



10.11.19 09:43

# Neural Machine Translation<sup>2</sup>

## CMPT 413/825, Fall 2019



Jetic Gū

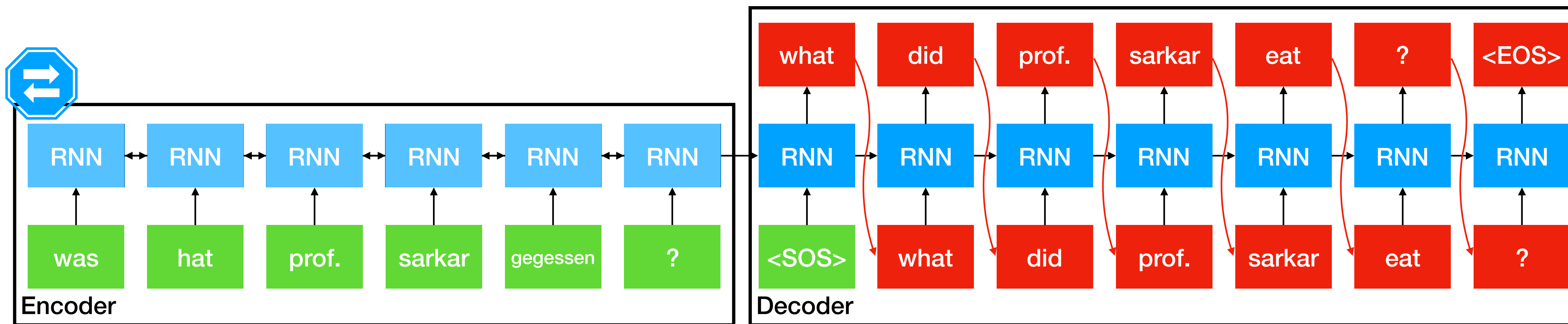
School of Computing Science

12 Nov. 2019

# Overview

- Focus: Neural Machine Translation
- Architecture: Encoder-Decoder Neural Network
- Main Story:
  - Extension to Seq2Seq: Copy Mechanism
  - Extension to Seq2Seq: Ensemble
  - Extension to Seq2Seq: BeamSearch
  - [Extra] Beyond Seq2Seq: Attention is all you need
  - [Extra] Beyond NMT

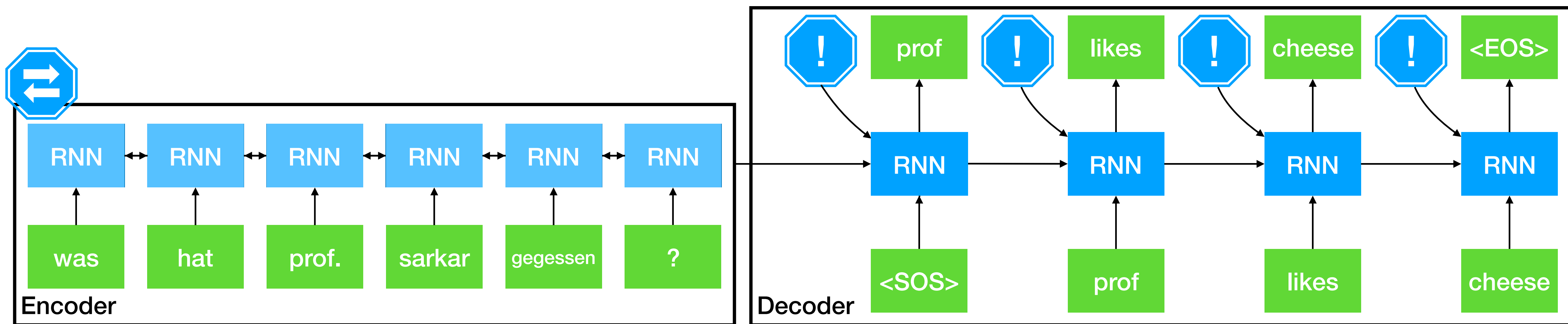
# Sequence-to-Sequence



$$\underline{Pr(E | F) = \prod_t Pr(e' = e_t | F, e_{<t})}$$

CLM

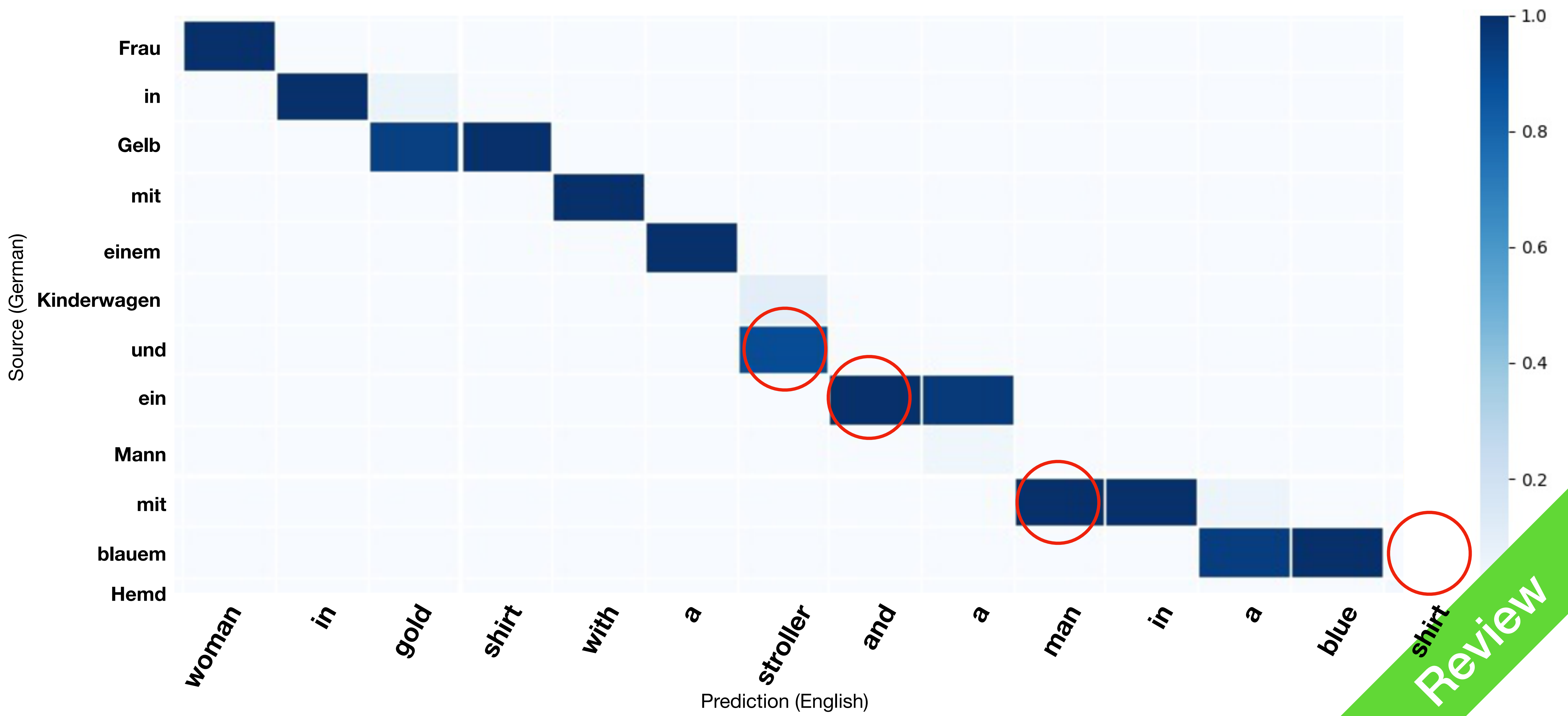
# Attention



# P0 Review



# Attention





# Attention

- Why does Attention work?
- What is in the context vector?
- What is in  $\alpha$ ?
  - ~~It's alignment~~<sup>1</sup> !
- learns to refer to useful information in src<sup>2</sup>
- similar to human attention: we pay attention to whatever is needed

$$\begin{aligned} score_{t,i} &= f(h_i^{enc}, h_t^{dec}) \\ \alpha_t &= \text{softmax}(score_{t,i}) \\ context_t &= \sum_i \alpha_{t,i} h_i^{enc} \end{aligned}$$

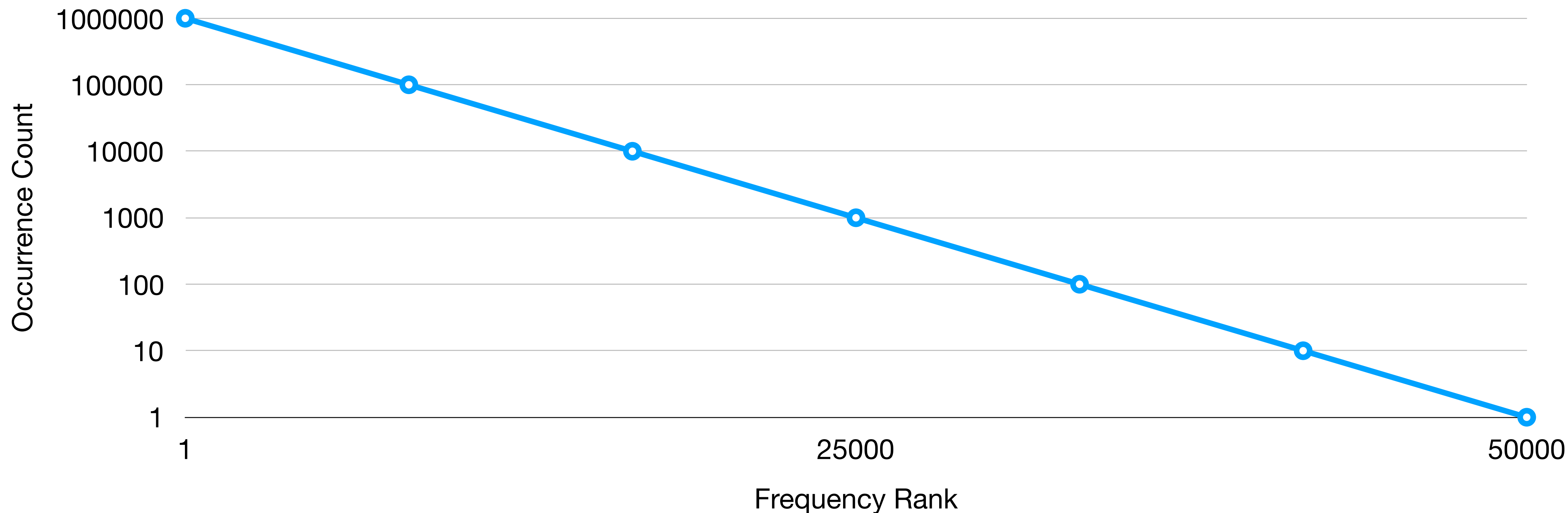
1. CL2015008 [Bahdanau et al.] Neural Machine Translation by Jointly Learning to Align and Translate  
2. CL2017342 [Ghader et Monz] What does Attention in Neural Machine Translation Pay Attention to?

# Common Problems of NMT

- Out-of-Vocabulary (OOV) Problem; Rare word problem
  - Frequent word are translated correctly, rare words are not
  - In any corpus, word frequencies are exponentially unbalanced

# OOV

## Zipf's Law

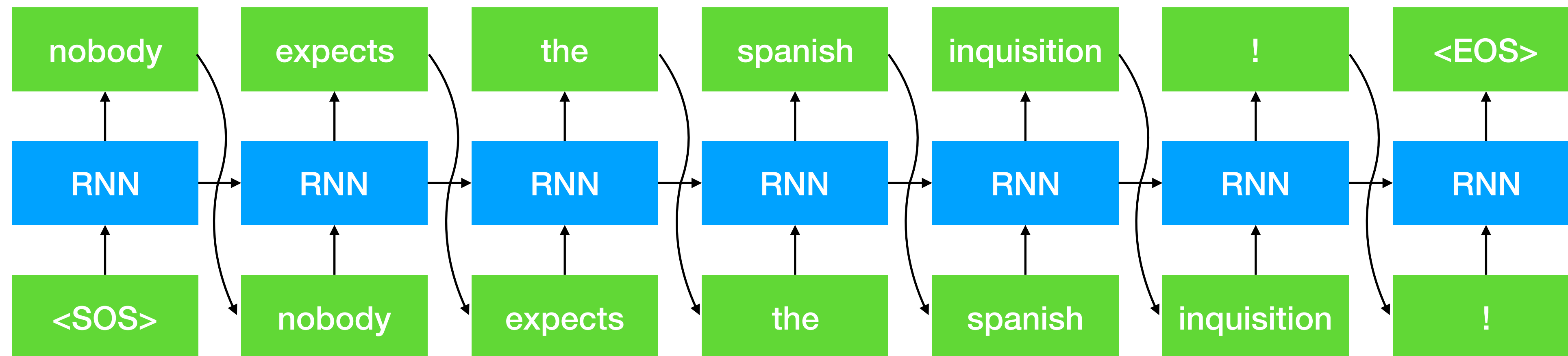


- rare word are exponentially less frequent than frequent words
- e.g. in HW4, 45%|65% (src|tgt) of the unique words occur once

# Common Problems of NMT

- Out-of-Vocabulary (OOV) Problem; Rare word problem
  - Frequent word are translated correctly, rare words are not
  - In any corpus, word frequencies are exponentially unbalanced
- Under translation
  - Crucial information are left untranslated; premature `<EOS>` generation

# Under-trans<EOS>

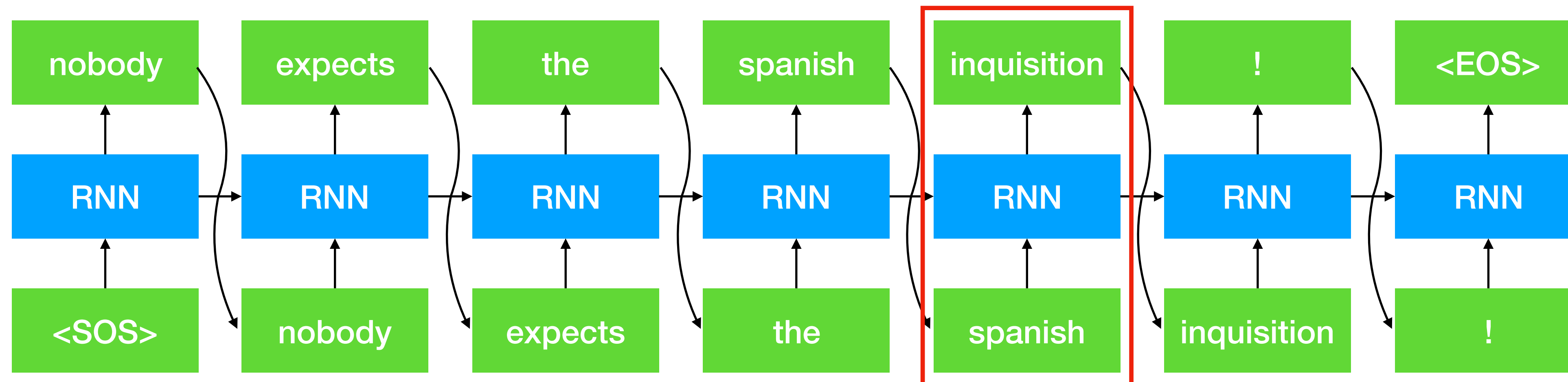


- Under translation
- Crucial information are left untranslated; premature <EOS> generation

# Under-trans<EOS>



'inquisition': 0.49  
<EOS>: 0.51

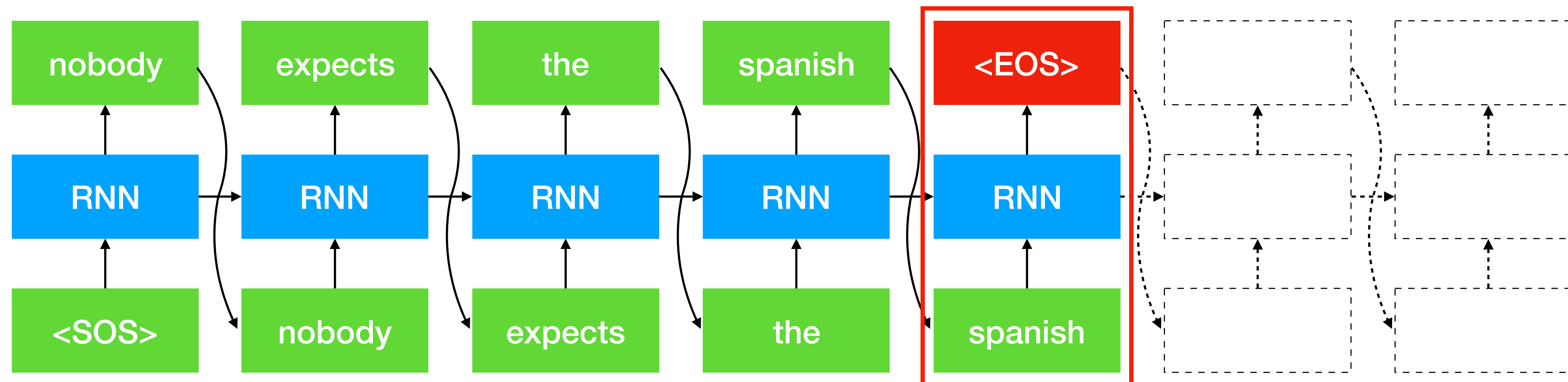


- Under translation
- Crucial information are left untranslated; premature <EOS> generation

# Under-trans<EOS>



'inquisition': 0.49  
<EOS>: 0.51



- Under translation
- Crucial information are left untranslated; premature <EOS> generation

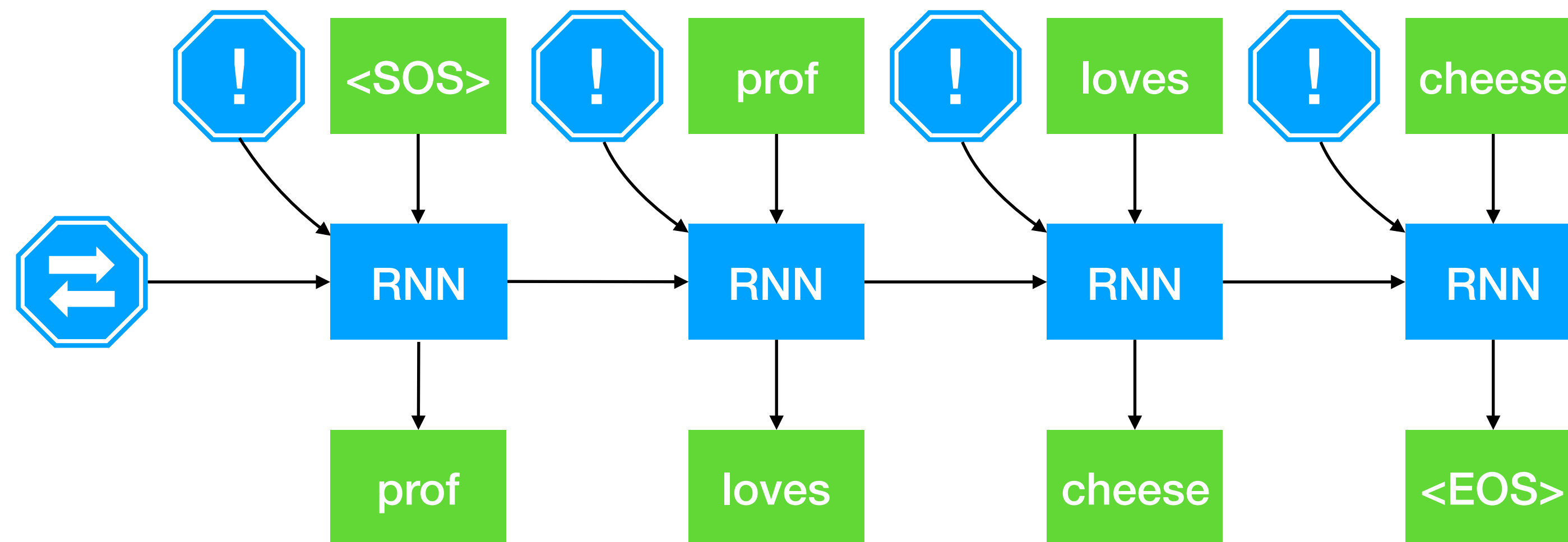
# Common Problems of NMT

- Out-of-Vocabulary (OOV) Problem; Rare word problem
  - Frequent word are translated correctly, rare words are not
  - In any corpus, word frequencies are exponentially unbalanced
- Under translation
  - Crucial information are left untranslated; premature `<EOS>` generation

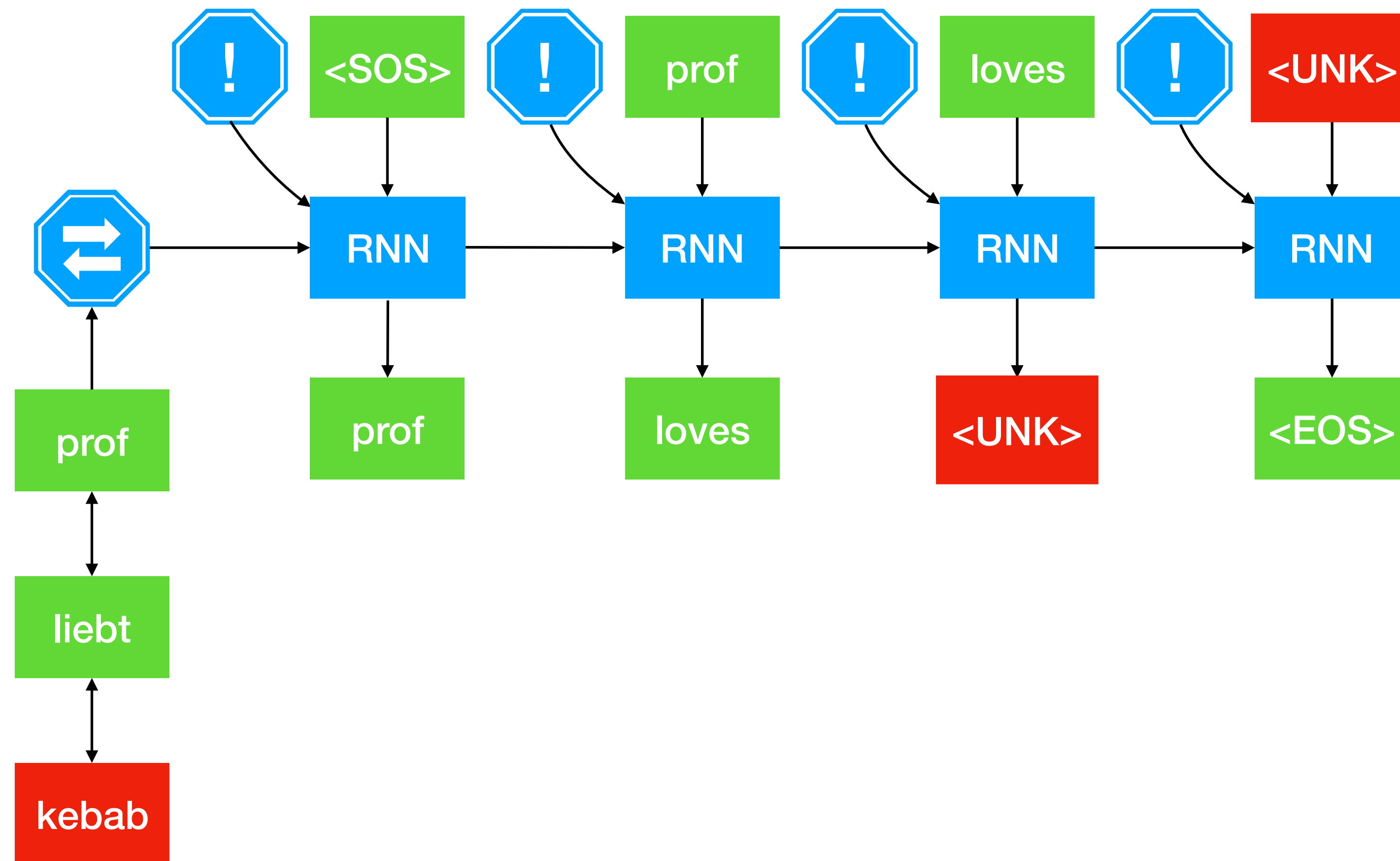
# Common Problems of NMT

- Out-of-Vocabulary (OOV) Problem; Rare word problem
  - Copy Mechanisms
  - Char-level Encoder (oops for logogram, e.g. Chinese)
- Under translation
  - Ensemble
  - Beam search
  - Coverage models

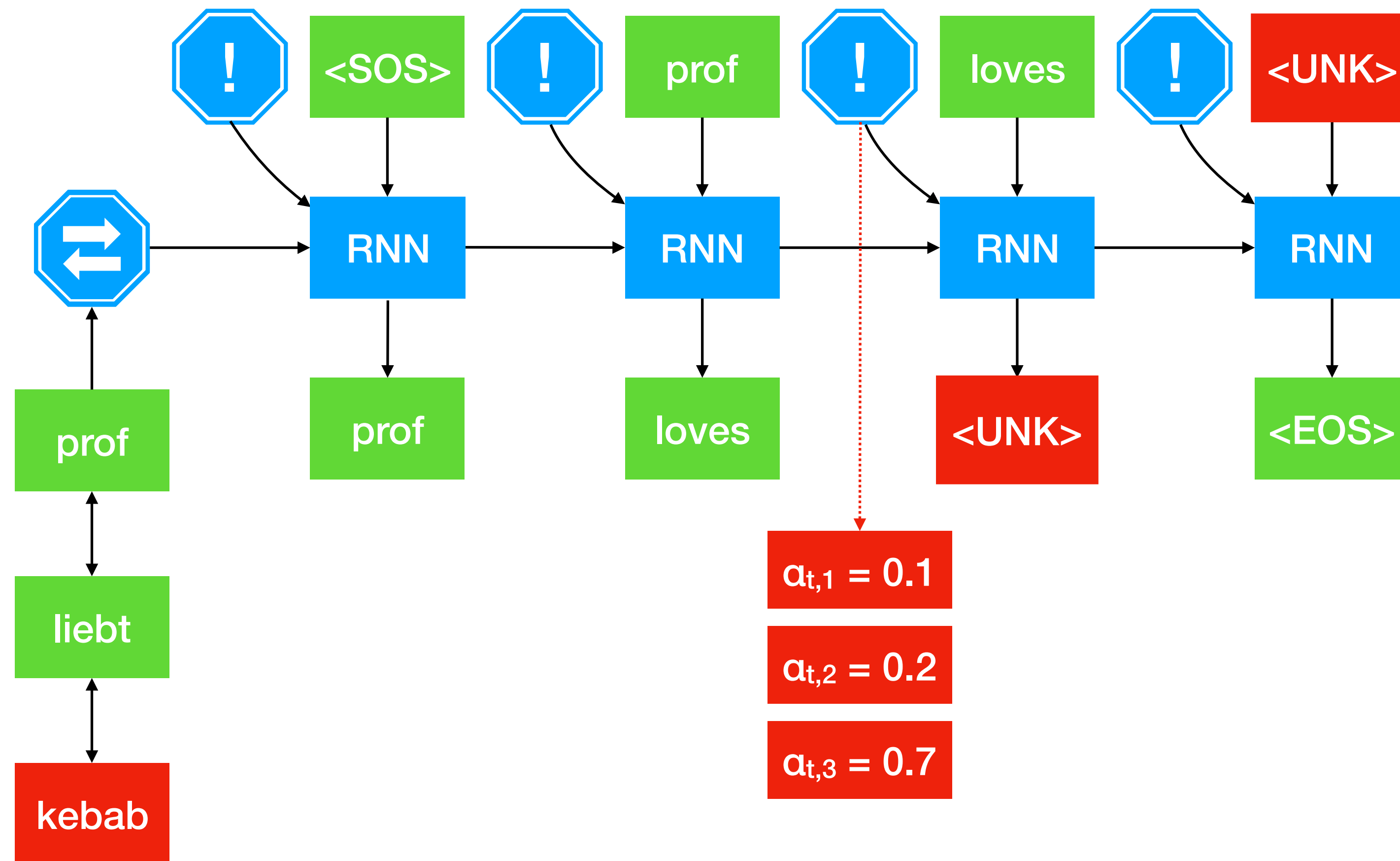
# Copy Mechanism



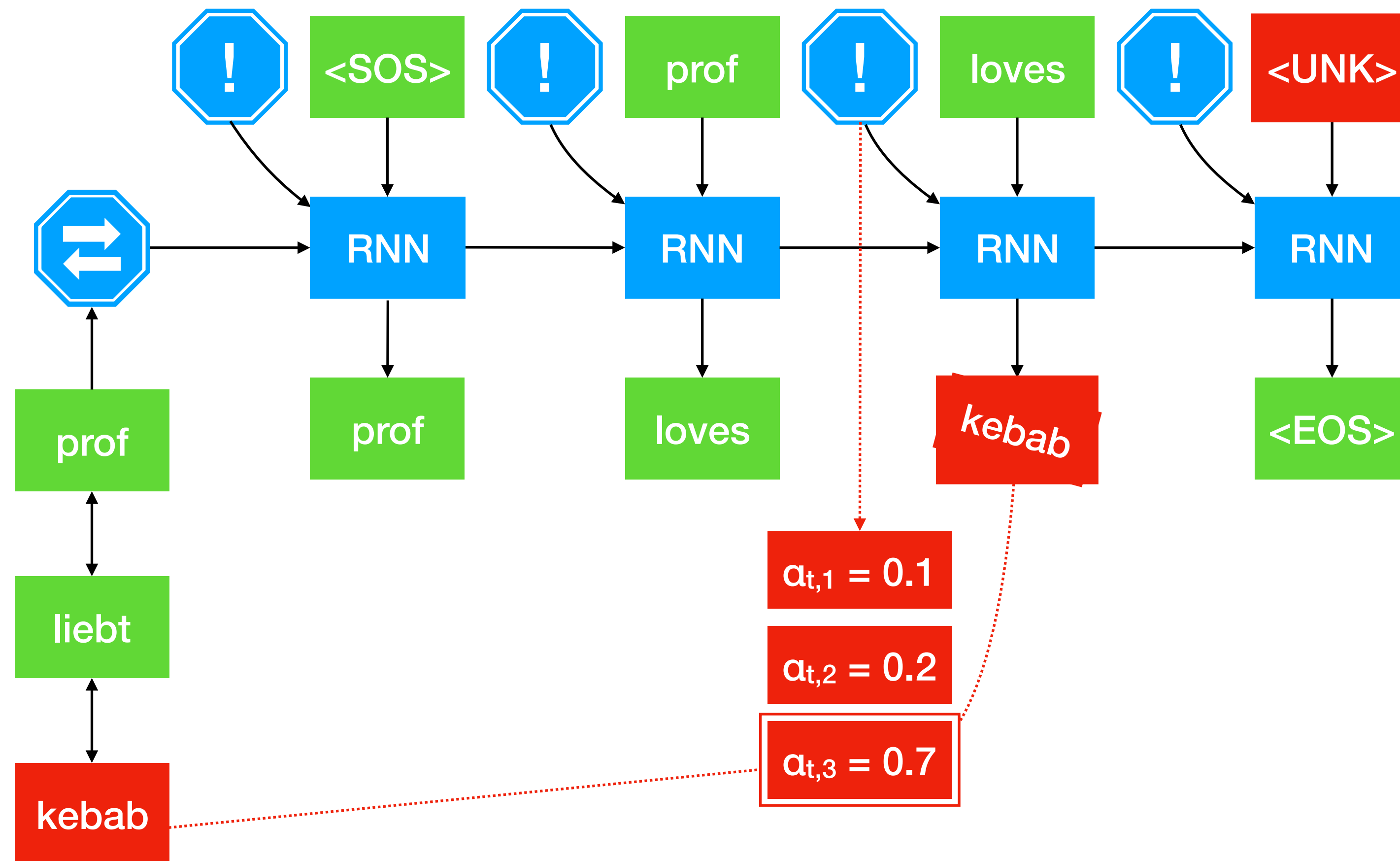
# Copy Mechanism



# Copy Mechanism

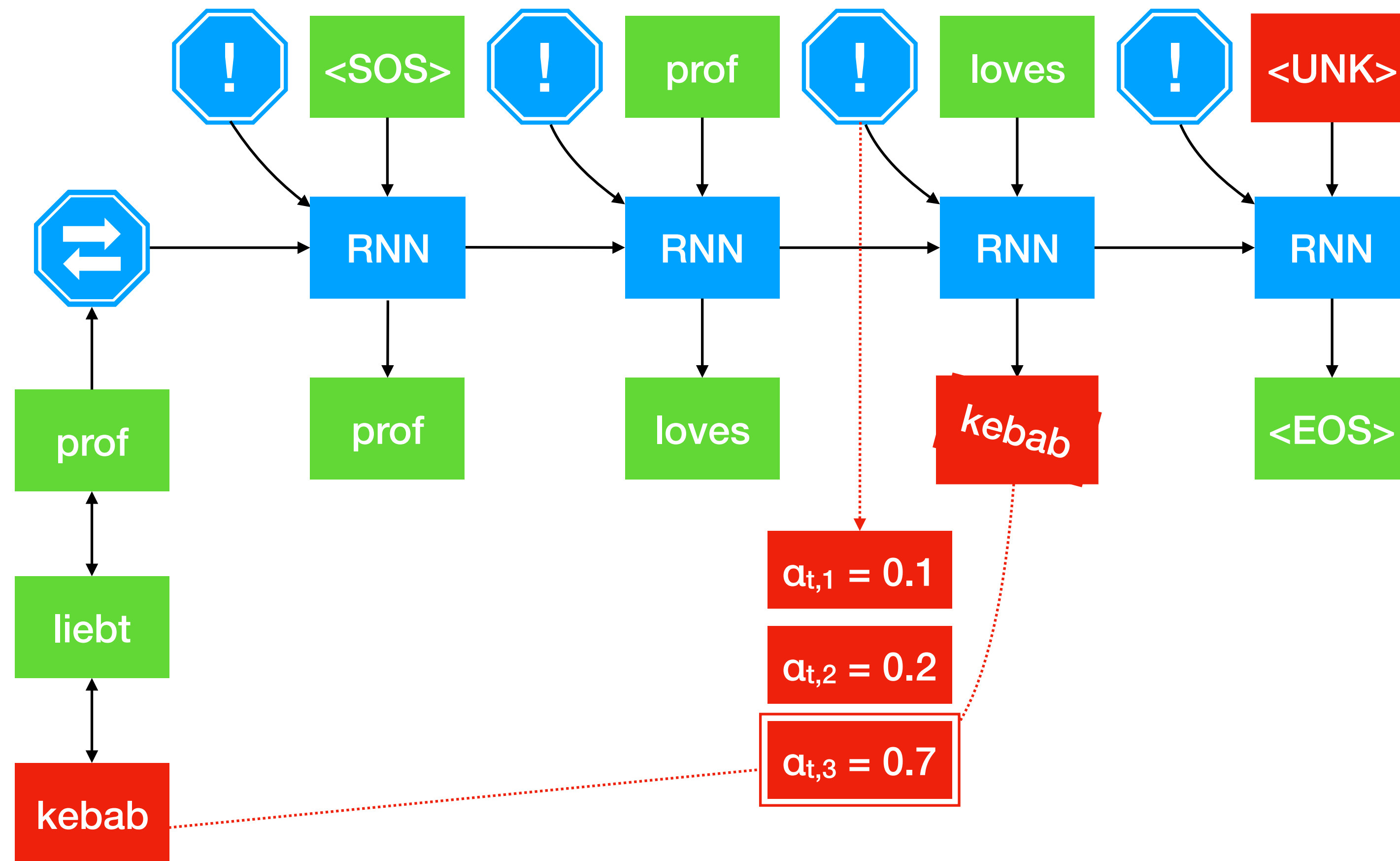


# Copy Mechanism

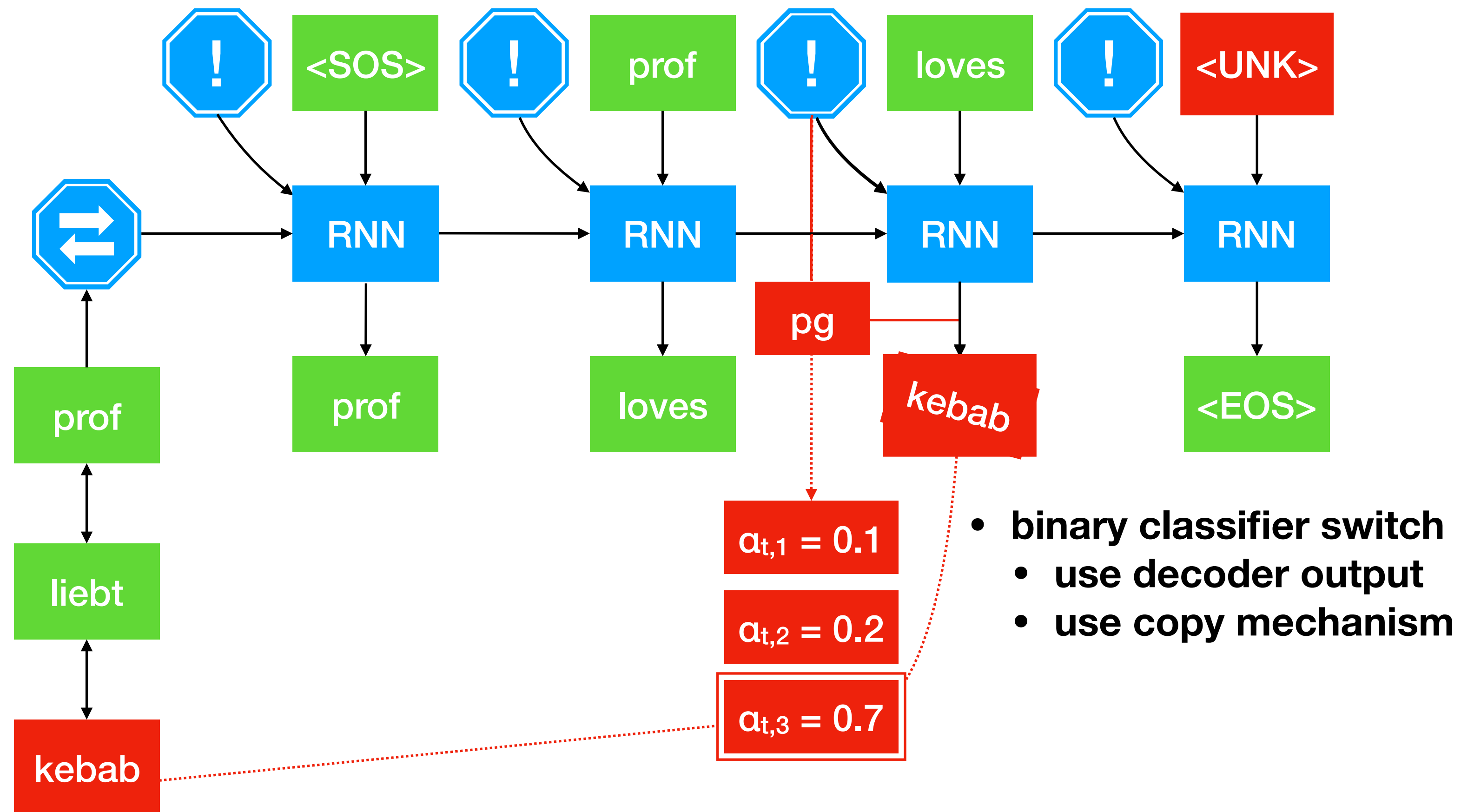


- 
- The diagram illustrates an unrolled Recurrent Neural Network (RNN) with three time steps. The input sequence is "prof", "loves", "kebab". The hidden state is passed through three "RNN" blocks. The output sequence is "<SOS>", "prof", "loves". The third time step shows a red "X" over the output "kebab" and a red box with  $\alpha_{t,3} = 0.7$ , indicating a high probability of a word boundary.

# \*Pointer-Generator Copy Mechanism

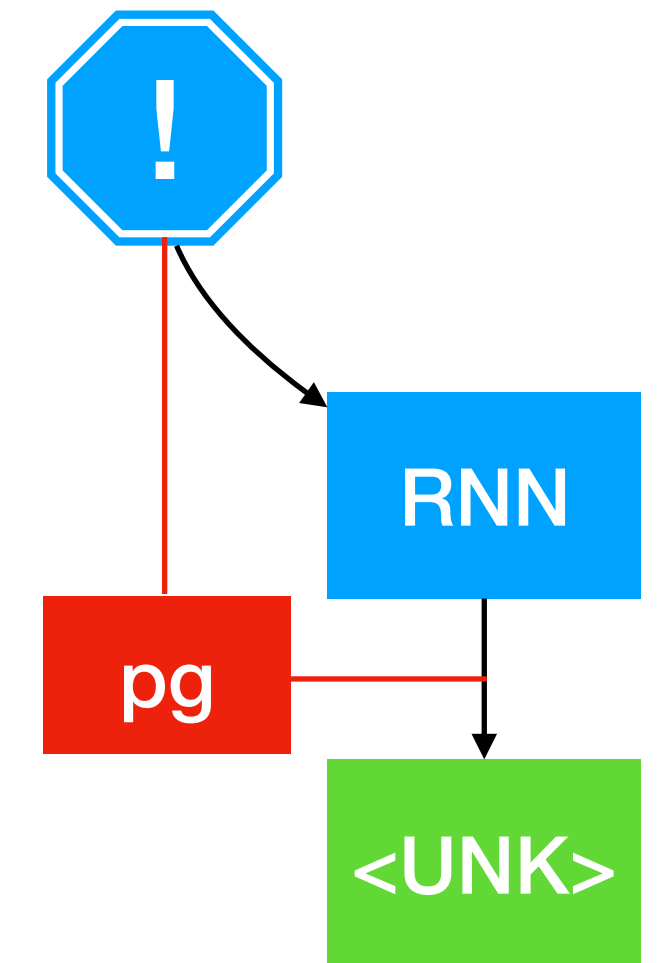


# \*Pointer-Generator Copy Mechanism



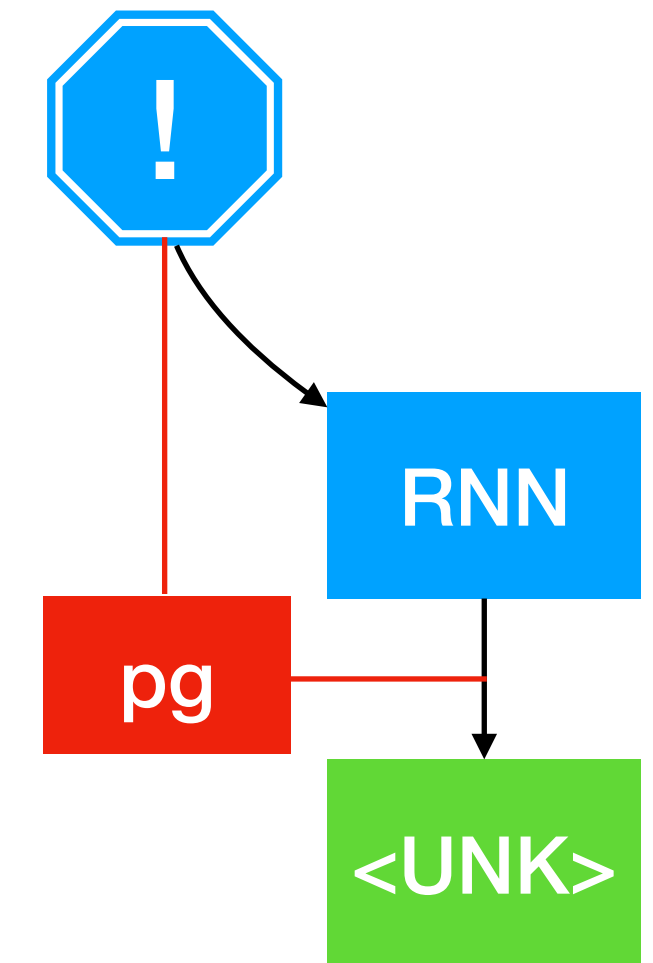
# \*Pointer-Generator Copy Mechanism

- Sees  $\langle \text{UNK} \rangle$  at step  $t$ , or  $pg([h_t^{dec}; context_t]) \leq 0.5$
- looks at attention weight  $\alpha_t$
- replace  $\langle \text{UNK} \rangle$  with the source word  $f_{\arg\max_i \alpha_{t,i}}$



# \*Pointer-Based Dictionary Fusion

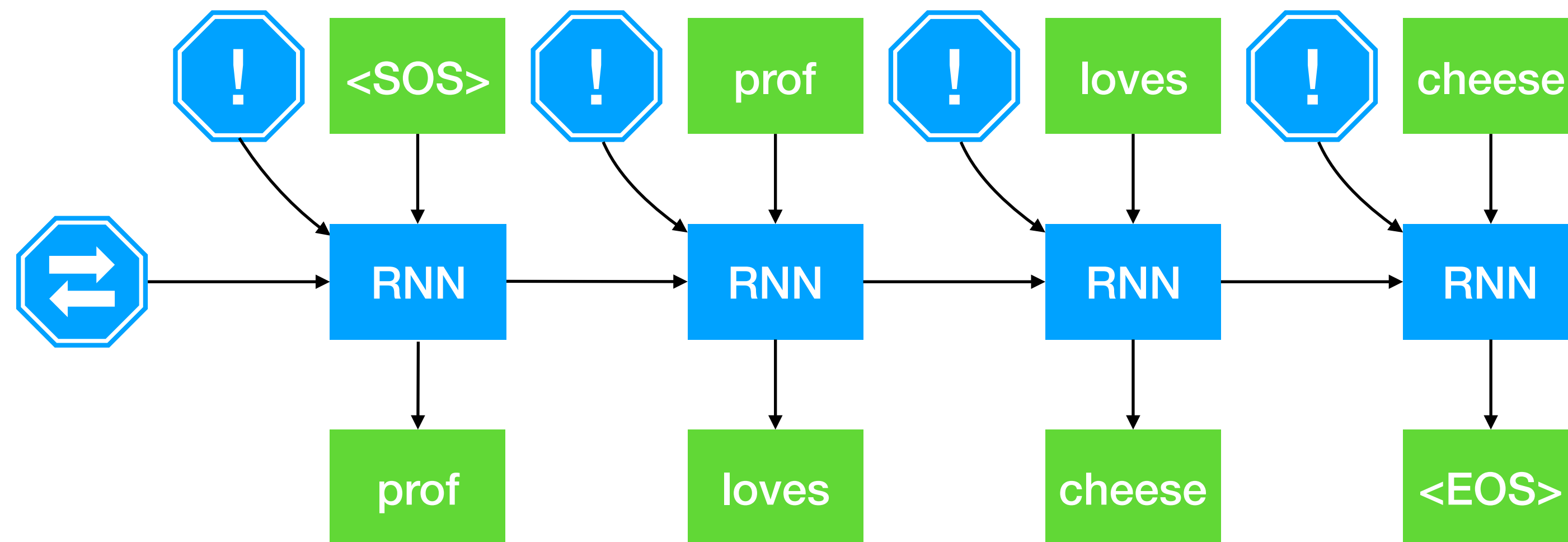
- Sees  $\langle \text{UNK} \rangle$  at step  $t$ , or  $pg([h_t^{dec}; context_t]) \leq 0.5$
- looks at attention weight  $\alpha_t$
- replace  $\langle \text{UNK} \rangle$  with **translation of** the source word  
 $\text{dict}(f_{\text{argmax}_i \alpha_{t,i}})$



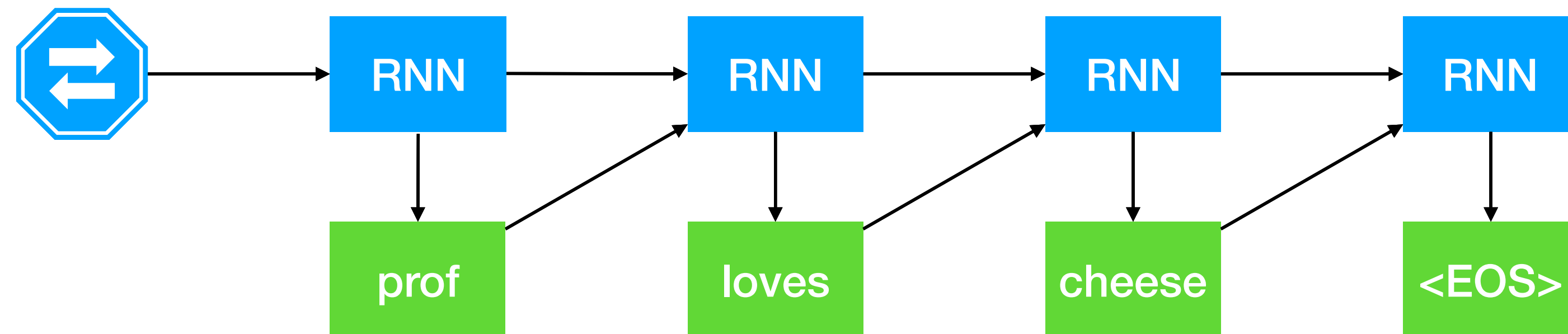
# Ensemble

- Similar to voting mechanism, but with probabilities
- multiple models of different parameters (usually from different checkpoints of the same training instance)
- use the output with the highest probability across all models

# Ensemble

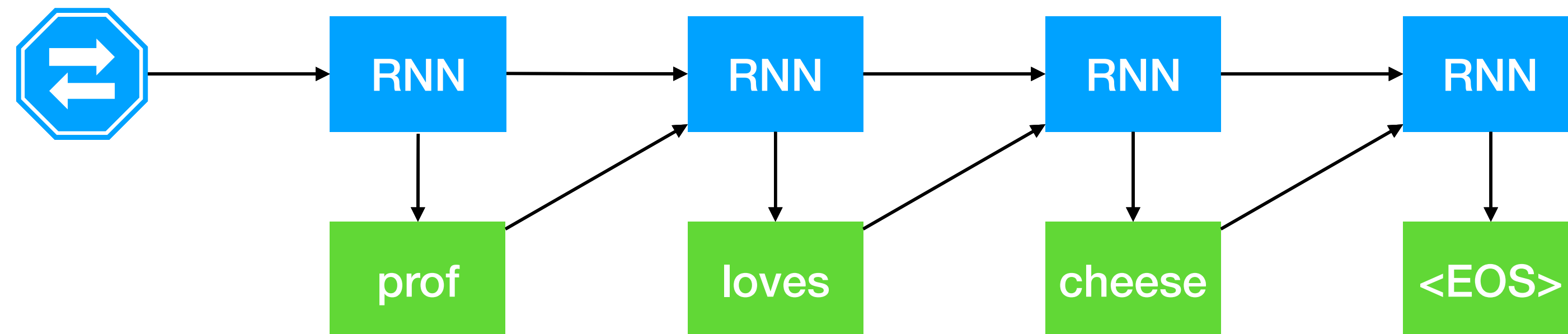


# Ensemble



# Ensemble

seq2seq\_E049.pt

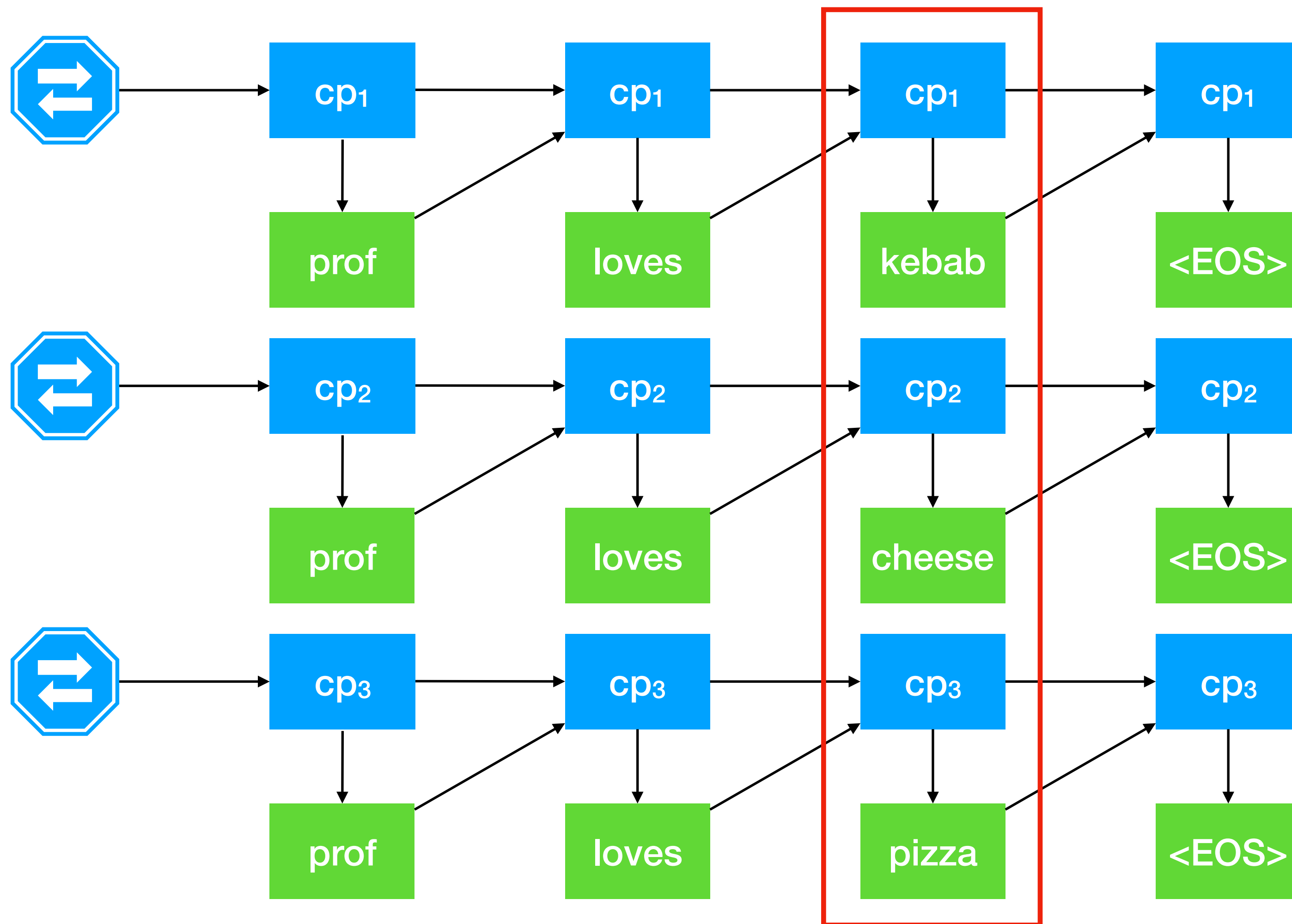


# Ensemble

seq2seq\_E049.pt

seq2seq\_E048.pt

seq2seq\_E047.pt

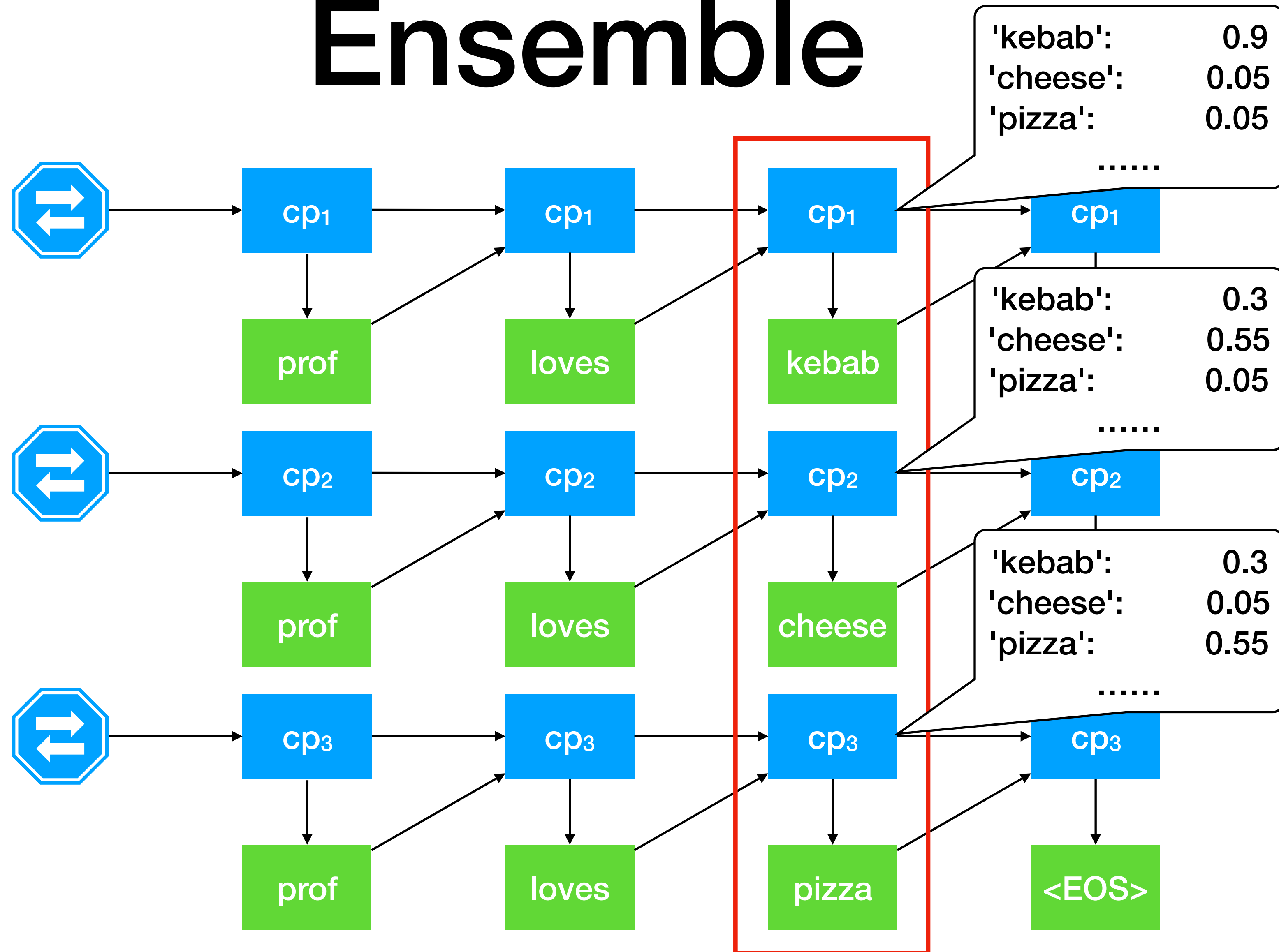


# Ensemble

seq2seq\_E049.pt

seq2seq\_E048.pt

seq2seq\_E047.pt

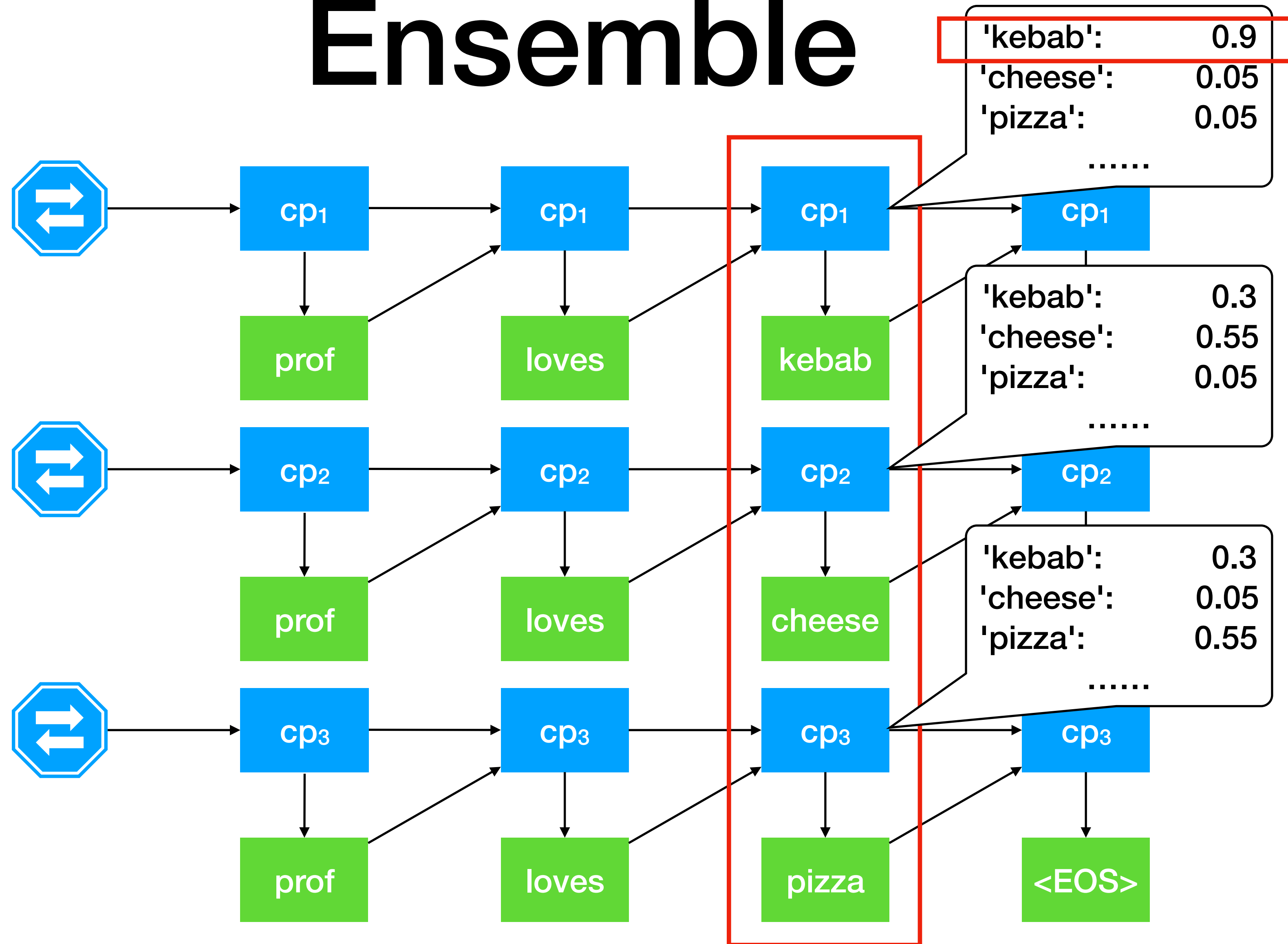


# Ensemble

seq2seq\_E049.pt

seq2seq\_E048.pt

seq2seq\_E047.pt

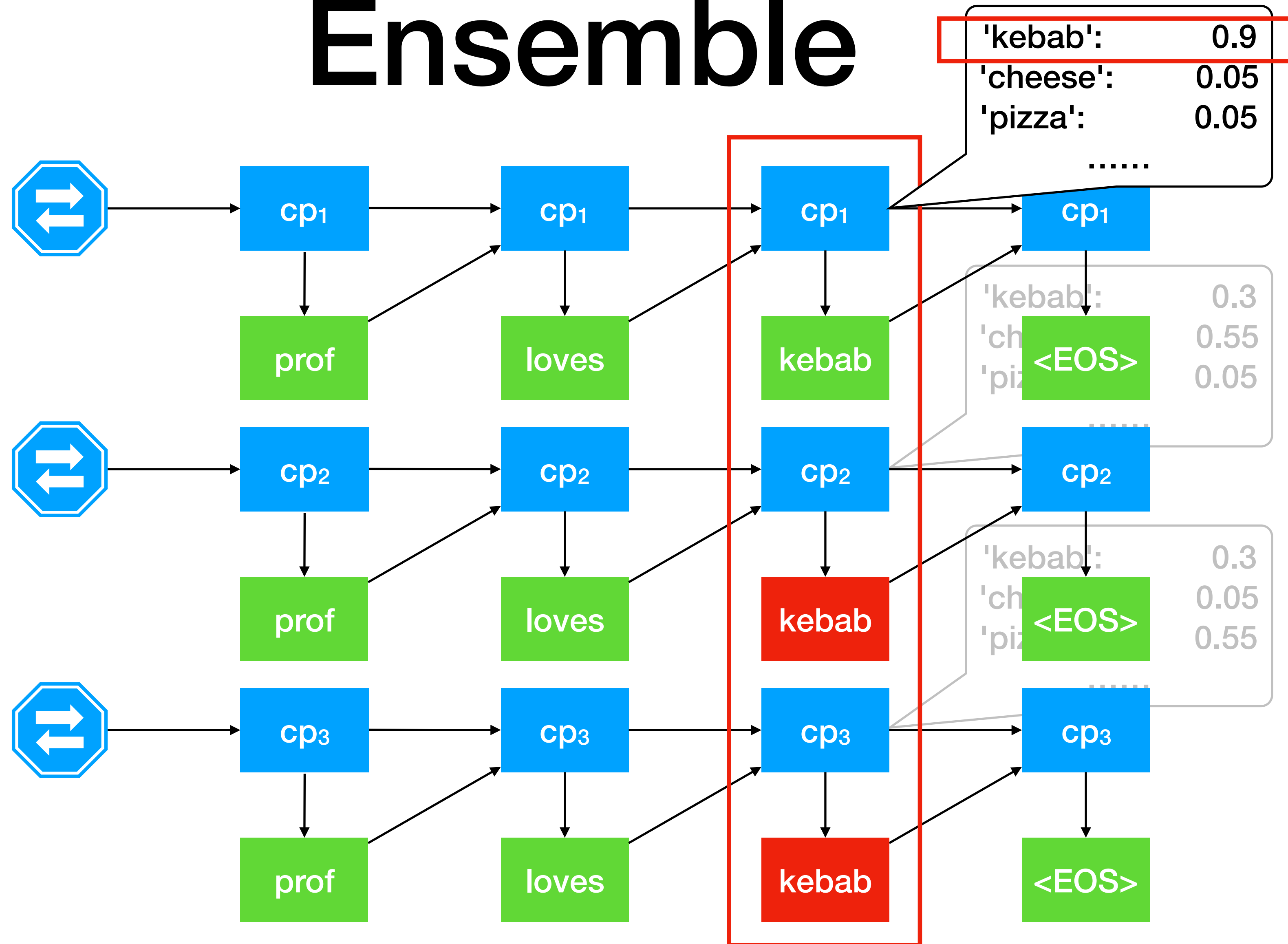


# Ensemble

seq2seq\_E049.pt

seq2seq\_E048.pt

seq2seq\_E047.pt

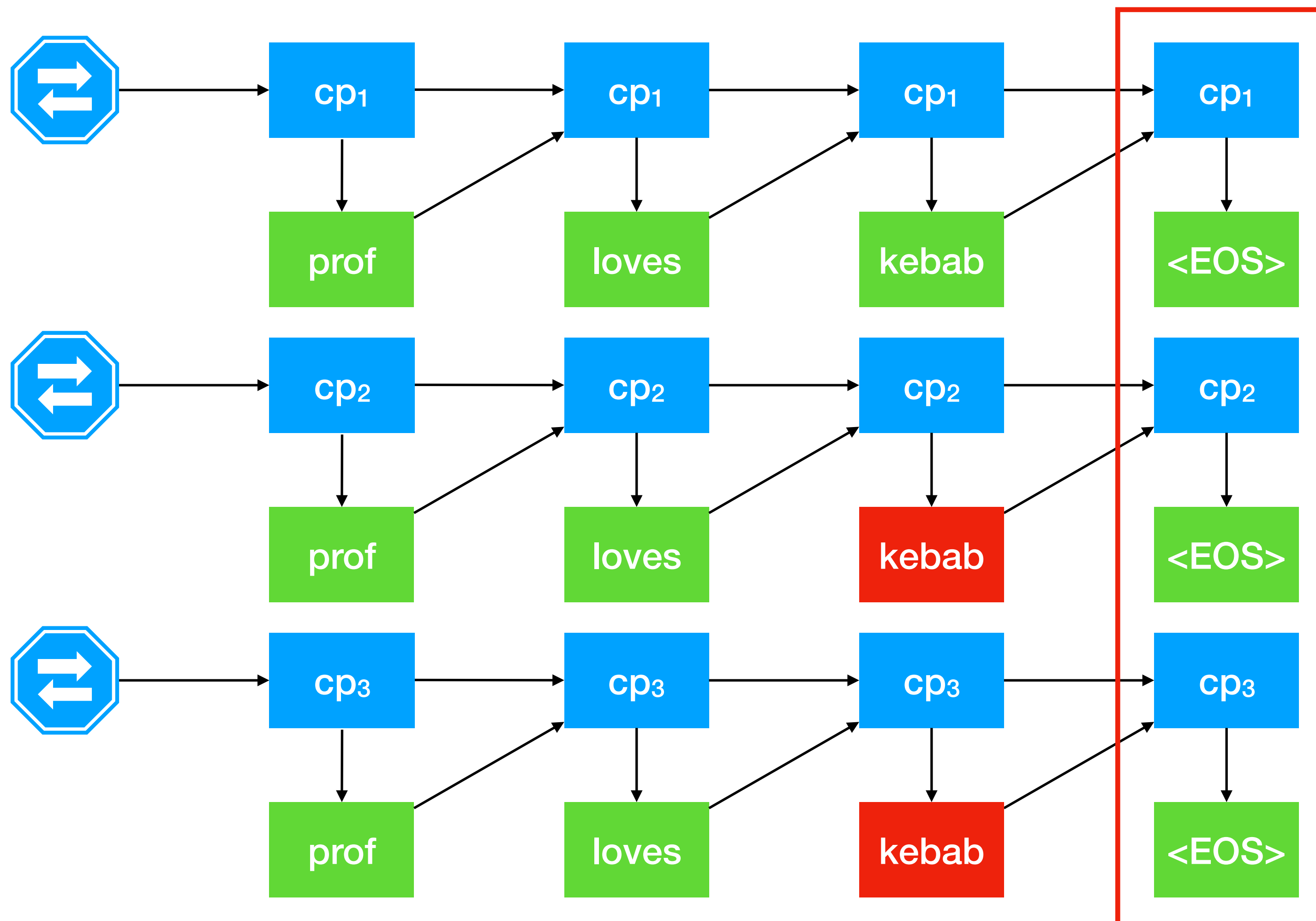


# Ensemble

seq2seq\_E049.pt

seq2seq\_E048.pt

seq2seq\_E047.pt

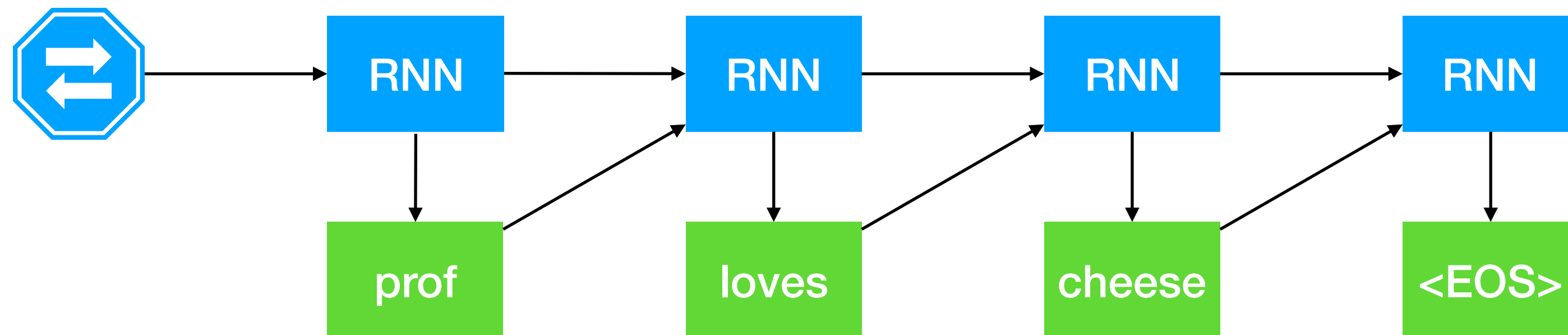


Demo

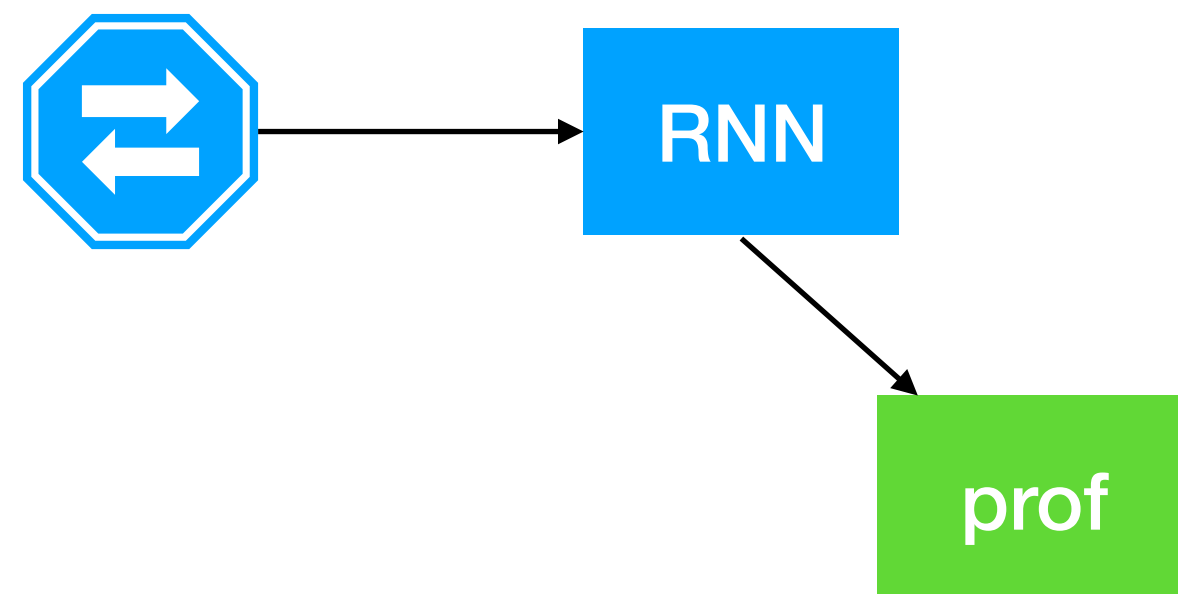
# BeamSearch

- Consider multiple hypothesis at each step  $t$ , formulate decoding as a search programme on a tree

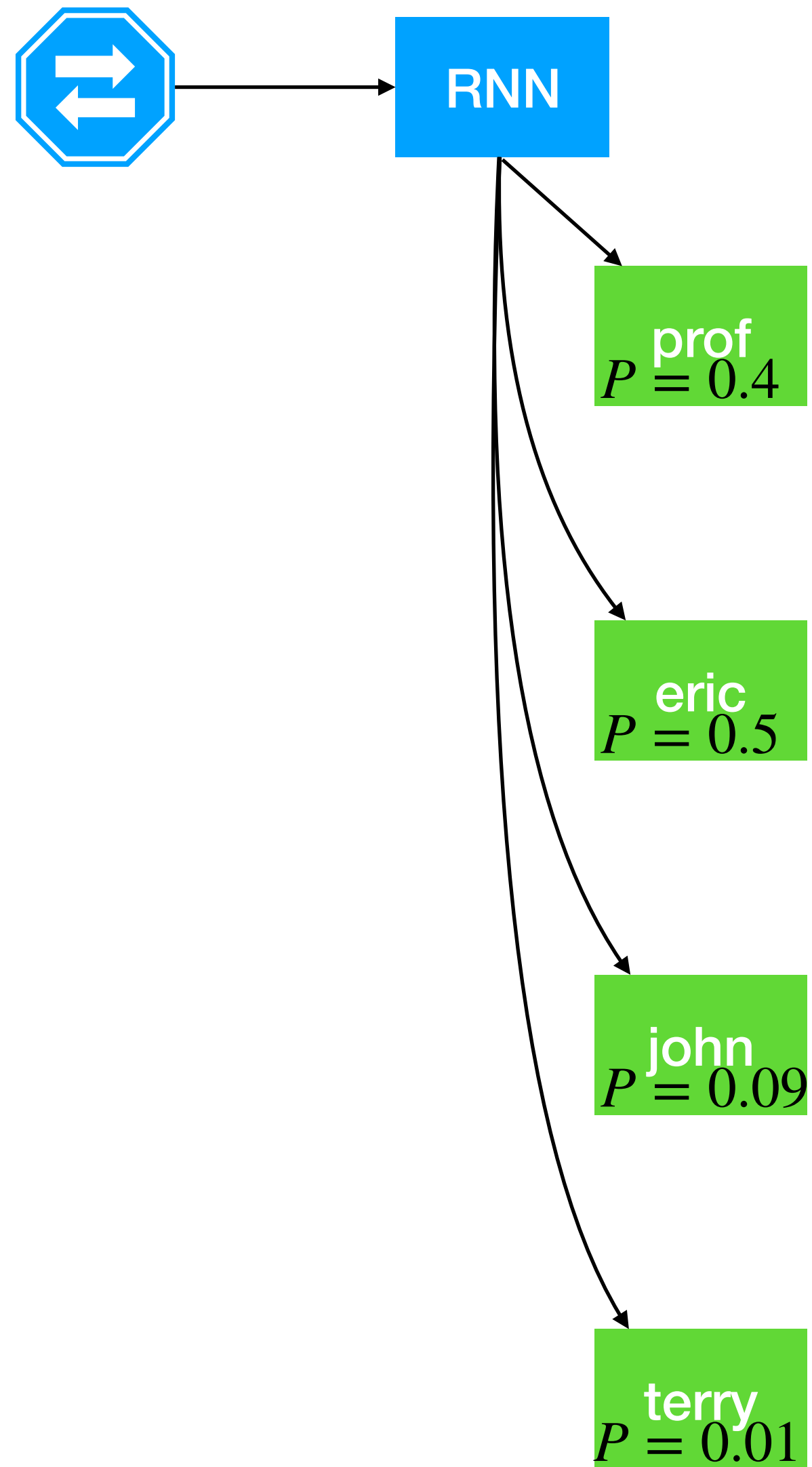
# BeamSearch



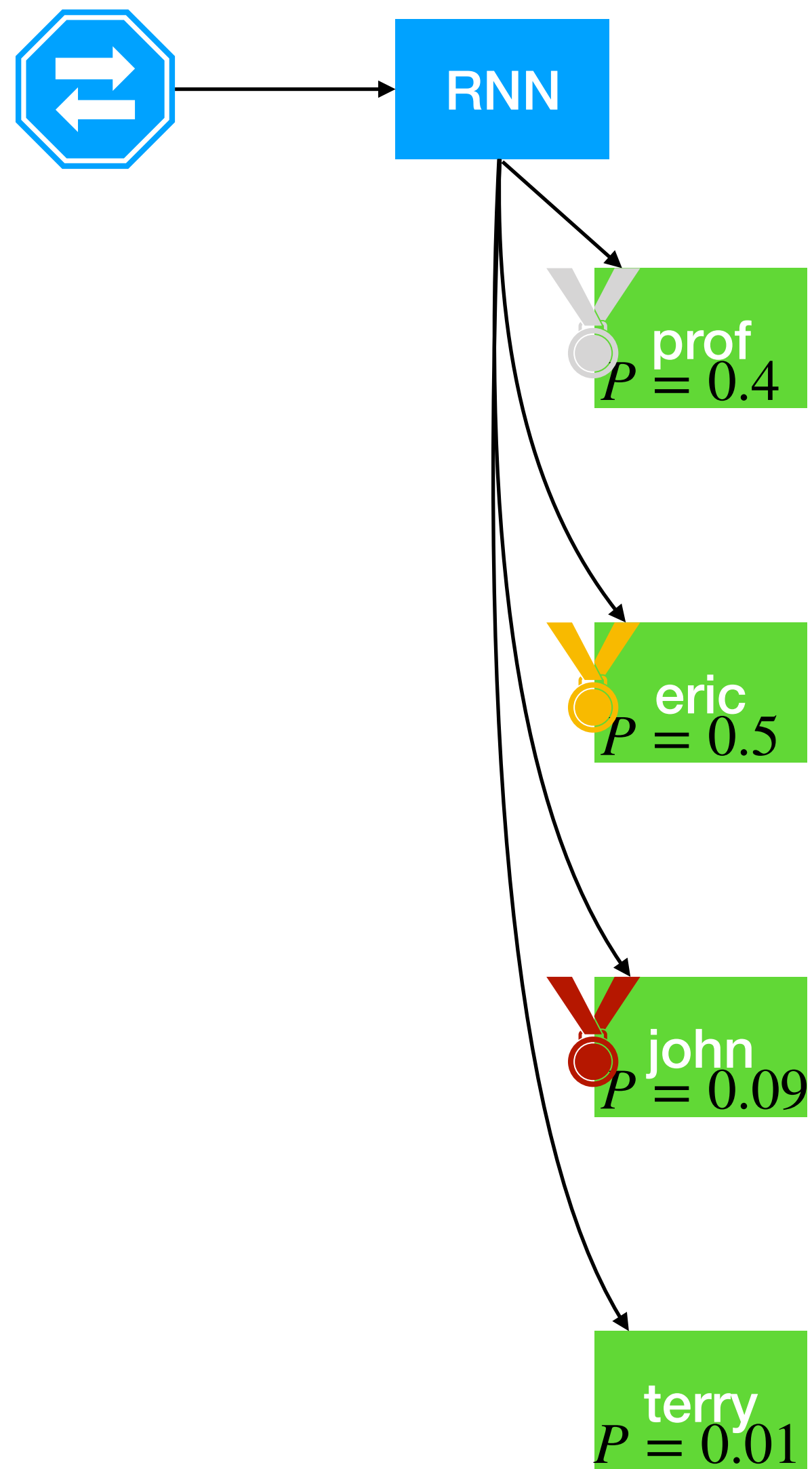
# BeamSearch (Size=3)



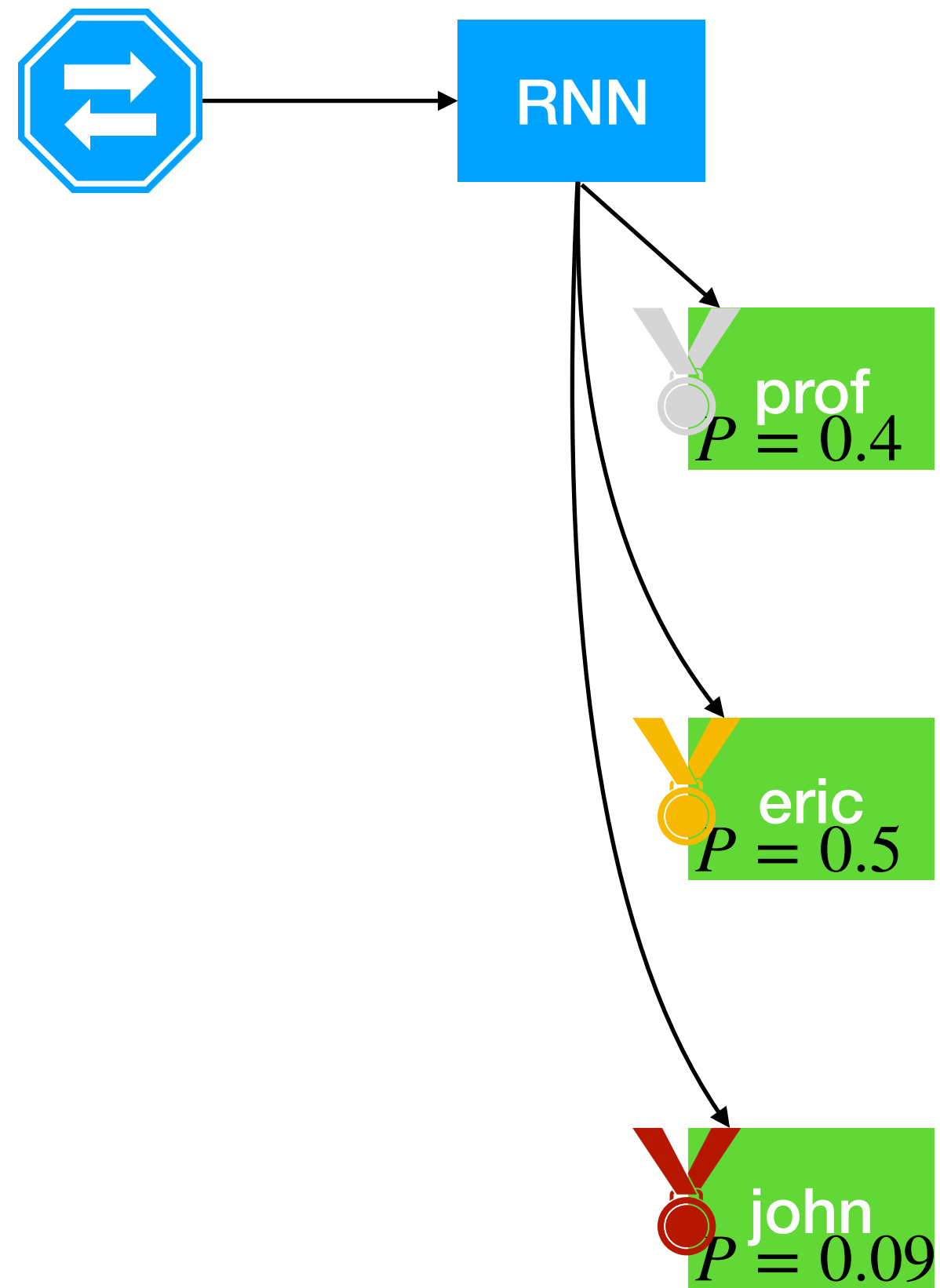
# BeamSearch (Size=3)



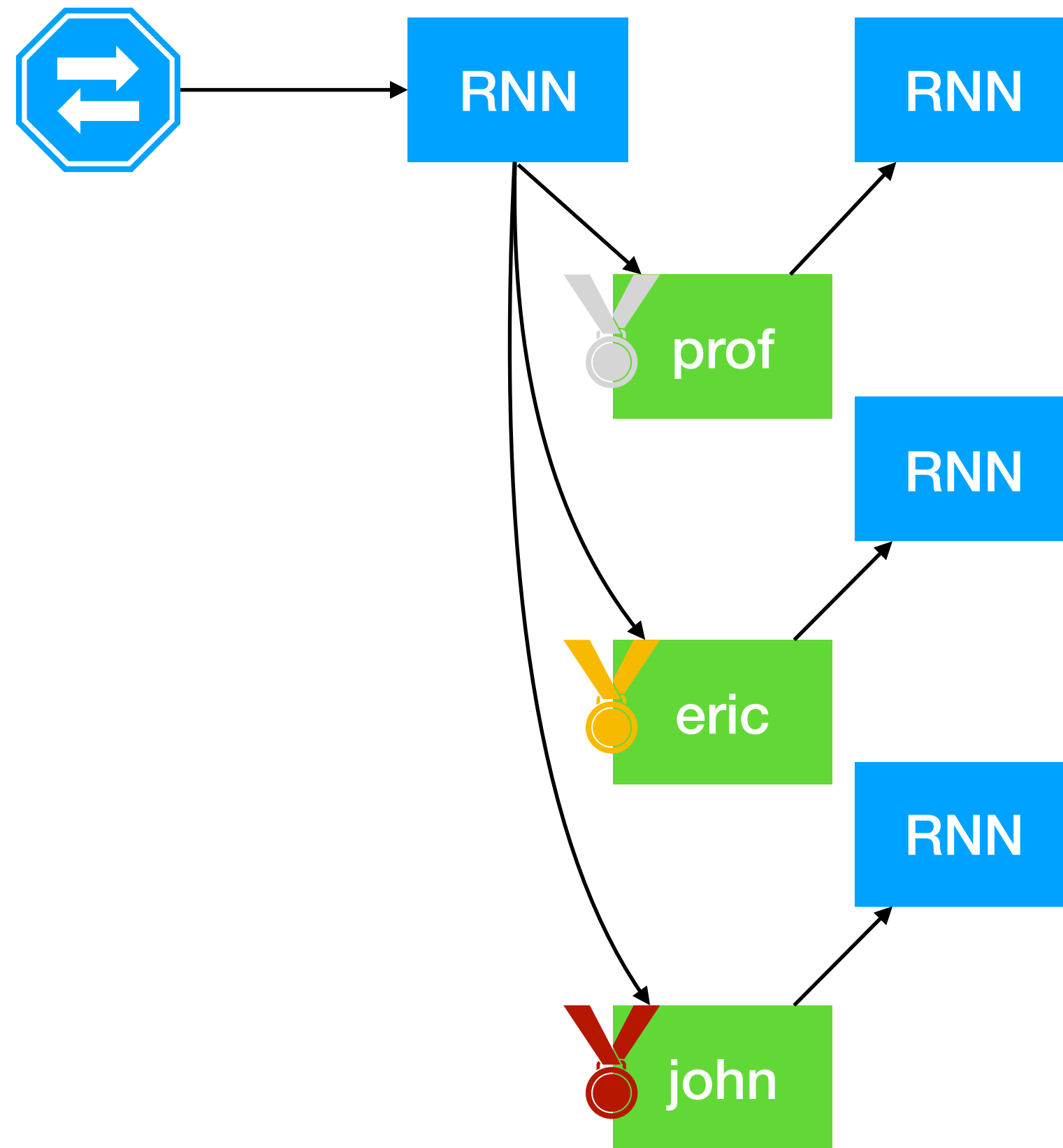
# BeamSearch (Size=3)



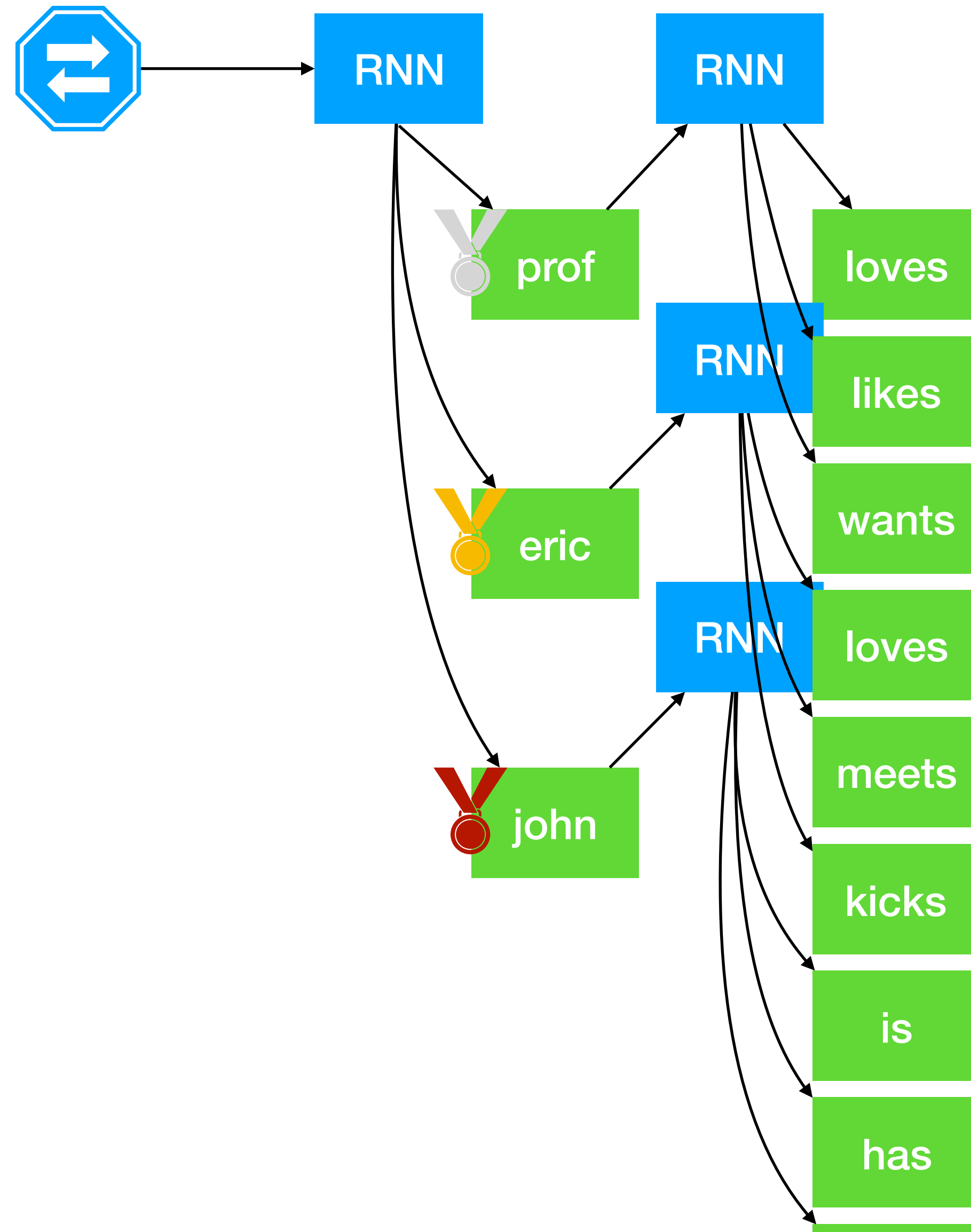
# BeamSearch (Size=3)



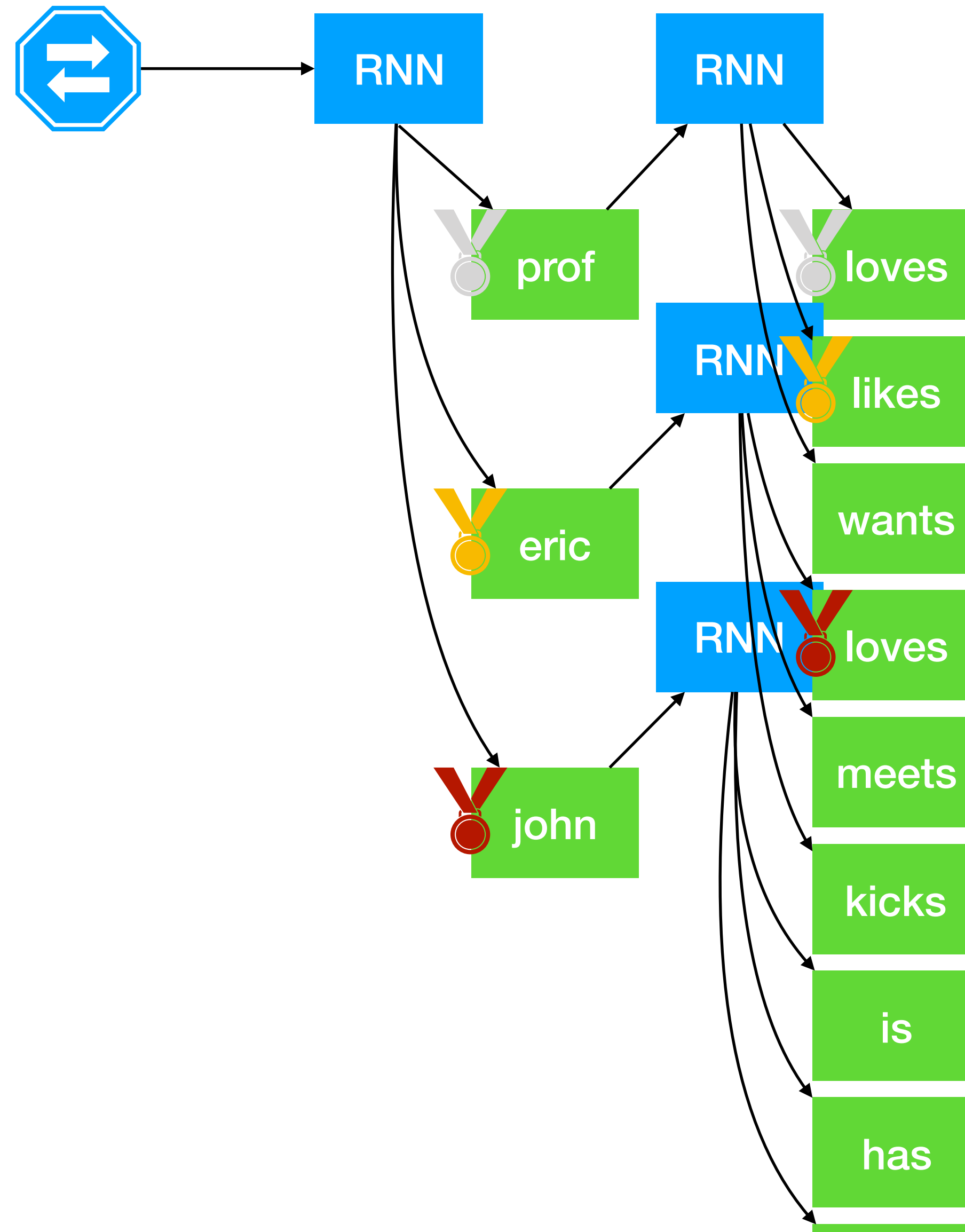
# BeamSearch (Size=3)



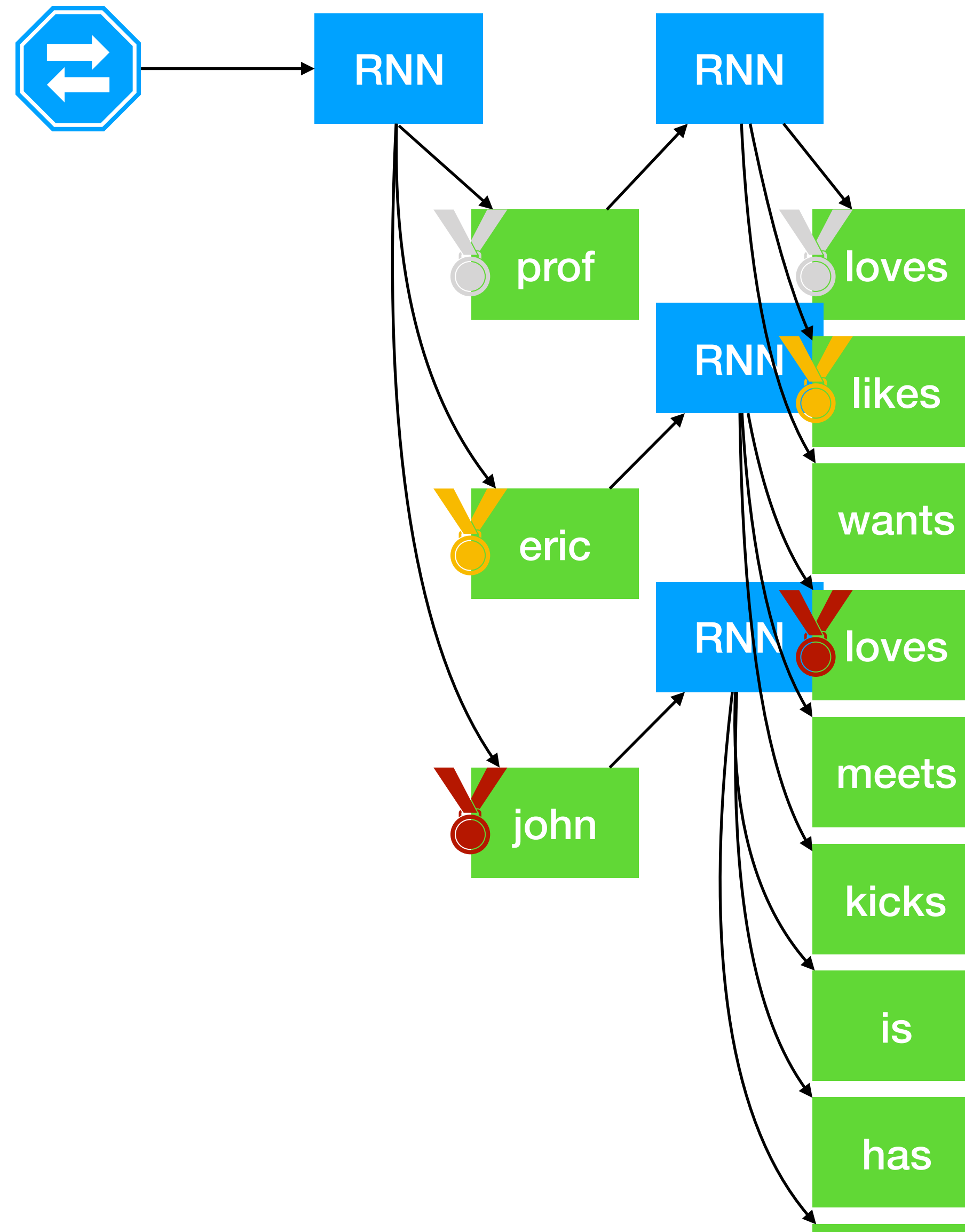
# BeamSearch (Size=3)



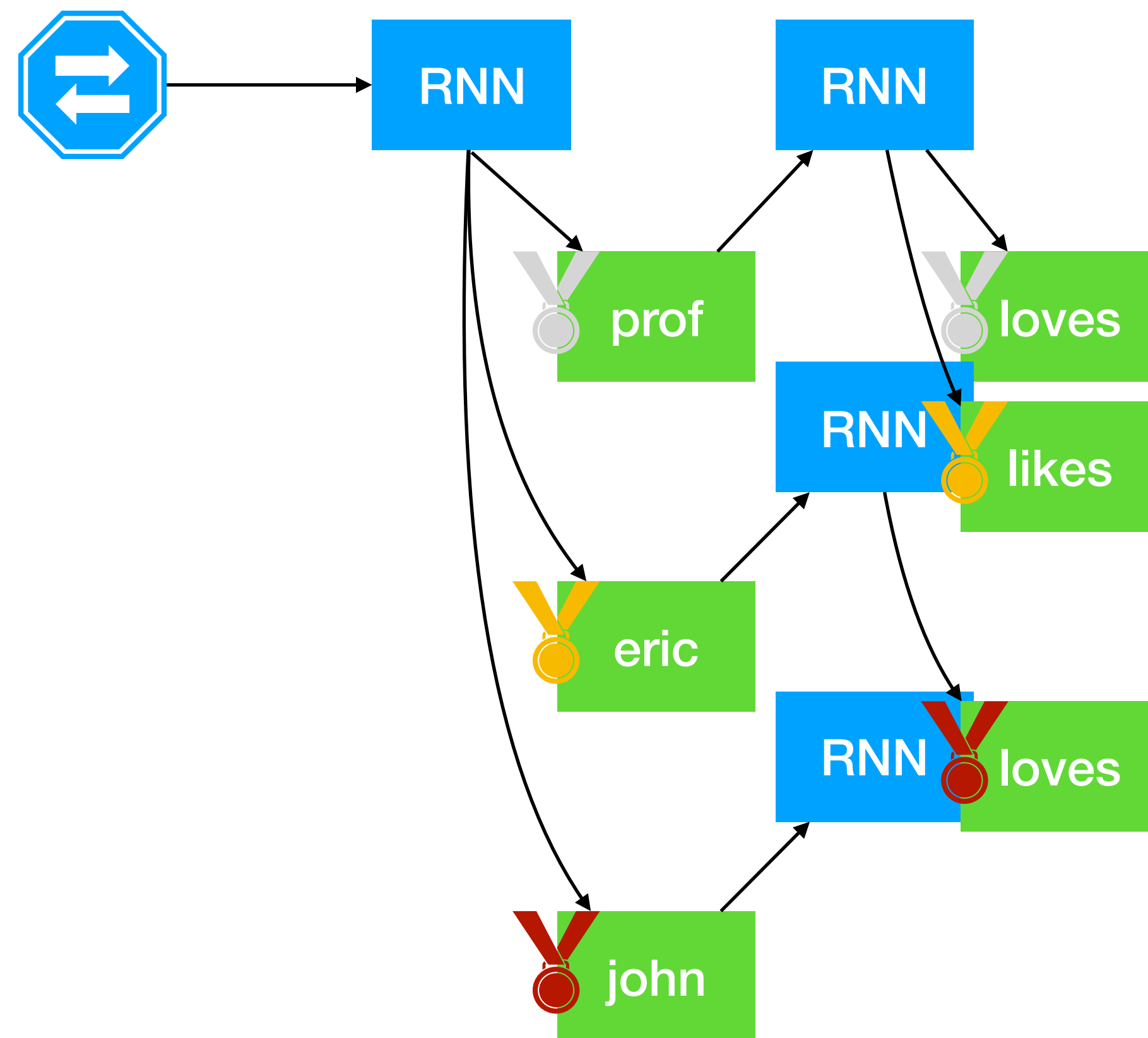
# BeamSearch (Size=3)



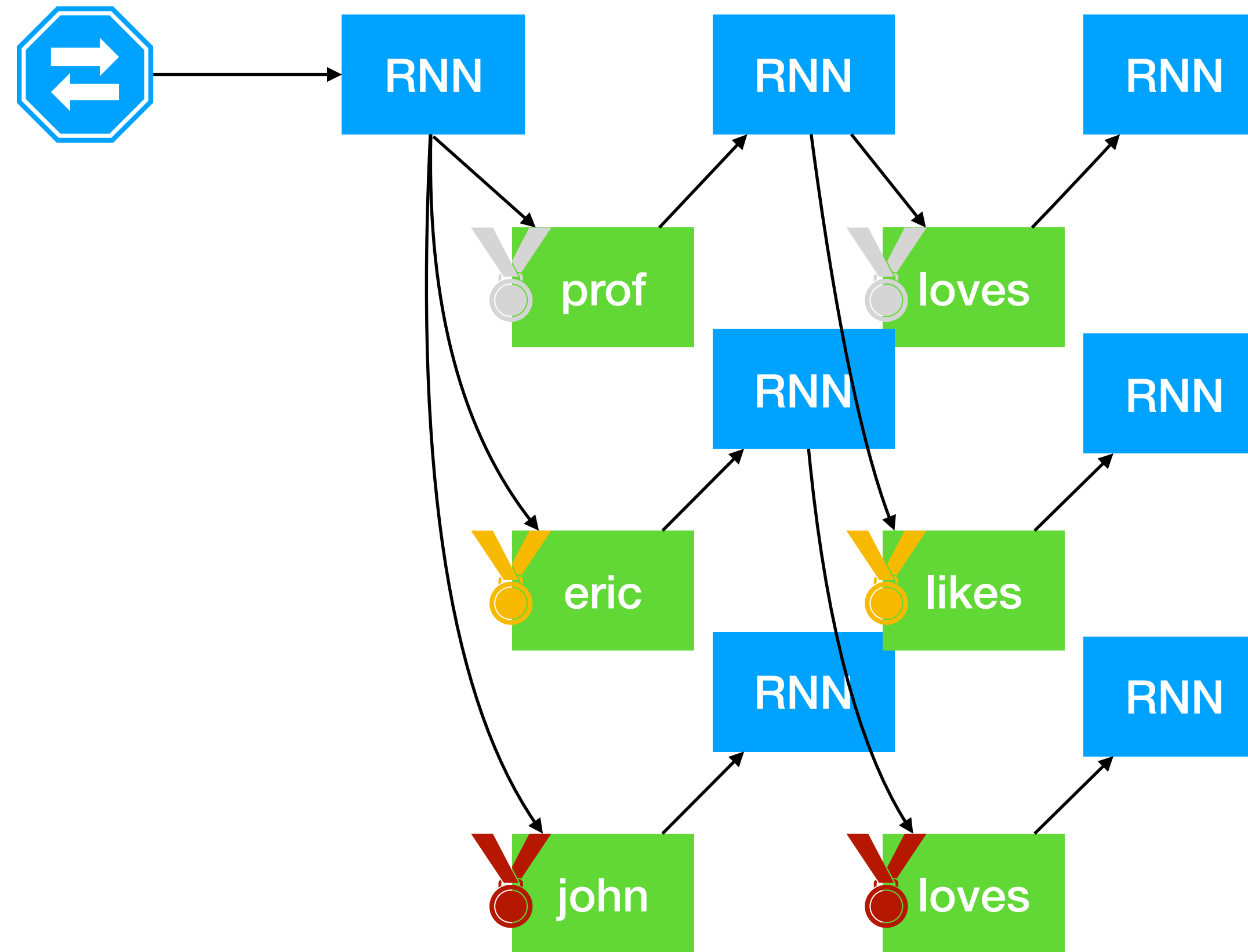
# BeamSearch (Size=3)



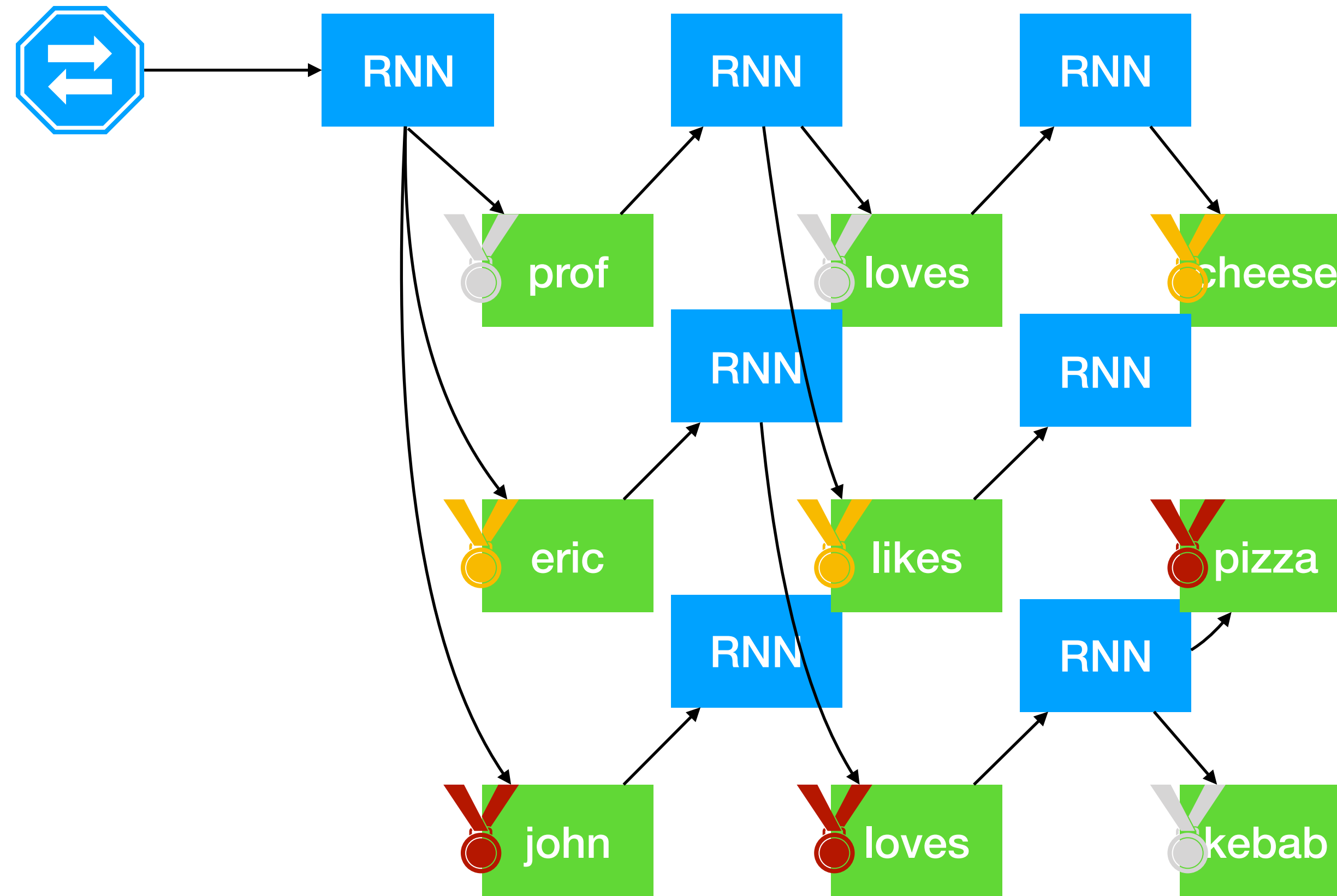
# BeamSearch (Size=3)



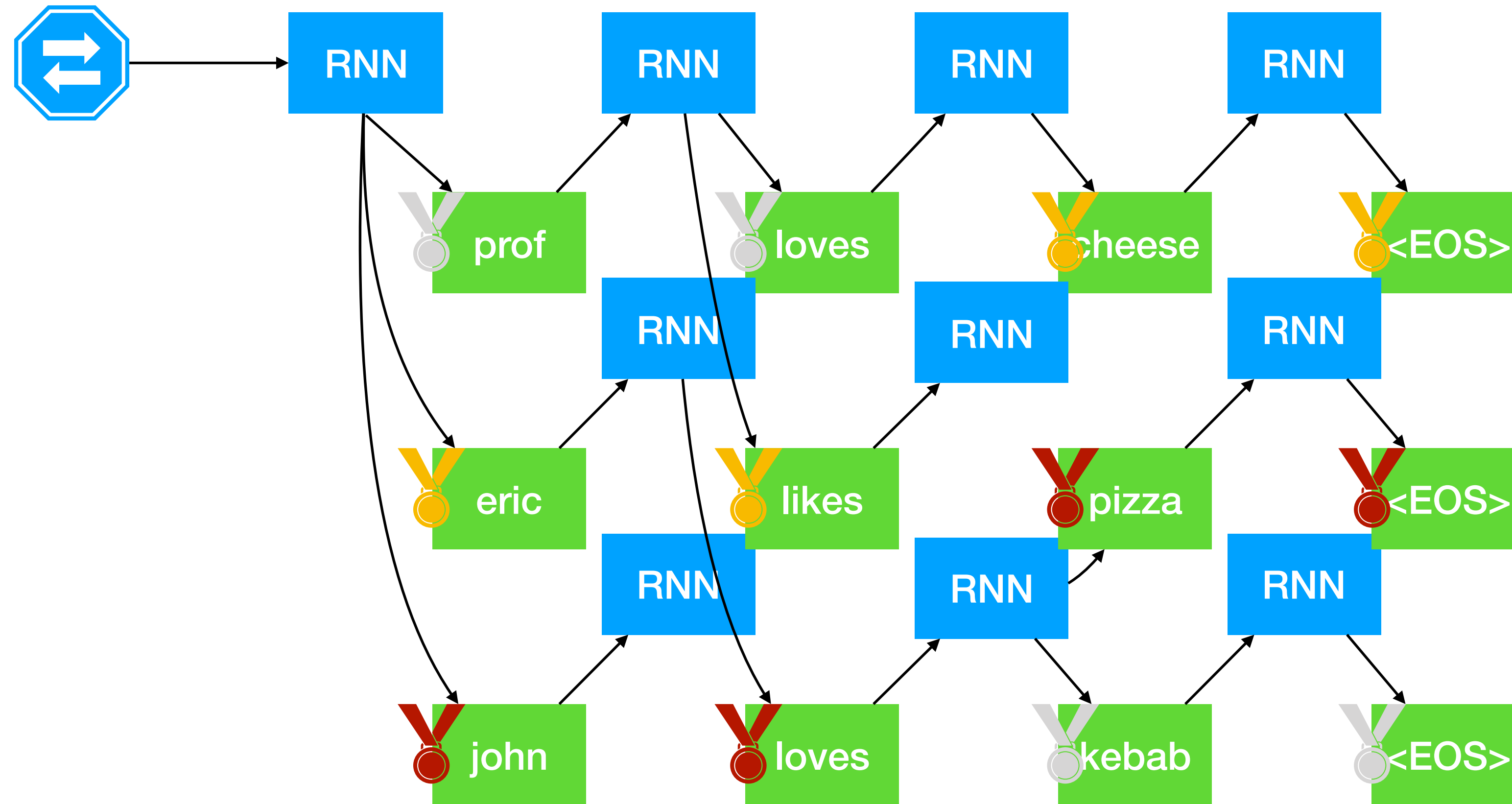
# BeamSearch (Size=3)



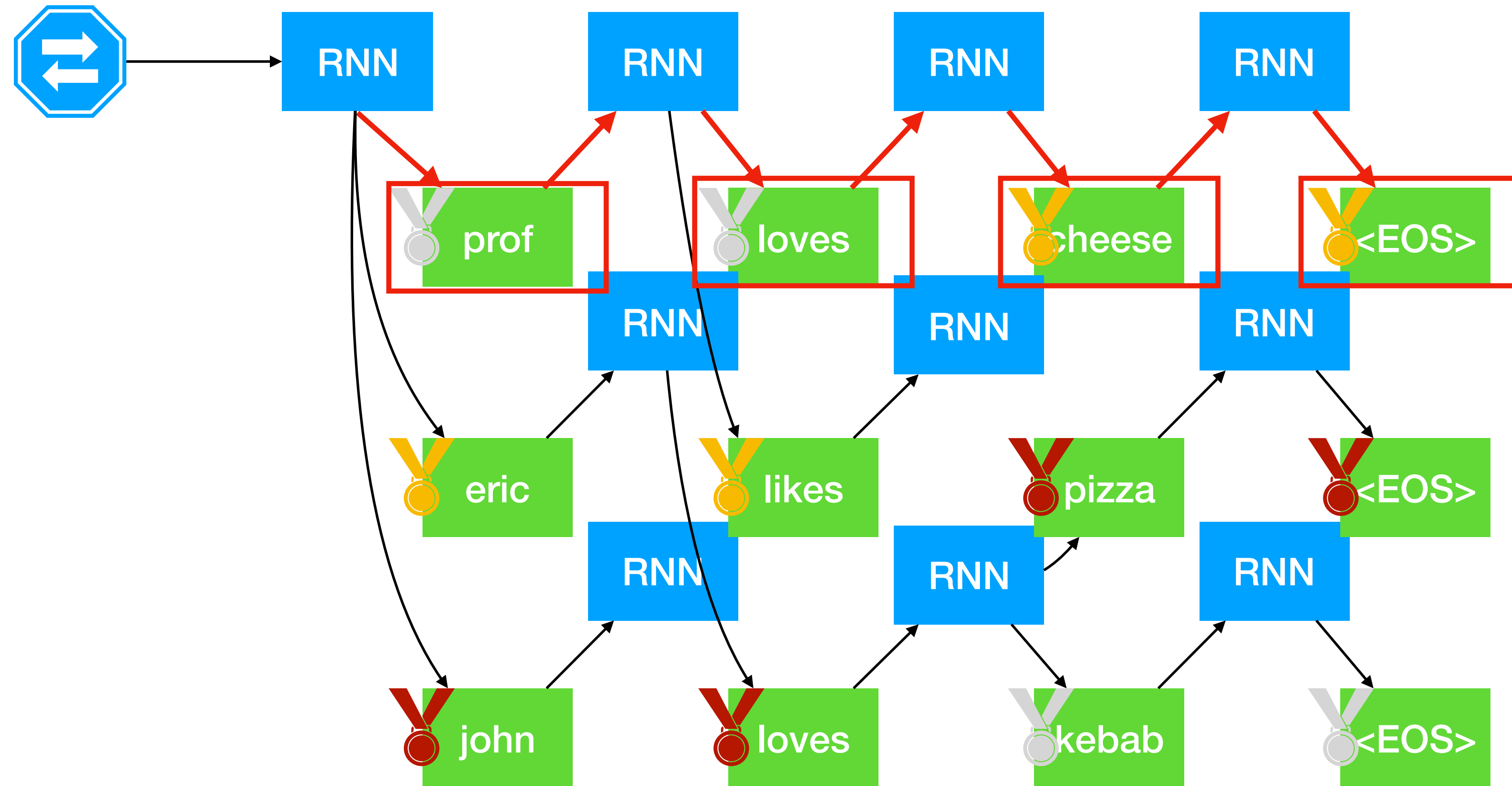
# BeamSearch (Size=3)



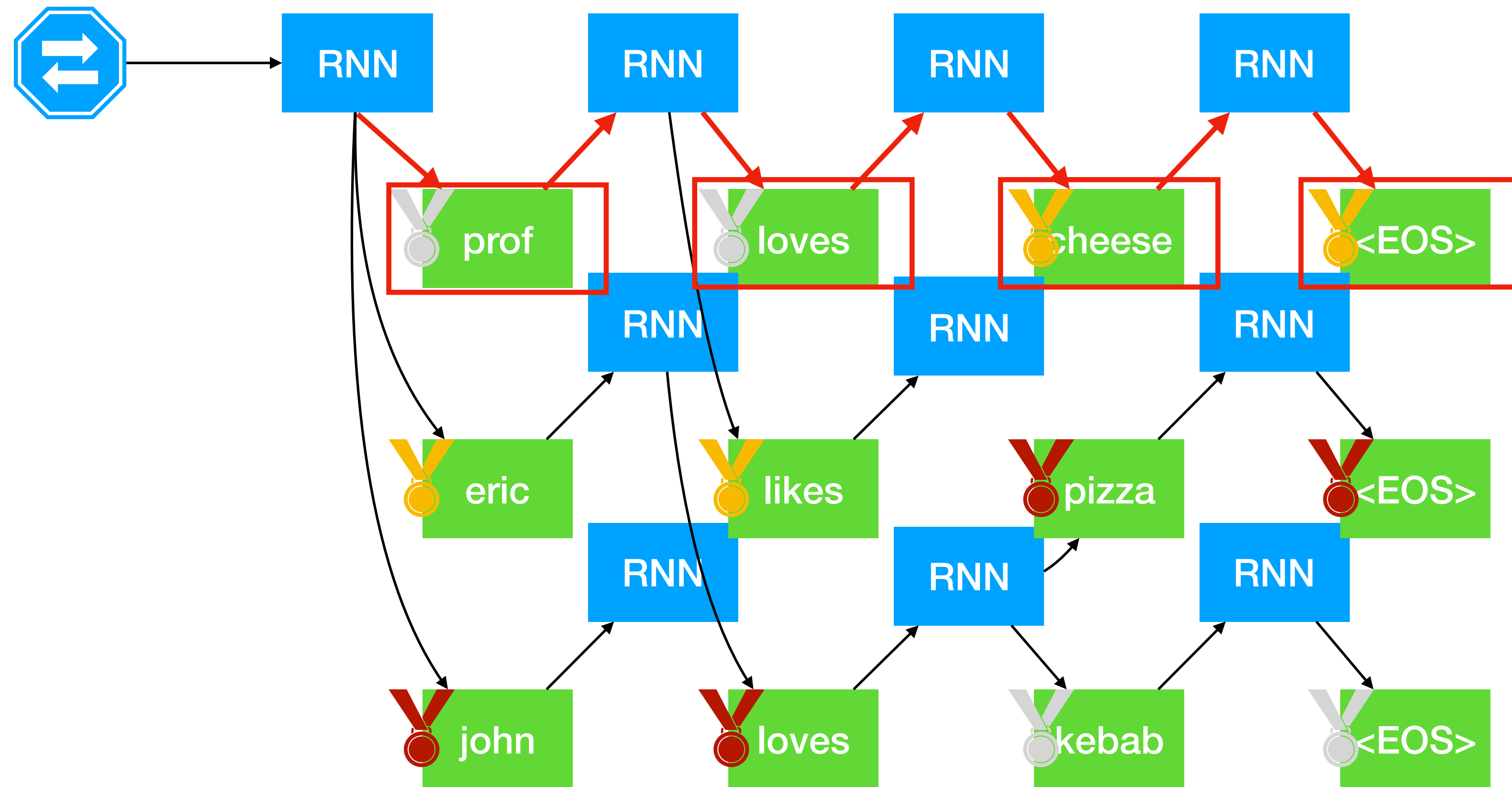
# BeamSearch (Size=3)



# BeamSearch (Size=3)

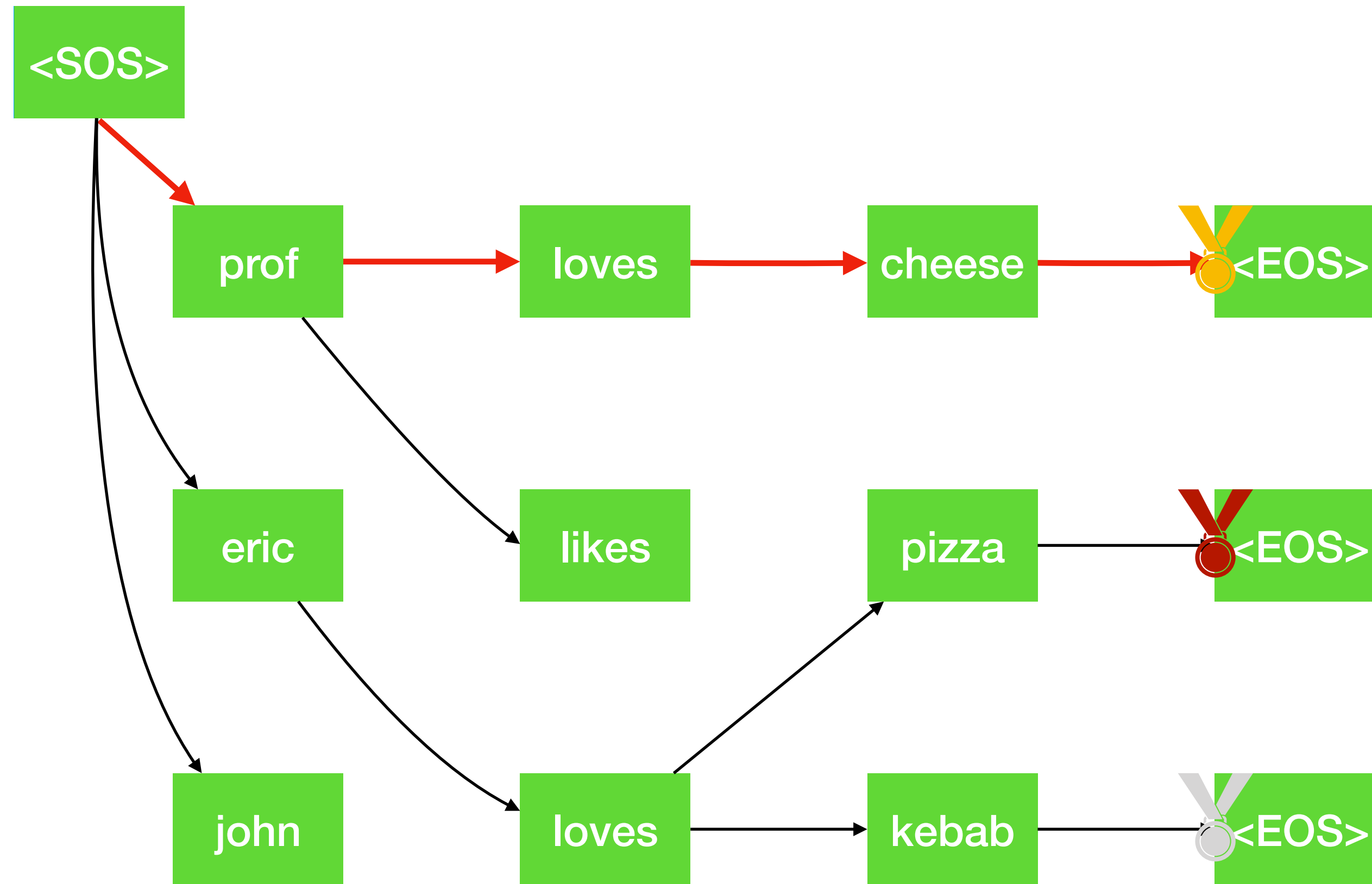


# BeamSearch (Size=3)



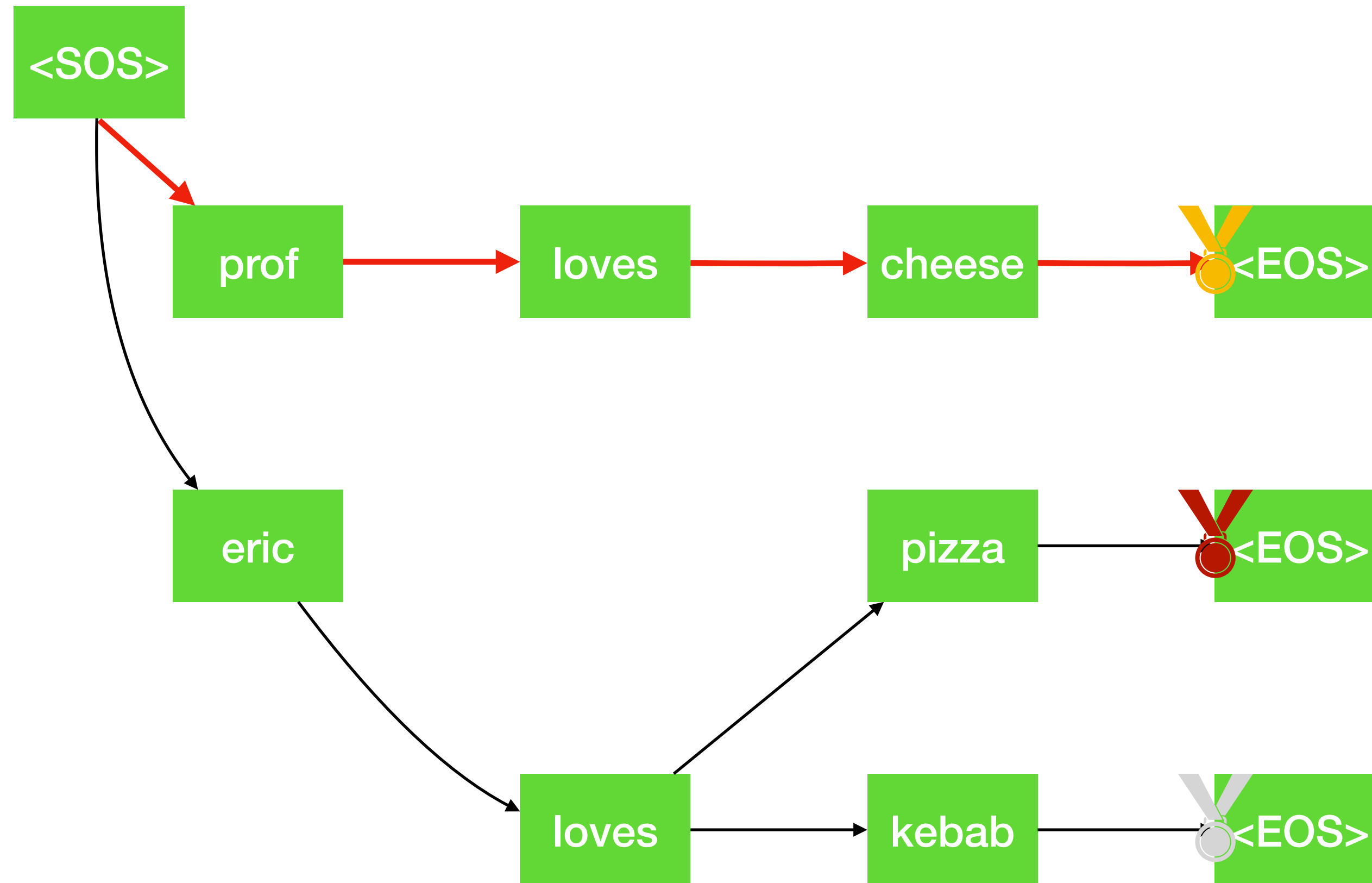
Output: "prof loves cheese"

# BeamSearch (Size=3)



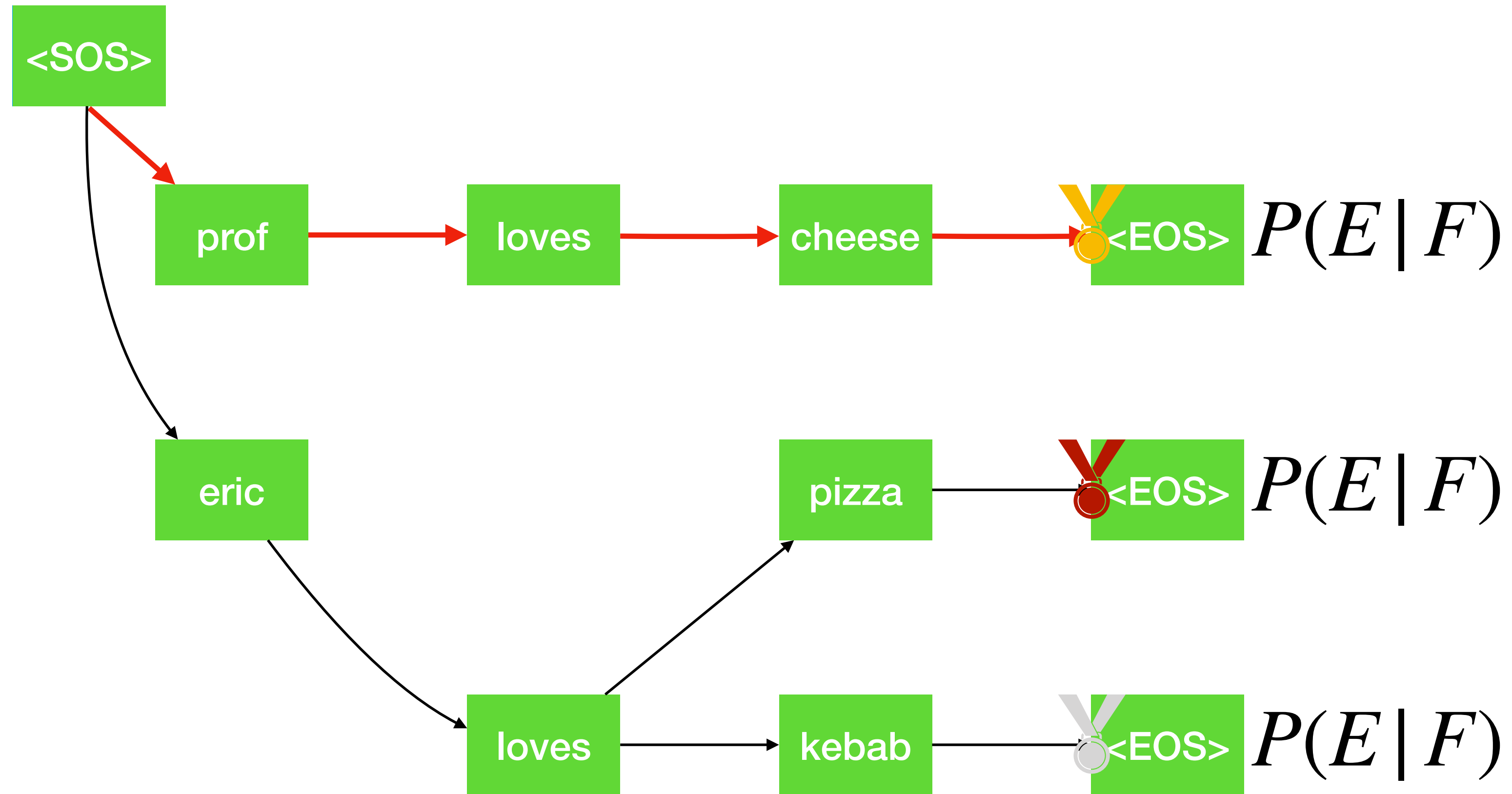
**Output: "prof loves cheese"**

# BeamSearch (Size=3)



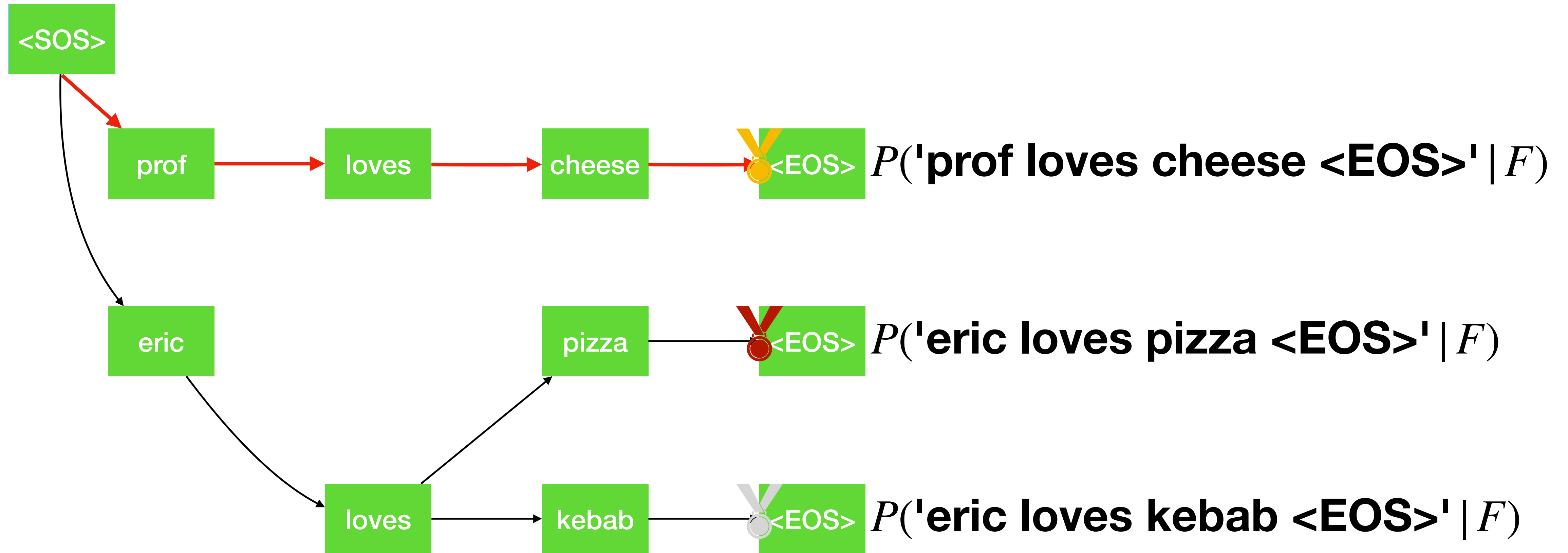
**Output: "prof loves cheese"**

# BeamSearch (Size=3)



**Output: "prof loves cheese"**

# BeamSearch (Size=3)



**Output: "prof loves cheese"**

# BeamSearch

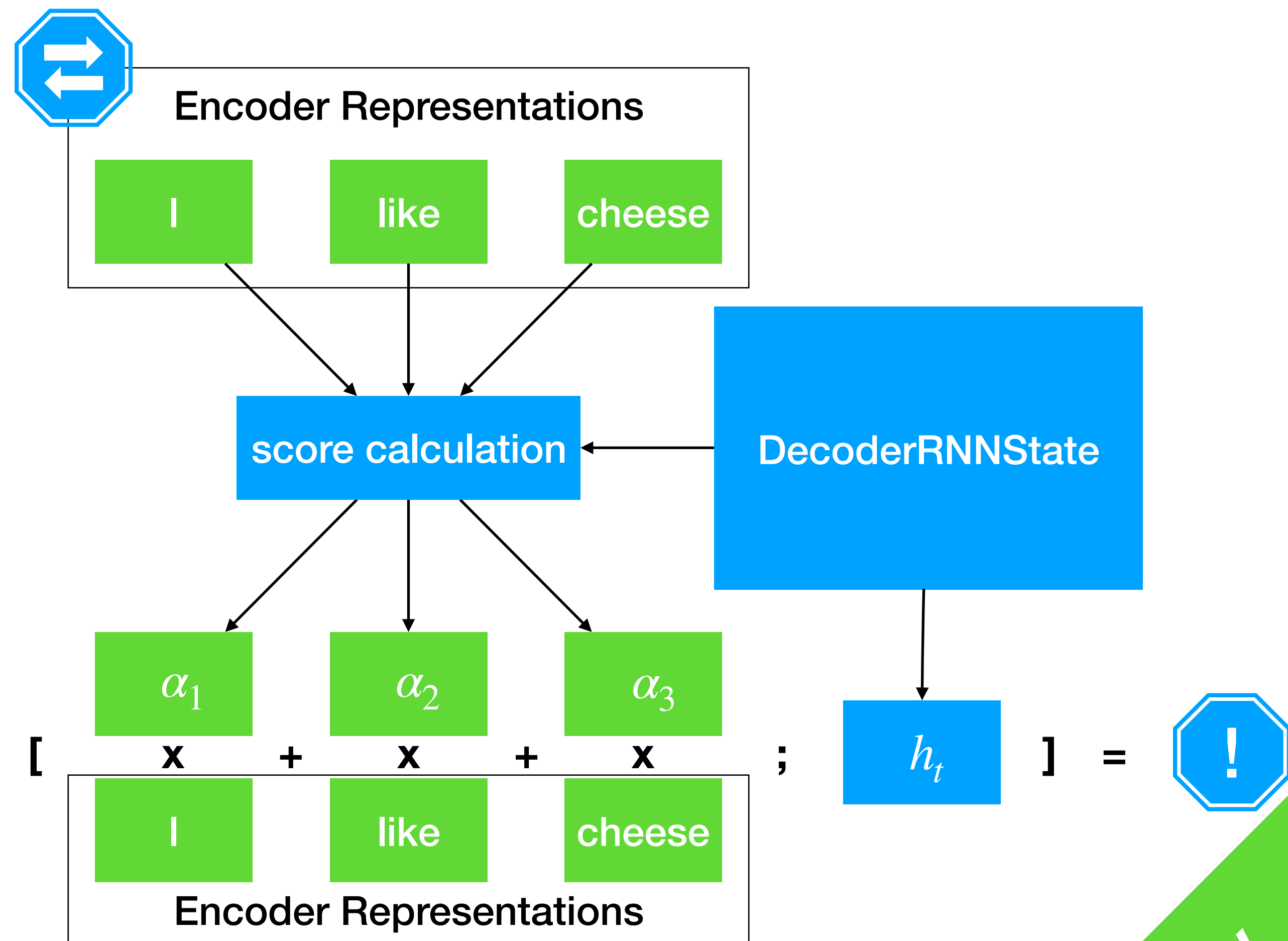
- Consider multiple hypotheses at each step  $t$ , formulate decoding as a search programme on a tree
- Consider multiple complete translation hypotheses, but reduces the total search space from  $|V_E|^m$  to

# BeamSearch

- Consider multiple hypotheses at each step  $t$ , formulate decoding as a search programme on a tree
- Consider multiple complete translation hypotheses, but reduces the total search space from  $|V_E|^m$  to  $\text{batchSize} \times m$

# Self-Attentive NMT

- Aggregating src information  $h_{1:|F|}^{enc}$ ,  
with  $h_t^{dec}$  as query



# Self-Attentive NMT

- Aggregating src information  $h_{1:|F|}^{enc}$ , with  $h_t^{dec}$  as query
- $\vec{h}^{enc}$ : aggregating information  $f_{<i}$ , with  $f_i$  as query
- $\overleftarrow{h}^{enc}$ : aggregating information  $f_{>i}$ , with  $f_i$  as query
- Can we replace the RNN with attention?

# Self-Attentive NMT

- Transformer<sup>1</sup>
  - Encoder-Decoder
  - No RNN -> all RNNs in Seq2Seq replaced by identical self-attention blocks
  - Multi-headed attention blocks\*
    - read the paper for more information

# Self-Attentive NMT

- Does everything Seq2Seq does
- BERT: Bidirectional Encoder Representation of Transformer

# Beyond NMT

- NMT weaknesses
  - Longer src sentences -> SMT does it better, augment attention?
  - Growing lexicon/Terminologies -> Pointer-based Dictionary Fusion
  - Sensitive domain -> offline computing, human-involved system to ensure accuracy
  - Mobile Platform deployment -> model compression

# Beyond NMT

- NMT weaknesses
  - Longer src sentences -> SMT does it better, augment attention?
  - Growing lexicon/Terminology **HYBRID** based Dictionary Fusion
  - Sensitive domain -> offline computing, human-involved system to ensure accuracy
  - Mobile Platform deployment -> model compression

# Beyond NMT

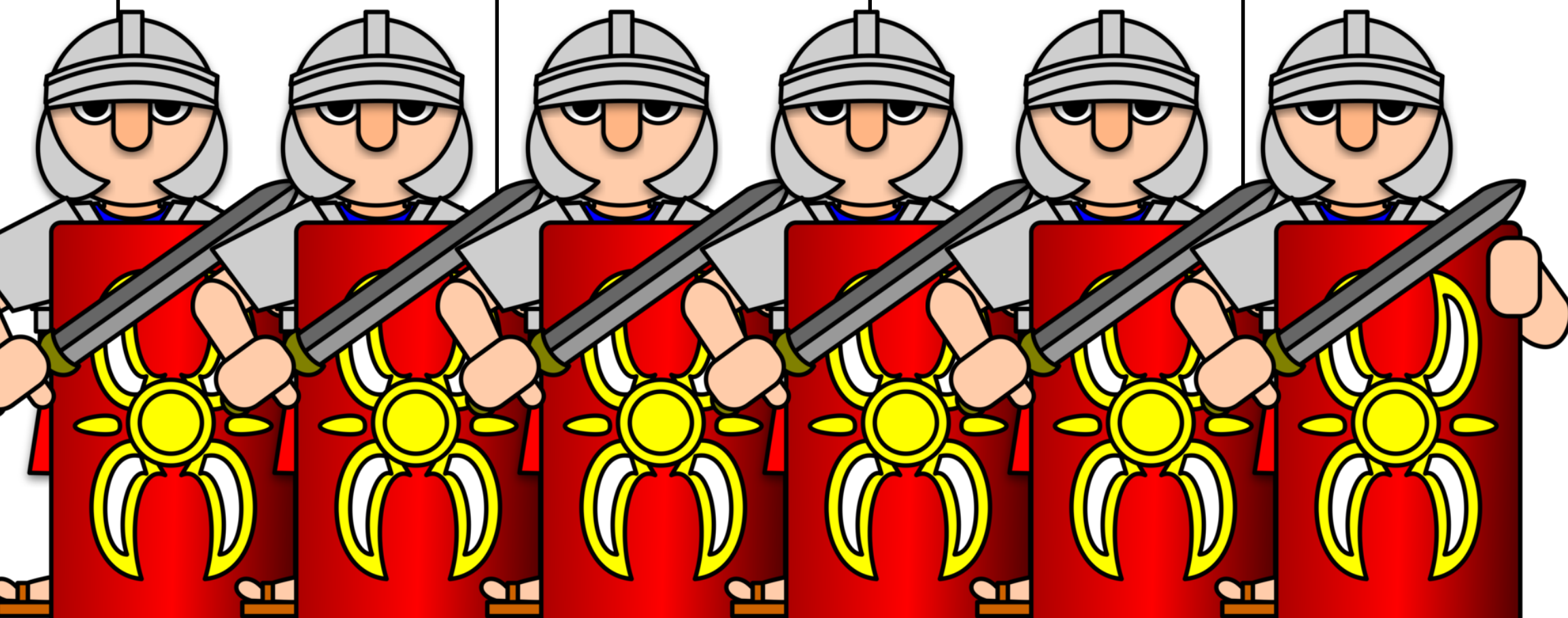
**HYBRID**

HYBRIDS  
ARE DUMB

PURISM  
FOREVER!

MAKE NMT  
GREAT AGAIN

ATTENTION IS  
ALL YOU  
NEED



Concept

**My work here is done, thank you.**