# Faster Decoding for Phrases and Syntax

Kenneth Heafield

# Translation is Expensive

"speed-up in tuning time but affects the performance"
"18 days using 12 cores"
[Williams et al WMT 2014]

"Time-sensitive BLEU score"
[Chung and Galley, 2012]

"Due to time constraints, this procedure was not used"
[Servan et al, WMT 2012]

$\implies$ Routine Quality Compromises

蘭州國中室內溫水游泳池

In-Room Lukewarm Water Swimming Pool

# Blame the Language Model

"LM queries often account for more than 50% of the CPU"
[Green et al, WMT 2014]

# Blame the Language Model

"LM queries often account for more than 50% of the CPU"
[Green et al, WMT 2014]

Faster queries (KenLM)
**More effective queries**

# Cube pruning

- Widely used for phrase-based and syntax-based MT

- May be applied in conjunction with a bottom-up decoder, or as a second "rescoring" pass

  - Nodes may also be grouped together (for example, all nodes corresponding to a certain source span)

- Requirement for topological ordering means translation hypergraph may not have cycles
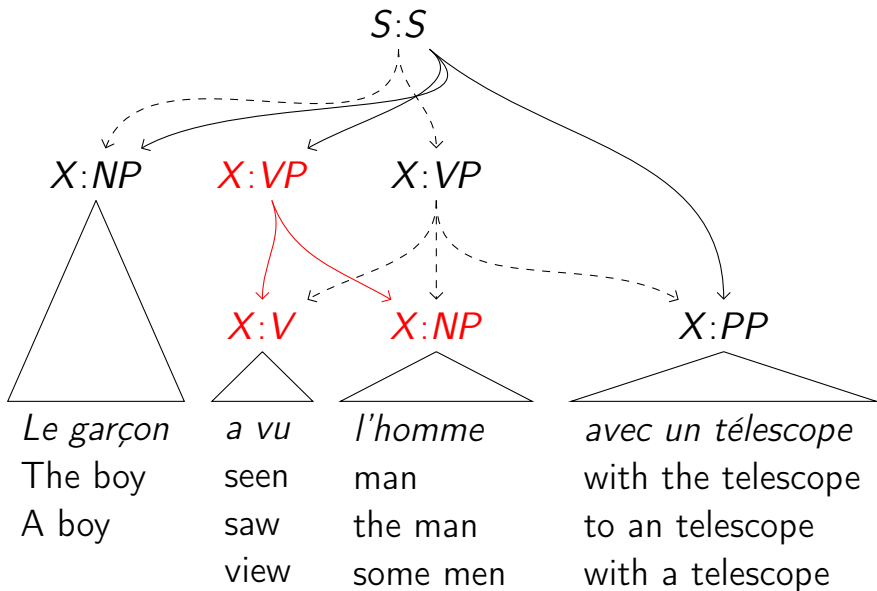
# Cube pruning

- Widely used for phrase-based and syntax-based MT

- May be applied in conjunction with a bottom-up decoder, or as a second "rescoring" pass

  - Nodes may also be grouped together (for example all nodes corresponding to a certain source span)

- Requirement for topological ordering means translation hypergraph may not have cycles

# Decoding Example: Input

*Le garçon     a vu     l'homme     avec un télescope*

# Decoding Example: Parse with SCFG

# Decoding Example: Read Target Side

# Decoding Example: One Constituent



S:S

X:NP          X:VP          X:VP

X:V          X:NP          X:PP

| Le garçon | a vu | l'homme | avec un télescope |
| The boy | seen | man | with the telescope |
| A boy | saw | the man | to an telescope |
|  | view | some men | with a telescope |

$X{:}VP$

$X{:}V$  $X{:}NP$

*a vu*  *l'homme*

**Hyp**  **Hyp**

seen  man

saw  the man

view  some men

$X$:VP

$X$:V        $X$:NP

*a vu*        *l'homme*

**Hyp**

seen

saw

view

**Hyp**

man

the man

some men

$X$:VP

*a vu l'homme*

**Hypothesis**

seen man

seen the man

seen some men

saw man

saw the man

saw some men

view man

view the man

view some men

$X{:}VP$

$X{:}V$     $X{:}NP$

$X{:}VP$

*a vu l'homme*

| Hypothesis | Score |
|---|---|
| seen man | -8.8 |
| seen the man | -7.6 |
| seen some men | -9.5 |
| saw man | -8.3 |
| saw the man | -6.9 |
| saw some men | -8.5 |
| view man | -8.5 |
| view the man | -8.9 |
| view some men | -10.8 |

*a vu*

| Hyp | Score |
|---|---|
| seen | -3.8 |
| saw | -4.0 |
| view | -4.0 |

*l'homme*

| Hyp | Score |
|---|---|
| man | -3.6 |
| the man | -4.3 |
| some men | -6.3 |

$X{:}VP$

$X{:}V$   $X{:}NP$

*a vu*

| Hyp | Score |
|------|-------|
| seen | -3.8 |
| saw | -4.0 |
| view | -4.0 |

*l'homme*

| Hyp | Score |
|------|-------|
| man | -3.6 |
| the man | -4.3 |
| some men | -6.3 |

$X{:}VP$

*a vu l'homme*

| Hypothesis | Score |
|------------|-------|
| saw the man | -6.9 |
| seen the man | -7.6 |
| saw man | -8.3 |
| saw some men | -8.5 |
| view man | -8.5 |
| seen man | -8.8 |
| view the man | -8.9 |
| seen some men | -9.5 |
| view some men | -10.8 |

X:VP

X:V      X:NP

*a vu*      *l'homme*

| Hyp | Score |
|---|---|
| seen | -3.8 |
| saw | -4.0 |
| view | -4.0 |

| Hyp | Score |
|---|---|
| man | -3.6 |
| the man | -4.3 |
| some men | -6.3 |

X:VP

*a vu l'homme*

| Hypothesis | Score |
|---|---|
| saw the man | -6.9 |
| seen the man | -7.6 |
| saw man | -8.3 |
| saw some men | -8.5 |
| view man | -8.5 |
| seen man | -8.8 |
| view the man | -8.9 |
| seen some men | -9.5 |
| view some men | -10.8 |

## Scores do not sum

X:VP

X:V — a vu

| Hyp | Score |
|-----|-------|
| seen | -3.8 |
| saw | -4.0 |
| view | -4.0 |

X:NP — l'homme

| Hyp | Score |
|-----|-------|
| man | -3.6 |
| the man | -4.3 |
| some men | -6.3 |

X:VP — a vu l'homme

| Hypothesis | Score |
|------------|-------|
| saw the man | -6.9 |
| seen the man | -7.6 |
| saw man | -8.3 |
| saw some men | -8.5 |
| view man | -8.5 |
| seen man | -8.8 |
| view the man | -8.9 |
| seen some men | -9.5 |
| view some men | -10.8 |

Pruning is Approximate

# Appending Strings

Hypotheses are built by string concatenation.
Language model probability changes when this is done:

$$\frac{p(\text{saw the man})}{p(\text{saw})p(\text{the man})} = \frac{p(\text{the} \mid \text{saw})p(\text{man} \mid \text{saw the})}{p(\text{the})\quad p(\text{man} \mid \text{the})}$$

# Appending Strings

Hypotheses are built by string concatenation.
Language model probability changes when this is done:

$$\frac{p(\text{saw the man})}{p(\text{saw})p(\text{the man})} = \frac{p(\text{the} \mid \text{saw})p(\text{man} \mid \text{saw the})}{p(\text{the})} \frac{}{p(\text{man} \mid \text{the})}$$

Log probability is part of the score
$\implies$ Scores do not sum
$\implies$ Local decisions may not be globally optimal
$\implies$ Search is hard.

# Beam Search

|  | **man** $-3.6$ | **the man** $-4.3$ | **some men** $-6.3$ |
|---|---|---|---|
| **seen** $-3.8$ | seen man $-8.8$ | seen the man $-7.6$ | seen some men $-9.5$ |
| **saw** $-4.0$ | saw man $-8.3$ | saw the man $-6.9$ | saw some men $-8.5$ |
| **view** $-4.0$ | view man $-8.5$ | view the man $-8.9$ | view some men $-10.8$ |

[Lowerre, 1976; Chiang, 2005]

# Cube Pruning

man  −3.6   the man −4.3   some men −6.3

seen −3.8   Queue
saw  −4.0
view −4.0

**Queue**

| Hypothesis | Sum |
|---|---|
| → seen man | $-3.8-3.6=-7.4$ |

[Chiang, 2007]

# Cube Pruning

|  | man −3.6 | the man −4.3 | some men −6.3 |
|---|---|---|---|
| seen −3.8 | seen man −8.8 | Queue | |
| saw −4.0 | Queue | | |
| view −4.0 | | | |

**Queue**

| Hypothesis | Sum |
|---|---|
| → saw man | −4.0−3.6=−7.6 |
| seen the man | −3.8−4.3=−8.1 |

[Chiang, 2007]

# Cube Pruning

|          | man       | the man  | some men |
|----------|-----------|----------|----------|
| **seen** −3.8 | seen man −8.8 | Queue    |          |
| **saw** −4.0  | saw man −8.3  | Queue    |          |
| **view** −4.0 | Queue     |          |          |

|          | **man** −3.6 | **the man** −4.3 | **some men** −6.3 |
|----------|-----------|----------|----------|

### Queue

| Hypothesis   | Sum              |
|--------------|------------------|
| → view man   | −4.0−3.6=−7.6    |
| seen the man | −3.8−4.3=−8.1    |
| saw the man  | −4.0−4.3=−8.3    |

[Chiang, 2007]

# Cube Pruning

|           |          | **man**   | $-3.6$ | **the man** | $-4.3$ | **some men** | $-6.3$ |
|-----------|----------|-----------|--------|-------------|--------|--------------|--------|
| **seen**  | $-3.8$   | seen man  | $-8.8$ | Queue       |        |              |        |
| **saw**   | $-4.0$   | saw man   | $-8.3$ | Queue       |        |              |        |
| **view**  | $-4.0$   | view man  | $-8.5$ | <span style="color:red">Queue</span> |        |              |        |

### Queue

| Hypothesis | Sum |
|------------|-----|
| → seen the man | $-3.8-4.3=-8.1$ |
| saw the man | $-4.0-4.3=-8.3$ |
| view the man | $-4.0-4.3=-8.3$ |

[Chiang, 2007]

### Beam Search
Make every dish. Keep the best $k$, throw the rest out.

### Cube pruning
Combine the best ingredients. Only make $k$ dishes.

# Cube Pruning Hypotheses are Atomic

**String**
is a
are a

**String**
countries that
countries which
country

**String**
is a countries that
are a countries that
are a countries which
⋮

No notion that "a countries" is bad.

### Beam Search
Make every dish. Keep the best $k$, throw the rest out.

### Cube pruning
Combine the best ingredients. Only make $k$ dishes.

### Coarse-to-Fine
Make small portions, taste, and order the best ones.

# Coarse-to-Fine

Decode multiple times, adding detail each time:

Increased LM order, words instead of classes

Detect and prune "a countries" with a bigram LM.

[Zhang et al, 2008; Petrov et al, 2008]

# Coarse-to-Fine

Decode multiple times, adding detail each time:

Increased LM order, words instead of classes

Detect and prune "a countries" with a bigram LM.

[Zhang et al, 2008; Petrov et al, 2008]

Requires tuning each pruning pass.
Operates in lock step.

# Coarse-to-Fine

Decode multiple times, adding detail each time:

Increased LM order, words instead of classes

Detect and prune "a countries" with a bigram LM.

[Zhang et al, 2008; Petrov et al, 2008]

Requires tuning each pruning pass.
Operates in lock step.

# Can coarse-to-fine be done on the fly?

# Observations

Competing translations have words in common:

is a, are a

# Observations

Competing translations have words in common:
is a, are a

Words at the boundary matter most:
a + country, a + countries

# Observations

Competing translations have words in common:
is a, are a

Words at the boundary matter most:
a + country, a + countries

<span style="color:red">Emphasize boundary words</span>

### Beam Search
Make every dish. Keep the best $k$, throw the rest out.

### Cube pruning
Combine the best ingredients. Only make $k$ dishes.

### Coarse-to-Fine
Make small portions, taste, and order the best ones.

### Incremental
Taste during cooking. Share ingredients.

# Boundary Words

1 Left-to-right phrase-based: one side
2 Bottom-up syntax: both sides

# Partial Translations

## Plain text

The United Kingdom is a $+$ ...
Scotland and Wales are a $+$ ...

## Tree

The United Kingdom $\longleftarrow$ is

$\qquad\qquad\qquad\qquad\qquad$ a $\longleftarrow \epsilon$

Scotland and Wales $\longleftarrow$ are

# Phrase Continuations

## Plain text
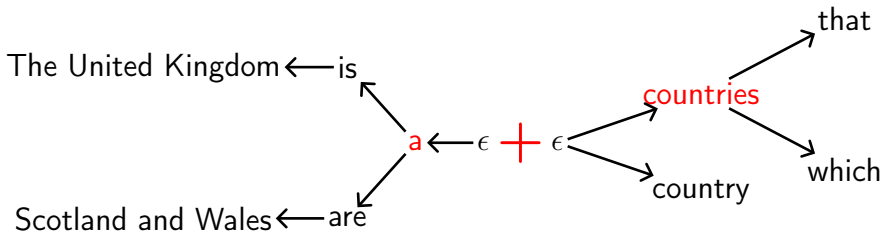
... + countries that

... + countries which

... + country

## Tree



```
                                    →that
              →countries
  ε                                 ↘which
              →country
```

The United Kingdom ← is

a ← $\epsilon$ + $\epsilon$ → countries → that

Scotland and Wales ← are

countries → which

country

The United Kingdom ← is

a ← $\epsilon$ + $\epsilon$ → countries → that

countries → which

country

Scotland and Wales ← are

Does the model like "a + countries"?

# Exploring and Backtracking

Does the model like "a + countries"?

Yes Try more detail.

No Consider alternatives.

# Exploring and Backtracking

Does the model like "a + countries"?

Yes Try more detail.

No Consider alternatives.

Formally: best-first search with a priority queue.

# The queue entry

"a + $\epsilon$"

## splits into

Best Child "a + countries"
Other Children "a + country"

# Scores come from the best descendant:

$$\text{Score}(a) = \max\{\text{Score}(\text{is } a), \text{Score}(\text{are } a)\}$$

Scores come from the best descendant:

$$\text{Score}(a) = \max\{\text{Score}(\text{is a}), \text{Score}(\text{are a})\}$$

The language model updates scores:

$$\text{Score}(a + \text{countries}) < \text{Score}(a) + \text{Score}(\text{countries})$$

Scores come from the best descendant:

$$\text{Score}(a) = \max\{\text{Score}(\text{is } a), \text{Score}(\text{are } a)\}$$

The language model updates scores:

$$\text{Score}(a + \text{countries}) < \text{Score}(a) + \text{Score}(\text{countries})$$

Formally: $p(\text{countries} \mid a)$ replaces $p(\text{countries})$

# Best-First Algorithm Summary

Populate the queue with $\epsilon + \epsilon$

    Loop until $k$ complete options have been found:

        Split the top-scoring option

    Build a tree from the $k$ complete options

# Summary

Translations are assembled from left to right.

Partial translations often share suffixes.
Phrases often share prefixes.

Test suffixes and prefixes before full combinations.
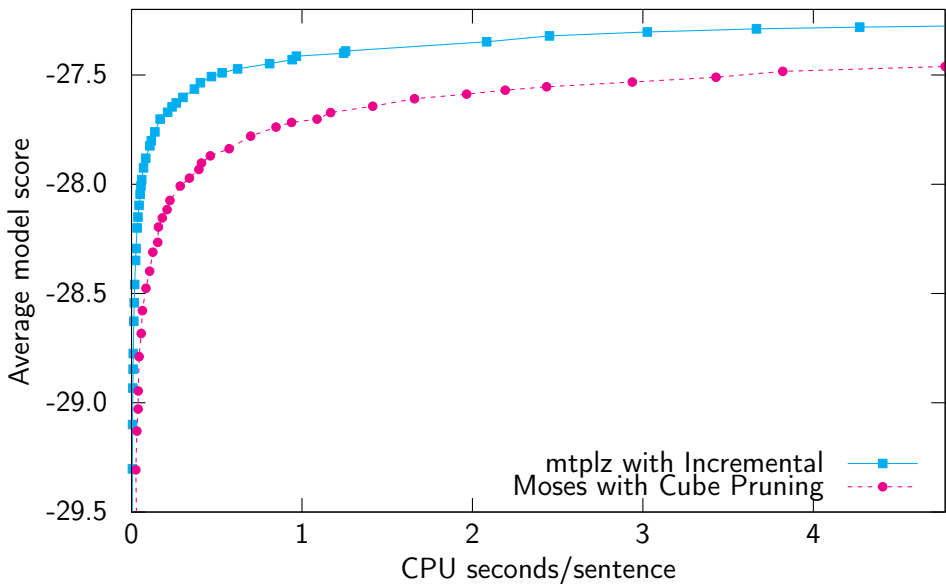
# Experiment

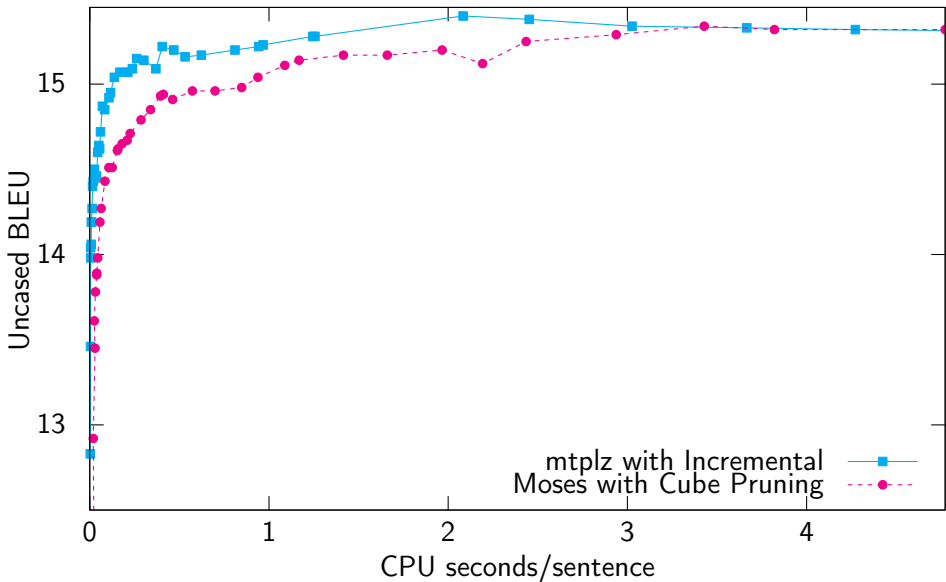Task Chinese–English

Source Stanford

Model Phrase-based

Software My own decoder, mtplz, versus Moses

# Phrase-Based Results



Plot of Average model score (y-axis, ranging from -29.5 to -27.5) versus CPU seconds/sentence (x-axis, 0 to 4). Two curves: "mtplz with Incremental" (blue squares) and "Moses with Cube Pruning" (magenta dashed circles).

# Phrase-Based Results



Legend:
- mtplz with Incremental
- Moses with Cube Pruning

X-axis: CPU seconds/sentence
Y-axis: Uncased BLEU

# Search

The language model cares most about adjacent words.

Test them first.

Share work for shared words.

# Boundary Words

1 Left-to-right phrase-based: one side
2 **Bottom-up syntax: both sides**

# Bottom-Up Syntax: Both Sides

is a $X$:$NP1$ $</s>$
is a $X$:$NP1$ that

How do we find the best value to substitute?
Manage words on both sides.

# Example Hypotheses

**Left State**                                              **Right State**

countries that maintain diplomatic relations with North Korea .
                                          ties

countries that have an embassy in DPR Korea .

country that maintains some diplomatic ties in North Korea .

nations which has some diplomatic ties with DPR Korea .

country that maintains some diplomatic ties with DPR Korea .

# Example Hypotheses

**Left State**                    **Right State**

(countries that      ⋄ with North Korea .)

(nations which has  ⋄ with DPR Korea .)

(countries that have ⋄      DPR Korea .)

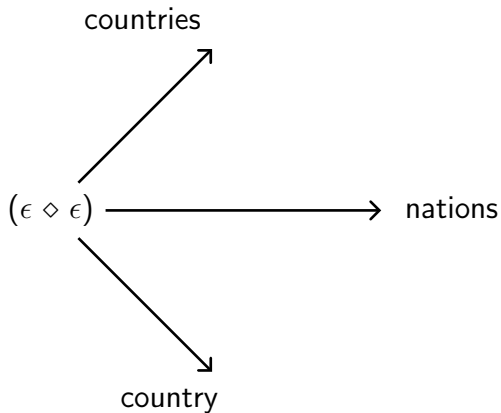(country             ⋄    in North Korea .)

(country             ⋄ with DPR Korea .)
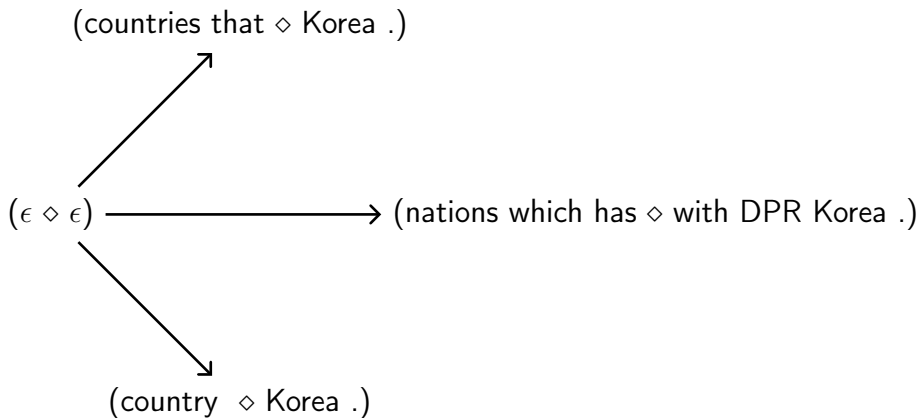
⋄ Words the language model does not care about
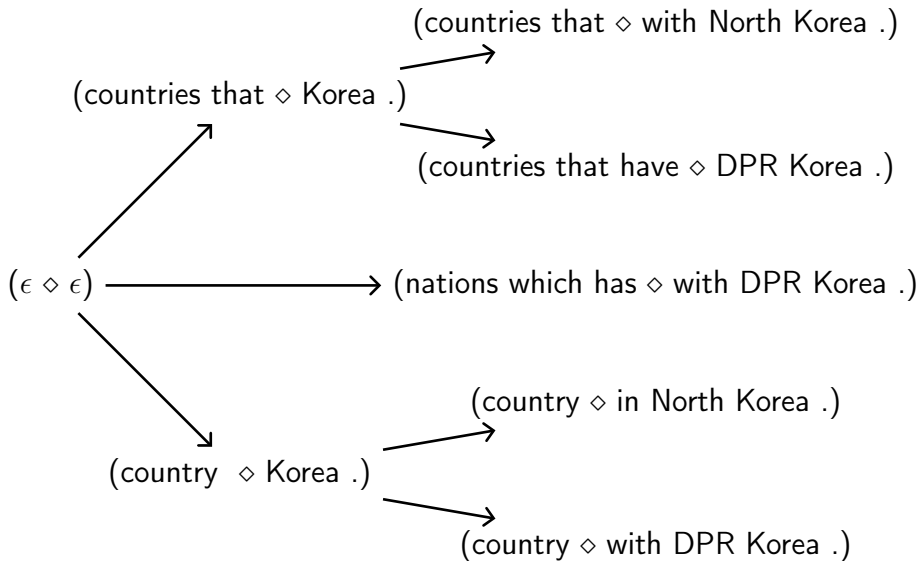
Idea: alternate between left and right side

# Group by Leftmost Word



countries

$(\epsilon \diamond \epsilon)$ ⟶ nations

country

# Reveal Common Words in Each Group

(countries that ◇ Korea .)

$(\epsilon \diamond \epsilon)$ ⟶ (nations which has ◇ with DPR Korea .)

(country ◇ Korea .)

# Alternate Sides Until Tree is Full

(countries that ◇ with North Korea .)

(countries that ◇ Korea .)

(countries that have ◇ DPR Korea .)

(ε ◇ ε) ⟶ (nations which has ◇ with DPR Korea .)

(country ◇ in North Korea .)

(country ◇ Korea .)

(country ◇ with DPR Korea .)

# Using Rules

is a $X{:}NP1$ </s> $\qquad$ $X{:}V1$ the $X{:}N2$

turns into $\qquad\qquad$ turns into

is a $(\epsilon \diamond \epsilon)$ </s> $\qquad$ $(\epsilon \diamond \epsilon)$ the $(\epsilon \diamond \epsilon)$

$$\underbrace{\qquad\quad}_{X{:}V1} \quad \underbrace{\qquad\quad}_{X{:}N2}$$

# Exploring and Backtracking

Does the LM like "is a (countries that $\diamond$ Korea .) </s>"?

Yes Try more detail.
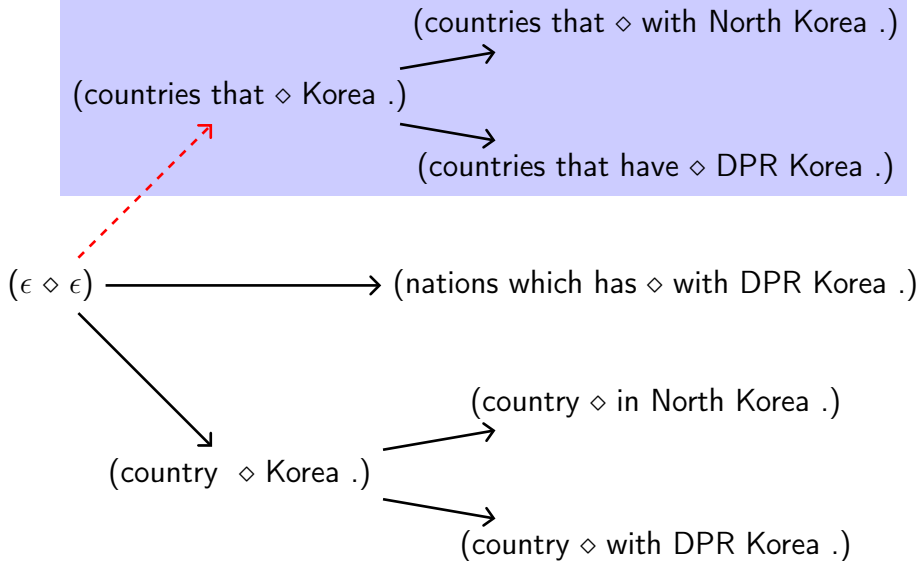
No Consider alternatives.

# Exploring and Backtracking

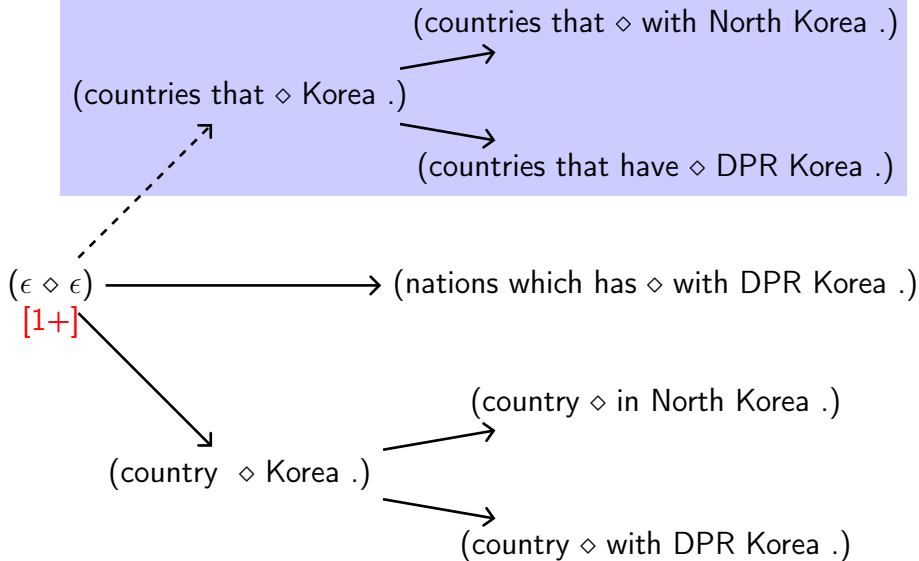Does the LM like "is a (countries that ◇ Korea .) </s>"?

Yes Try more detail.

No Consider alternatives.

Formally: priority queue containing breadcrumbs.

# Split and Leave Breadcrumbs

(countries that ⋄ with North Korea .)

(countries that ⋄ Korea .)

(countries that have ⋄ DPR Korea .)

$(\epsilon \diamond \epsilon)$ ⟶ (nations which has ⋄ with DPR Korea .)

(country ⋄ Korea .)

(country ⋄ in North Korea .)

(country ⋄ with DPR Korea .)

# Split and Leave Breadcrumbs



(countries that ◇ with North Korea .)

(countries that ◇ Korea .)

(countries that have ◇ DPR Korea .)

(ε ◇ ε)
[1+]

(nations which has ◇ with DPR Korea .)

(country ◇ Korea .)

(country ◇ in North Korea .)

(country ◇ with DPR Korea .)

# The queue entry

is a $(\epsilon \diamond \epsilon)$ </s>

## splits into

Zeroth Child "is a (countries that $\diamond$ Korea .) </s>"
Other Children "is a $(\epsilon \diamond \epsilon)$[1+] </s>"

Children except the zeroth.

A priority queue contains competing entries:

    is a (countries that $\diamond$ Korea .) $</s>$
    $(\epsilon \diamond \epsilon)$ the $(\epsilon \diamond \epsilon)$
    is a $(\epsilon \diamond \epsilon)[1+]$ $</s>$

    The algorithm pops the top entry,
    splits a non-terminal, and pushes.

# Best-First Algorithm

Populate the queue with rules like "is a $(\epsilon \diamond \epsilon)$ </s>"

Loop until $k$ complete options have been found:

Split the top-scoring option, leave a breadcrumb

Build a tree from the $k$ complete options

# Syntax

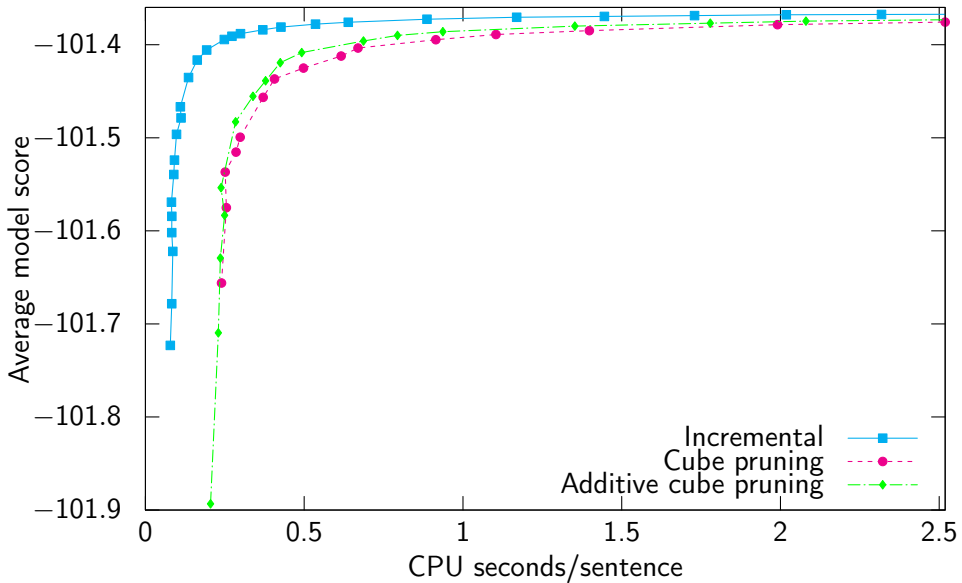Same as phrase-based, just concatenate on left and right.
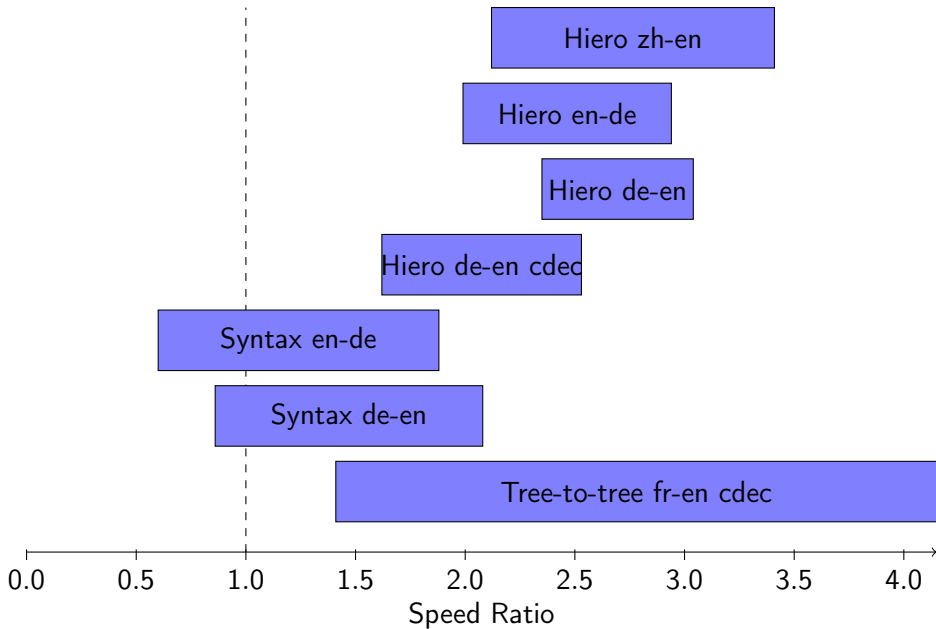
# Experiment

Task WMT 2011 German-English

Model Hierarchical

Decoder Moses

# Moses Hierarchical

# Moses Hierarchical

Speed Ratio

# Incremental

A series of coarse-to-fine estimates.

Continually taste the dish and adjust.

# Takeaway

Search limits what translation can do.

Long-distance models like gender and number are harder.

Open the black box.

Language models can produce intermediate scores.