

Machine Translation

COMP90042

Natural Language Processing
Lecture 17



Introduction

- **Machine translation (MT)** is the task of translating text from one **source language** to another **target language**

虽然北风呼啸，但天空依然十分清澈



However, the sky remained clear under the strong north wind

Why?

- Removes language barrier
- Makes information in any languages accessible to anyone
- But translation is a classic “AI-hard” challenge
 - ▶ Difficult to preserve the **meaning** and the **fluency** of the text after translation

MT is Difficult

- Not just simple word for word translation
- Structural changes, e.g., syntax and semantics
- Multiple word translations, idioms
- Inflections for gender, case etc *European languages.*
- Missing information (e.g., determiners)
Japanese.

虽然 北 风 呼啸 , 但 天空 依然 十分 清澈
although north wind howls , but sky still extremely limpid

Outline

- Statistical Machine Translation
 - ▶ Word-based approach
- Neural Machine Translation
 - ▶ Encoder-decoder architecture
- Attention Mechanism
 - ▶ Improvement to Neural MT
- Machine Translation Evaluation



Statistical MT

Early Machine Translation

- Started in early 1950s
- Motivated by the Cold War to translate Russian to English
- Rule-based system
 - ▶ Use bilingual dictionary to map Russian words to English words
- Goal: translate 1-2 million words an hour within 5 years

Statistical MT

probabilistic Model

- Given French sentence f , aim is to find the best English sentence e

$$\triangleright \operatorname{argmax}_e P(e | f) \xrightarrow{\text{French sentence}} \xleftarrow{\text{English sentence}}$$

- Use Baye's rule to decompose into two components

BR.

$$\triangleright \operatorname{argmax}_e P(f|e)P(e)$$

TM \downarrow *language M*

$$P(e|f) = \frac{P(f|e)P(e)}{P(f)}$$

Known.

Language vs Translation Model

- $\operatorname{argmax}_e P(f|e)P(e)$
- $P(e)$: language model
 - ▶ learns how to write fluent English text
- $P(f|e)$: translation model
 - ▶ learns how to translate words and phrases from English to French

learn English

learn + translation

How to Learn LM and TM?

- Language model:
 - ▶ Based on text frequency in large **monolingual corpora**
English.
- Translation model:
 - ▶ Based on word co-occurrences in **parallel corpora**
 - ▶ i.e. English-French sentence pairs

Parallel Corpora

- One text in multiple languages
- Produced by human translation
 - ▶ Bible, news articles, legal transcripts, literature, subtitles
 - ▶ Open parallel corpus:
<http://opus.nlpl.eu/>

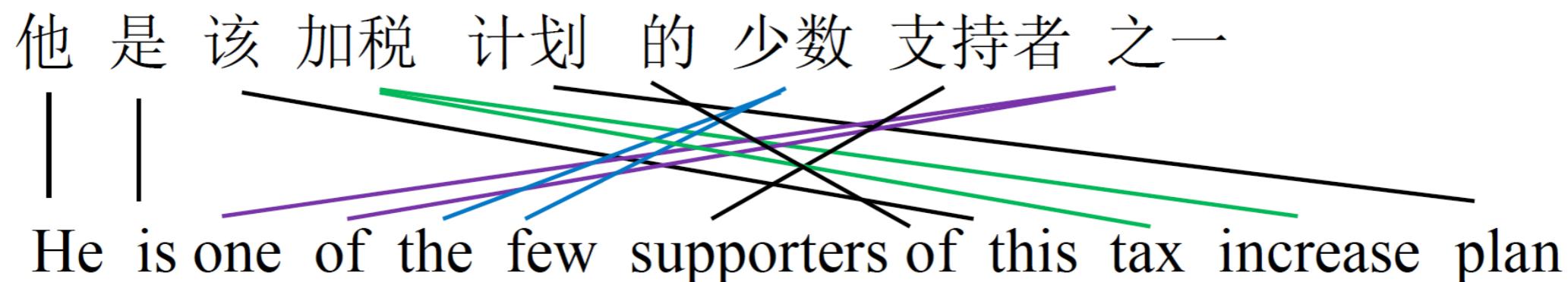


The Rosetta Stone

Models of Translation

Sentence pairs

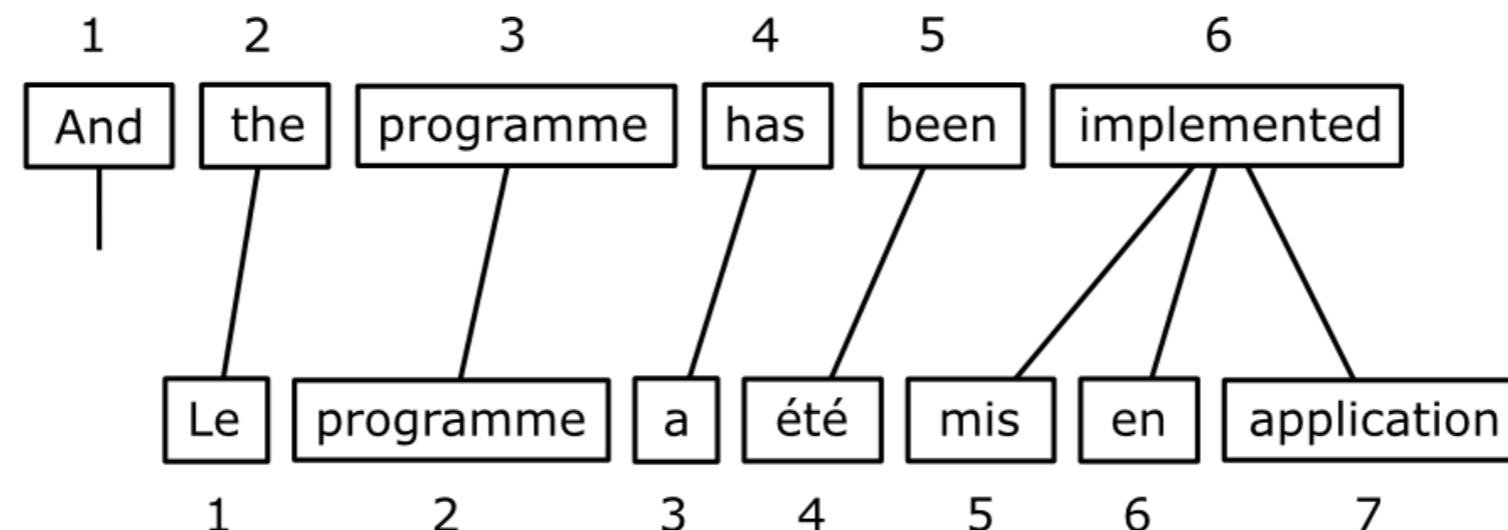
- How to learn $P(f|e)$ from parallel text?
- We only have sentence pairs; words are not aligned in the parallel text *no words alignment*
- Not word to word translation:
 - ▶ rearrangement of words



Alignment

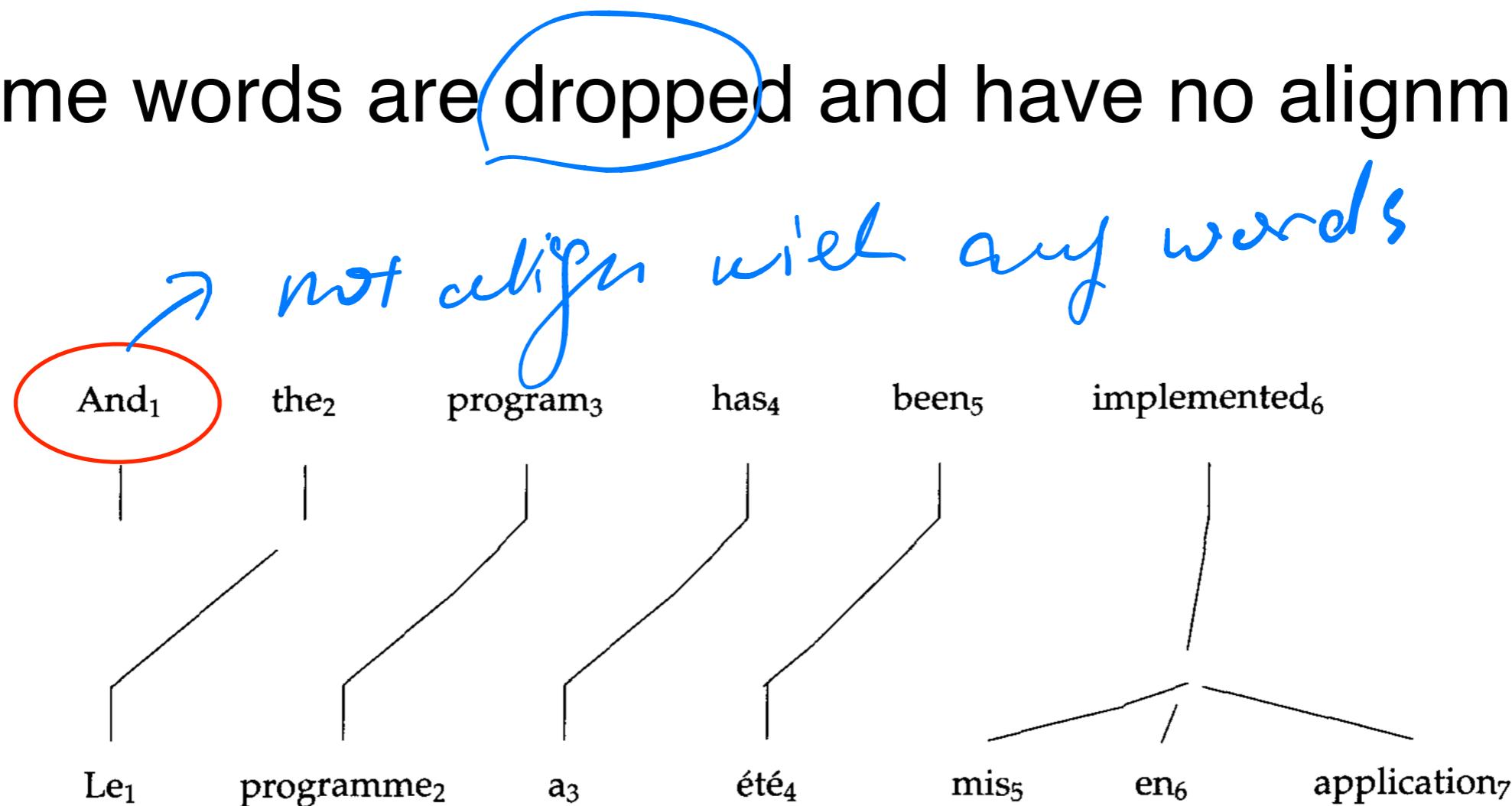
- Idea: introduce **word alignment** as a **latent variable** into the model

► $P(f, a | e)$



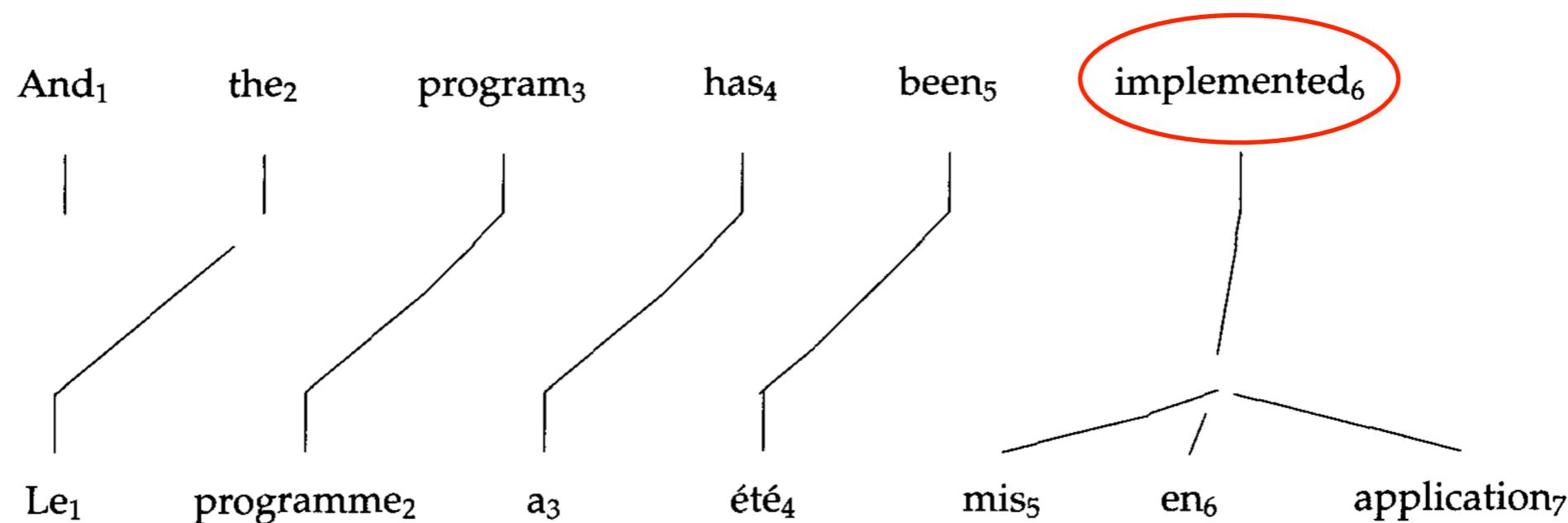
Complexity of Alignment

- Some words are **dropped** and have no alignment



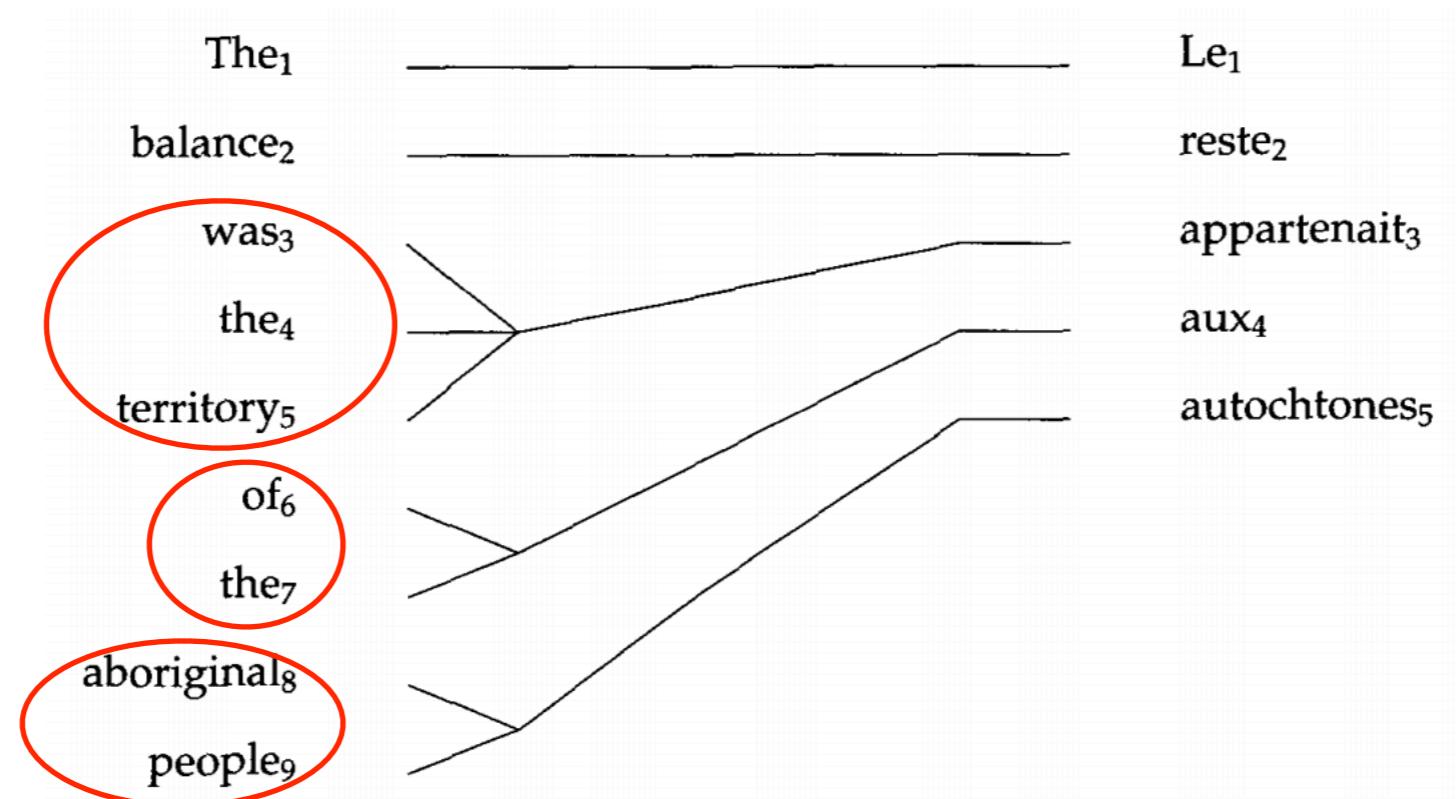
Complexity of Alignment

- One-to-many alignment



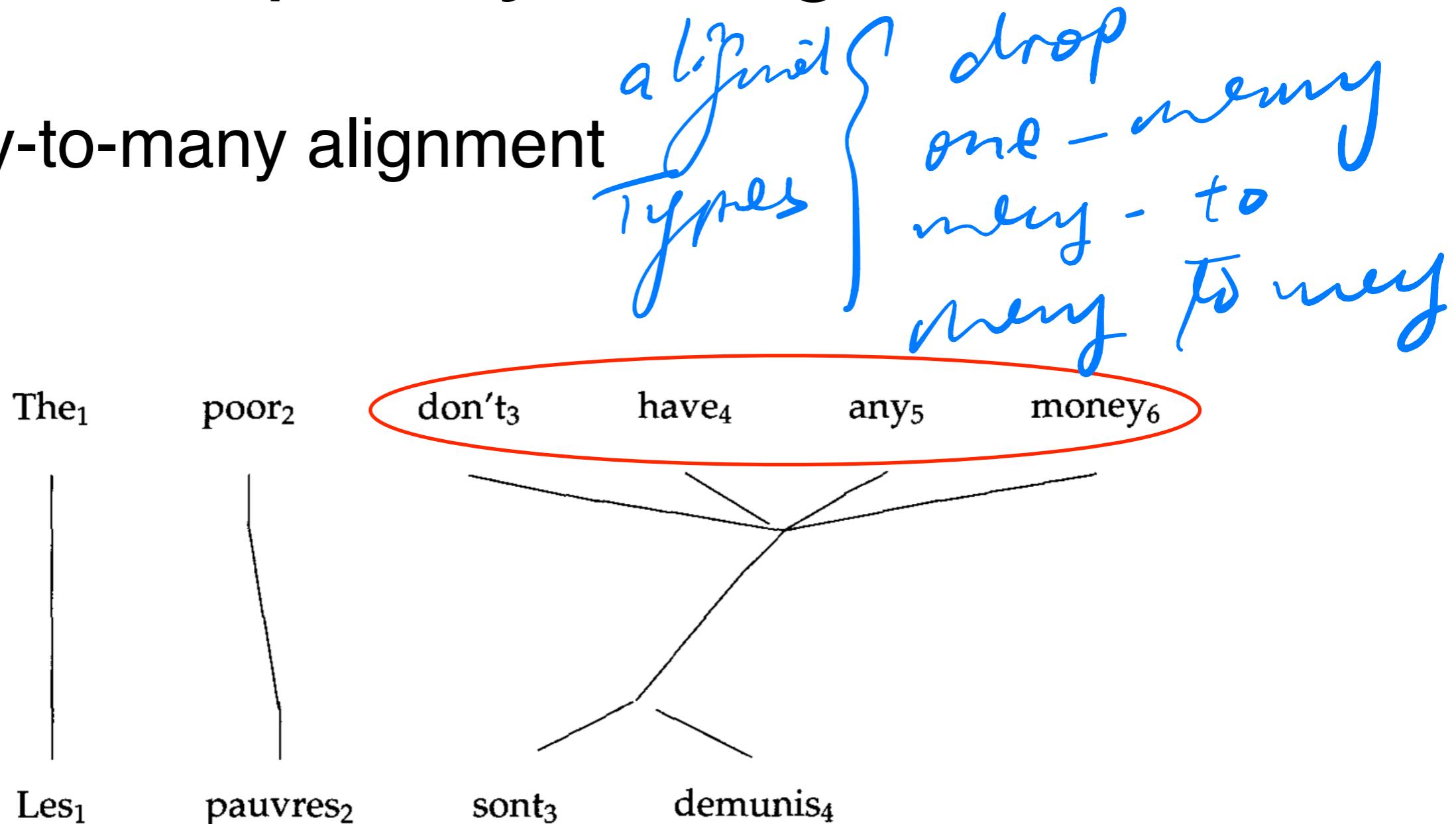
Complexity of Alignment

- Many-to-one alignment



Complexity of Alignment

- Many-to-many alignment



Learning the Alignment

- Word alignments are rarely provided in parallel corpora
 - ▶ As such alignment a introduced as a latent variable
- Use algorithms such as expectation maximisation (EM) to learn

Statistical MT: Summary

- A very popular field of research in NLP prior to 2010s
- Lots of feature engineering
- State-of-the-art systems are very complex
 - ▶ Difficult to maintain
 - ▶ Significant effort needed for new language pairs

Neural Machine Translation

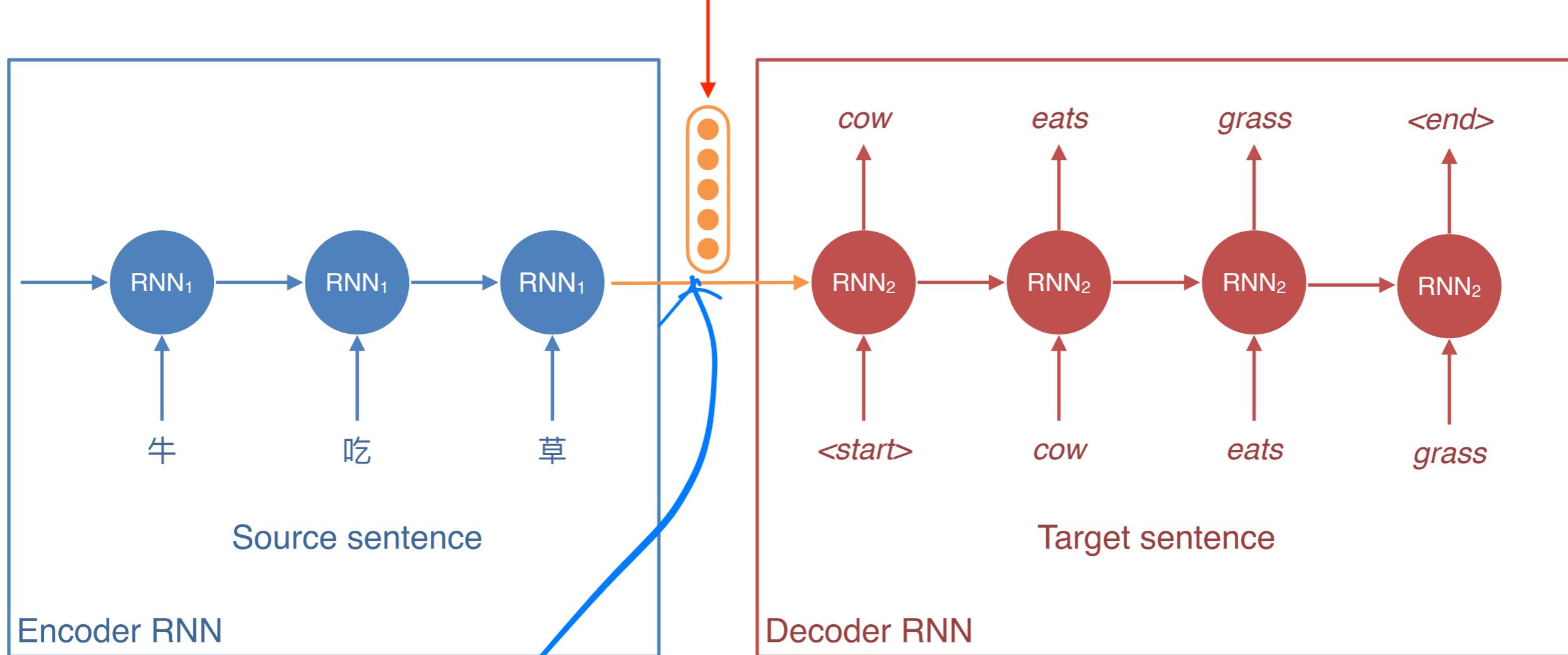
Introduction

Google Translations from Static \rightarrow NN

- Neural machine translation is a new approach to do machine translation
- Use a single neural model to directly translate from source to target *Seq \rightarrow Seq Model*
- Architecture: encoder-decoder model
 - ▶ 1st RNN to encode the source sentence
 - ▶ 2nd RNN to decode the target sentence

last hidden representation

This vector encodes the whole source sentence;
it is used as the initial state for decoder RNN



x denotes
lengths of
source sentences.

$$h_i = RNN_1(h_{i-1}, x_i)$$

$$s_t = RNN_2(s_{t-1}, y_t)$$

$$s_1 = h_{|x|}$$

words ·
↑ previous state

Neural MT

- The decoder RNN can be interpreted as a **conditional language model**
 - ▶ Language model: predicts the next word given previous words in target sentence y
 - ▶ Conditional: prediction is also conditioned on the source sentence x
- $P(y|x) = P(y_1|x)P(y_2|y_1, x)\dots P(y_t|y_1, \dots, y_{t-1}, x)$

Language Model with
add. final condition

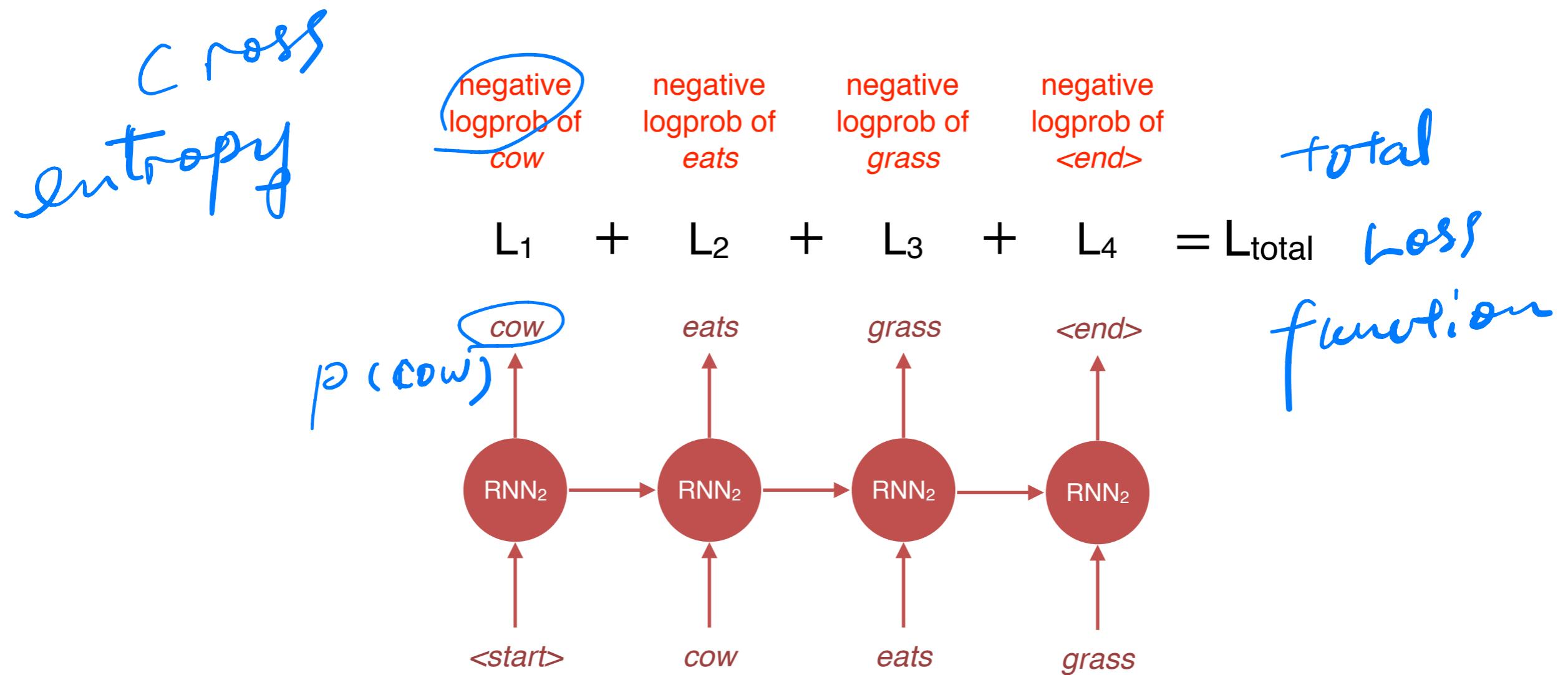
Conditioned
on source

Training Neural MT

- Requires parallel corpus just like statistical MT
- Trains with next word prediction, just like a language model!

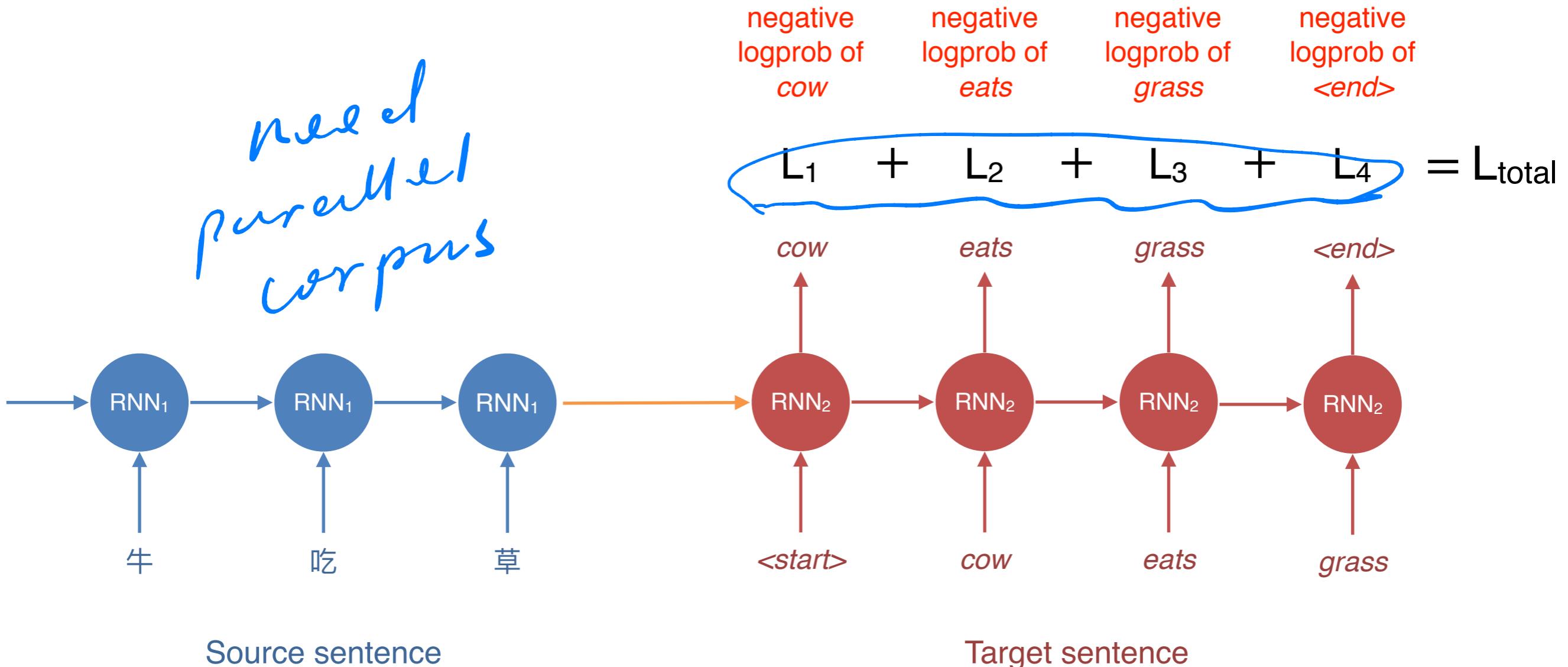
train like LM

Language Model Training Loss

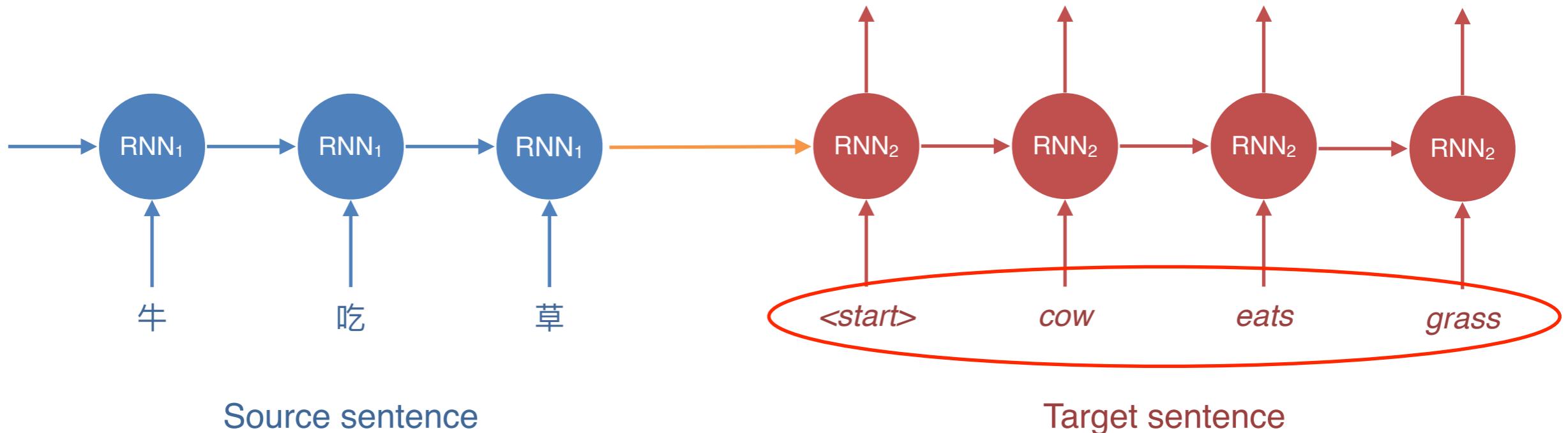


Neural MT Training Loss

*need
parallel
corpus*

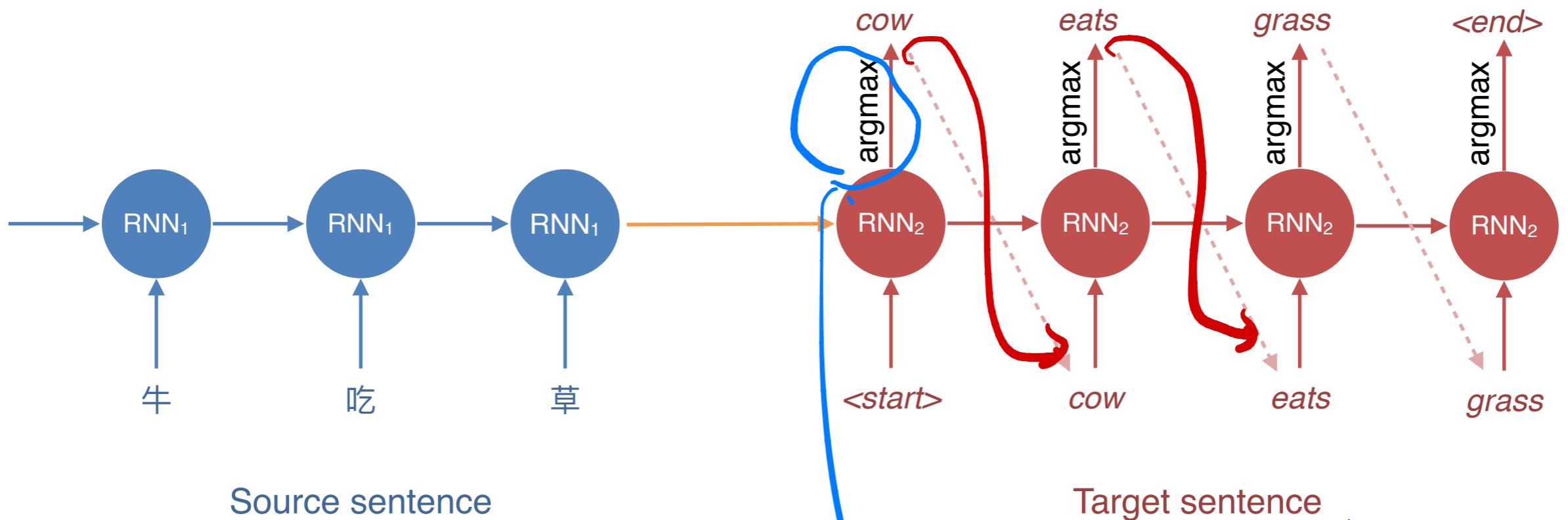


Training



- During training, we have the target sentence
- We can therefore feed the right word from target sentence, one step at a time

Decoding at Test Time



- But at test time, we don't have the target sentence (that's what we're trying to predict!)
- argmax: take the word with the highest probability at every step

Greedy Decoding

- argmax decoding is also called greedy decoding
- Issue: does not guarantee optimal probability

$$P(y|x)$$

=

suffer
exposure bias

local
opt

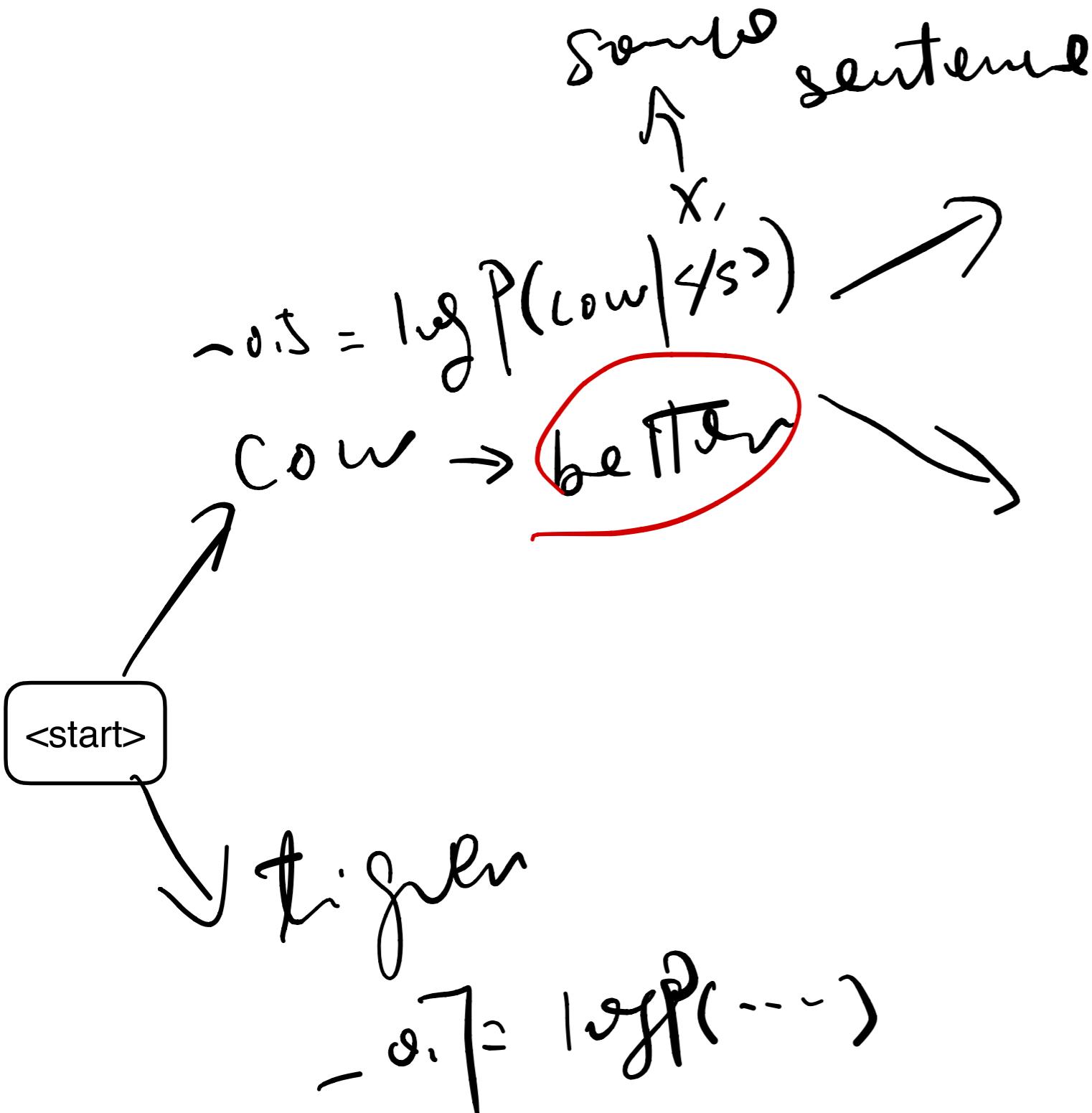
- $y_1, y_2, y_3, \dots, y_n$.

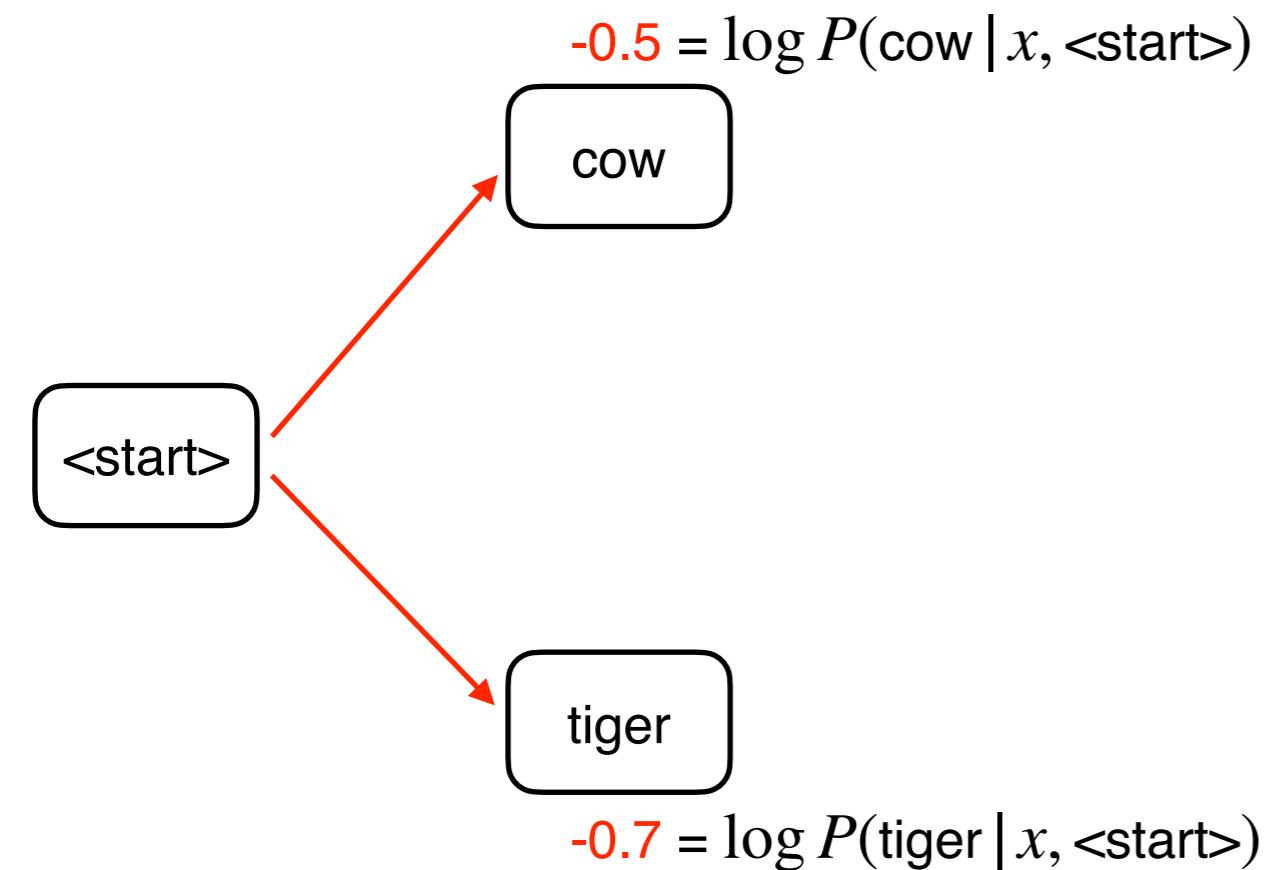
Exhaustive Search Decoding

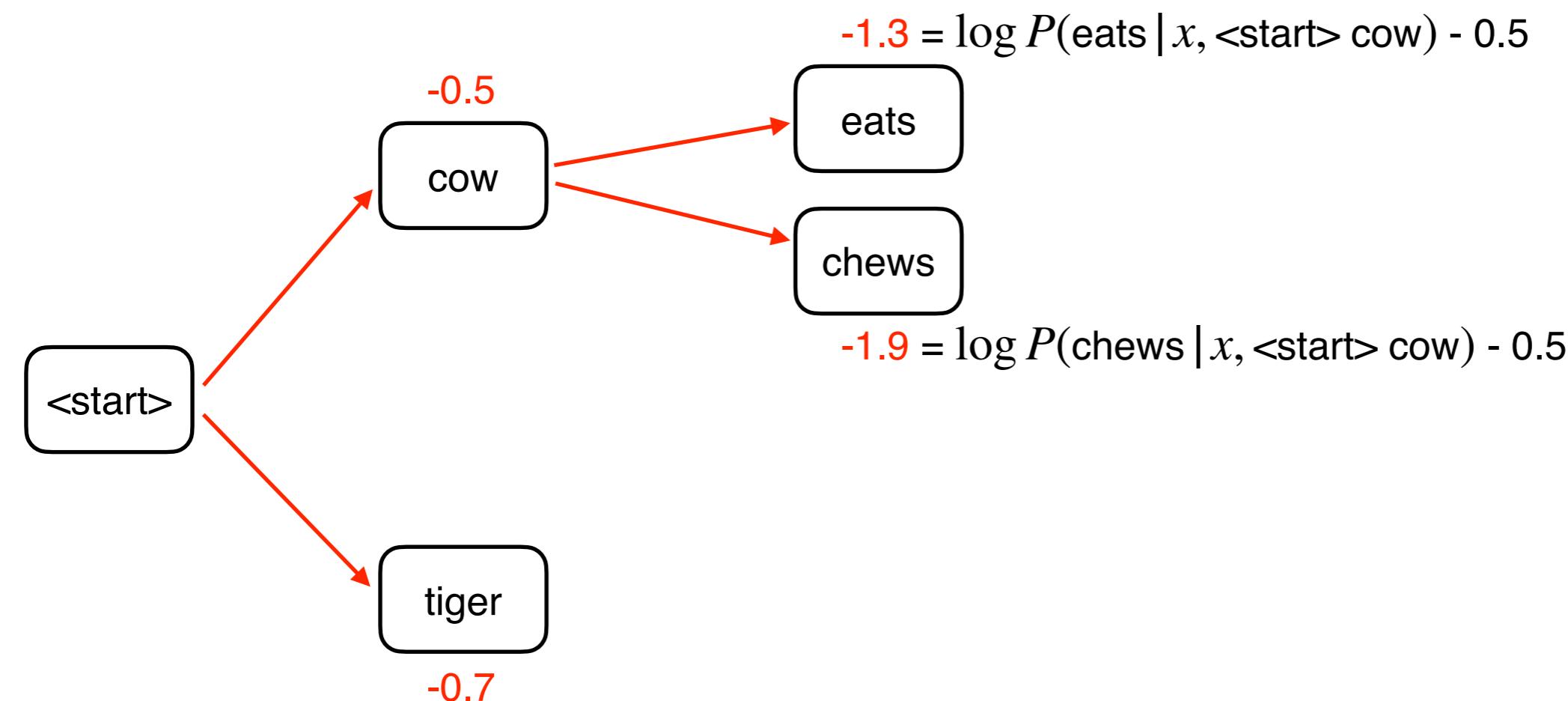
- To find optimal $P(y | x)$, we need to consider every word at every step to compute the probability of all possible sequences
- $O(V^n)$; $V = \text{vocab size}$; $n = \text{sentence length}$
- Far too expensive to be feasible

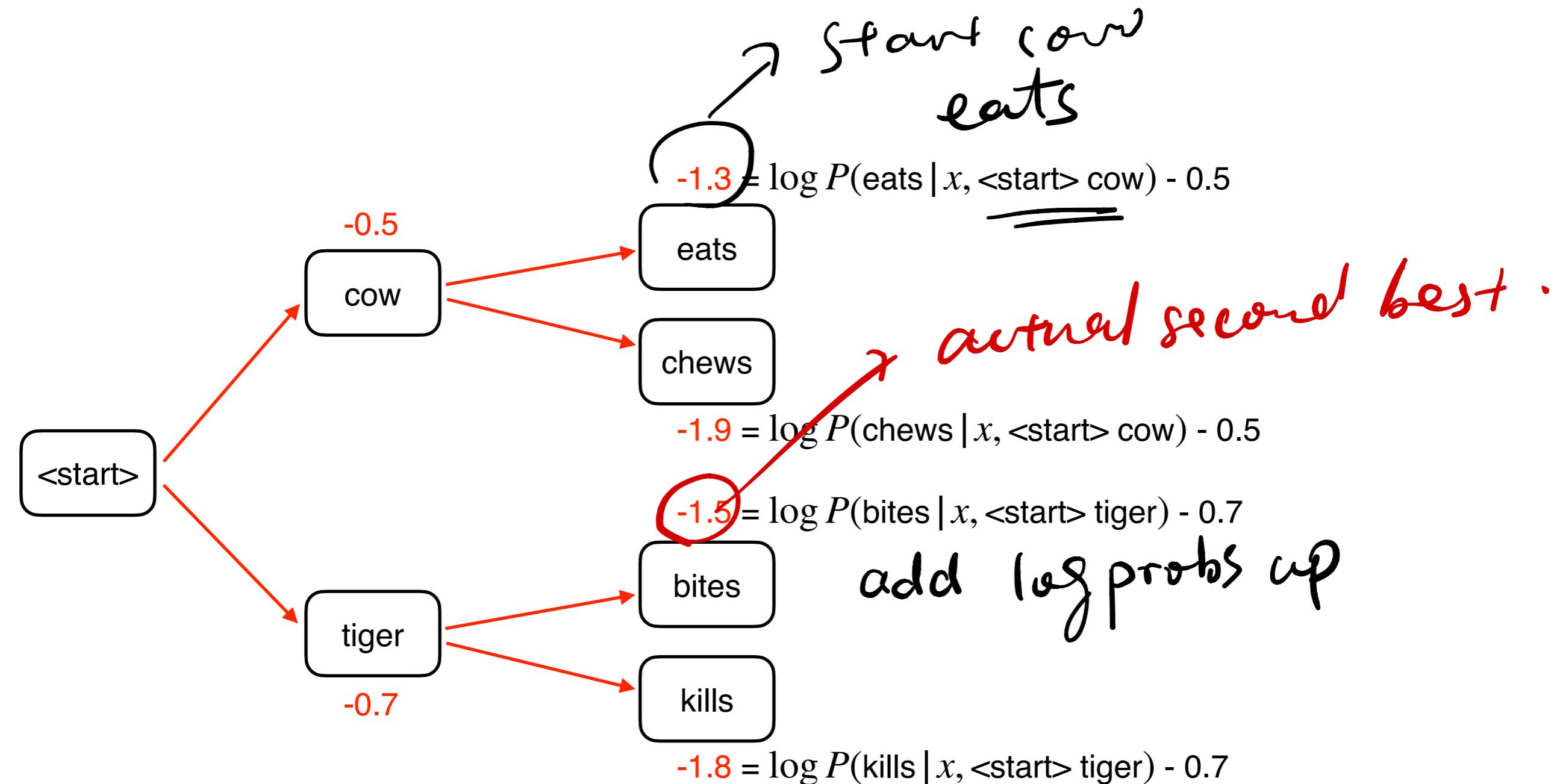
Beam Search Decoding

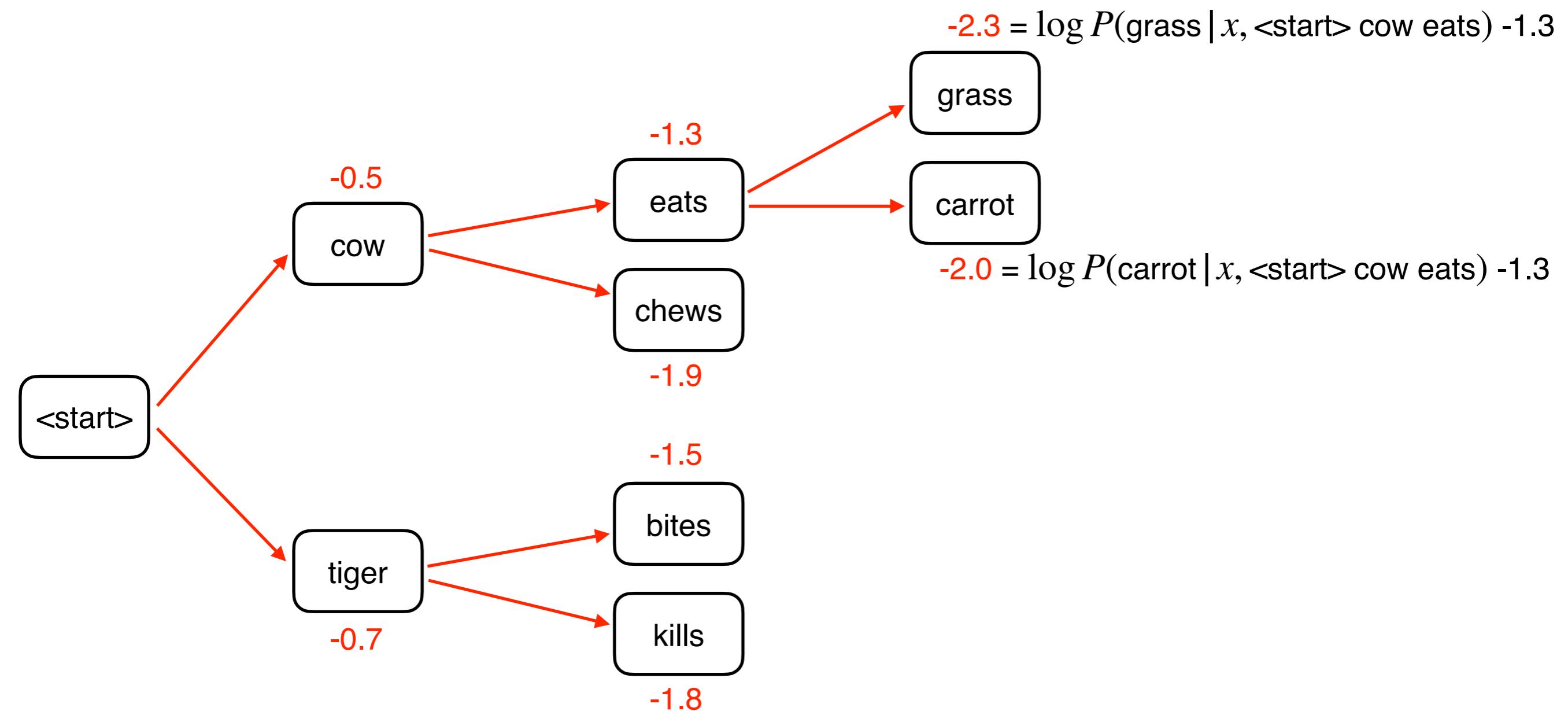
- Instead of considering all possible words at every step, consider **k best words**
→ greedy best hypotheses.
- That is, we keep track of the **top-k words** that produce the best partial translations (**hypotheses**) thus far
- $k = \text{beam width}$ (typically 5 to 10)
- $k = 1 = \text{greedy decoding} \rightarrow \text{argmax best 1.}$
- $k = V = \text{exhaustive search decoding}$

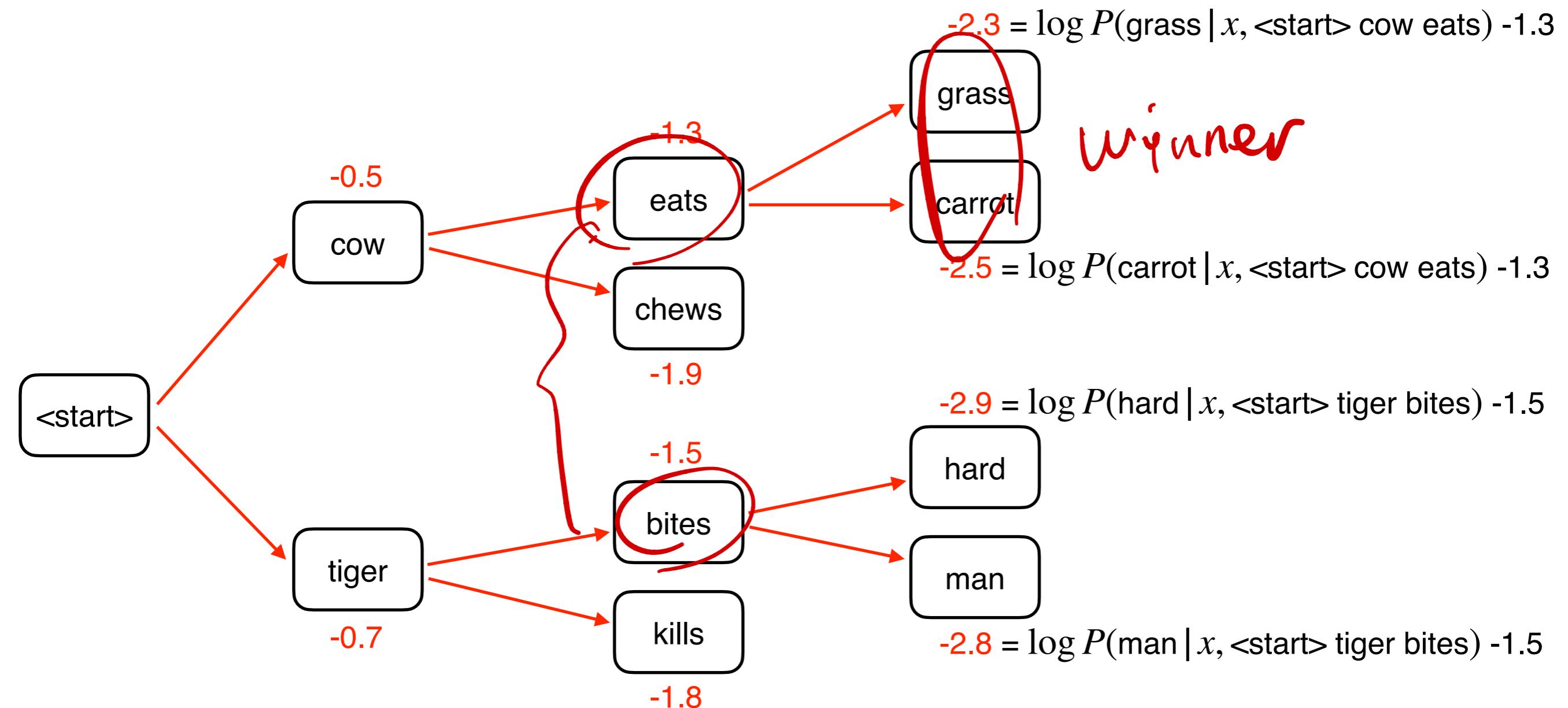


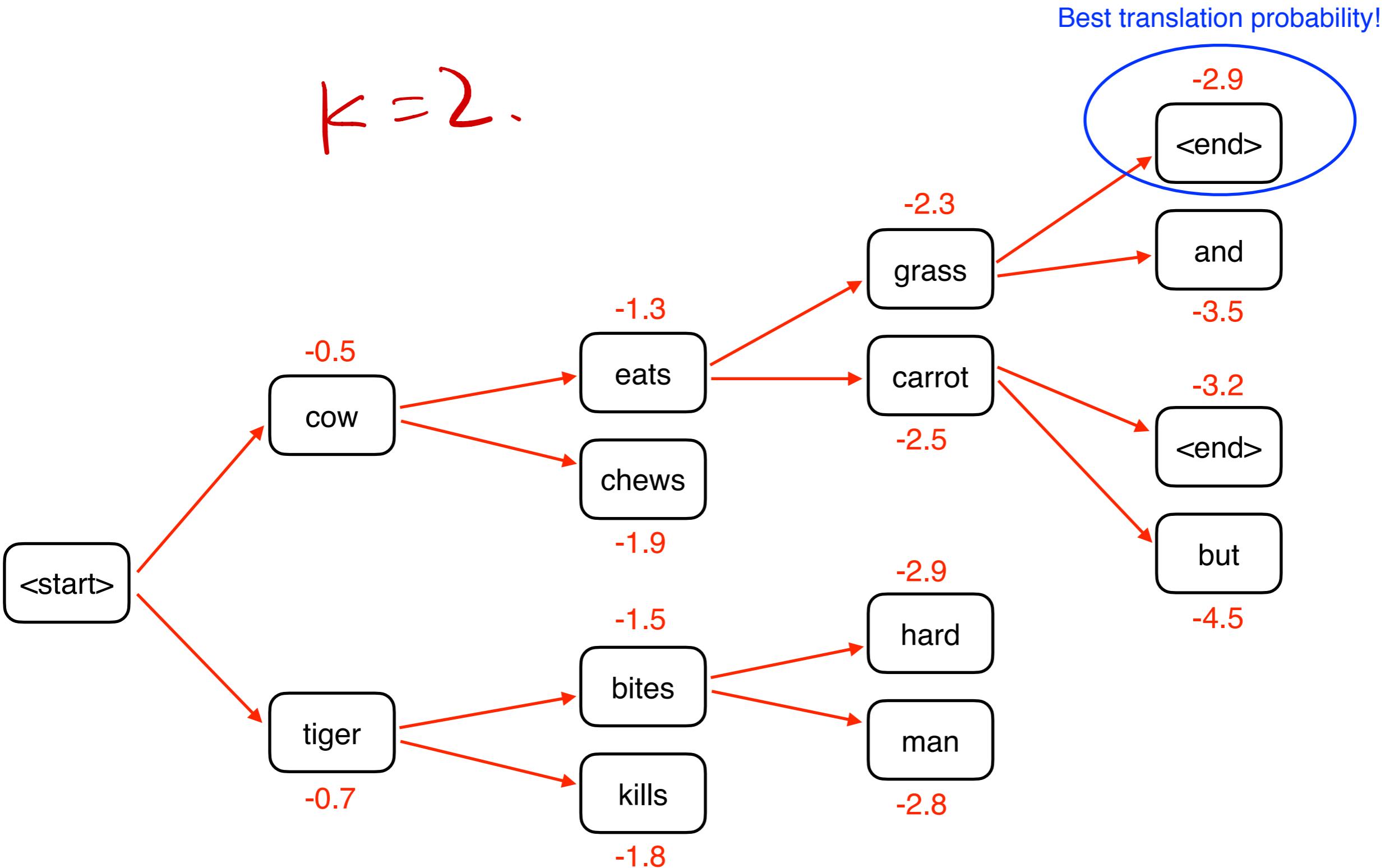




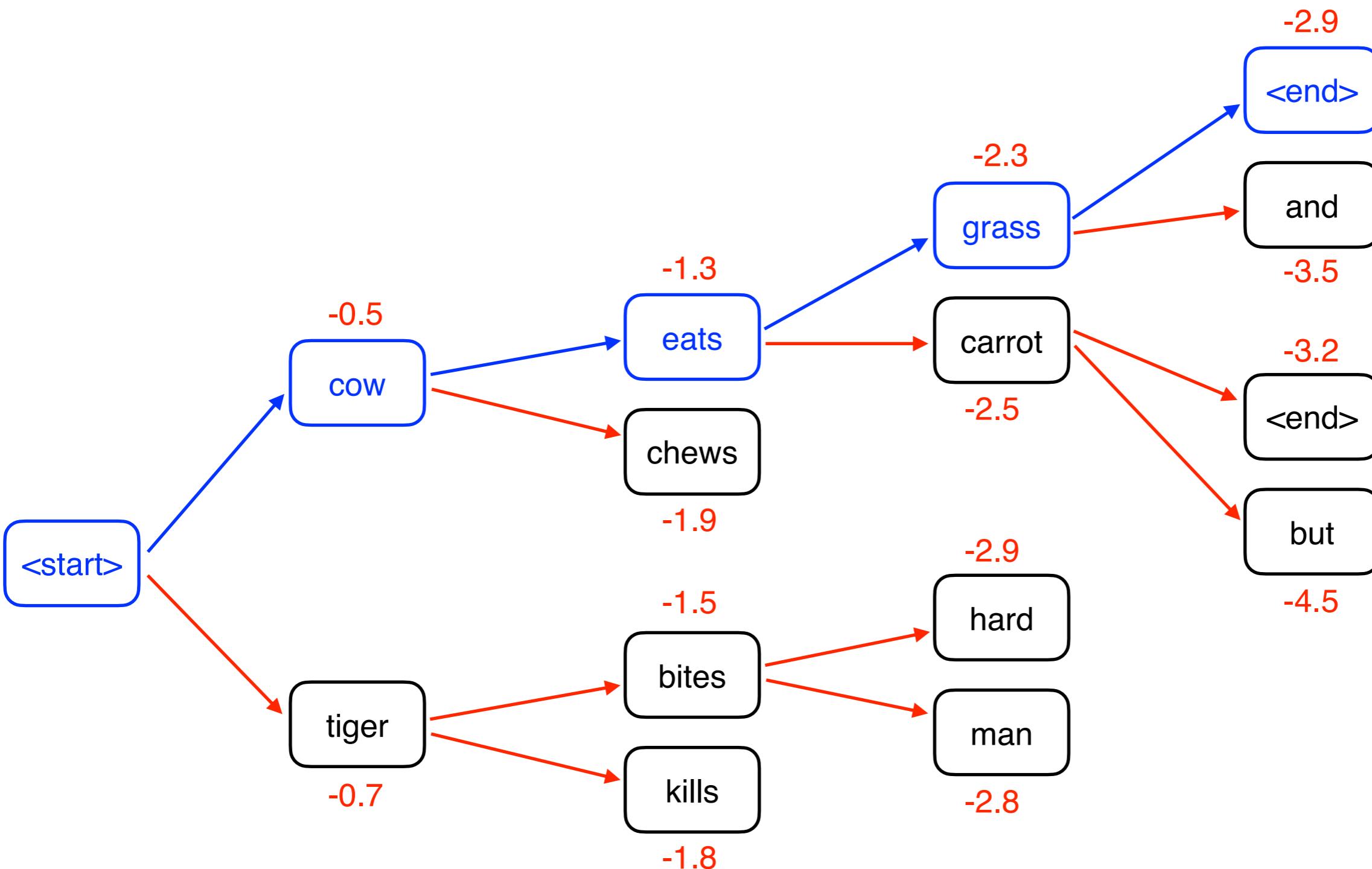






Beam search decoding with $k=2$

Follow pointers to get the target sentence



Beam search decoding with $k=2$

When to Stop?

- When decoding, we stop when we generate <end> token
- But multiple hypotheses may terminate their sentence at different time steps
- We store hypotheses that have terminated, and continue explore those that haven't
stop until all sentences terminated
- Typically we also set a maximum sentence length that can be generated (e.g. 50 words)

Neural MT: Summary

- Single **end-to-end** model
 - ▶ Statistical MT systems have multiple sub-components
- Less feature engineering
- Generated translation is more fluent
- But can produce new details that are not in the source sentence (**hallucination**)

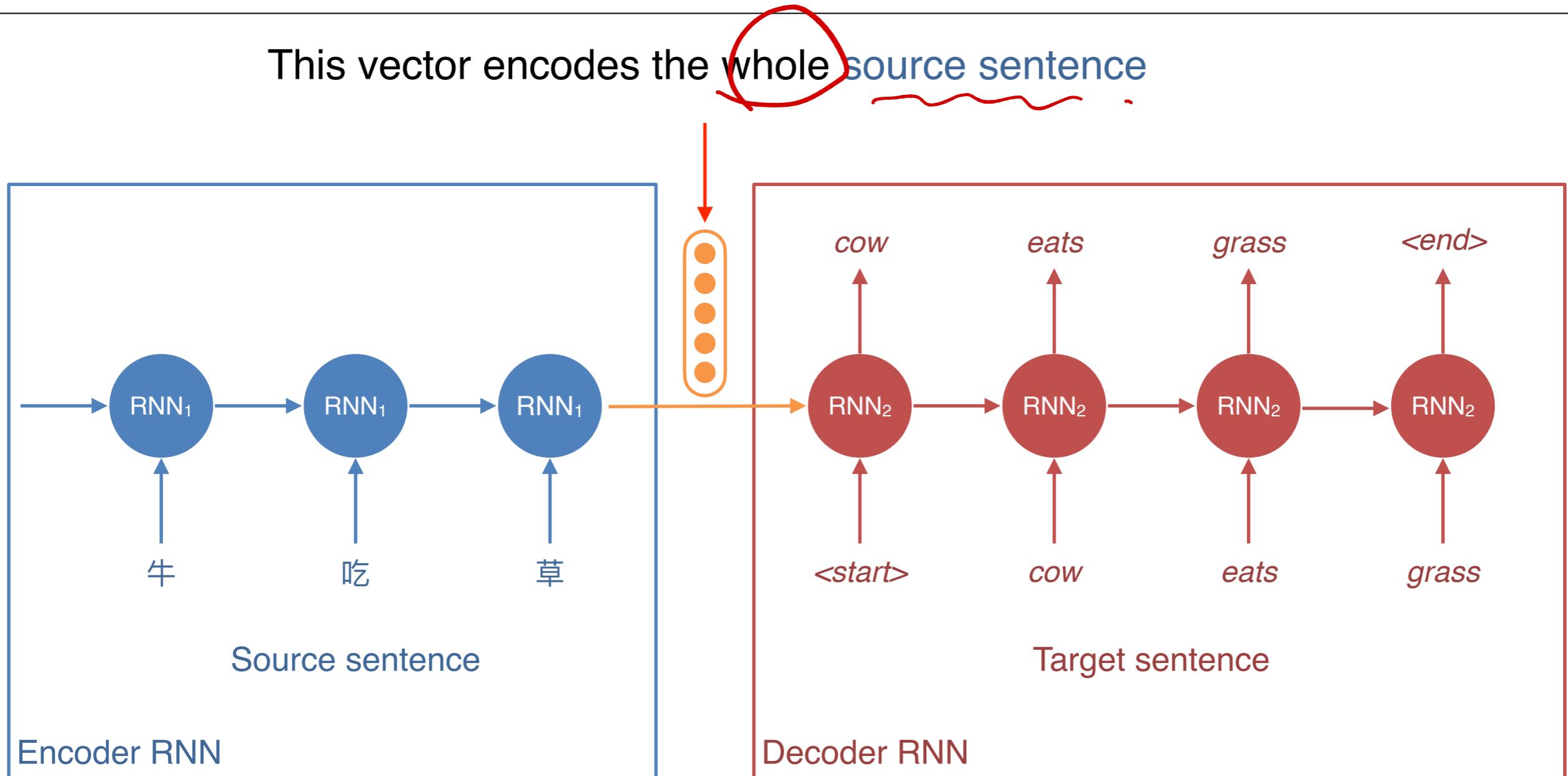
details not in source.

Source: Caldiero sprach mit E! Nachrichten nach dem hart erkämpften Sieg, noch immer unter dem Schock über den Gewinn des Großen Preises von 1 Million \$.

Reference: Caldiero spoke with E! News after the hard-fought victory, still in shock about winning the \$1 million grand prize.

NMT Translation: Caldiero spoke with E, after the hard won victory, still under the shock of the winning of the **Grand Prix** of 1 million \$.

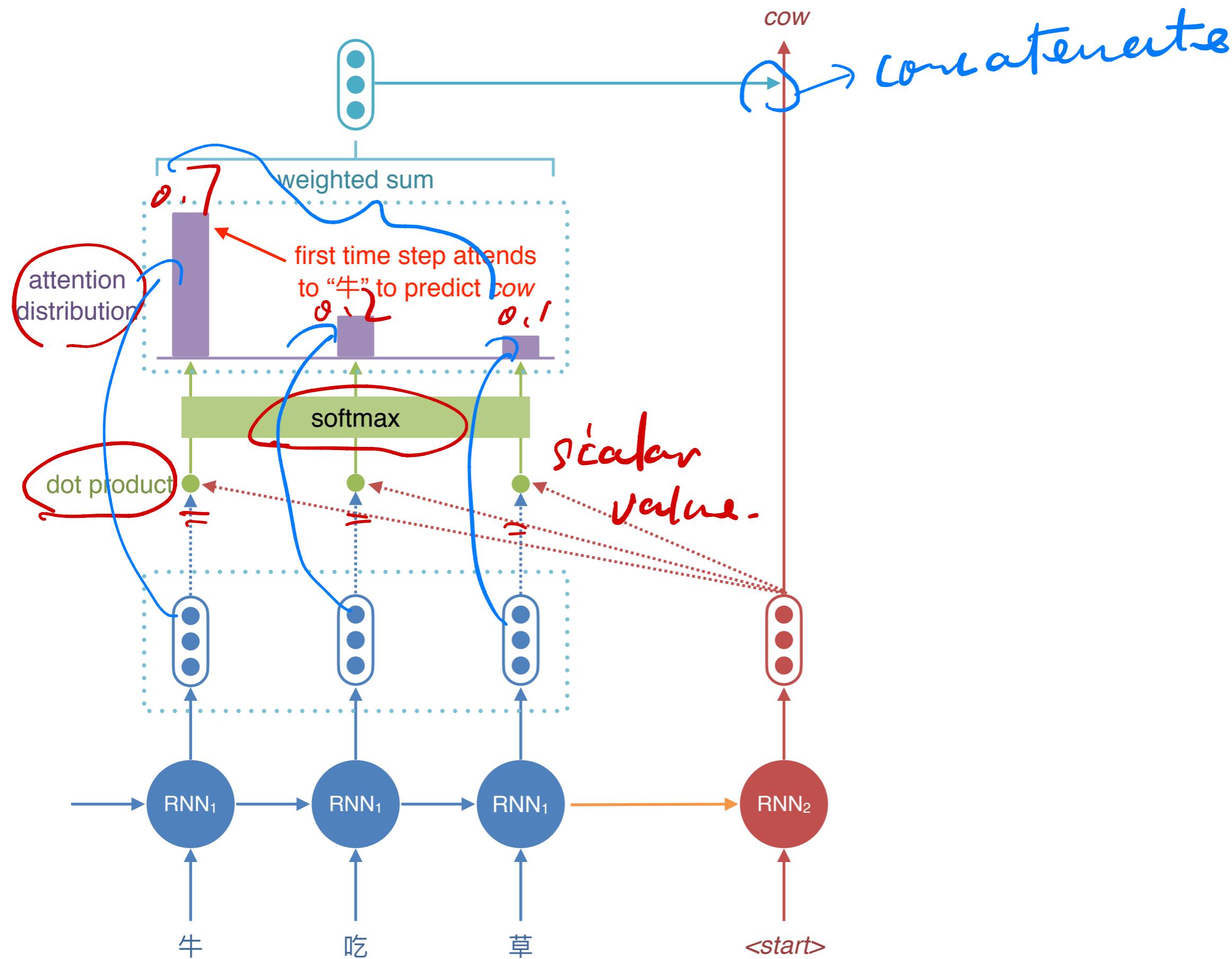
Attention (!)

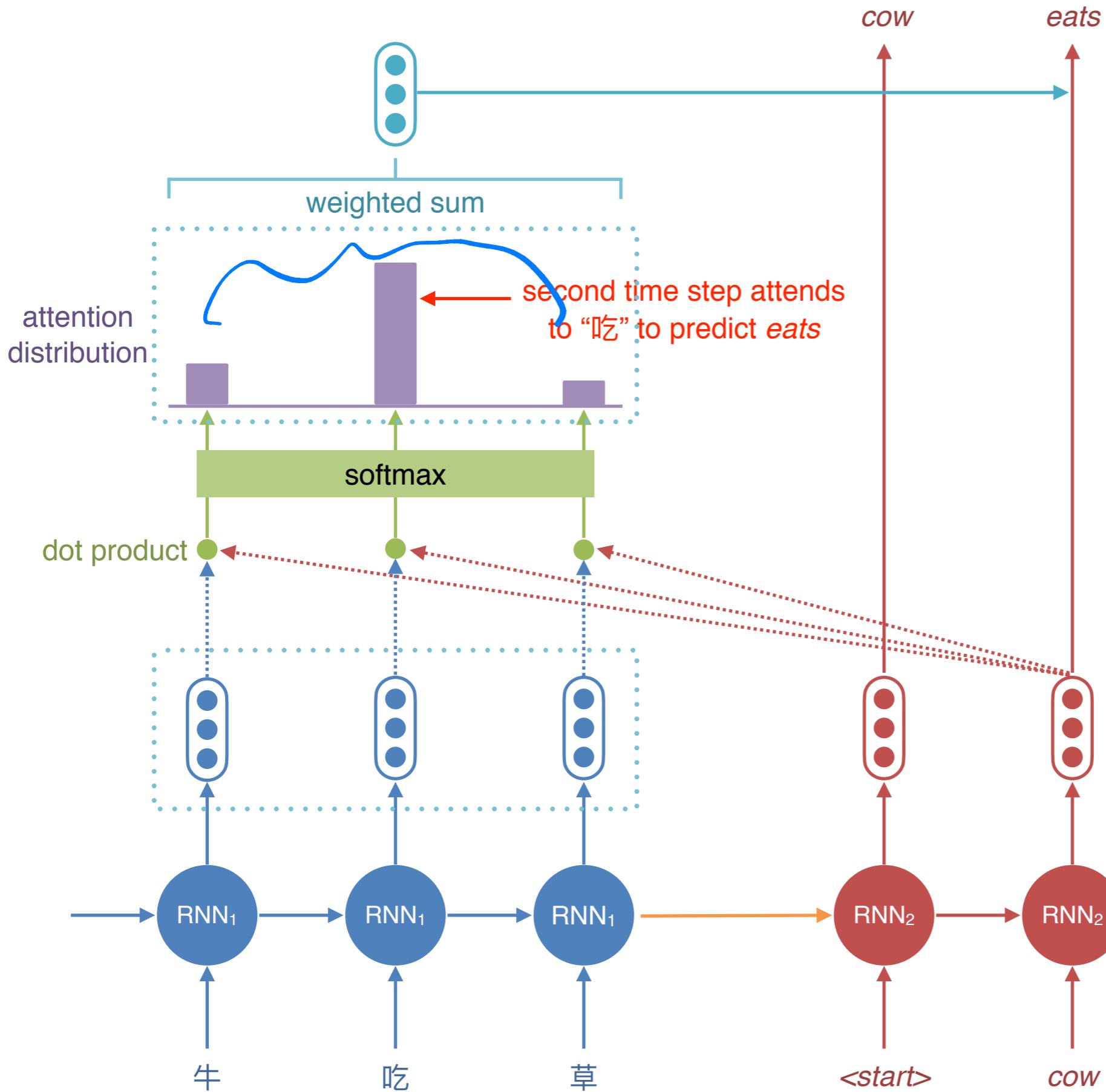


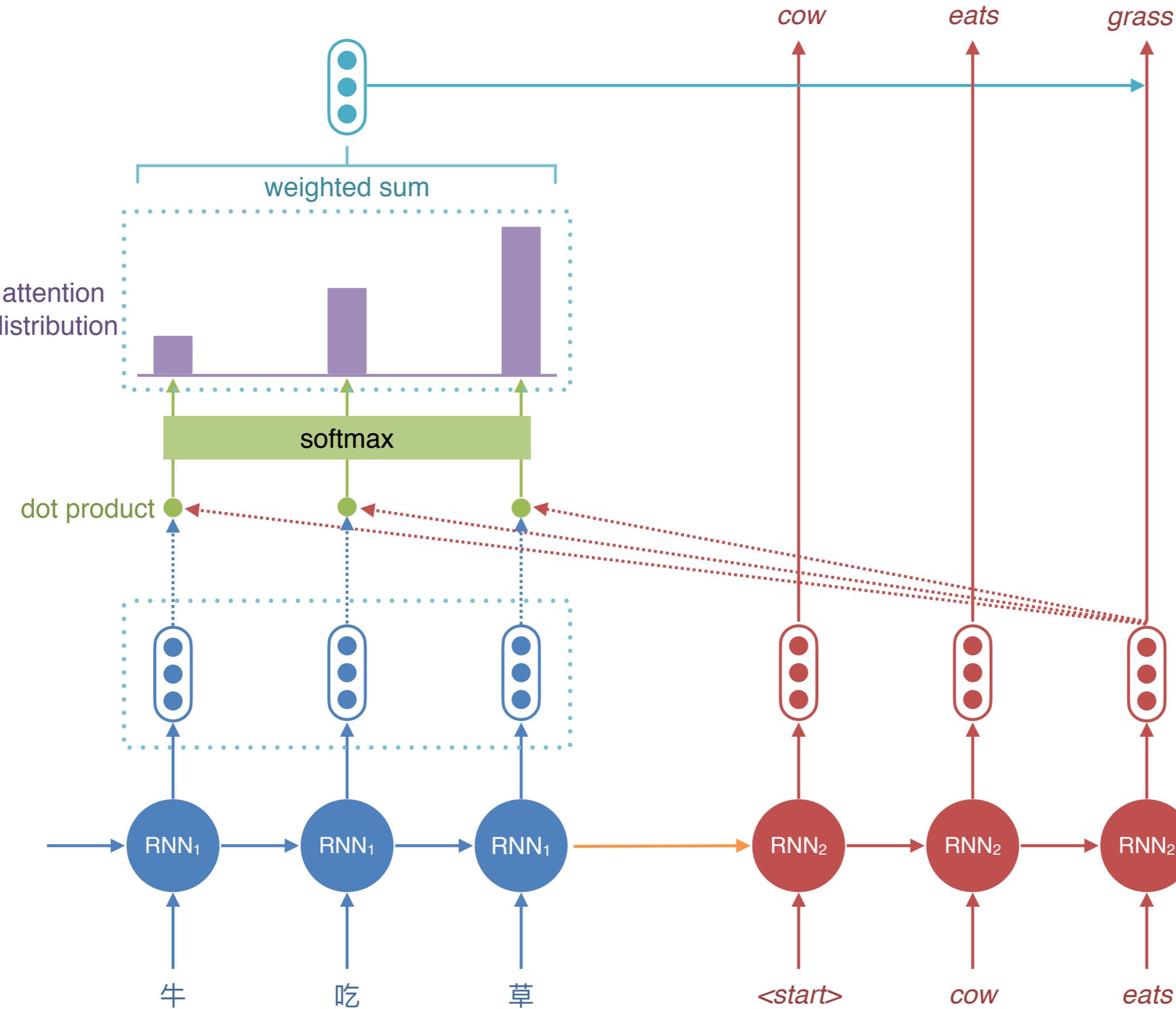
- With a long source sentence, the encoded vector is unlikely to capture all the information in the sentence
- This creates an **information bottleneck**

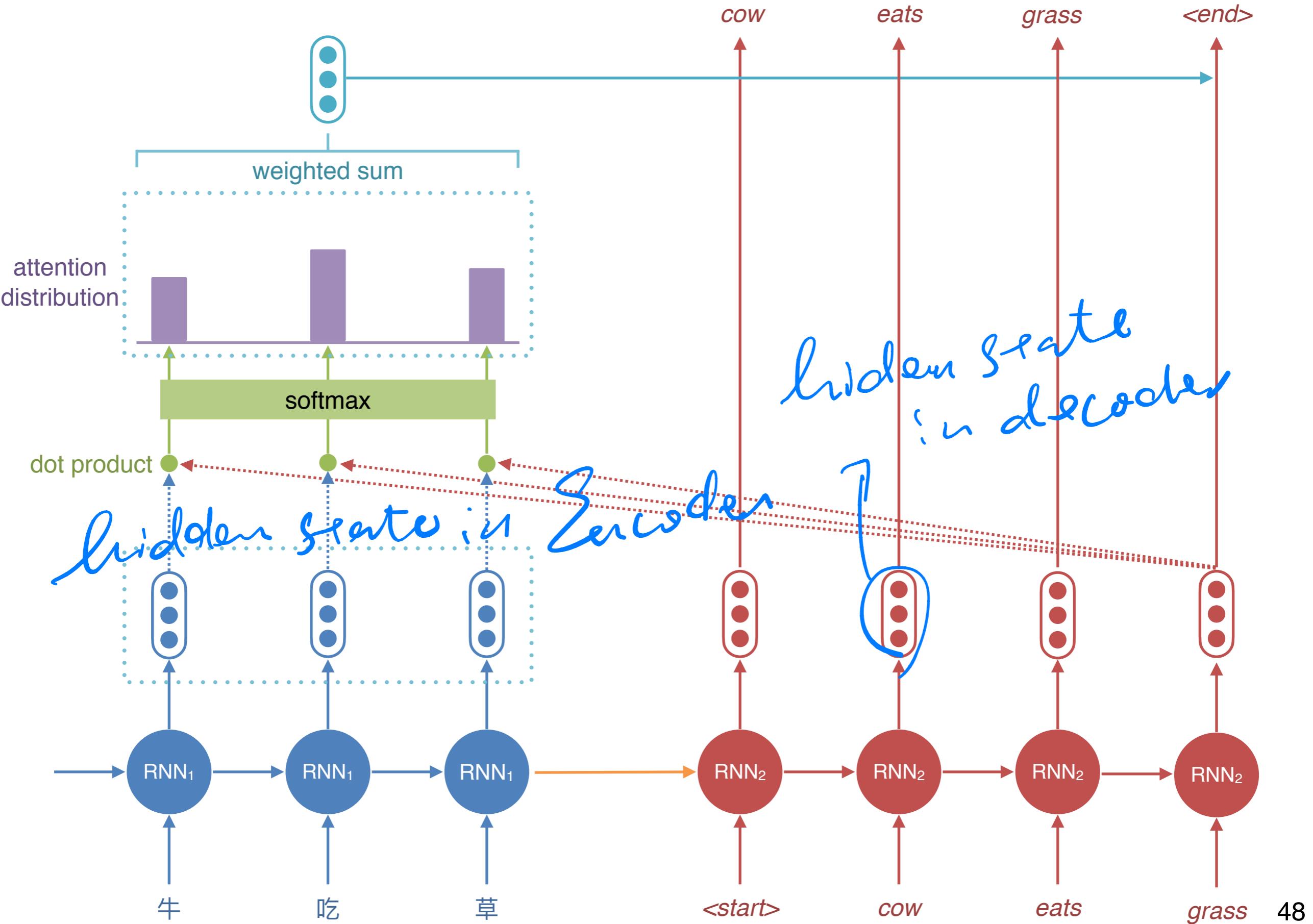
Attention

- For the decoder, at every time step allow it to '**attend**' to words in the source sentence









Encoder-Decoder with Attention

- Encoder hidden states: $h_i = RNN_1(h_{i-1}, x_i)$
- Decoder hidden states: $s_t = RNN_2(s_{t-1}, y_t)$
- For each time step in the decoder, attend to each of the hidden states to produce the attention weights:

→ $e_t = [s_t^T h_1, s_t^T h_2, \dots, s_t^T h_{|x|}]$

dec enc

- Apply softmax to the attention weights to get a valid probability distribution:

→ $\alpha_t = \underbrace{\text{softmax}(e_t)}$

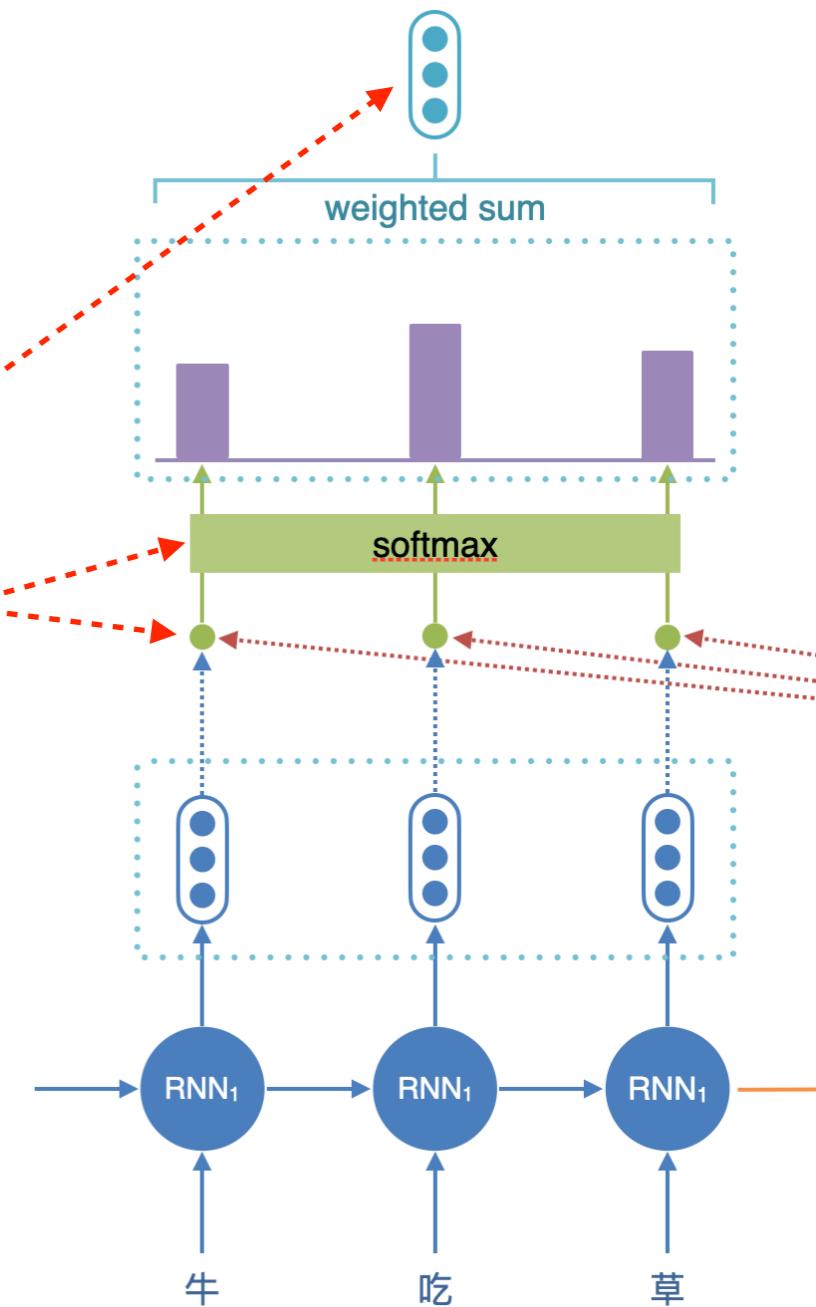
- Compute weighted sum of the encoder hidden states:

→ $c_t = \sum_{i=1}^{|x|} \alpha_t^i h_i$

weighted ell sum

- Concatenate c_t and s_t to predict the next word

context vector



Variants

- Attention:
 - ▶ Dot product: $s_t^\top h_i$
 - ▶ Bilinear: $s_t^\top W h_i$
 - ▶ Additive: $v^\top \tanh(W_s s_t + W_h h_i)$
- c_t can be injected to the current state (s_t) , or to the input word (y_t) .

Attention: Summary

- Solves the information bottleneck issue by allowing decoder to have access to the source sentence words directly *Avoid Bottleneck.*
- Provides some form of interpretability
 - ▶ Attention weights can be seen as word alignments
- Most state-of-the-art MT systems are based on this architecture
 - ▶ **Google Translate** (<https://slator.com/technology/google-facebook-amazon-neural-machine-translation-just-had-its-busiest-month-ever/>)

Evaluation

MT Evaluation

- **BLEU:** compute n-gram overlap between “reference” translation and generated translation

- Typically computed for 1 to 4-gram

Strong penalty for short output.

$$\text{BLEU} = \text{BP} \times \exp \left(\frac{1}{N} \sum_n^N \log p_n \right)$$

*n=2 bigram
n=3 trigram*

otherwise

it is easy to just

“Brevity Penalty” to penalise short outputs

generate short output.

$$p_n = \frac{\# \text{ correct n-grams}}{\# \text{ predicted n-grams}}$$

$$\text{BP} = \min \left(1, \frac{\text{output length}}{\text{reference length}} \right)$$

penalise short output

this one appears in human translation

A Final Word

- Statistical MT
- Neural MT
- Encoder-decoder with attention architecture is a general architecture that can be used for other tasks
 - ▶ Summarisation (lecture 21)
 - ▶ Dialogue generation

Reading

- https://web.stanford.edu/class/cs224n/readings/cs224n-2019-notes06-NMT_seq2seq_attention.pdf (Section 1-5)