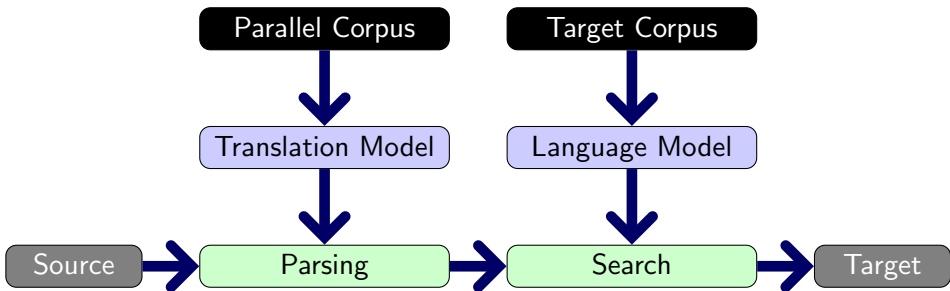


Efficient Language Modeling Algorithms with Applications to Statistical Machine Translation

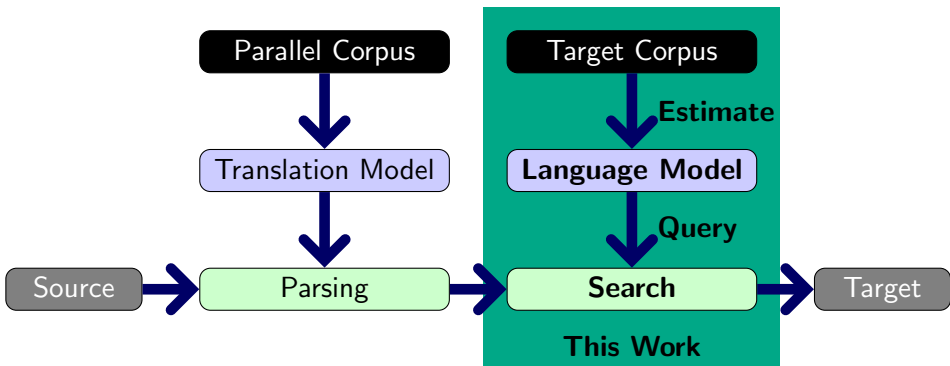
Kenneth Heafield

September 20, 2013

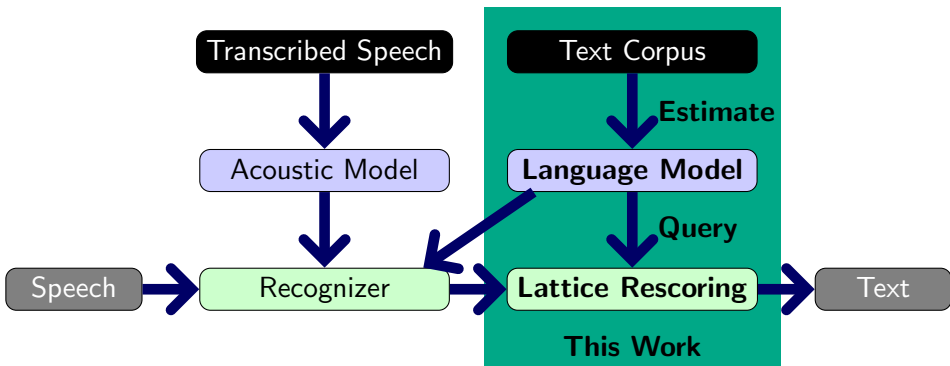
Application: Syntactic MT



Application: Syntactic MT

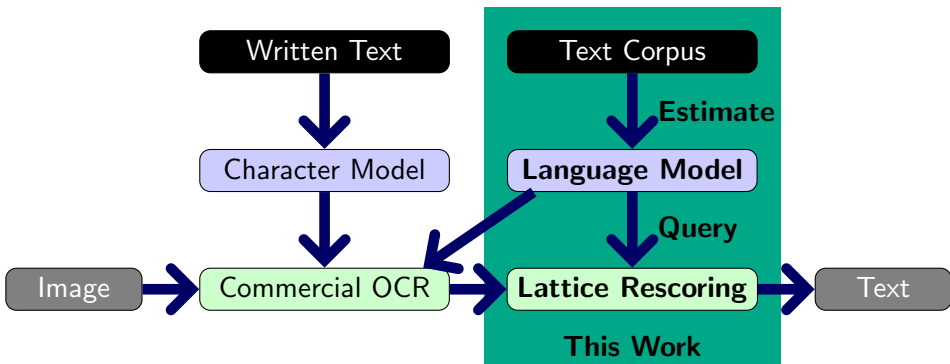


Speech Recognition



Some of the thesis is already used for speech [Kim et al, 2012; Si et al, 2013].

Optical Character Recognition



Numen (2013) is using most of this thesis for OCR.

Language Models Are Expensive

Store a sparse set of 121 billion n -grams

⇒ RAM

Language Models Are Expensive

Store a sparse set of 121 billion n -grams \implies RAM

Millions of probability queries per sentence \implies CPU

Language Models Are Expensive

Store a sparse set of 121 billion n -grams \implies RAM

Millions of probability queries per sentence \implies CPU

Probability does not multiply when strings are concatenated:

$$\begin{aligned} & p(\text{saw the man}) \\ & \neq \\ & p(\text{saw})p(\text{the man}) \end{aligned}$$

\implies Search is hard \implies CPU

Thesis Problem

Much of the CPU and RAM cost is due to the language model.
Researchers routinely compromise quality due to these costs.

Costs Due To Language Models

Estimation from text

Probability queries

Search when the objective includes log probability

Results Preview

Speed RAM

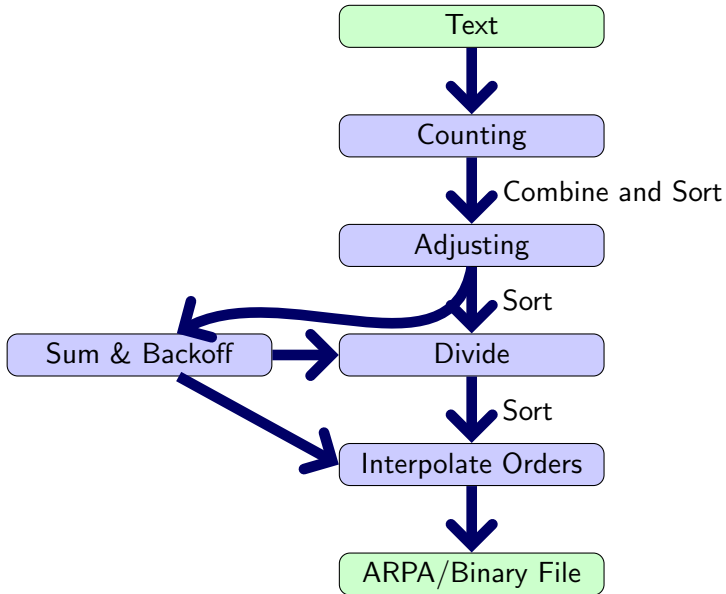
| | | |
|----------------------|-----------|-------|
| Estimation from Text | 7.1x | 0.07x |
| Raw Queries | 2.4x | 0.57x |
| Decoding | 3.2–10.0x | 0.85x |

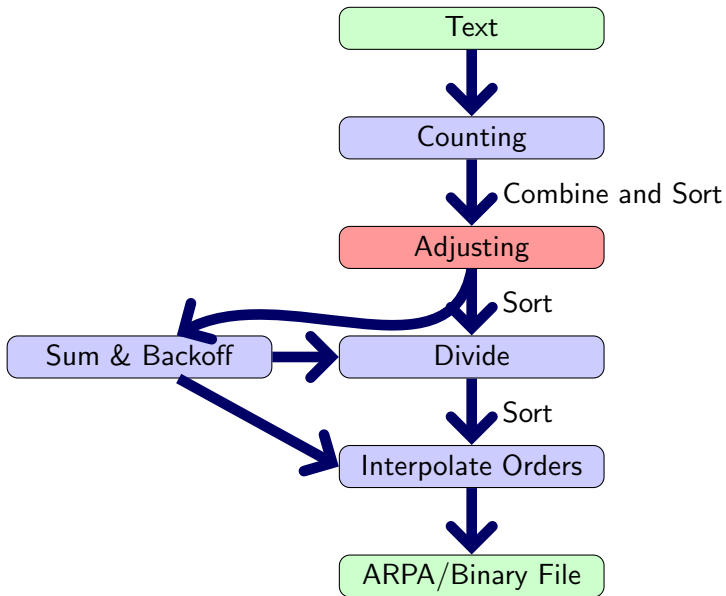
Decoding performance includes $\approx 1.15x$ speedup from raw queries.
Baseline: SRILM and cube pruning (more later).

Implz Features

- Disk-based streaming and sorting
- User-specified RAM
- Fast
- Interpolated modified Kneser-Ney

7.7% of SRI's RAM, 14% of SRI's wall time





Adjusting

Adjusted counts are:

Trigrams Same as counts.

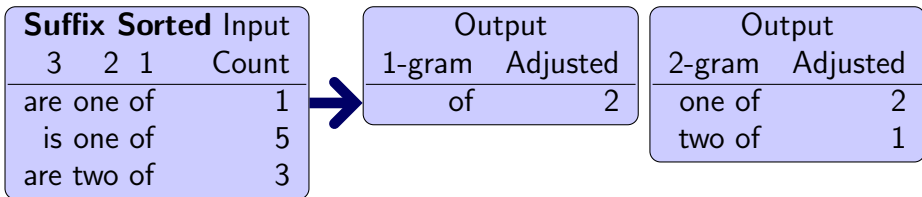
Others Number of unique words to the left.

Adjusting

Adjusted counts are:

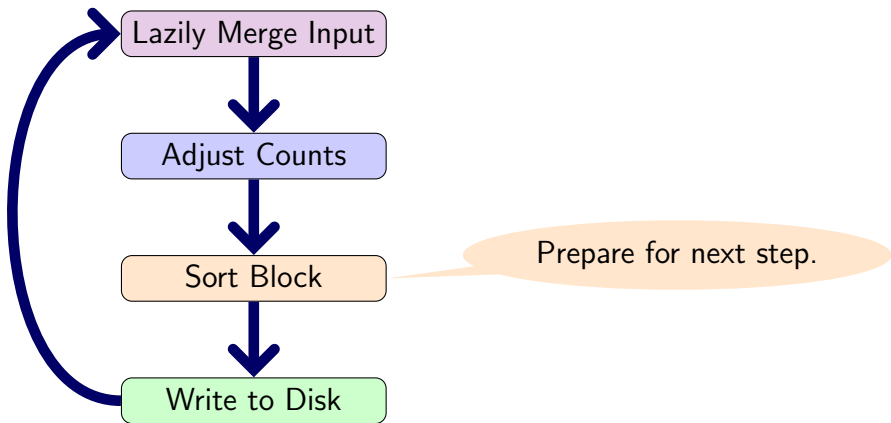
Trigrams Same as counts.

Others Number of unique words to the left.

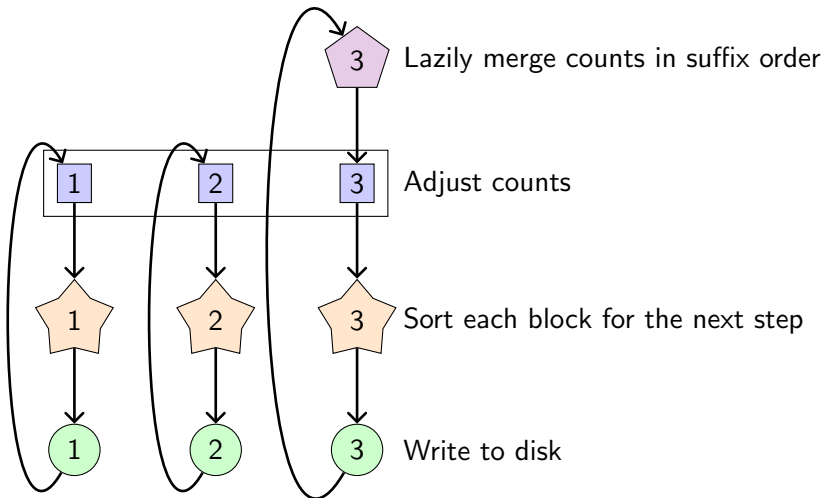


Streaming Framework

Memory is divided into blocks. Blocks are recycled.



Adjusted Counts Detail



Each vertex is a thread \implies Simultaneous disk and CPU.

Experiment: Toolkit Comparison

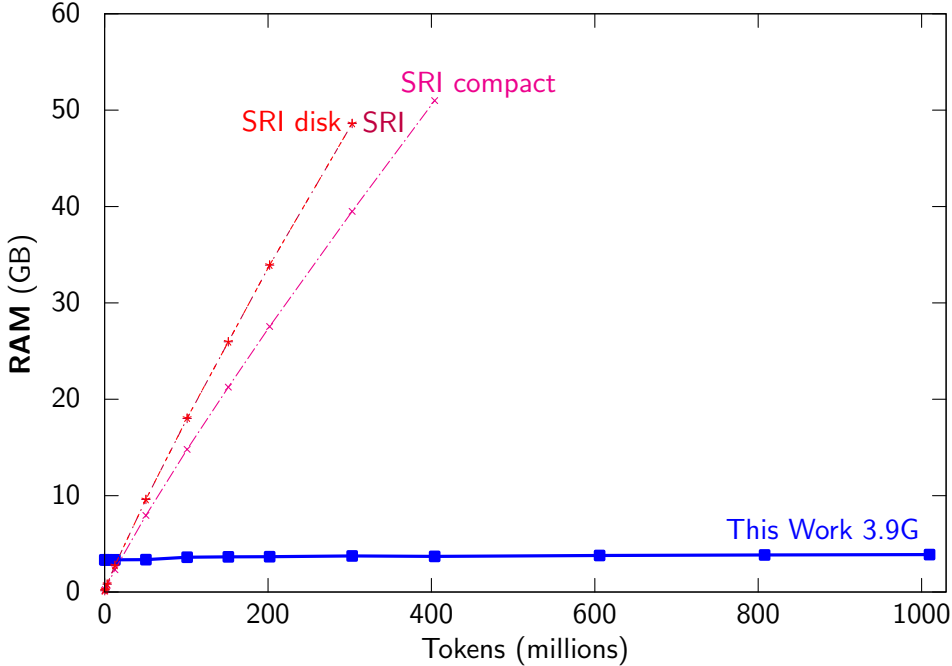
Task Build an unpruned 5-gram LM

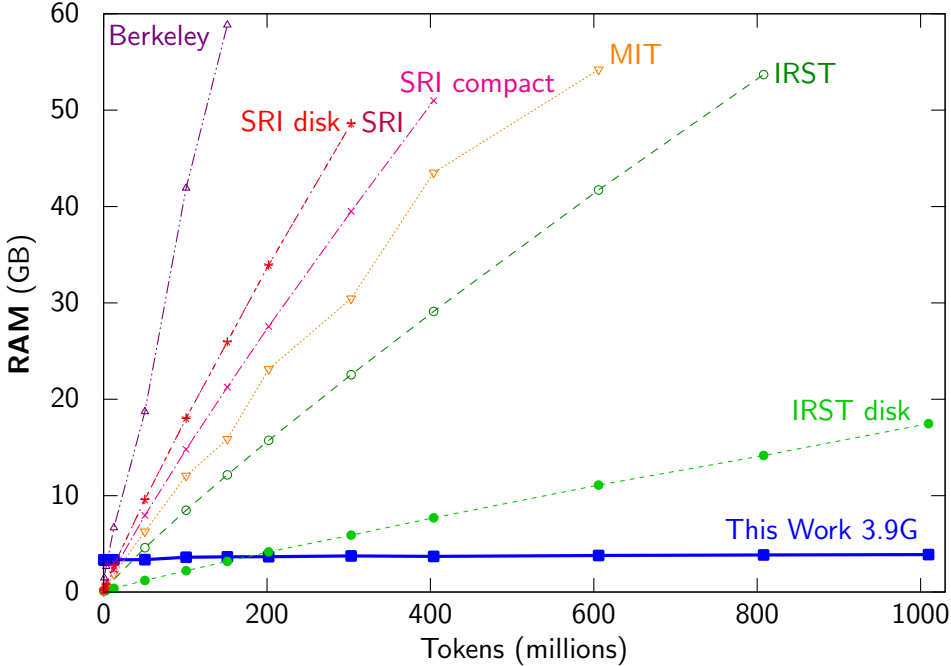
Data Subset of English ClueWeb09 (webpages)

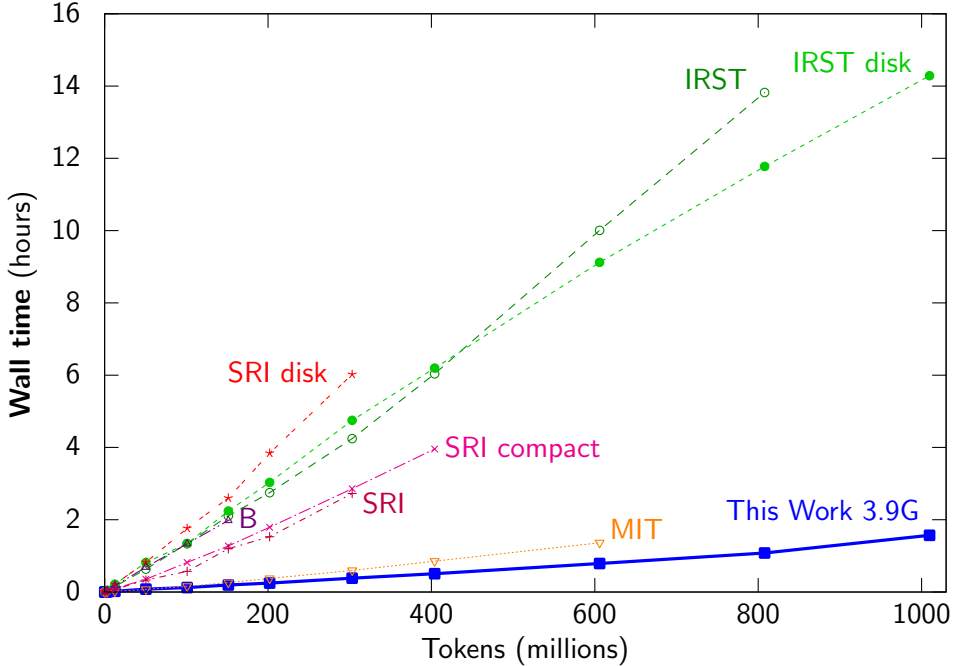
Machine 64 GB RAM

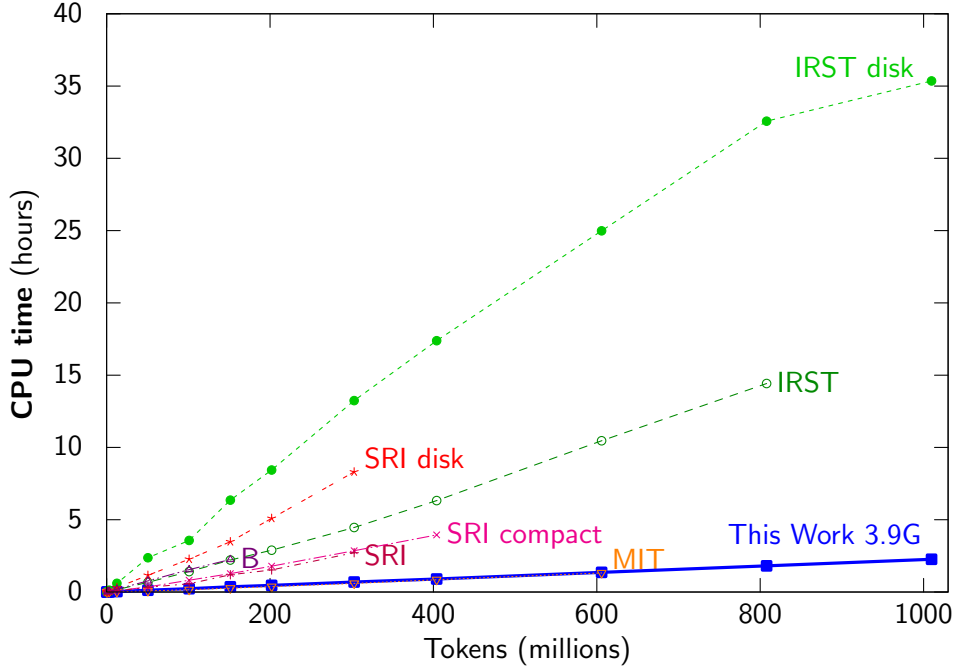
Output Format Binary (or ARPA when faster)

IRST disk: 3-way split. Peak RAM of any one process (as if run serially).
Berkeley: Binary search for minimum JVM memory.









Scaling

| | Tokens | Smoothing | Machines | Days |
|------------------|---------------|------------------|-----------------|-------------|
| This Work | 126 billion | Kneser-Ney | 1 | 2.8 |

Counts

| | 1 | 2 | 3 | 4 | 5 |
|-------------------------|----------|----------|----------|----------|----------|
| This Work 126B | 393m | 3,775m | 17,629m | 39,919m | 59,794m |
| Pruned Google 1T | 14m | 315m | 977m | 1,313m | 1,176m |

(This work used a machine with 140 GB RAM and a RAID5 array.)

Scaling

| | Tokens | Smoothing | Machines | Days | Year |
|------------------|---------------|------------------|-----------------|-------------|-------------|
| This Work | 126 billion | Kneser-Ney | 1 | 2.8 | 2013 |
| Google | 31 billion | Kneser-Ney | 400 | 2 | 2007 |
| Google | 230 billion | Kneser-Ney | ? | ? | 2013 |
| Google | 1800 billion | Stupid | 1500 | 1 | 2007 |

Counts

| | 1 | 2 | 3 | 4 | 5 |
|-------------------------|----------|----------|----------|----------|----------|
| This Work 126B | 393m | 3,775m | 17,629m | 39,919m | 59,794m |
| Pruned Google 1T | 14m | 315m | 977m | 1,313m | 1,176m |

(This work used a machine with 140 GB RAM and a RAID5 array.)

Workshop on Statistical MT Results

- 1 Compress the big LM to 676 GB
- 2 Decode with 1 TB RAM
- 3 Make three WMT 2013 submissions

| | Czech–English | | French–English | | Spanish–English | |
|------------------|----------------------|-------|-----------------------|-------|------------------------|-------|
| | Rank | BLEU | Rank | BLEU | Rank | BLEU |
| This Work | 1 | 28.16 | 1 | 33.37 | 1 | 32.55 |
| Google | 2–3 | 27.11 | 2–3 | 32.62 | 2 | 33.65 |
| Baseline | 3–5 | 27.38 | 2–3 | 32.57 | 3–5 | 31.76 |

Future Work on Estimation

- Pruning
- Linearly interpolate separately trained models
→ SRI's ARPA output is misleading.
- More smoothing methods
- Parallelization by data splitting

Outline

| | Speed | RAM | Published |
|------------------------|--------------|--------------|--|
| 1 Estimation from Text | 7.1x | 0.07x | ACL 2013 |
| 2 Raw Queries | 2.4x | 0.57x | WMT 2011 |
| 3 Decoding | 3.2–10.0x | 0.85x | { IWSLT 2011, EMNLP 2012, NAACL 2013 |

Raw Queries

Answer language model queries using less time and memory.

$$\log p(\text{iran} \mid \langle s \rangle) = -3.33437$$

$$\log p(\text{is} \mid \langle s \rangle \text{ iran}) = -1.05931$$

$$\log p(\text{one} \mid \langle s \rangle \text{ iran is}) = -1.80743$$

$$\log p(\text{of} \mid \langle s \rangle \text{ iran is one}) = -0.03705$$

$$\log p(\text{the} \mid \text{iran is one of}) = -0.08317$$

$$\log p(\text{few} \mid \text{is one of the}) = -1.20788$$

Example Language Model

Unigrams

| Words | log p | log b |
|-------|-----------|-------|
| <s> | $-\infty$ | -2.0 |
| iran | -4.1 | -0.8 |
| is | -2.5 | -1.4 |
| one | -3.3 | -0.9 |
| of | -2.5 | -1.1 |

Bigrams

| Words | log p | log b |
|----------|-------|-------|
| <s> iran | -3.3 | -1.2 |
| iran is | -1.7 | -0.4 |
| is one | -2.0 | -0.9 |
| one of | -1.4 | -0.6 |

Trigrams

| Words | log p |
|-------------|-------|
| <s> iran is | -1.1 |
| iran is one | -2.0 |
| is one of | -0.3 |

Example Queries

Unigrams

| Words | log p | log b |
|-------|-----------|-------|
| <s> | $-\infty$ | -2.0 |
| iran | -4.1 | -0.8 |
| is | -2.5 | -1.4 |
| one | -3.3 | -0.9 |
| of | -2.5 | -1.1 |

Bigrams

| Words | log p | log b |
|----------|-------|-------|
| <s> iran | -3.3 | -1.2 |
| iran is | -1.7 | -0.4 |
| is one | -2.0 | -0.9 |
| one of | -1.4 | -0.6 |

Trigrams

| Words | log p |
|-------------|-------|
| <s> iran is | -1.1 |
| iran is one | -2.0 |
| is one of | -0.3 |

Query: <s> iran is

$$\log p(\text{is} \mid \text{<s> iran}) = -1.1$$

Query: iran is of

$$\begin{array}{r} \log p(\text{of}) \quad \quad \quad -2.5 \\ \log b(\text{is}) \quad \quad \quad -1.4 \\ \log b(\text{iran is}) \quad \quad + -0.4 \\ \hline \log p(\text{of} \mid \text{iran is}) = -4.3 \end{array}$$

Trie Based

| | |
|---------------|---|
| CMU-Cambridge | Early implementation |
| SRI | Popular, considered fast, high-memory |
| IRST | Smaller than SRI, single-threaded |
| MIT | Batch querying |
| TPT | Memory locality |
| Joshua | Java, “not as scalable as the SRILM” [Li et al] |
| Berkeley | Java |

Trie Based

| | |
|---------------|---|
| CMU-Cambridge | Early implementation |
| SRI | Popular, considered fast, high-memory |
| IRST | Smaller than SRI, single-threaded |
| MIT | Batch querying |
| TPT | Memory locality |
| Joshua | Java, “not as scalable as the SRILM” [Li et al] |
| Berkeley | Java |

Lossy Low-Memory

| | |
|--------|---|
| Rand | Bloom maps |
| Shef | Minimal perfect hashing |
| Google | Minimal perfect hashing, larger than Shef |

KenLM Features

- Faster than all baselines
- Lowest lossless memory
- Multithreaded
- Quick loading via memory mapping
- Easier to compile

Data Structures

Probing Fast.

Trie Small. But still fast.

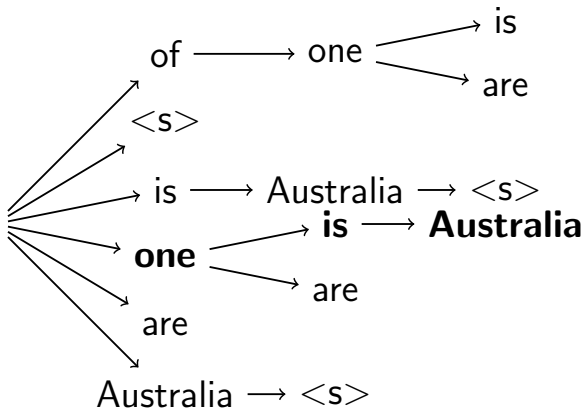
Probing

Hash every n -gram to a 64-bit integer. Ignore collisions.
Store n -grams in custom linear probing hash tables.

Fastest, 24 bytes/ n -gram (still less than SRI).

Trie

Reverse n -grams, arrange in a trie.



Smaller than most, faster than all but probing.

Optimizing the Trie

CPU Interpolation search instead of binary search [Yehoshua et al, 1978]

RAM Pack at the bit level i.e. $\log p$ has no sign bit

Options to Save More Memory

Cluster floats into 2^f bins, store f bits/float.



Chop bits from integer sequences, store offsets.

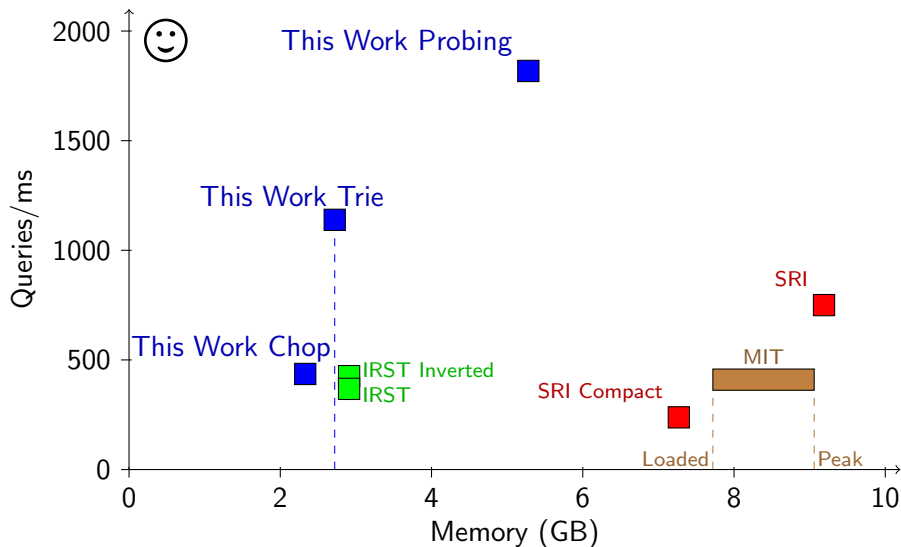
[Whittaker and Raj, 2001; Raj and Whittaker, 2003]

Experiment: Raw Queries

Task Score the English Gigaword corpus
Model 5-gram Europarl + deduped news crawl 2011

Queries/ms Excludes loading and file reading time
Loaded RAM Resident after loading
Peak RAM Peak virtual after scoring

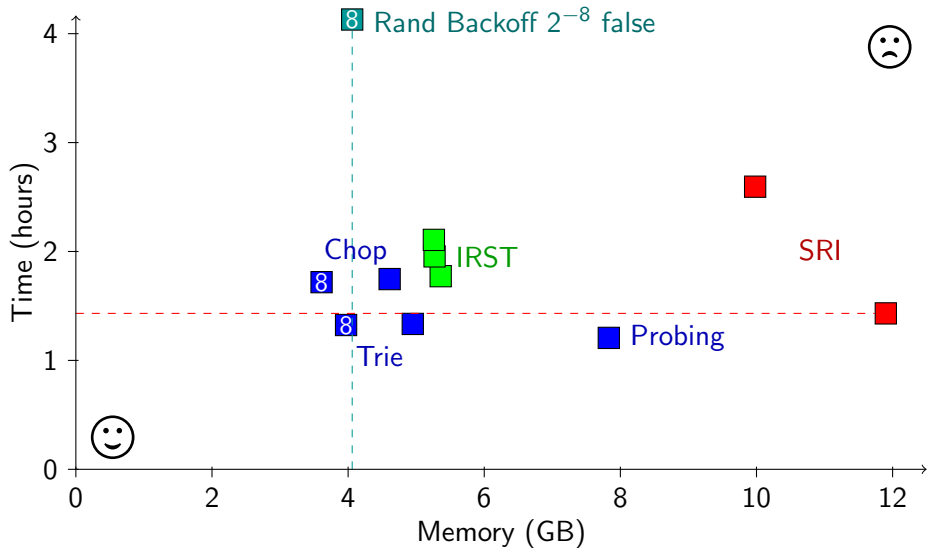
Raw Queries: Exact Variants



Experiment: Translate 3003 Sentences

| | |
|-----------|--|
| Task | WMT 2011 French–English baseline |
| Decoder | Moses |
| Model | 5-gram Europarl+News LM (same as before) |
| Formalism | Phrase-based from Europarl |
| Time | Total wall time, including loading |
| Memory | Total resident memory after decoding |

Moses Benchmarks: Single Threaded



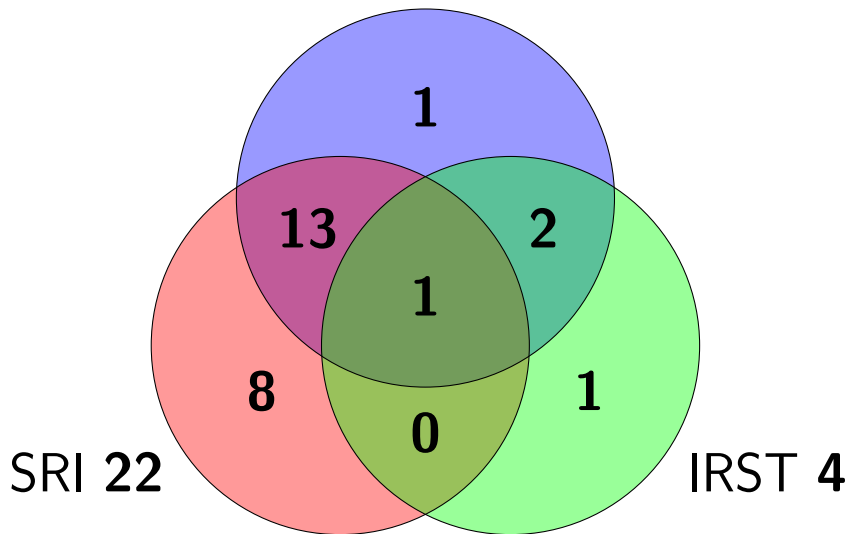
New: Yasuhara et al, EMNLP 2013

“An Efficient Language Model Using Double-Array Structures”

- 19% less RAM and 4-9% faster than this work’s probing method.
- More RAM than the trie method.
- 4 days to build a data structure with 936 million n -grams.

WMT 2013 Adoption: Any Task

This Work **17**



Outline

| | Speed | RAM | Published |
|------------------------|------------------|--------------|---|
| 1 Estimation from Text | 7.1x | 0.07x | ACL 2013 |
| 2 Raw Queries | 2.4x | 0.57x | WMT 2011 |
| 3 Decoding | 3.2–10.0x | 0.85x | { IWSLT 2011, EMNLP 2012, NAACL 2013 |

Decoding performance includes $\approx 1.15x$ speedup from raw queries.
Baseline: SRILM and cube pruning (more later).

Parsing-Based MT is Slow

26 CPU hours to translate 3000 sentences

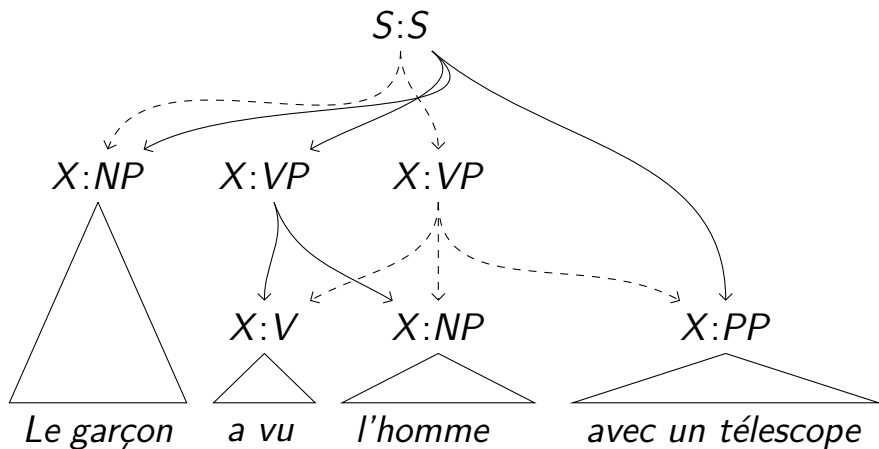


French–English system from Ammar et al [2013] using cdec, a 4-gram LM, and cube pruning with beam size 200.

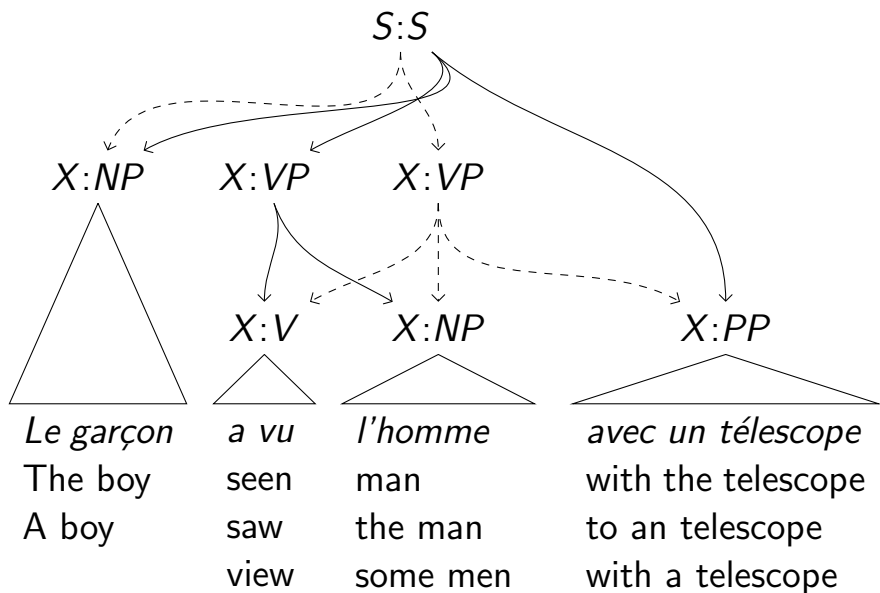
Decoding Example: Input

Le garçon a vu l'homme avec un télescope

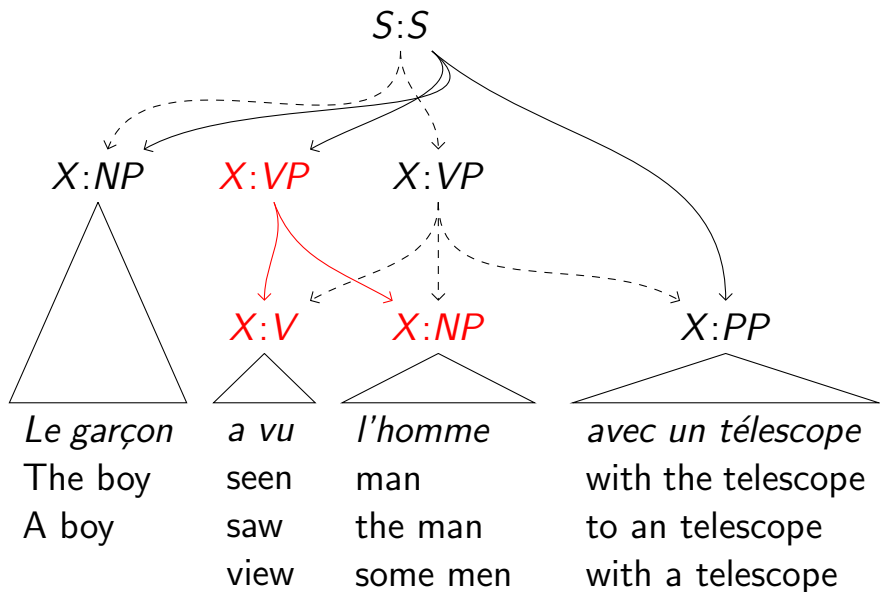
Decoding Example: Parse with SCFG



Decoding Example: Read Target Side



Decoding Example: One Constituent



$X:VP$



$X:V$

$X:NP$

a vu

l'homme

Hyp

seen

saw

view

Hyp

man

the man

some men

$X:VP$

$X:VP$

$X:V$

$X:NP$

a vu l'homme

a vu

l'homme

Hyp

seen

saw

view

Hyp

man

the man

some men

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

$X:VP$

$X:VP$

$X:V$

$X:NP$

a vu l'homme

a vu

l'homme

Hyp Score

Hyp

Score

seen -3.8

man -3.6

saw -4.0

the man -4.3

view -4.0

some men -6.3



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

$X:VP$

$X:VP$

$X:V$

$X:NP$

a vu l'homme

a vu

l'homme

Hyp Score

Hyp

Score

seen -3.8

man

-3.6

saw -4.0

the man

-4.3

view -4.0

some men

-6.3



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

$X:VP$

$X:VP$

$X:V$

$X:NP$

a vu l'homme

a vu

l'homme

Hyp Score

Hyp

Score

seen -3.8

man -3.6

saw -4.0

the man -4.3

view -4.0

some men -6.3

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

$X:VP$

$X:VP$

$X:V$

$X:NP$

a vu l'homme

a vu

l'homme



Hyp Score

Hyp

Score

seen -3.8

man

-3.6

saw -4.0

the man

-4.3

view -4.0

some men

-6.3

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~~

~~-10.2~~

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:

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

Log probability is part of the score

- ⇒ Scores do not sum
- ⇒ Local decisions may not be globally optimal
- ⇒ Search is hard.

Outline

| | Speed | RAM | Published |
|---------------------------|--------------|------------|------------------|
| 1 Estimation from Text | 7.1x | 0.07x | ACL 2013 |
| 2 Raw Queries | 2.4x | 0.57x | WMT 2011 |
| 3 Decoding | 3.2–10.0x | 0.85x | |
| 1 State and Recombination | | | IWSLT 2011 |
| 2 Score Estimates | | | EMNLP 2012 |
| 3 New Search Algorithm | | | NAACL 2013 |

Appending Strings

Hypotheses are built by string concatenation.

Language model probability changes when this is done:

$$c(\text{saw} \bullet \text{the man}) = \frac{p(\text{saw the man})}{p(\text{saw})p(\text{the man})} = \frac{p(\text{the} | \text{saw})p(\text{man} | \text{saw the})}{p(\text{the}) p(\text{man} | \text{the})}$$

What words does correction c examine?

Markov Assumption

A 5-gram language model uses up to 4 words of context:
 $p(\text{man} \mid \langle s \rangle \text{ the boy saw the}) = p(\text{man} \mid \text{the boy saw the})$



Correction c examines up to 4 words from each string:

$c(\langle s \rangle \boxed{\text{the boy saw the}} \bullet \boxed{\text{man with a telescope}} \text{.})$

Right State

Left State

End of state

End of state

Hypotheses have State

$X:X$

Left State

countries that maintain diplomatic \neg

countries that maintain diplomatic \neg

countries that maintain some \neg

relations

ties

diplomatic ties

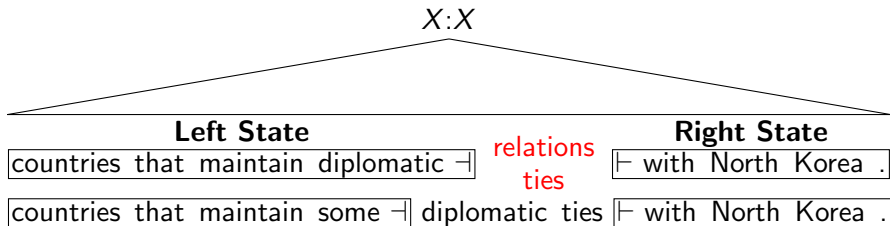
Right State

\vdash with North Korea .

\vdash with North Korea .

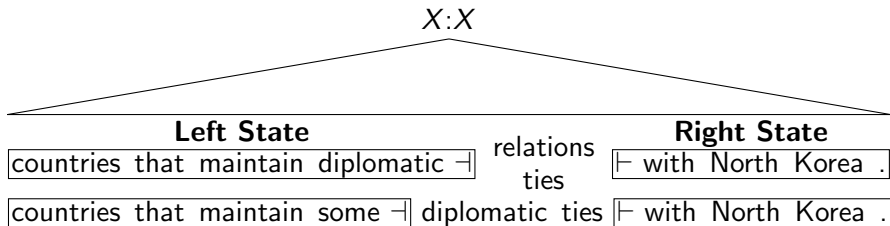
\vdash with North Korea .

State Controls Recombination



The decoder may recombine hypotheses with equal state.

State Controls Recombination



The decoder may recombine hypotheses with equal state.

Smaller state

- \Rightarrow More recombination
- \Rightarrow Reason over more hypotheses at once
- \Rightarrow **Improved time-accuracy tradeoff.**

Efficiently Minimizing State

Li et al [2008] Criterion for state minimization.
“Inefficient implementation”

This Work (IWSLT 2011) Repurpose log probability sign bit.
Use existing lookups.
Encode state to make queries faster.

Efficiently Minimizing State

Li et al [2008] Criterion for state minimization.
“Inefficient implementation”

This Work (IWSLT 2011) Repurpose log probability sign bit.
Use existing lookups.
Encode state to make queries faster.

11% faster than right state minimization alone.
On hierarchical Chinese–English with beam size 1000.

Baseline: How to Score a Fragment

$$\begin{array}{rcl} \log p_5(\text{is}) & = & -2.63 \\ \log p_5(\text{one} \mid \text{is}) & = & -2.03 \\ \log p_5(\text{of} \mid \text{is one}) & = & -0.24 \\ \log p_5(\text{the} \mid \text{is one of}) & = & -0.47 \\ + \log p_5(\text{few} \mid \text{is one of the}) & = & -1.26 \\ \hline = \log p_5(\text{is one of the few}) & = & -6.62 \end{array}$$

The Problem: Lower Order Entries

5-Gram Model: $\log p_5(is) = -2.63$

Unigram Model: $\log p_1(is) = -2.30$

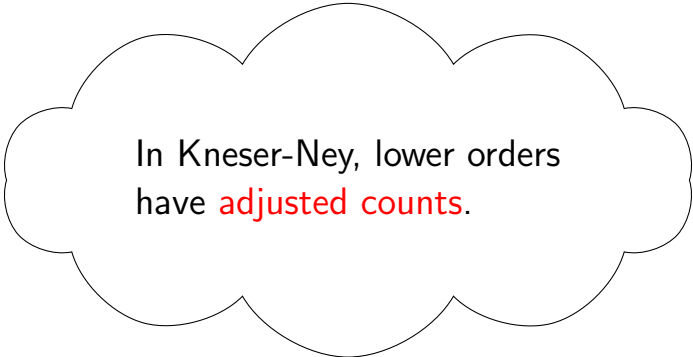
Same training data.

The Problem: Lower Order Entries

5-Gram Model: $\log p_5(is) = -2.63$

Unigram Model: $\log p_1(is) = -2.30$

Same training data.



In Kneser-Ney, lower orders
have **adjusted counts**.

Build One Model For Each Order

| | Baseline | Lower |
|--|-----------------|-------------------------------|
| $\log p_5(\text{is})$ | $= -2.63$ | $-2.30 = \log p_1$ |
| $\log p_5(\text{one} \mid \text{is})$ | $= -2.03$ | $-1.92 = \log p_2$ |
| $\log p_5(\text{of} \mid \text{is one})$ | $= -0.24$ | $-0.08 = \log p_3$ |
| $\log p_5(\text{the} \mid \text{is one of})$ | $= -0.47$ | $-0.21 = \log p_4$ |
| $+ \log p_5(\text{few} \mid \text{is one of the})$ | $= -1.26$ | $-1.26 = \log p_5$ |
| <hr/> $= \log p_5(\text{is one of the few})$ | $= -6.62$ | $-5.77 = \log p_{\text{Low}}$ |

Storing Lower Order Models

One extra float per entry, except for longest order.

Unigrams

| Words | $\log p_5$ | $\log b_5$ | $\log p_1$ |
|-----------|------------|------------|------------|
| australia | -3.9 | -0.6 | -3.6 |
| is | -2.6 | -1.5 | -2.3 |
| one | -3.4 | -1.0 | -2.9 |
| of | -2.5 | -1.1 | -1.7 |

No need for backoff b_1

If backoff occurs, use of p_5 is appropriate.

Storing Lower Order Models

One extra float per entry, except for longest order.

Unigrams

| Words | $\log p_5$ | $\log b_5$ | $\log p_1$ |
|-----------|------------|------------|------------|
| australia | -3.9 | -0.6 | -3.6 |
| is | -2.6 | -1.5 | -2.3 |
| one | -3.4 | -1.0 | -2.9 |
| of | -2.5 | -1.1 | -1.7 |

No need for backoff b_1

If backoff occurs, use of p_5 is appropriate.

Related: store upper bounds [Carter et al, also EMNLP 2012].

Pessimism

Assume backoff all the way to unigrams.

$$q(\text{is one of}) = p(\text{is one of})b(\text{is one of})b(\text{one of})b(\text{of})$$

Sentence Scores Are Unchanged

$$q(\langle s \rangle \dots \langle /s \rangle) = p(\langle s \rangle \dots \langle /s \rangle)$$

$$\text{because } b(\dots \langle /s \rangle) = 1$$

Pessimism

Assume backoff all the way to unigrams.

$$q(\text{is one of}) = p(\text{is one of})b(\text{is one of})b(\text{one of})b(\text{of})$$

Sentence Scores Are Unchanged

$$q(\langle s \rangle \dots \langle /s \rangle) = p(\langle s \rangle \dots \langle /s \rangle)$$

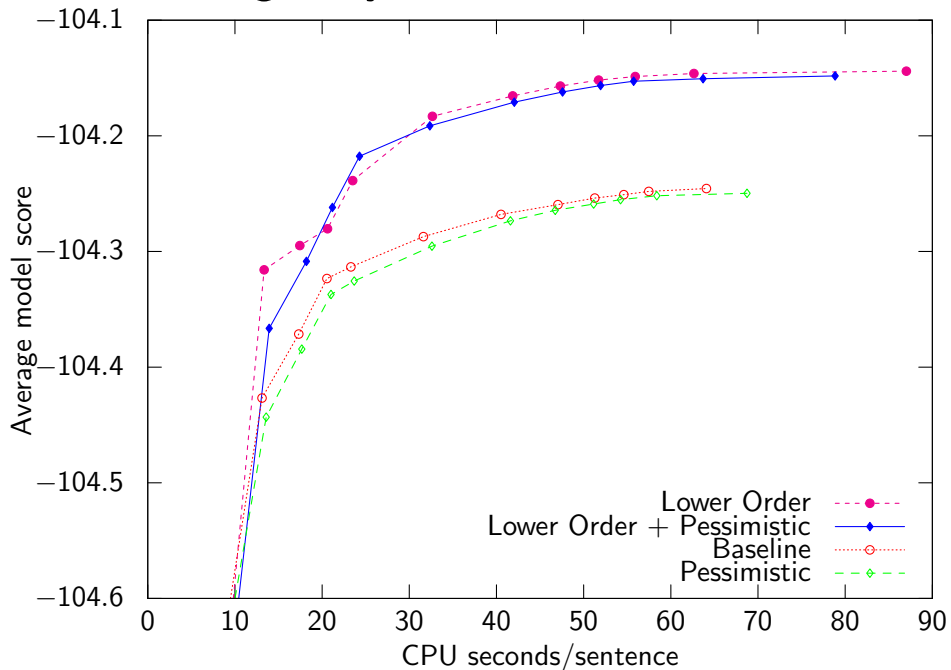
$$\text{because } b(\dots \langle /s \rangle) = 1$$

Telescoping

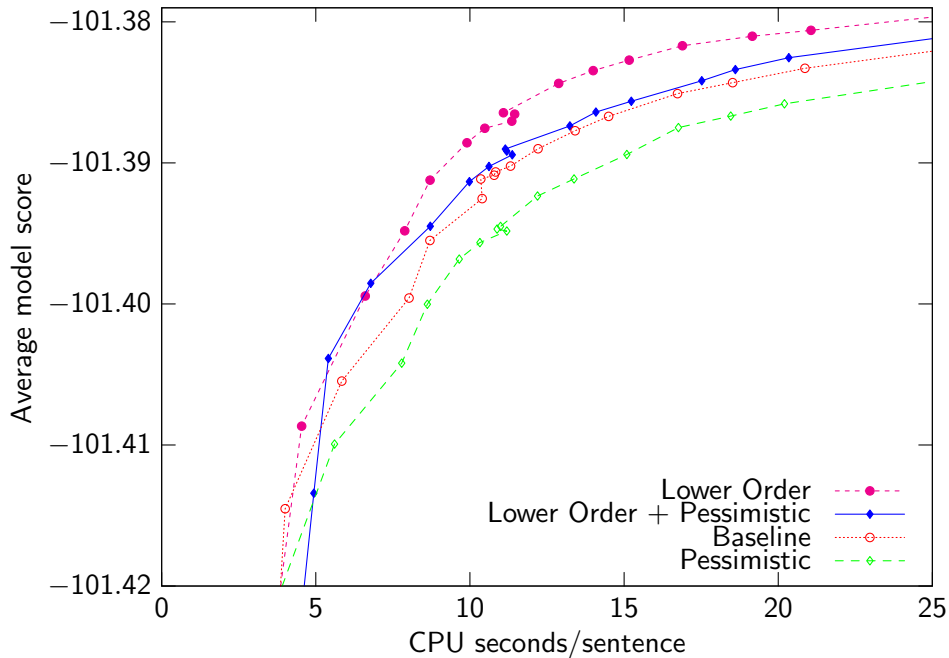
$$q(\text{is}) = p(\text{is})b(\text{is})$$

$$q(\text{one} \mid \text{is}) = p(\text{one} \mid \text{is}) \frac{b(\text{is one})b(\text{one})}{b(\text{is})}$$

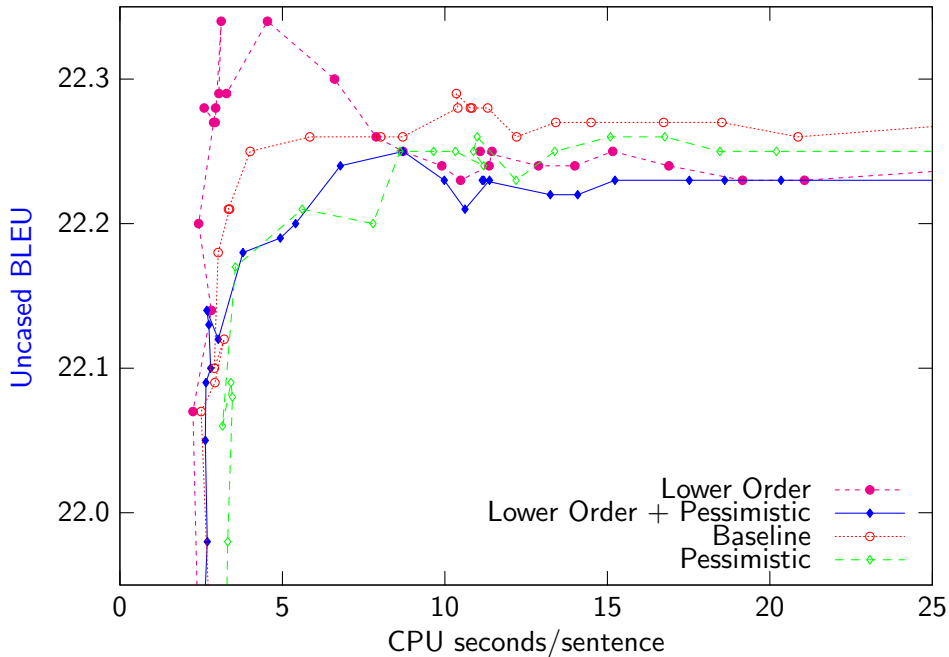
Target-Syntax Model Score



Hierarchical Model Score



Hierarchical BLEU



Memory

Effect of adding or removing a float per entry.

| Structure | Baseline (MB) | Change (MB) | % |
|----------------------------|----------------------|--------------------|----------|
| Probing | 4,072 | 517 | 13% |
| Trie | 2,647 | 506 | 19% |
| 8-bit quantized trie | 1,236 | 140 | 11% |
| 8-bit minimal perfect hash | 540 | 140 | 26% |

$X:VP$

$X:VP$

$X:V$

$X:NP$

a vu l'homme

a vu

l'homme

Hyp Score

Hyp

Score

seen -3.8

man

-3.6

saw -4.0

the man

-4.3

view -4.0

some men

-6.3



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 -10.2~~

Pruning is Approximate

Baseline: Cube Pruning [Chiang, 2007]

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 |

Baseline: Cube Pruning [Chiang, 2007]

| | | | | | | |
|-------------|------------|--------------|----------------|--------------|-----------------|-------------|
| | 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 | $-3.8 - 4.3 = -8.1$ |

Baseline: Cube Pruning [Chiang, 2007]

| | | | | | | |
|-------------|------------|-------------|----------------|-------------|-----------------|-------------|
| | 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 | Queue | | | | |

Queue

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

Problem With Cube Pruning

Hyp
is a
are a

Hypothesis
countries that
countries which
country

Hypothesis
is a countries that
are a countries that
are a countries which
:

No notion that “a countries” is bad.

Problem With Cube Pruning

Hyp
is a
are a

Hypothesis
countries that
countries which
country

Hypothesis
is a countries that
are a countries that
are a countries which
:

No notion that “a countries” is bad.

Idea: group by outermost words.

Example Hypotheses

Left State

countries that \neg

maintain diplomatic

relations
ties

countries that have \neg

an embassy in

\vdash DPR Korea .

country \neg

that maintains some diplomatic ties

\vdash in North Korea .

nations which has \neg

some diplomatic ties

\vdash with DPR Korea .

country \neg

that maintains some diplomatic ties

\vdash with DPR Korea .

Right State

\vdash with North Korea .

Example Hypotheses

Left State

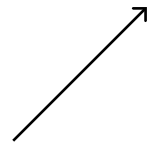
Right State

| | | | | |
|----------------------|---|---|---|---------------------|
| (countries that | ┆ | ◇ | ┆ | with North Korea .) |
| (nations which has | ┆ | ◇ | ┆ | with DPR Korea .) |
| (countries that have | ┆ | ◇ | ┆ | DPR Korea .) |
| (country | ┆ | ◇ | ┆ | in North Korea .) |
| (country | ┆ | ◇ | ┆ | with DPR Korea .) |

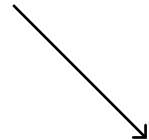
- ┆ Left state is completely present.
- ◇ Stands for elided words
- ┆ Right state is completely present.

Reveal Common Words in Each Group

(countries that \diamond Korea .)

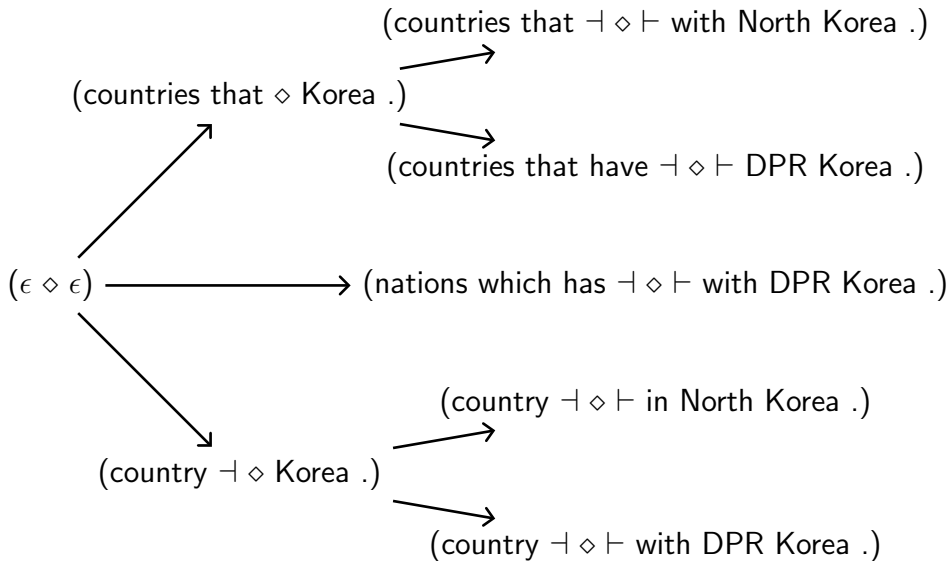


$(\epsilon \diamond \epsilon)$ \rightarrow (nations which has $\dashv \diamond \vdash$ with DPR Korea .)



(country $\dashv \diamond$ Korea .)

Alternate Sides Until Tree is Full



Using Rules

is a $X:NP1$ $\langle /s \rangle$

turns into

is a $(\epsilon \diamond \epsilon)$ $\langle /s \rangle$

$X:V1$ the $X:N2$

turns into

$(\epsilon \diamond \epsilon)$ the $(\epsilon \diamond \epsilon)$

$\underbrace{\hspace{1.5cm}}$

$X:V1$

$\underbrace{\hspace{1.5cm}}$

$X:N2$

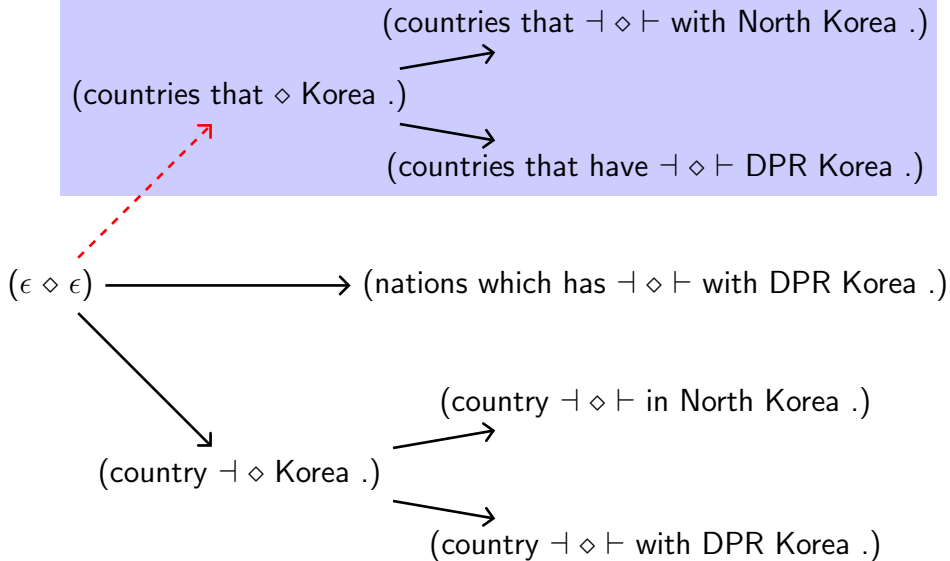
Exploring and Backtracking

Does the LM like “is a (countries that ◊ Korea .) </s>”?

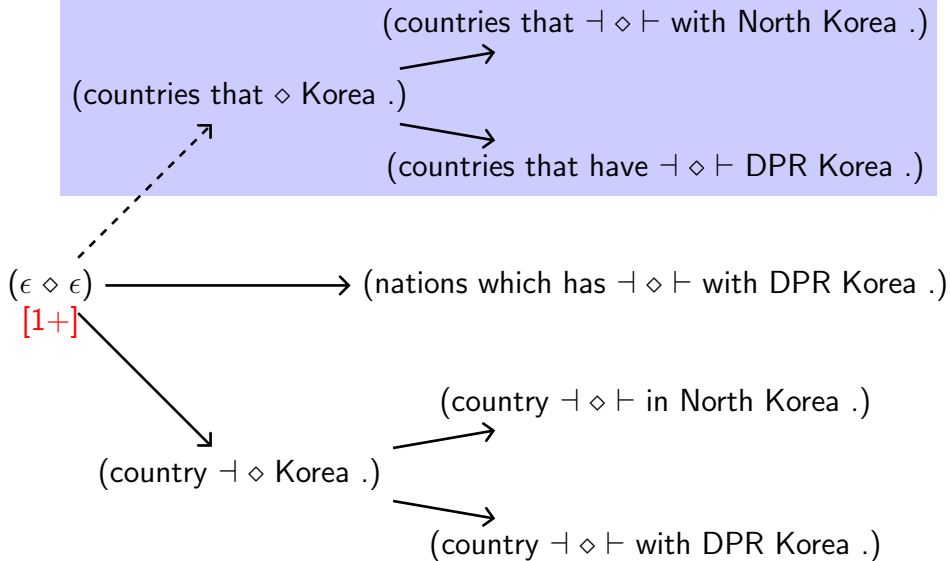
Yes Try more detail.

No Consider alternatives.

Split and Leave Breadcrumbs



Split and Leave Breadcrumbs



The queue entry

is a $(\epsilon \diamond \epsilon) \langle /s \rangle$

splits into

Zeroth Child “is a (countries that \diamond Korea .) $\langle /s \rangle$ ”

Other Children “is a $(\epsilon \diamond \epsilon)[1+]$ $\langle /s \rangle$ ”

Children except the zeroth.

Summary So Far

A priority queue contains competing entries:

is a (countries that \diamond Korea .) $\langle /s \rangle$

$(\epsilon \diamond \epsilon)$ the $(\epsilon \diamond \epsilon)$

is a $(\epsilon \diamond \epsilon)[1+]$ $\langle /s \rangle$

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

Summary So Far

A priority queue contains competing entries:

is a (countries that \diamond Korea .) $\langle /s \rangle$

$(\epsilon \diamond \epsilon)$ the $(\epsilon \diamond \epsilon)$

is a $(\epsilon \diamond \epsilon)[1+]$ $\langle /s \rangle$

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

Next: Scoring queue entries


Scores come from the best descendant:

$$\text{Score}(\epsilon \diamond \epsilon) = \text{Score}(\text{countries that } \vdash \diamond \vdash \text{ with North Korea .})$$

$$\geq$$

$$\text{Score}(\epsilon \diamond \epsilon)[1+] = \text{Score}(\text{nations which has } \vdash \diamond \vdash \text{ with DPR Korea .})$$

Estimates Update as Words are Revealed

is a ($\epsilon \diamond \epsilon$) $\langle /s \rangle$  is a (countries that \diamond Korea .) $\langle /s \rangle$

$p(\text{is})$

$p(\text{a} \mid \text{is})$

$p(\text{countries})$

$p(\text{that} \mid \text{countries})$

$p(\langle /s \rangle)$

$p(\text{is})$

$p(\text{a} \mid \text{is})$

$p(\text{countries} \mid \text{is a})$

$p(\text{that} \mid \text{is a countries})$

$p(\langle /s \rangle \mid \text{Korea .})$

Summary: Processing a Constituent

- 1 **Initialize:** Push rules onto a priority queue.
- 2 **Best-First Loop:**
 - 1 Pop the top entry.
 - 2 If it's complete, add to the beam. Otherwise, split and push.
- 3 **Finalize:** Convert the beam to a tree (lazily).

Summary: Processing a Constituent

- 1 **Initialize:** Push rules onto a priority queue.
- 2 **Best-First Loop:**
 - 1 Pop the top entry.
 - 2 If it's complete, add to the beam. Otherwise, split and push.
- 3 **Finalize:** Convert the beam to a tree (lazily).

Process constituents in bottom-up order (like cube pruning).

Coarse-to-Fine [Zhang et al, 2008; Petrov et al, 2008]

Decode multiple times, each with more detail:

- LM order
- Word classes

Key Difference

| | |
|----------------|--------------------------|
| Coarse-to-Fine | Lock-step refinement |
| This Work | Locally refine on demand |

Future Work

- Use this work for each decoding pass
- Word classes for this work

Exact Algorithms

- Weighted finite state transducers [Iglesias et al, 2011]
- Integer linear programming [Rush et al, 2011]
- Later: upper bounds and LM refinement [Aziz et al, WMT 2013]

Currently intractable for large MT (7 hours for a 7-word sentence)

⇒ Used as first pass of *approximate* coarse-to-fine.

Exact Algorithms

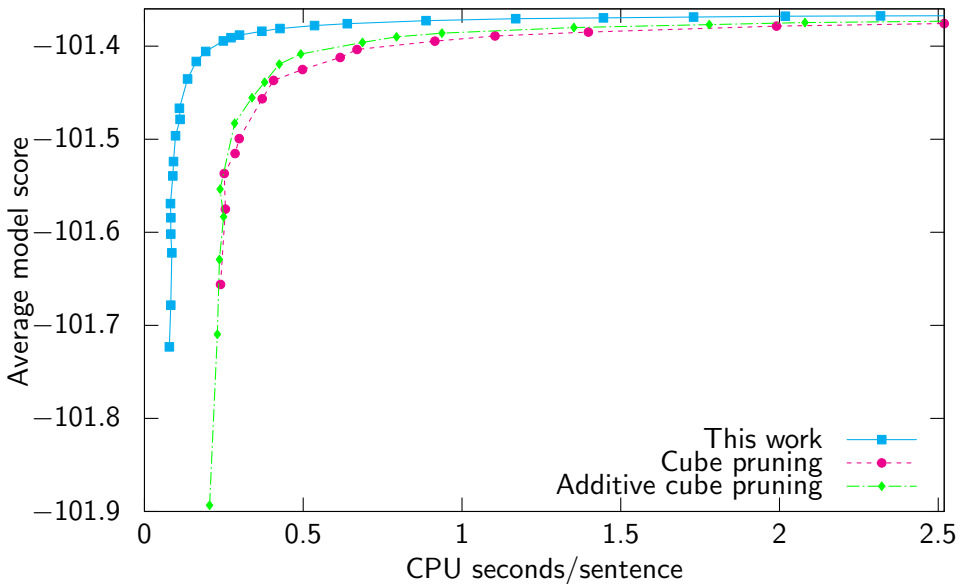
- Weighted finite state transducers [Iglesias et al, 2011]
- Integer linear programming [Rush et al, 2011]
- Later: upper bounds and LM refinement [Aziz et al, WMT 2013]

Currently intractable for large MT (7 hours for a 7-word sentence)
⇒ Used as first pass of *approximate* coarse-to-fine.

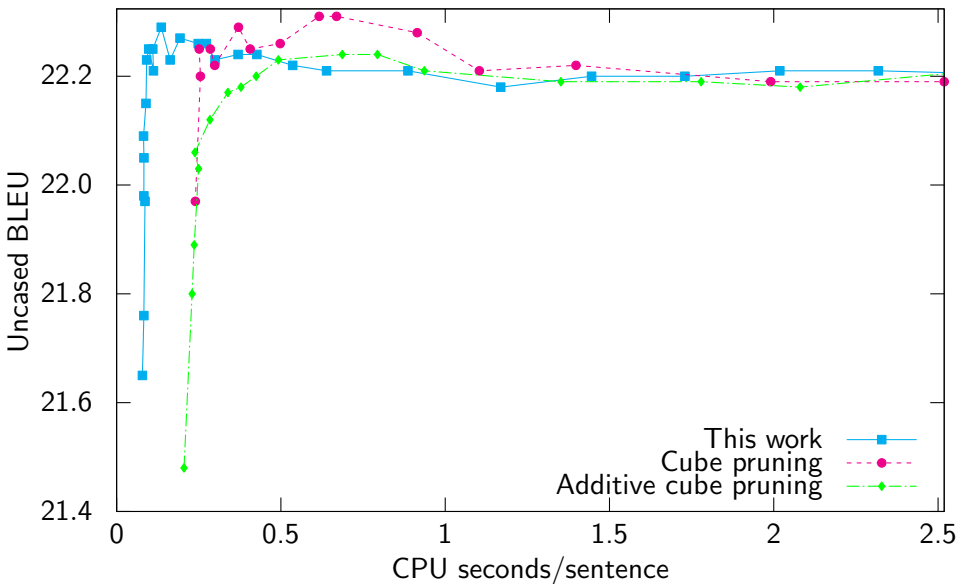
Key Difference

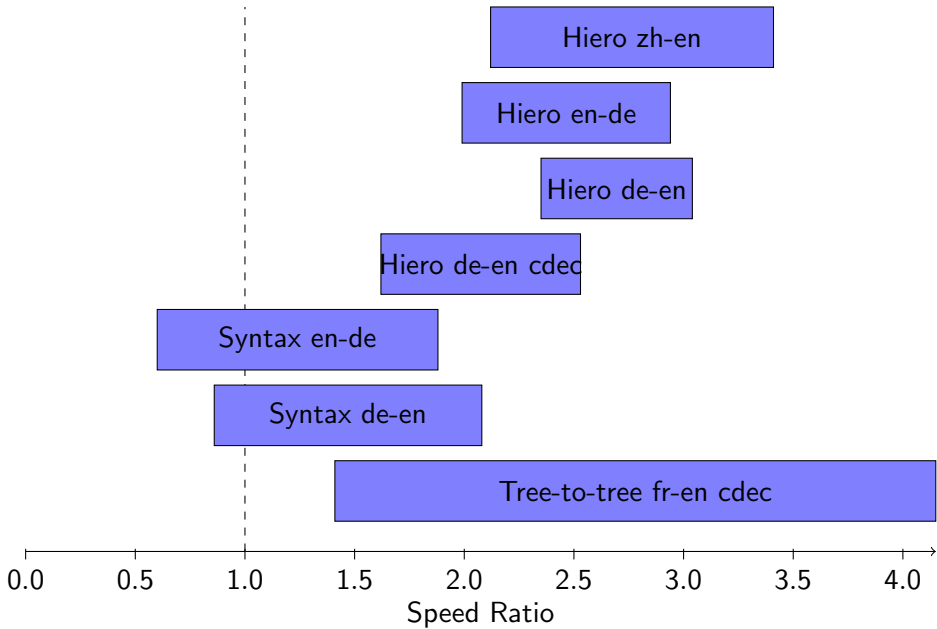
Approximation based on average-case scores *before* expanding hypotheses.

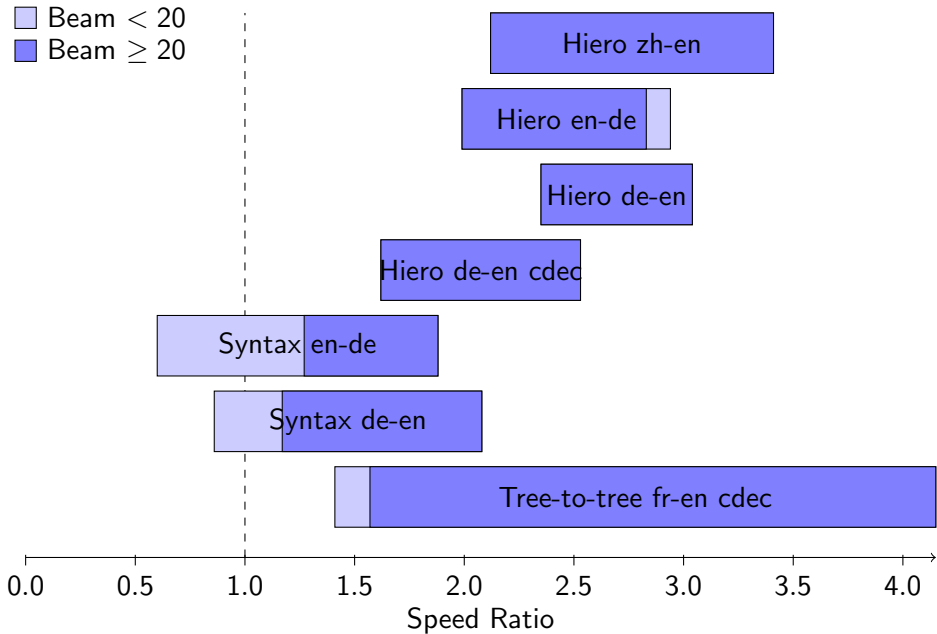
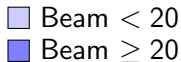
Moses Hierarchical



Moses Hierarchical







Summary

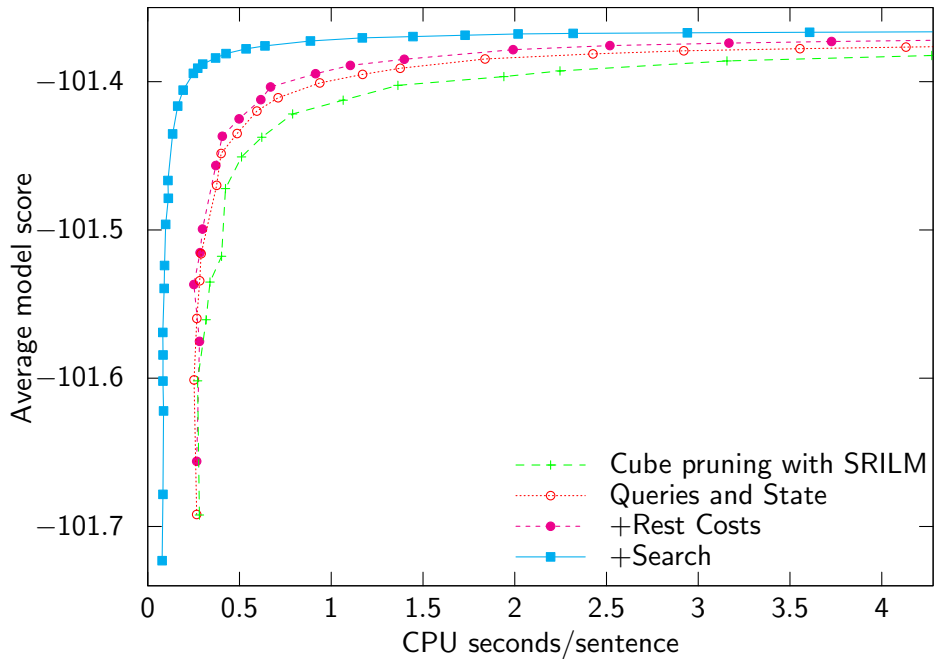
Optimized the entire LM pipeline from estimation to search.

Summary

Optimized the entire LM pipeline from estimation to search.

Comparison

| | |
|---------|-------------------------|
| Task | WMT 2011 German-English |
| Built | [Koehn et al, 2011] |
| Model | Hierarchical |
| Decoder | Moses |



Decoder Support

| | Moses | cdec | Joshua | Jane |
|------------|-------|------|--------|------|
| Estimation | ✓ | ✓ | ✓ | |
| Queries | ✓ | ✓ | ✓ | ✓ |
| State | ✓ | ✓ | ✓ | ✓ |
| Rest Costs | ✓ | ✓ | ✓ | ✓ |
| Search | ✓ | ✓ | | |

`kheafield.com/code`

Questions?

Full State Minimization

Keep only words that might form cross-hypothesis n -grams.

Left State [Joshua]

For any word w , does the model contain:

- w countries that maintain some ✗
- w countries that maintain ✗
- w countries that ✓

⇒ Left state minimizes to “countries that” ∪.

Full State Minimization

Keep only words that might form cross-hypothesis n -grams.

Left State [Joshua]

For any word w , does the model contain:

- w countries that maintain some ✗
- w countries that maintain ✗
- w countries that ✓

⇒ Left state minimizes to “countries that” \neg .

Right State [SRI, Rand, Joshua]

For any word w , does the model contain:

with North Korea . w ✓

⇒ Right state minimizes to \vdash “with North Korea .”

Related Work on State

Joshua Left and right but “inefficient implementation”
SRI Right only, additional lookups

This Work Repurpose memory, existing lookups
Also: encode state to make queries faster

Experimental Setup

| | |
|---------|--|
| Task | NIST Chinese–English [Koehn, 2011] |
| LM | Xinhua and AFP from English Gigaword 4 + Parallel Data |
| Grammar | Hierarchical |
| Decoder | Moses with cube pruning and faster raw queries |

Experimental Setup

| | |
|---------|--|
| Task | NIST Chinese–English [Koehn, 2011] |
| LM | Xinhua and AFP from English Gigaword 4 + Parallel Data |
| Grammar | Hierarchical |
| Decoder | Moses with cube pruning and faster raw queries |

11% faster

Experiments: Systems

Hierarchical with Moses [Koehn, 2012]

- German–English also ported to cdec, Joshua, and Jane
- English–German
- Chinese–English

Target Syntax with Moses [Koehn, 2012]

- German–English
- English–German

Tree-to-Tree with cdec [Ammar et al, 2013]

- French–English

Experiments: Systems and Scenarios

Hierarchical with Moses [Koehn, 2012]

- German–English also ported to cdec, Joshua, and Jane
- English–German
- Chinese–English

Target Syntax with Moses [Koehn, 2012]

- German–English
- English–German

Tree-to-Tree with cdec [Ammar et al, 2013]

- French–English

Baseline and improved rest costs, 2–3 flavors of cube pruning.