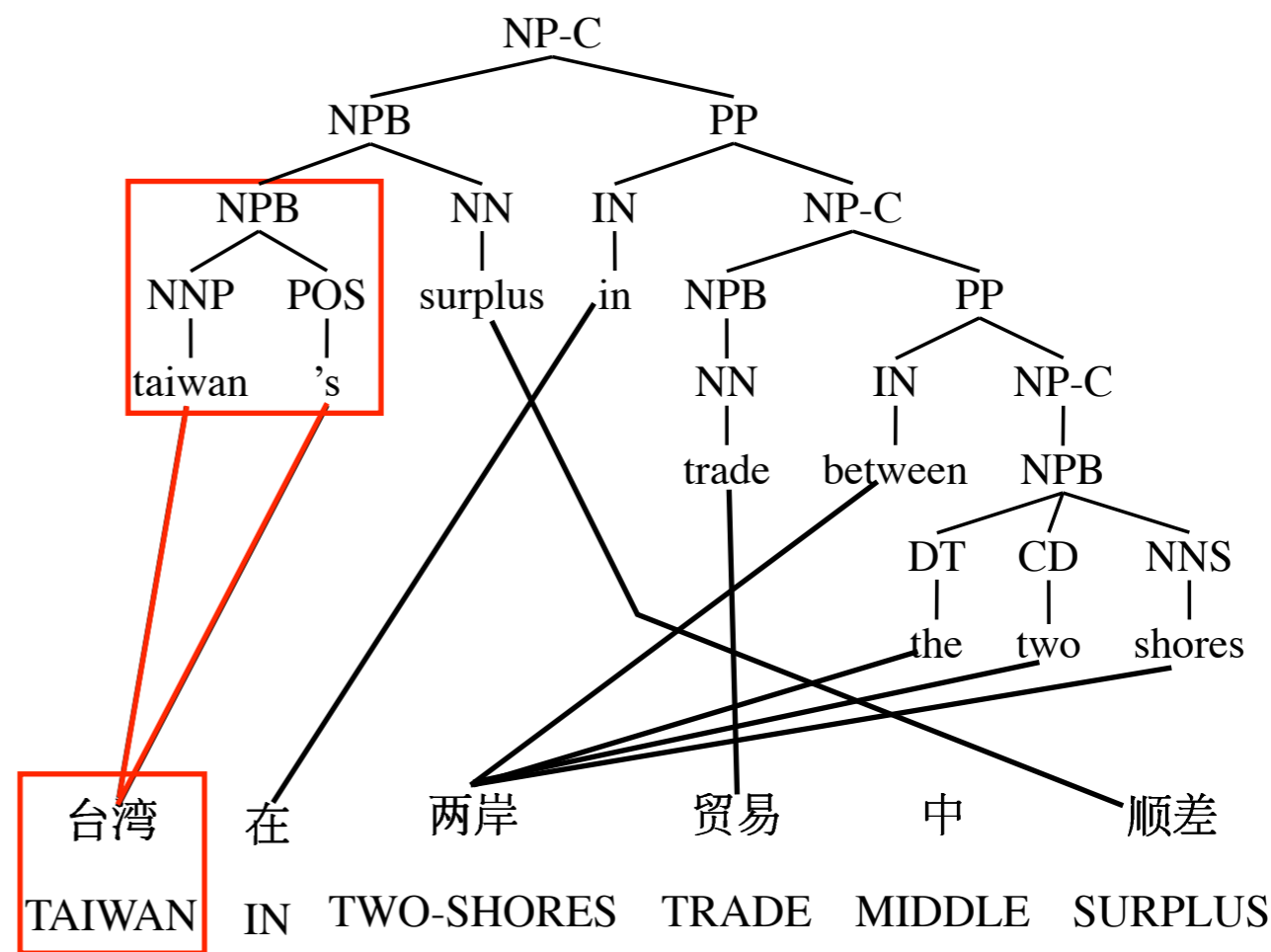
(Galley et al. '04, '06)

151
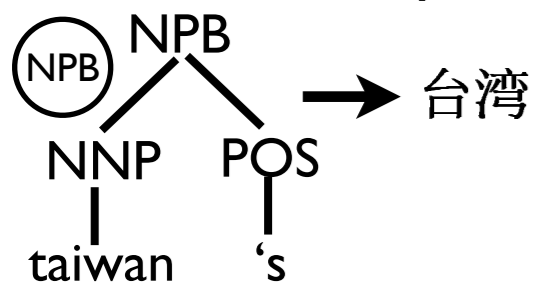
# Extracting syntactic rules

## 3) Extract rules



(Galley et al. '04, '06)

151

# Extracting syntactic rules

3) Extract rules



(Galley et al. '04, '06)
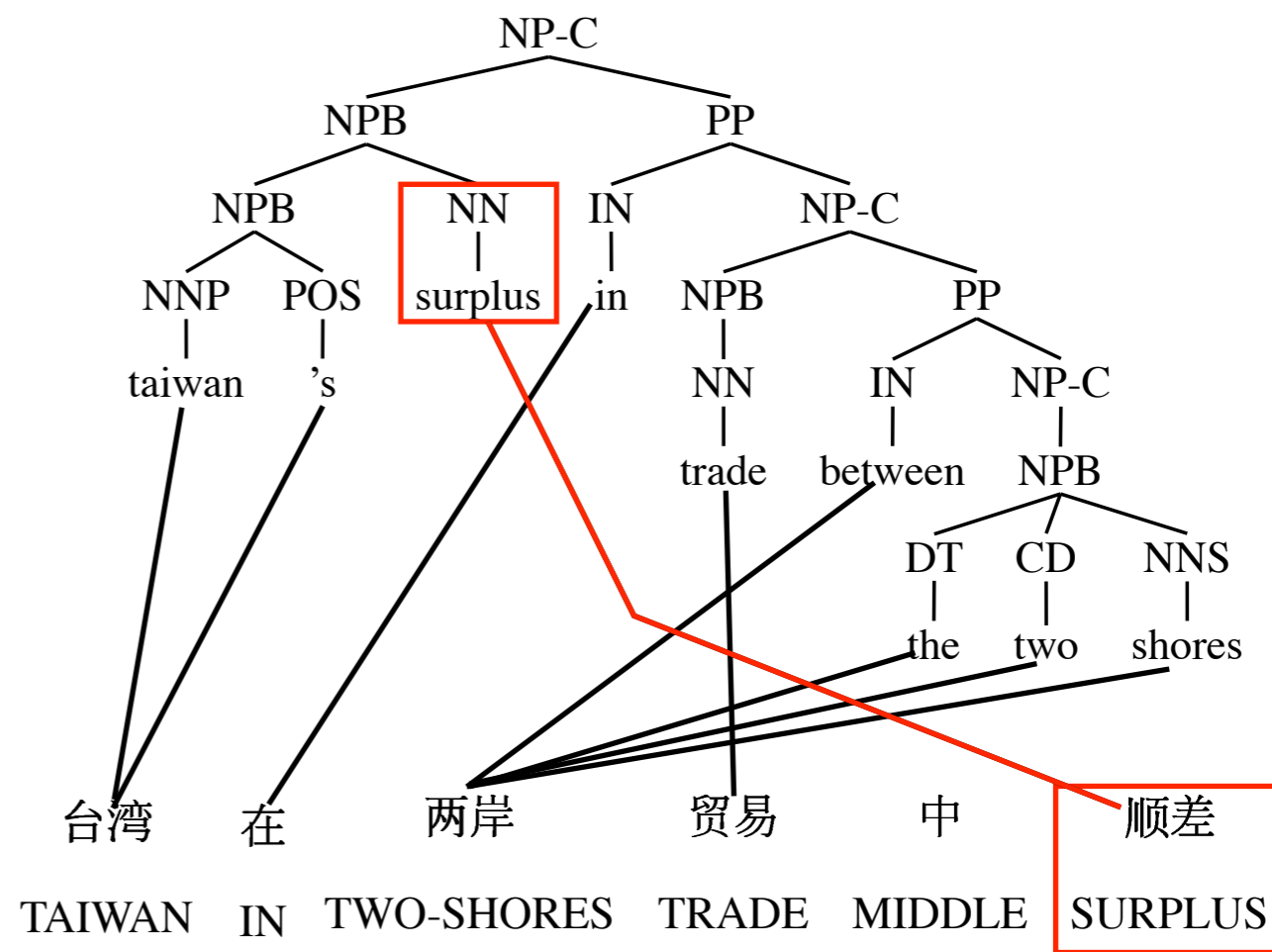
151

# Extracting syntactic rules
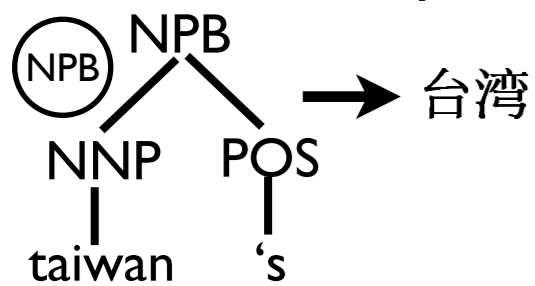
## 3) Extract rules



(Galley et al. '04, '06)

151

# Extracting syntactic rules

## 3) Extract rules



(Galley et al. '04, '06)

151

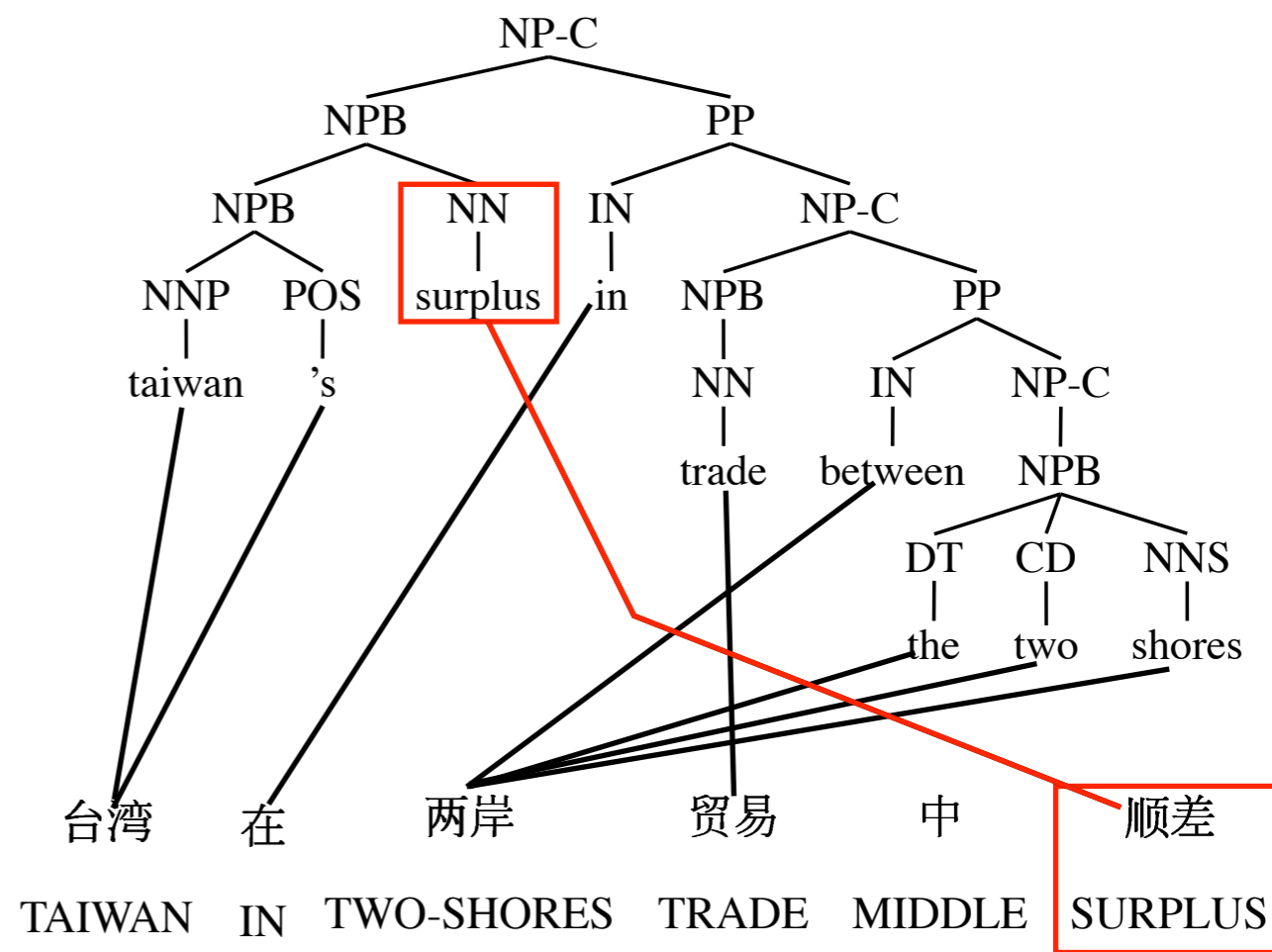# Extracting syntactic rules

## 3) Extract rules



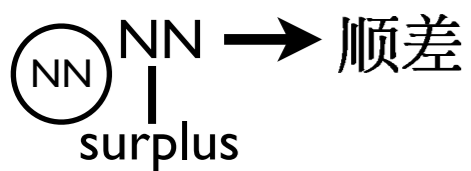(Galley et al. '04, '06)

151

# Extracting syntactic rules

## 3) Extract rules



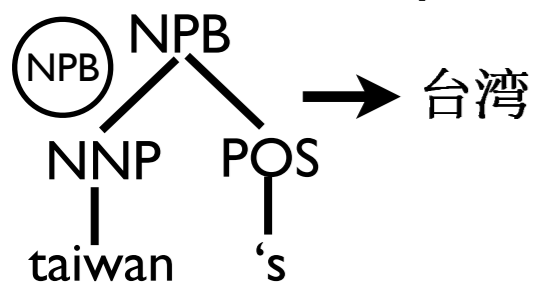(Galley et al. '04, '06)

151

# Extracting syntactic rules

## 3) Extract rules



(Galley et al. '04, '06)

151

# Extracting syntactic rules

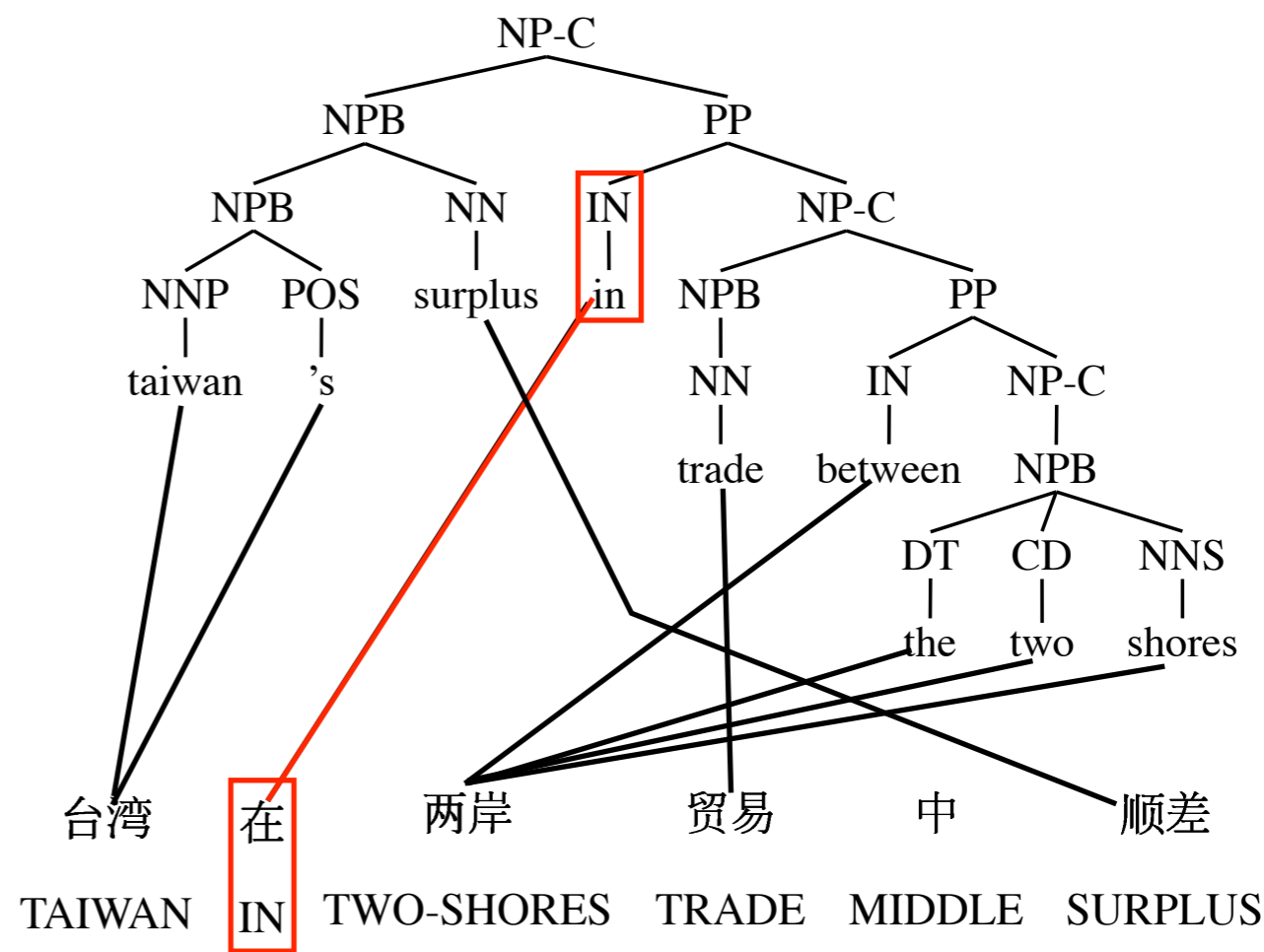## 3) Extract rules



(Galley et al. '04, '06)

151

# Extracting syntactic rules

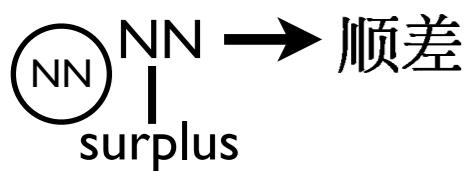## 3) Extract rules



(Galley et al. '04, '06)

151

# Extracting syntactic rules

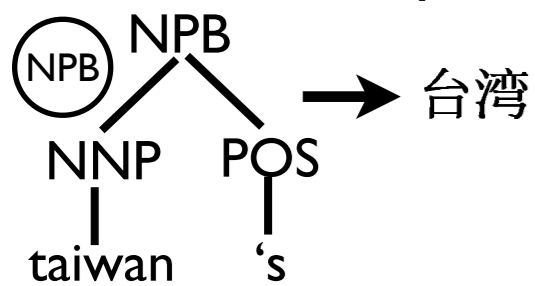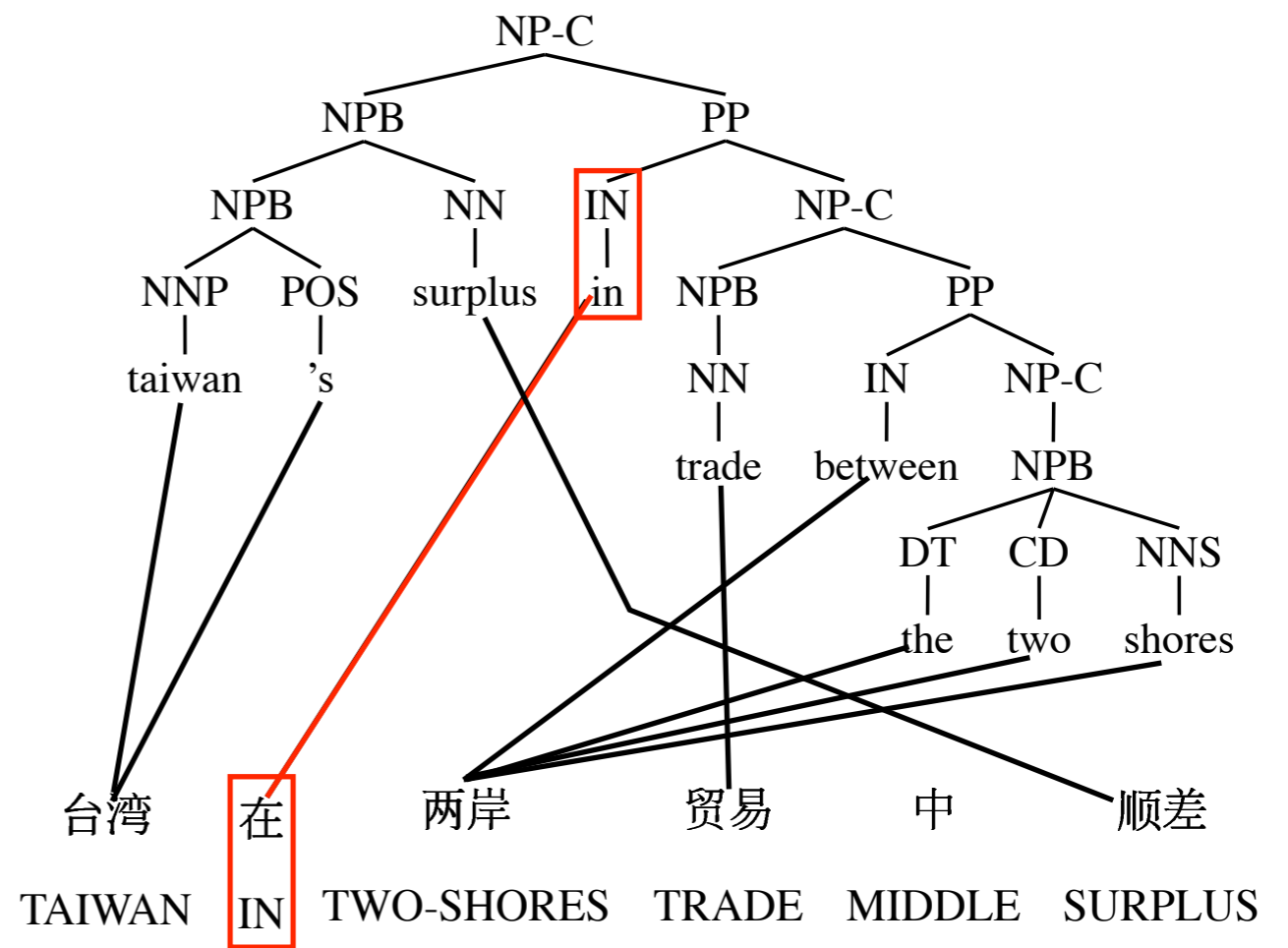## 3) Extract rules



(Galley et al. '04, '06)

151

# Extracting syntactic rules

## 3) Extract rules



(Galley et al. '04, '06)

151

# Extracting syntactic rules

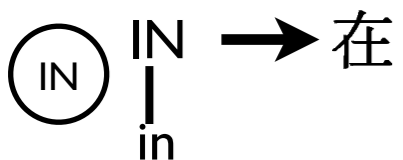## 3) Extract rules



(Galley et al. '04, '06)

151

# Extracting syntactic rules

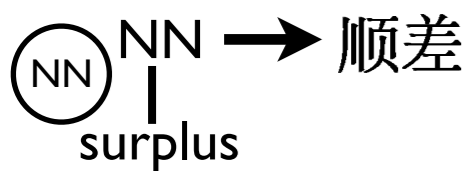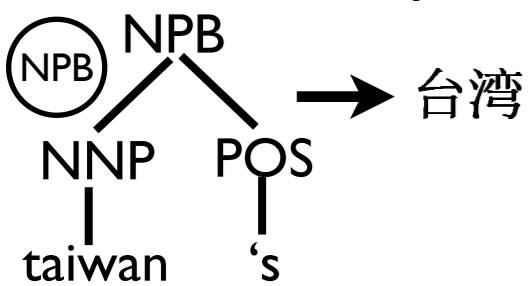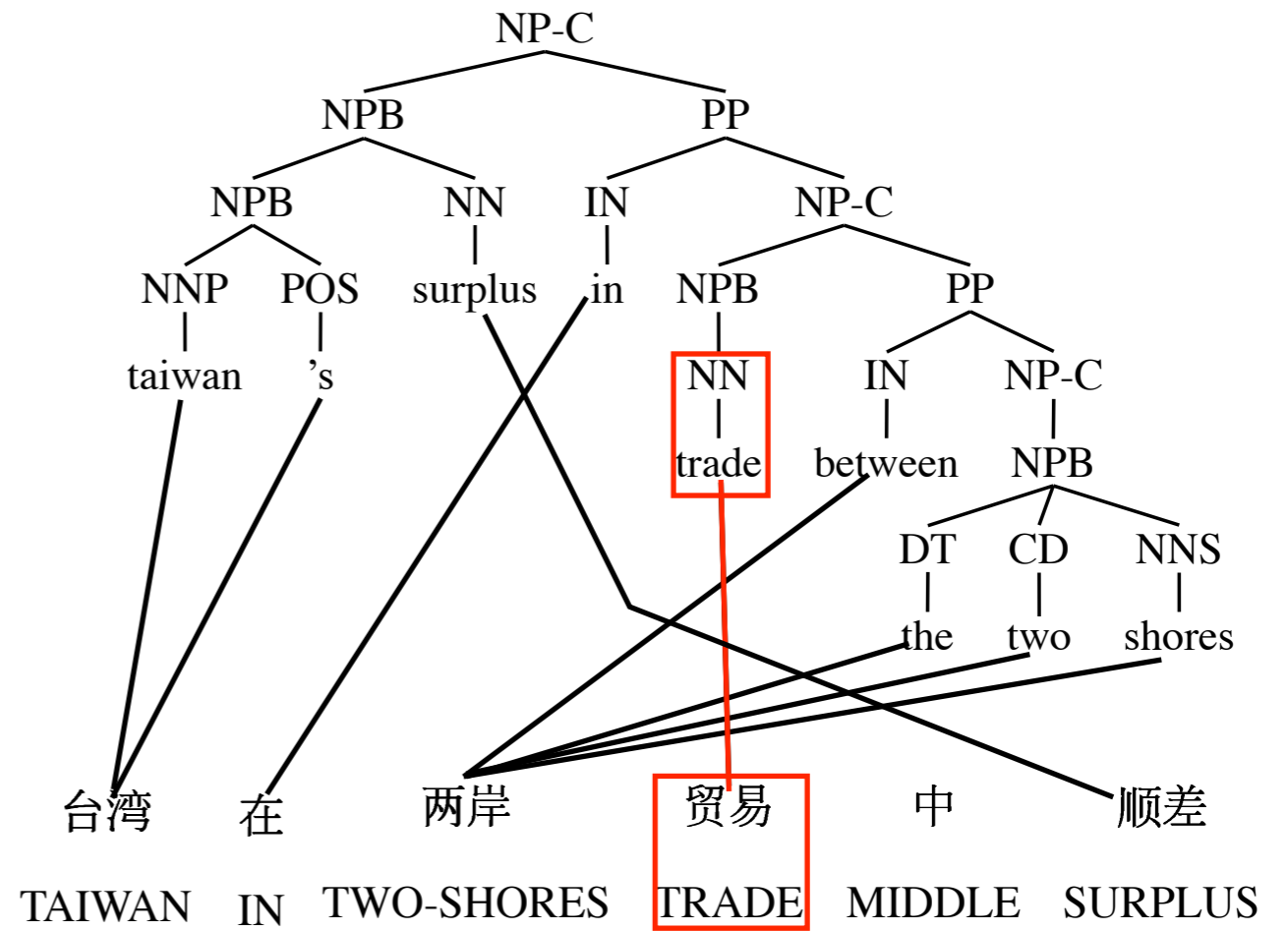## 3) Extract rules



(Galley et al. '04, '06)

151

# Extracting syntactic rules

## 3) Extract rules



(Galley et al. '04, '06)

151

# Extracting syntactic rules



## 3) Extract rules

(Galley et al. '04, '06)

151

# Extracting syntactic rules

## 3) Extract rules



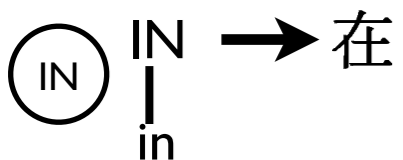(Galley et al. '04, '06)

151

# Extracting syntactic rules

## 3) Extract rules



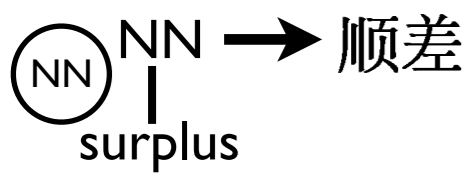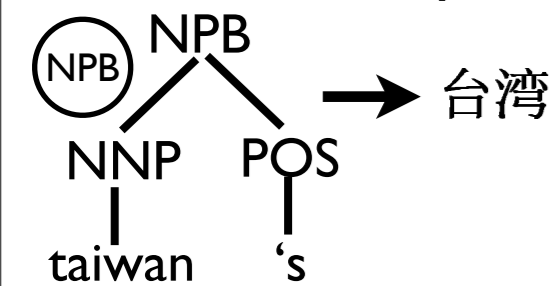(Galley et al. '04, '06)

151

# Extracting syntactic rules
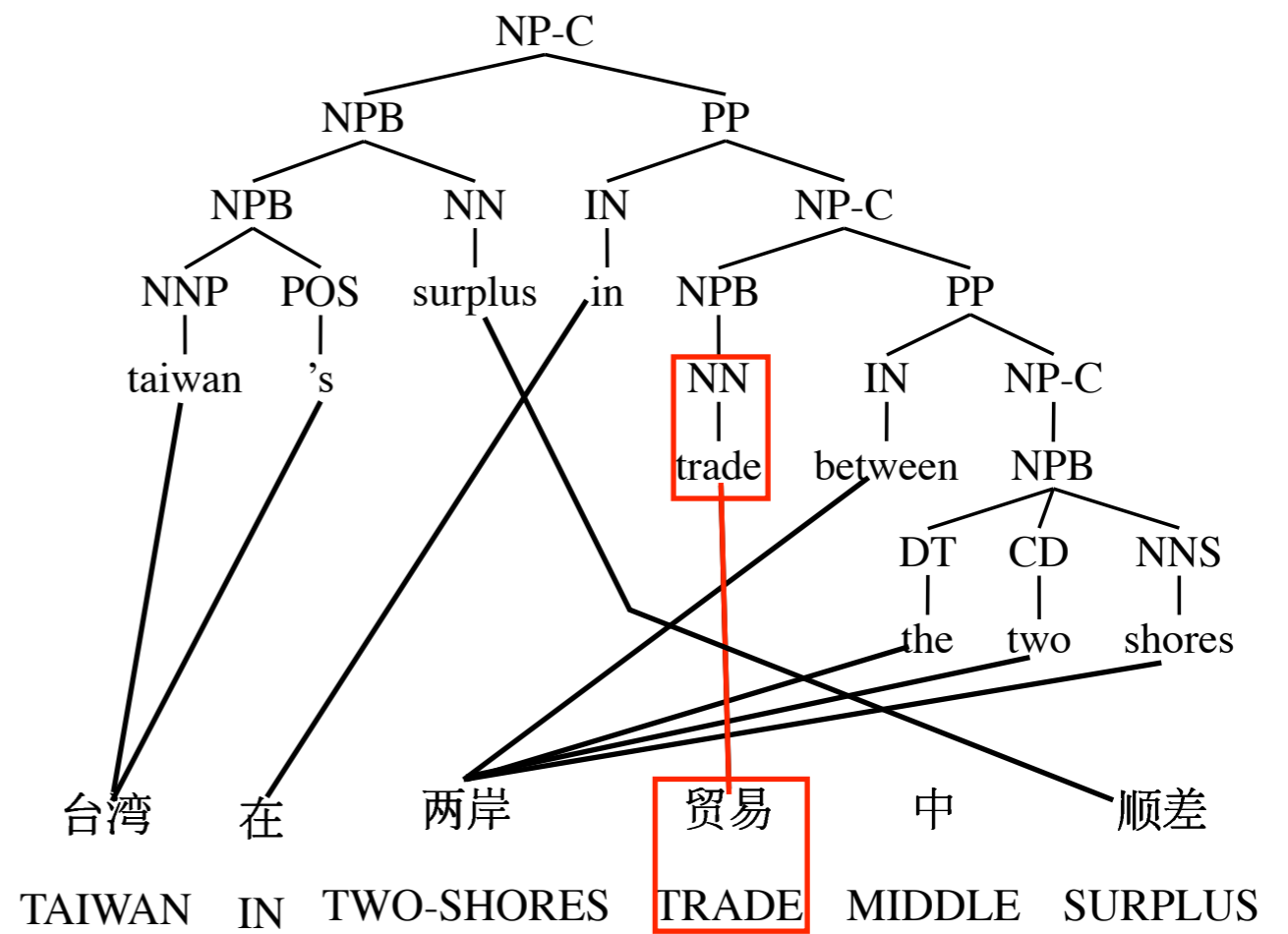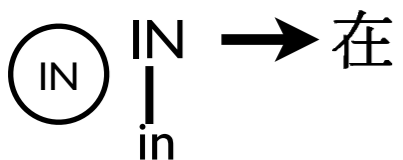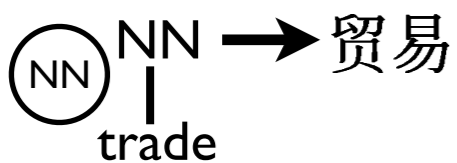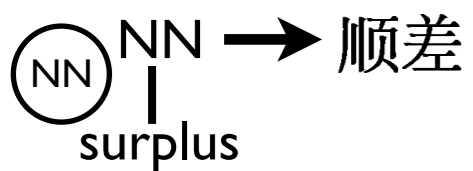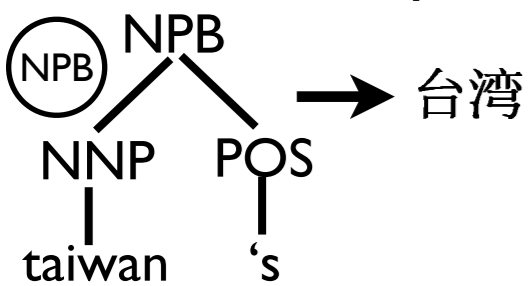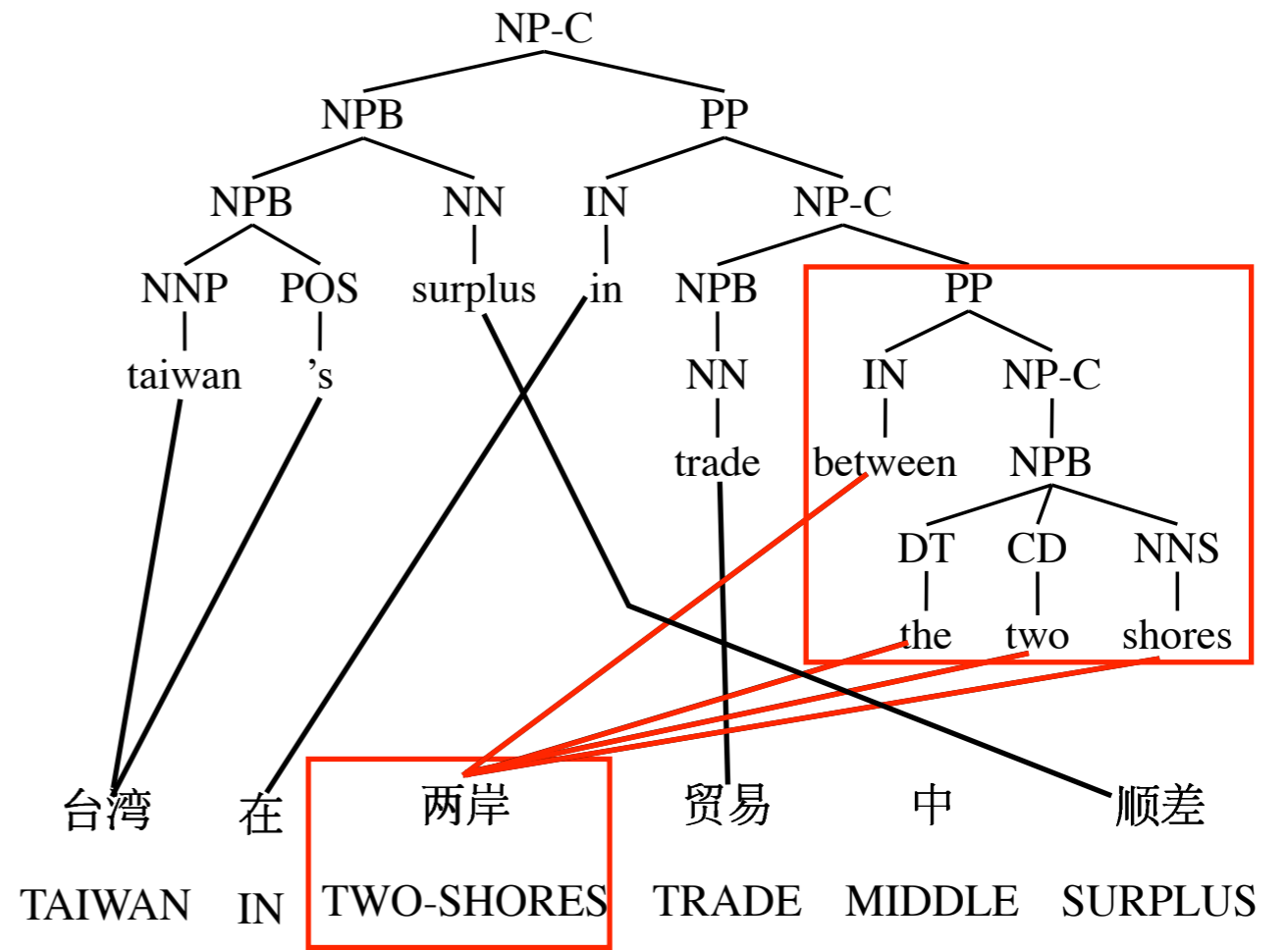
## 3) Extract rules



(Galley et al. '04, '06)

151

# Bad alignments make bad rules



One bad link makes a totally unusable syntax rule!

152

# Bad alignments make bad rules



One bad link makes a totally unusable syntax rule!

152

# Where do the alignments come from?

(word pairs)

(e-string, f-string)

seed data

sentence pairs

unsupervised learning

Viterbi alignments

generative model

(GIZA++)
(Och and Ney, '03)

(IBM model 4)
(Brown et al., '93)

Notice, nothing about syntax!

153

# Let's add syntax!

(tree-string
syntax rules)

(e-**tree**, f-string)

```
seed data
```

```
sentence
pairs
```

```
unsupervised
learning
```

```
Viterbi
alignments
```

```
generative
model
```

(Training Tree Transducers)
(Graehl, Knight, May '08)

(syntax model)

154

# Let's add syntax!

(tree-string
syntax rules)

(e-**tree**, f-string)

```
seed data
```

```
sentence
pairs
```

```
unsupervised
learning
```

```
Viterbi
alignments
```

```
generative
model
```

(Training Tree Transducers)
(Graehl, Knight, May '08)

(syntax model)

tree-to-string transducer

154

# Let's add syntax!

(tree-string
syntax rules)

(e-**tree**, f-string)

seed data

sentence
pairs

tiburon -t ...

unsupervised
learning

Viterbi
alignments

generative
model

(Training Tree Transducers)
(Graehl, Knight, May '08)

(syntax model)

tree-to-string transducer

154

# How we learn

NPB

NNP  POS   = ?

taiwan   's

台湾

中
台湾

中

**RULE SET**

R2: (NPB) NPB → 台湾　在　　两岸　　　贸易

NNP POS

taiwan   's

R11: (NPB) NPB → 台湾

(NPB) ● POS

在　　两岸

R12: (NNP) NNP → 台灣　在　　两岸

taiwan

R17: (NPB) NPB → (POS) ●

(NPB) NNP 两岸

在

taiwan

R18: (POS) POS → 台湾

's

- For each training sentence, build *derivation forest* containing each possible tree of rules that satisfies the sentence pair

- EM iterations set highest probability to most useful rules

- Viterbi derivation has syntax-aware alignments and bad rules are not extracted

R2 | R11 | R17
    | R12 | R18

(Graehl, Knight, May, '08)

# Experiments

GIZA

Vs.

GIZA

Syntactic
Re-Alignment

Syntax
MT

Syntax
MT

- Build a bootstrap alignment with GIZA

- Obtain new alignments with syntactic re-alignment

- Compare syntax MT system performance on rules extracted from each alignment

156

# Results

| source language | original alignments | type | MT system rules (millions) | NIST 2003 BLEU | Δ |
|---|---|---|---|---|---|
| Arabic | weak | baseline | 2.3 | 47.3 | +.6 |
| | | re-alignment | 2.5 | **47.9** | |
| | strong | baseline | 3.2 | 49.6 | +.4 |
| | | re-alignment | 3.6 | **50.0** | |
| Chinese | weak | baseline | 19.1 | 37.8 | +.9 |
| | | re-alignment | 26.0 | **38.7** | |
| | strong | baseline | 23.4 | 38.9 | +1.1 |
| | | re-alignment | 33.4 | **40.0** | |

157

# Conclusions and future work

- Algorithmic contributions

  - Determinization of weighted tree automata

  - Efficient inference through weighted tree transducer cascades

- Practical contributions

  - Weighted tree automata and transducer toolkit

  - Improvements in SMT using tree transducer EM

158

# Future work

- More algorithms!
    - approximate linear k-best
    - on-the-fly tree-to-string inference
- More applications!
    - financial systems
    - gene sequencing
- More formalisms!
    - unranked automata
    - tree-adjoining grammars

159

# Conclusions

- Tiburon makes it easy to use tree transducers in NLP

- (known) Theses using Tiburon:

  - Alexander Radzievskiy -- Masters on parsing with semantic role labels

  - Joseph Tepperman -- PhD on pronunciation evaluation

  - Victoria Fossum -- PhD on machine translation and parsing

- July 2010: ATANLP in Uppsala!

160

# Thanks!

Erika Barragan-Nunez, Rahul Bhagat, Marlynn Block, Matthias Büchse, Gully Burns, Marco Carbone, David Chiang, Hal Daumé III, Steve DeNeefe, John DeNero, Jason Eisner, Victoria Fossum, Alex Fraser, Jonathan Graehl, Erica Greene, Carmen Heger, Ulf Hermjakob, Johanna Högberg, Dirk Hovy, Ed Hovy, Liang Huang, David Kempe, Kevin Knight, Sven Koenig, Zornitsa Kozareva, Lorelei Laird, Kary Lau, Jerry Levine, Andreas Maletti, Daniel Marcu, Mitch Marcus, Howard May, Irena May, Rutu Mehta, Alma Nava, Adam Pauls, Fernando Pereira, Ben Plantan, Oana Postolache, Michael Pust, David Pynadath, Sujith Ravi, Deepak Ravichandran, Jason Riesa, Bill Rounds, Lee Rowland, Tom Russ, Shri Narayanan, Radu Soricut, Magnus Steinby, Shang-Hua Teng, Cătălin Tîrnăucă, Ashish Vaswani, Jens Vöckler, Heiko Vogler, David Foster Wallace, Wei Wang, Ralph Weischedel, Kenji Yamada

161

# Backup Slides

# Non-deterministic and nonterminal?

# MT Details

- Decoded 116 short Chinese sentences using the string-to-tree MT model based on (Galley et al. 2004)

  - No language model

  - No reranking

- Counted number of trees in each forest before and after determinization

- 86.3% trees in forest are duplicates on average

  - $1.4\text{x}10^{12}$ median per forest pre-determ

  - $2.0\text{x}10^{11}$ median per forest post-determ

- Determinization changes top tree 77.6% of the time

- Crunching matches determinization 50.6% of the time

165

# xLNT not closed!

could be arbitrarily long!

(Maletti, Graehl, Hopkins, Knight, '09)

# Closure Under Composition and Recognizability Preservation

| closed | forward recog | backward recog |
|--------|---------------|----------------|
| wLNT | wxLNT | xT |
| | | wxLT |

# Where do the rules come from?

```
        A
       / \
      B   C
     / \  |
    D   E c
    |   |
    d   e
    f g h
```

=

103 possible rules

- Ideally we would add all possible rules

- To avoid overflow, we bootstrap with a previous (syntax-free) alignment model

- This follows a rich history in MT (Och & Ney '00, Fraser & Marcu '06)

# Other approaches to this problem

- Cherry and Lin '06: Discriminatively train ITG-based alignment model influenced by dependency graph

- DeNero and Klein '07: HMM model modified to incorporate syntax penalty into distortion

- Fossum et al. '08: Identify troublesome links and remove them

169

# Where do the rules come from?

# Where do the rules come from?



(Galley et al. '04)

# Where do the rules come from?



(Galley et al. '04)

# EM size bias



RULE CORPUS

R2: (NPB) NPB → 台湾    在    两岸       贸易    中    顺差
         NNP POS
         taiwan  's

R11: (NPB) NPB →  (NNP) ●
              ●    POS

R12: (NNP) NNP → 台湾    在    两岸       贸易    中    顺差
              taiwan

R17: (NPB) NPB → (POS) ●
              NNP    两岸
R18: (POS) POS → 台湾    在    两岸       贸易    中    顺差
              's         taiwan

R2 | R11 | R17
      |      |
     R12   R18

more useful

favored

- EM attempts to learn derivations with highest probability.

- Shorter derivations have fewer chances to take a probability "hit" and are thus biased to be favored.

- This, then, tends to favor larger rules, generally the opposite of what we want.

171

# Correcting size bias

**RULE CORPUS**

R2: (NPB) NPB → 台湾　　在　　　　两岸　　　　贸易　　中　　　顺差

    NNP POS    R11: (NPB) NPB → (NNP)

    taiwan 's       POS

R12: (NNP) NNP → 台湾　　在　　　　两岸　　　　贸易　　中　　　顺差

    taiwan

    R17: (NPB) NPB → (POS)

R18: (POS) POS → 台湾　　在　　NNP　两岸　　　贸易　　中　　　顺差

    's      taiwan

$S_3$ — $S_3$ — $R2$

$S_2$ — $R11$ — $R12$

$S_2$ — $R17$ — $R18$

- When using a rule with $n$ non-leaf nodes, prepend $n-1$ copies of a special *size rule* $S_n$

- Each competing derivation now has the same number of rules

- Size rules are built into the derivation forests and weights are learned by the same EM procedure

172

# Complexity Analysis

| | | |
|---|---|---|
| k-best (H&C) | $O(P + D_{max}k \log k)$ | P = rtg rules <br> $D_{max}$ = max deriv |
| determinization | $O(Ak^{zL})$ | A = alph size <br> k = max rank <br> z = max tree size <br> L = lang size |
| rtg+xLNT | $O(RP^l)$ | R = trans rules <br> P = rtg rules <br> l = max trans lhs |
| xT+LNT | $O(R_A R_B^r)$ | $R_A$ = xT rules <br> $R_B$ = LNT rules <br> r = max $R_A$ rhs |

# Dramatic use of size rules

$R_{bad}$:

广西　　　　对外　　　　开放

→ x0 对外 开放

```
              S-C
          /        \
       NP-C         VP
        |         /  |    \
       NPB     VBG  PRT    PP
      /   \     |    |    /   \
  x0:NNP  POS opening RP TO   NP-C
           |         |   |     |
           's        up  to   NPB
                             / |  \
                            DT JJ  NN
                            |  |   |
                          the outside world
```

$S_{15}$ 广西

$S_{15}$

$\vdots$

$\left.\begin{array}{c} \\ \\ \\ \\ \end{array}\right\}$ 14 times

$\vdots$

$S_{15}$　　　　开放

$R_{bad}$

对外　　　　　　广西

174

# Approximate Algorithms

- linear-time approximate *k*-best

- polynomial time determinization that fails to recognize some trees in the input

- weighted domain projection with relative ordering, but not exact weights, preserved

- mildly incorrect fast composition

- on-the-fly tree-to-string backward application

175

# Engineering

- Battle-test Tiburon implementations and bring it up to production level

- Make greater use of system on biological sequencing and financial systems analysis -- leads to more interesting algorithmic questions, different types of transducers

176

# Explore the limits of Tree Transducers

- Weighting scheme of Collins' parsing model[1] doesn't fit well

- Very large tree transducers needed in syntax MT[2]

- Can these models be simplified and still retain their power? Or should different formalisms be used?

1: Collins, 1997          2: DeNeefe and Knight, 2009

177