# CS230: Lecture 10
# Sequence models II

Kian Katanforoosh, Andrew Ng, Younes Bensouda Mourri

# Today's outline

We will learn how to:

- Automatically **score an NLP model**

- **Improve Machine Translation** results with Beam search

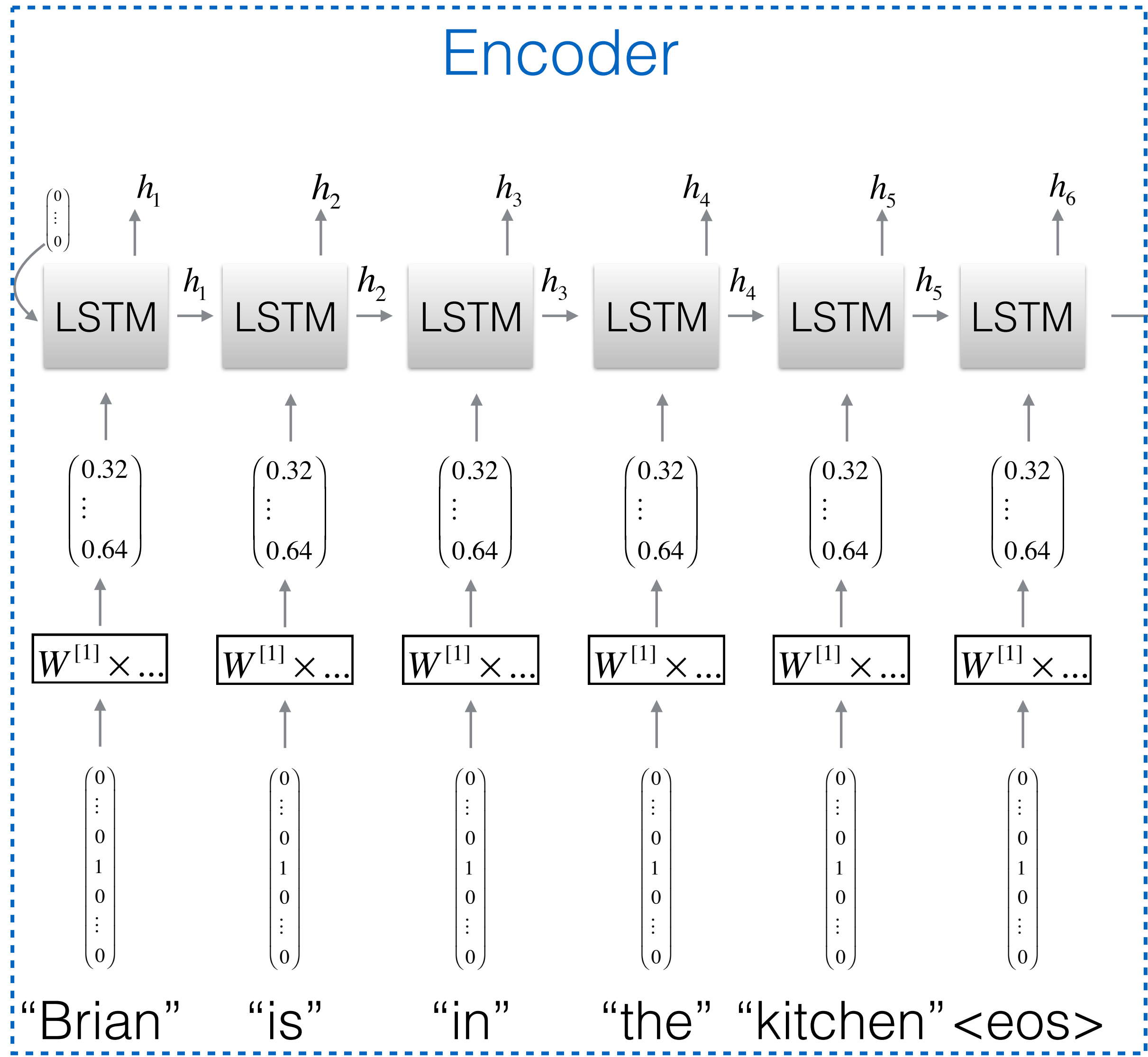- Build a **speech recognition** application
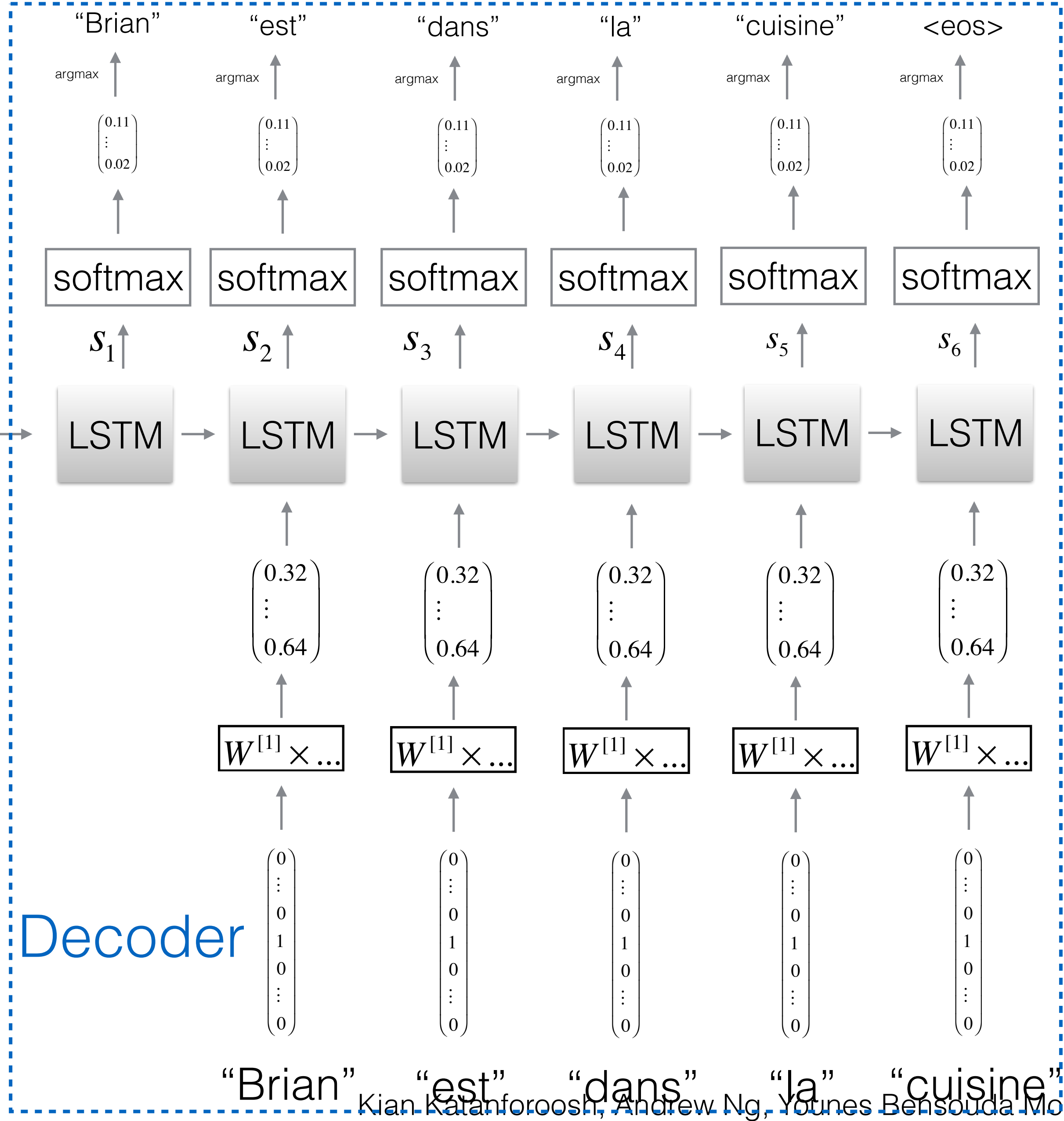
I. BLEU score

II. Beam Search

III. Speech Recognition

Kian Katanforoosh, Andrew Ng, Younes Bensouda Mourri

# BLEU score

## Neural Machine Translation



Encoder

Decoder

$h_1$ $h_2$ $h_3$ $h_4$ $h_5$ $h_6$

$s_1$ $s_2$ $s_3$ $s_4$ $s_5$ $s_6$

"Brian" "est" "dans" "la" "cuisine" <eos>

argmax argmax argmax argmax argmax argmax

softmax softmax softmax softmax softmax softmax

"Brian" "is" "in" "the" "kitchen" <eos>

"Brian" "est" "dans" "la" "cuisine"

Kian Katanforoosh, Andrew Ng, Younes Bensouda Mourri

# BLEU score

*"The baby to be walking by him"*

*"The baby walks by himself"*

<span style="color:green">Motivation</span>:
Human Evaluation of MT are extensive but expensive.

<span style="color:green">Goal</span>:
Construct a quick, inexpensive, language independent (and correlates highly with human evaluation) method to automatically evaluate Machine Translation models.

<span style="color:green">Centrale idea:</span>
*"The closer a machine translation is to a professional human translation, the better it is."*

Kian Katanforoosh, Andrew Ng, Younes Bensouda Mourri

Needs two ingredients:
- a numerical "translation closeness" metric
- a corpus of good quality human reference translations

In speech recognition:
a successful metric is *word error rate.* BLEU's closeness metric was built after it.

Kian Katanforoosh, Andrew Ng, Younes Bensouda Mourri

# BLEU score

## Machine Translations

Candidate 1    It is a guide to action which ensures that the military always obeys the commands of the party

Candidate 2        It is to insure the troops forever hearing the activity guidebook that party directs.

## Human Translations

Reference 1        It is a guide to action that ensures that the military will forever heed Party commands

Reference 2        It is the guiding principle which guarantees the military forces always being under the command of the Party

Reference 3        It is the practical guide for the army always to heed the directions of the party

Kian Katanforoosh, Andrew Ng, Younes Bensouda Mourri

# BLEU score

## Unigram count as precision metric:

| Machine Translation |
|---|
| Candidate                                the the the the the the the |

| Human Translations |
|---|
| Reference 1                               The cat is on the mat. |
| Reference 2                               There is a cat on the mat. |

Standard Unigram Precision = $\dfrac{\text{\# MT words occurring in any reference HT}}{\text{\# MT words}}$ = 100%

Modified Unigram Precision = $\dfrac{\text{\# MT words occurring in any reference HT (clipped)}}{\text{\# MT words}}$ = 2/6 = 33.3%

Kishore Papineni et al., BLEU: a Method for Automatic Evaluation of Machine Translation, 2002              Kian Katanforoosh, Andrew Ng, Younes Bensouda Mourri

# BLEU score

| Machine Translations | |
|---|---|
| Candidate 1 | It is a guide to action which ensures that the military always obeys the commands of the party |
| Candidate 2 | It is to insure the troops forever hearing the activity guidebook that party directs. |

| Human Translations | |
|---|---|
| Reference 1 | It is a guide to action that ensures that the military will forever heed Party commands |
| Reference 2 | It is the guiding principle which guarantees the military forces always being under the command of the Party |
| Reference 3 | It is the practical guide for the army always to heed the directions of the party |

$$\text{Modified Unigram Precision (1)} = \frac{\text{\# MT words occurring in any reference HT (clipped)}}{\text{\# MT words}} = 17/18 = 94\%$$

$$\text{Modified Unigram Precision (2)} = \frac{\text{\# MT words occurring in any reference HT (clipped)}}{\text{\# MT words}} = 8/14 = 57\%$$

Kishore Papineni et al., BLEU: a Method for Automatic Evaluation of Machine Translation, 2002                    Kian Katanforoosh, Andrew Ng, Younes Bensouda Mourri

## Generalizing to modified n-gram precision metric:

Modified n-gram precision = $\dfrac{\text{\# MT n-grams occurring in any reference HT (clipped)}}{\text{\# MT n-grams}}$

| Machine Translations |
|---|
| Candidate 1    It is a guide to action which ensures that the military always obeys the commands of the party |
| Candidate 2          It is to insure the troops forever hearing the activity guidebook that party directs. |

| Human Translations |
|---|
| Reference 1          It is a guide to action that ensures that the military will forever heed Party commands |
| Reference 2          It is the guiding principle which guarantees the military forces always being under the command of the Party |
| Reference 3          It is the practical guide for the army always to heed the directions of the party |

Modified bi-gram precision (Candidate 1) =  10/17
Modified bi-gram precision (Candidate 2) =  1/13

Kishore Papineni et al., BLEU: a Method for Automatic Evaluation of Machine Translation, 2002

Kian Katanforoosh, Andrew Ng, Younes Bensouda Mourri

Generalizing from a sentence precision, to a corpus precision:

$$p_n = \frac{\displaystyle\sum_{C\in\{Candidates\}} \sum_{ngram\in C} Count_{clip}(ngram)}{\displaystyle\sum_{C'\in\{Candidates\}} \sum_{ngram'\in C} Count(ngram')}$$

Kishore Papineni et al., BLEU: a Method for Automatic Evaluation of Machine Translation, 2002
Kian Katanforoosh, Andrew Ng, Younes Bensouda Mourri

# BLEU Score



Figure 1: Distinguishing Human from Machine
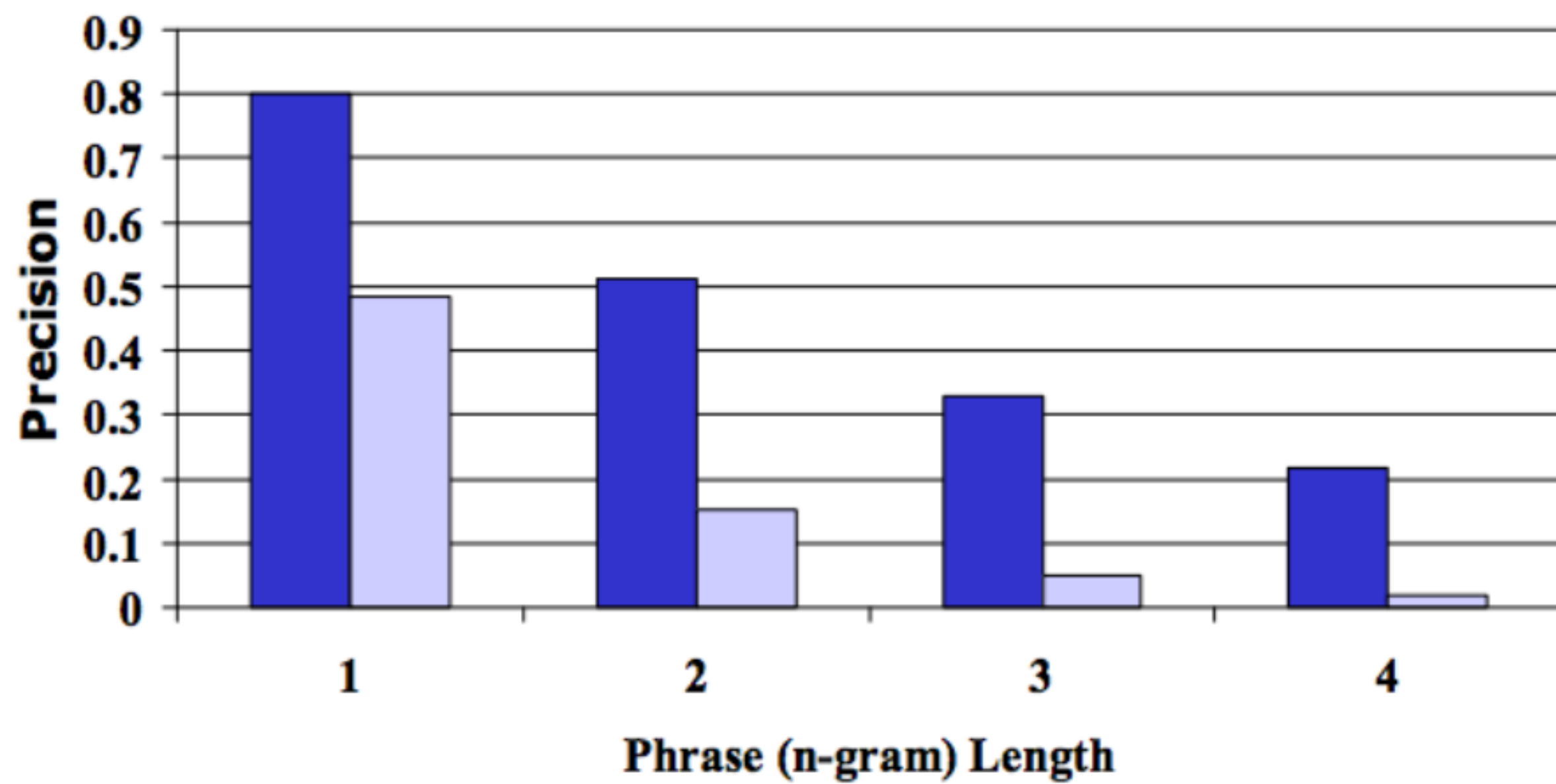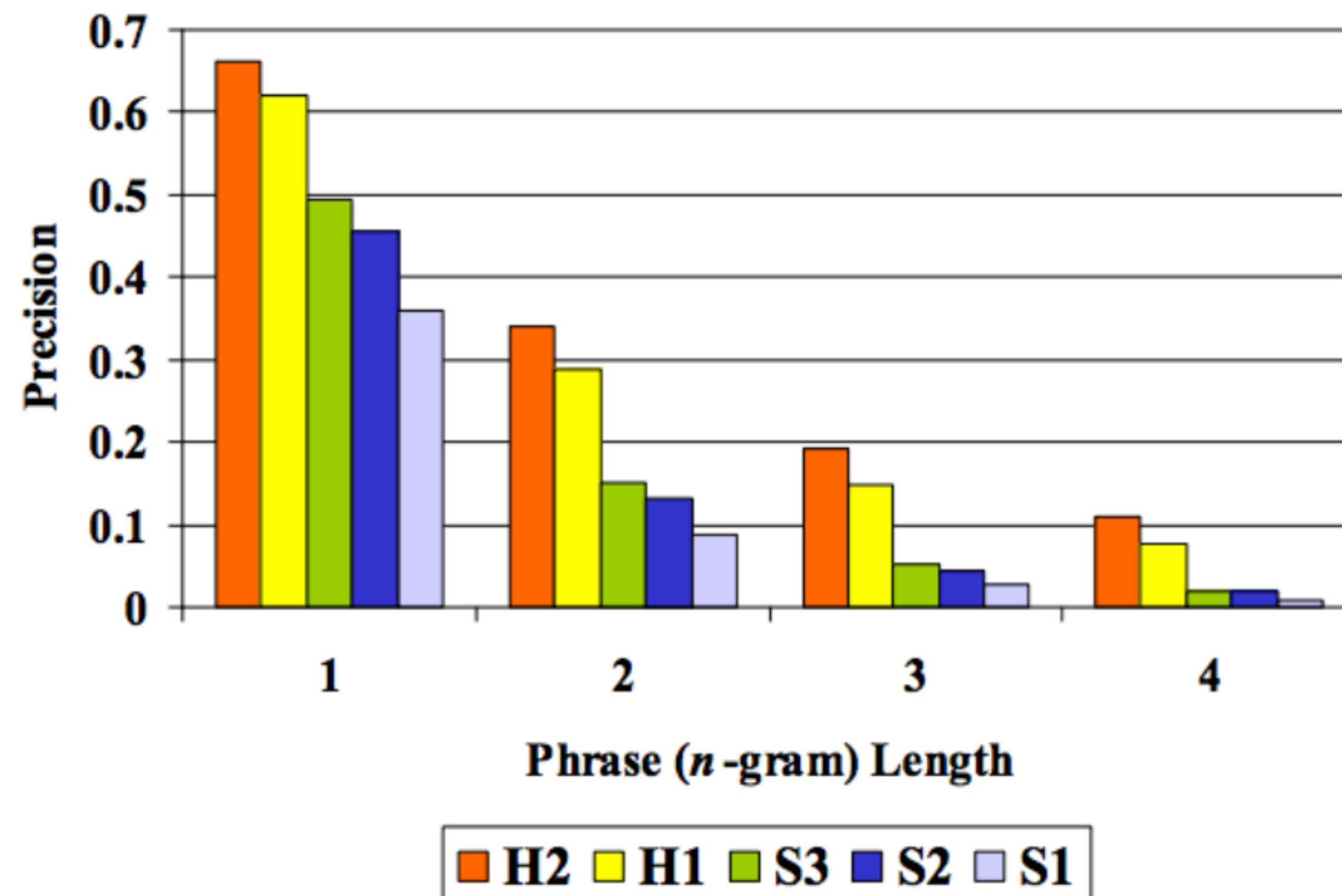


Figure 2: Machine and Human Translations

H2  H1  S3  S2  S1

Geometric average of modified n-grams precision measures = $\exp(\sum_{n=1}^{N} w_n \log(p_n)) = p_1^{w_1} p_2^{w_2} ... p_N^{w_N}$

Kishore Papineni et al., BLEU: a Method for Automatic Evaluation of Machine Translation, 2002

Kian Katanforoosh, Andrew Ng, Younes Bensouda Mourri

# BLEU Score

## Sentence length

### too short translation

| Machine Translations | | |
| --- | --- | --- |
| Candidate | | of the |

| Human Translations | |
| --- | --- |
| Reference 1 | It is a guide to action that ensures that the military will forever heed Party commands |
| Reference 2 | It is the guiding principle which guarantees the military forces always being under the command of the Party |
| Reference 3 | It is the practical guide for the army always to heed the directions of the party |

### too long translation

| Machine Translations | |
| --- | --- |
| Candidate 1 | I always invariably perpetually do |
| Candidate 2 | I always do |

| Human Translations | |
| --- | --- |
| Reference 1 | I always do |
| Reference 2 | I invariably do |
| Reference 3 | I perpetually do |

**Conclusion:** *longer translations are already penalized by the modified n-gram precision measure, not the too short translations*

$$\text{BP} = \text{Brevity penalty} = \text{decaying exponential} \sim \frac{\text{Total length of the machine's translation}}{\text{Total length of the candidate translation corpus}}$$

Kishore Papineni et al., BLEU: a Method for Automatic Evaluation of Machine Translation, 2002

Kian Katanforoosh, Andrew Ng, Younes Bensouda Mourri

BLEU definition:

$$BLEU = (p_1^{w_1} p_2^{w_2} ... p_N^{w_N}) \cdot BP \quad \text{where} \quad BP = \begin{cases} 1 & \text{if c > r} \\ e^{1-\frac{r}{c}} & \text{if c <= r} \end{cases}$$

log(BLEU) definition:

$$\log(BLEU) = \min(1 - \frac{r}{c}, 0) + \sum_{n=1}^{N} w_n \log(p_n)$$

Kishore Papineni et al., BLEU: a Method for Automatic Evaluation of Machine Translation, 2002

Kian Katanforoosh, Andrew Ng, Younes Bensouda Mourri

## Questions

- What is the main advantage of Beam search compared to other search algorithms?

*It is fast, and requires less computations.*

- What is the main disadvantage of Beam search compared to other search algorithms?

*It may not result in the optimal solution in terms of probability.*

- What is the time and memory complexity of Beam search?

*It is O(b\*Tx) in memory and O(b\*Tx) in time.*

# Speech Recognition Pipeline

Audio Data:

X? → model? → Y?

frequency

"Hello"

time

Kian Katanforoosh, Andrew Ng, Younes Bensouda Mourri

cross-entropy

H            E            L            L            O

argmax       argmax       argmax       argmax       argmax

probability
distribution
shape = (28,1)

$\begin{pmatrix} 0.31 \\ \vdots \\ 0.02 \end{pmatrix}$   $\begin{pmatrix} 0.03 \\ \vdots \\ 0.01 \end{pmatrix}$   $\begin{pmatrix} 0.10 \\ \vdots \\ 0.20 \end{pmatrix}$   $\begin{pmatrix} 0.08 \\ \vdots \\ 0.23 \end{pmatrix}$   $\begin{pmatrix} 0.11 \\ \vdots \\ 0.02 \end{pmatrix}$

| softmax | softmax | softmax | softmax | softmax |

$\begin{pmatrix} 0 \\ \vdots \\ 0 \end{pmatrix}$  $h_1$        $h_2$        $h_3$        $h_4$        $h_5$

LSTM → LSTM → LSTM → LSTM → LSTM

Spectrogram  $\begin{pmatrix} 0.32 \\ \vdots \\ 0.64 \end{pmatrix}$   $\begin{pmatrix} 0.21 \\ \vdots \\ 0.43 \end{pmatrix}$   $\begin{pmatrix} 0.33 \\ \vdots \\ 0.15 \end{pmatrix}$   $\begin{pmatrix} 0.12 \\ \vdots \\ 0.14 \end{pmatrix}$   $\begin{pmatrix} 0.22 \\ \vdots \\ 0.03 \end{pmatrix}$

**Raw Audio**

This never happens in practice because:

input length ≠ output length

Kian Katanforoosh, Andrew Ng, Younes Bensouda Mourri

Graves et al., 2006, Connectionist Temporal Classification: Labelling Unsegmented Sequence Data with Recurrent Neural Networks

Kian Katanforoosh, Andrew Ng, Younes Bensouda Mourri

## Examples

$$\beta(HH\_EEEE\_LL\_LOO) = "HELLO"$$

$$\beta(H\_E\_L\_LOO) = "HELLO"$$

$$\beta(H\_\_LLL\_OO) = "HELO"$$

$$\beta(BBAA\_NA\_NN\_\_AA\_A) = "BANANAA"$$

# Speech Recognition

<u>Independence assumption</u>

$$P(c_1|x) = \prod_{t=1}^{T_x} P(c_1^{\langle t \rangle}|x)$$

$$P(c|x) = \begin{bmatrix} 0.13 & HH\_E\_\_L\_LL\_OO \\ 0.04 & H\_EE\_L\_LL\_HO \\ 0.03 & HH\_E\_\_L\_LL\_OO \\ 0.01 & H\_EE\_\_L\_LL\_OO \\ 0.001 & H\_I\_ILL\_L\_OOOO \\ \vdots & \vdots \end{bmatrix}$$

$$P(y|x) = \sum_{c:\beta(c)=y} P(c|x)$$

$$P("HELLO") =$$

# Speech Recognition

Loss?

$\beta(..) \longrightarrow \hat{y} = "HELLO"$

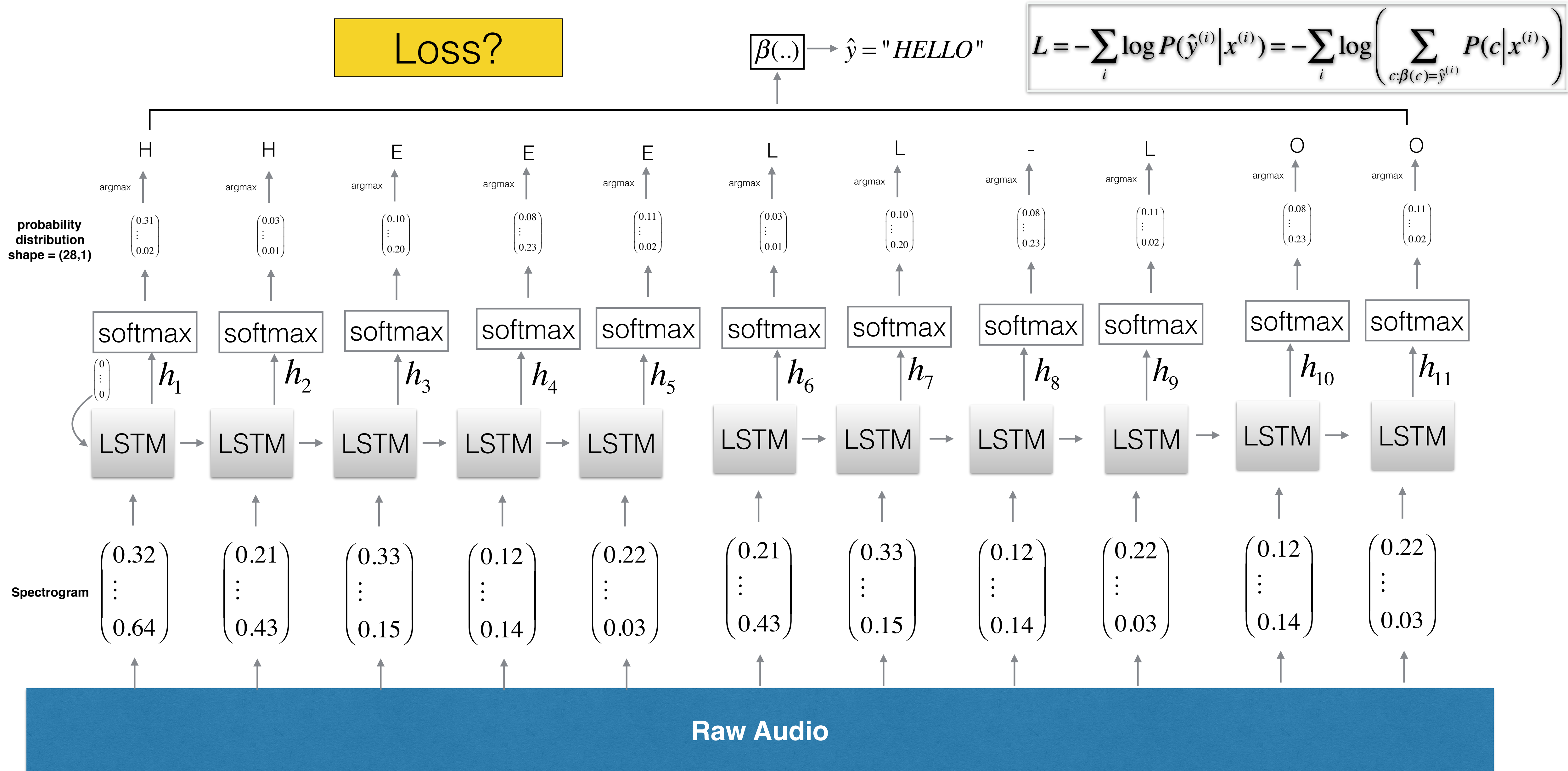$$L = -\sum_i \log P(\hat{y}^{(i)} | x^{(i)}) = -\sum_i \log \left( \sum_{c:\beta(c)=\hat{y}^{(i)}} P(c | x^{(i)}) \right)$$

| H | H | E | E | E | L | L | - | L | O | O |

argmax (×11)

probability
distribution
shape = (28,1)

$\begin{pmatrix} 0.31 \\ \vdots \\ 0.02 \end{pmatrix}$ $\begin{pmatrix} 0.03 \\ \vdots \\ 0.01 \end{pmatrix}$ $\begin{pmatrix} 0.10 \\ \vdots \\ 0.20 \end{pmatrix}$ $\begin{pmatrix} 0.08 \\ \vdots \\ 0.23 \end{pmatrix}$ $\begin{pmatrix} 0.11 \\ \vdots \\ 0.02 \end{pmatrix}$ $\begin{pmatrix} 0.03 \\ \vdots \\ 0.01 \end{pmatrix}$ $\begin{pmatrix} 0.10 \\ \vdots \\ 0.20 \end{pmatrix}$ $\begin{pmatrix} 0.08 \\ \vdots \\ 0.23 \end{pmatrix}$ $\begin{pmatrix} 0.11 \\ \vdots \\ 0.02 \end{pmatrix}$ $\begin{pmatrix} 0.08 \\ \vdots \\ 0.23 \end{pmatrix}$ $\begin{pmatrix} 0.11 \\ \vdots \\ 0.02 \end{pmatrix}$

softmax (×11)

$\begin{pmatrix} 0 \\ \vdots \\ 0 \end{pmatrix}$ $h_1$ $h_2$ $h_3$ $h_4$ $h_5$ $h_6$ $h_7$ $h_8$ $h_9$ $h_{10}$ $h_{11}$

LSTM → LSTM → LSTM → LSTM → LSTM → LSTM → LSTM → LSTM → LSTM → LSTM → LSTM

**Spectrogram**

$\begin{pmatrix} 0.32 \\ \vdots \\ 0.64 \end{pmatrix}$ $\begin{pmatrix} 0.21 \\ \vdots \\ 0.43 \end{pmatrix}$ $\begin{pmatrix} 0.33 \\ \vdots \\ 0.15 \end{pmatrix}$ $\begin{pmatrix} 0.12 \\ \vdots \\ 0.14 \end{pmatrix}$ $\begin{pmatrix} 0.22 \\ \vdots \\ 0.03 \end{pmatrix}$ $\begin{pmatrix} 0.21 \\ \vdots \\ 0.43 \end{pmatrix}$ $\begin{pmatrix} 0.33 \\ \vdots \\ 0.15 \end{pmatrix}$ $\begin{pmatrix} 0.12 \\ \vdots \\ 0.14 \end{pmatrix}$ $\begin{pmatrix} 0.22 \\ \vdots \\ 0.03 \end{pmatrix}$ $\begin{pmatrix} 0.12 \\ \vdots \\ 0.14 \end{pmatrix}$ $\begin{pmatrix} 0.22 \\ \vdots \\ 0.03 \end{pmatrix}$

**Raw Audio**

Kian Katanforoosh, Andrew Ng, Younes Bensouda Mourri

## Inference?

BEAM SEARCH > MAX DECODING :)

Kian Katanforoosh, Andrew Ng, Younes Bensouda Mourri

$$L = -\sum_i \log P(\hat{y}^{(i)} | x^{(i)}) = -\sum_i \log \left( \sum_{c:\beta(c)=\hat{y}^{(i)}} P(c|x^{(i)}) \right)$$

Implementations of CTC loss

tf.nn.ctc_loss(…)
Keras -> Custom loss

Kian Katanforoosh, Andrew Ng, Younes Bensouda Mourri

# Speech Recognition

- How to incorporate information about the future?

*A good way to efficiently incorporate future information in speech recognition is still an open problem.*

- What's the consequence of $P(c_1|x) = \prod_{t=1}^{T_x} P(c_1^{\langle t \rangle}|x)$ (conditional independence)?

*A model like CTC may have trouble producing such diverse transcripts for the same utterance because of conditional independence assumptions between frames.*

*But, on the other hand, it makes the model more robust to a change of settings.*

- What is the problem with our output $\hat{y}$ ?

*CTC model makes a lot of spelling and linguistic mistakes because P(y\x) directly models audio data. Some words are hard to spell based on their audios.*

- Can you think of any practical applications leveraging this model?

*Lipreading.*

Assael et al., LipNet: end-to-end sentence-level lipreading, 2016
Hannun, "Sequence Modeling with CTC", Distill, 2017.
Chan et al.,2015, Listen, Attend and Spell

Kian Katanforoosh, Andrew Ng, Younes Bensouda Mourri