

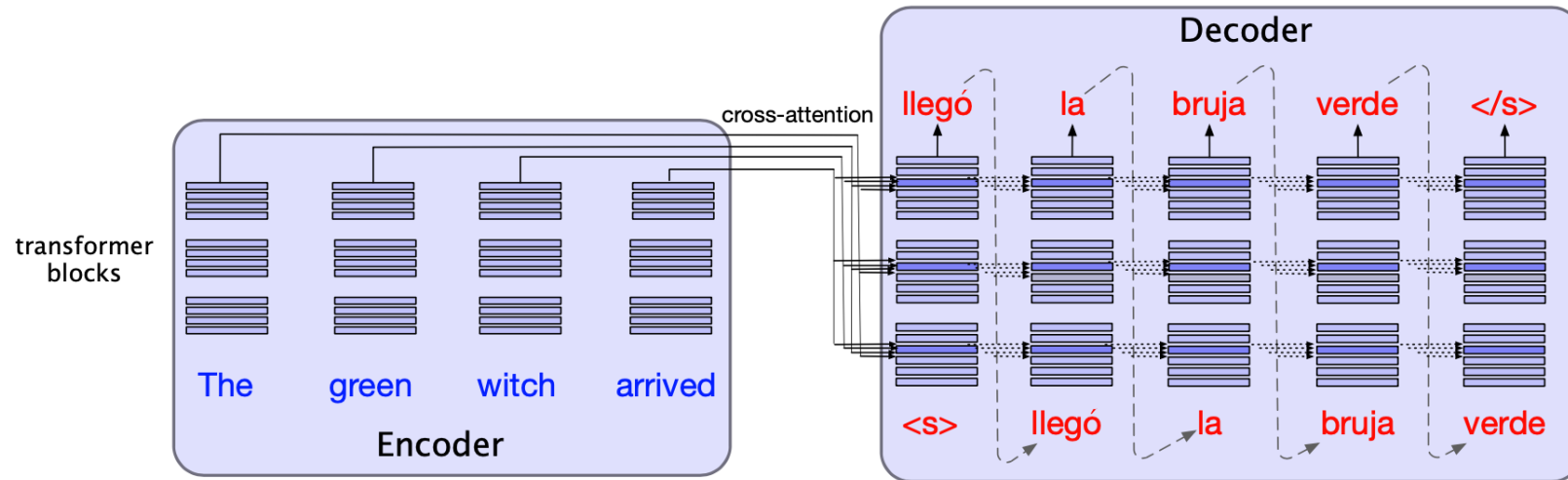
Machine  
Translation,  
RAG, and  
Evaluation

Machine Translation

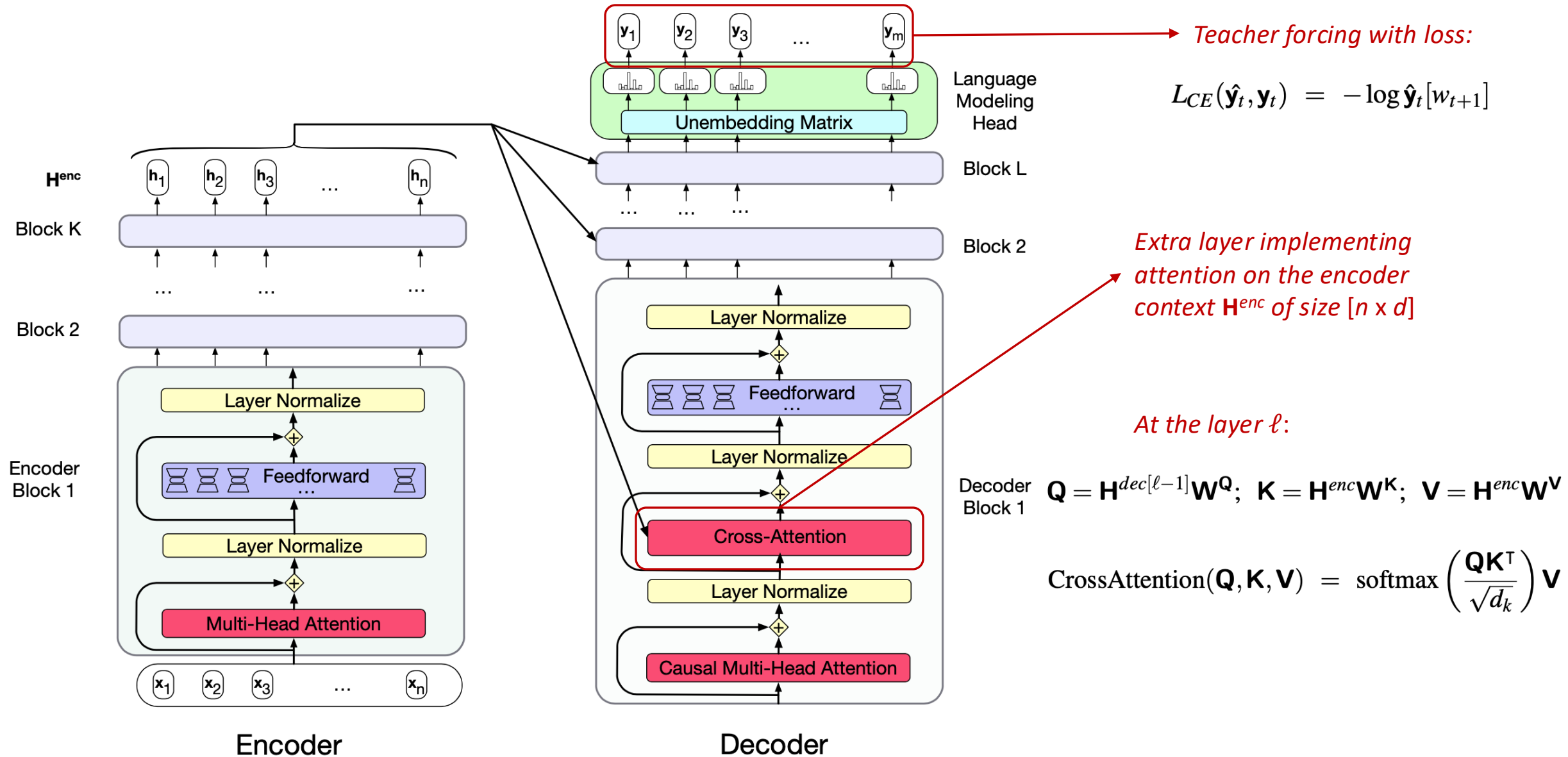
# Linguistic issues in Machine Translation

- Languages have some *universal* traits (there are words for referring to people, eating, being polite, and so on)
- Still they have many differences
  - Word Order Typology (SVO, SOV, VSO and use of prepositions or postpositions)
  - Lexical Divergences (one word in  $L_S \rightarrow$  different words in  $L_T$ , where more context is needed)
  - Morphological Typology (from isolated to fused morphemes)
  - Referential Density (things the language tends to omit)
    - *I talked to Juliet yesterday. She was very sad.  $\rightarrow$  Ieri ho parlato con Juliet. Era molto triste.*

# MT with the encoder-decoder model



# MT with the encoder-decoder model

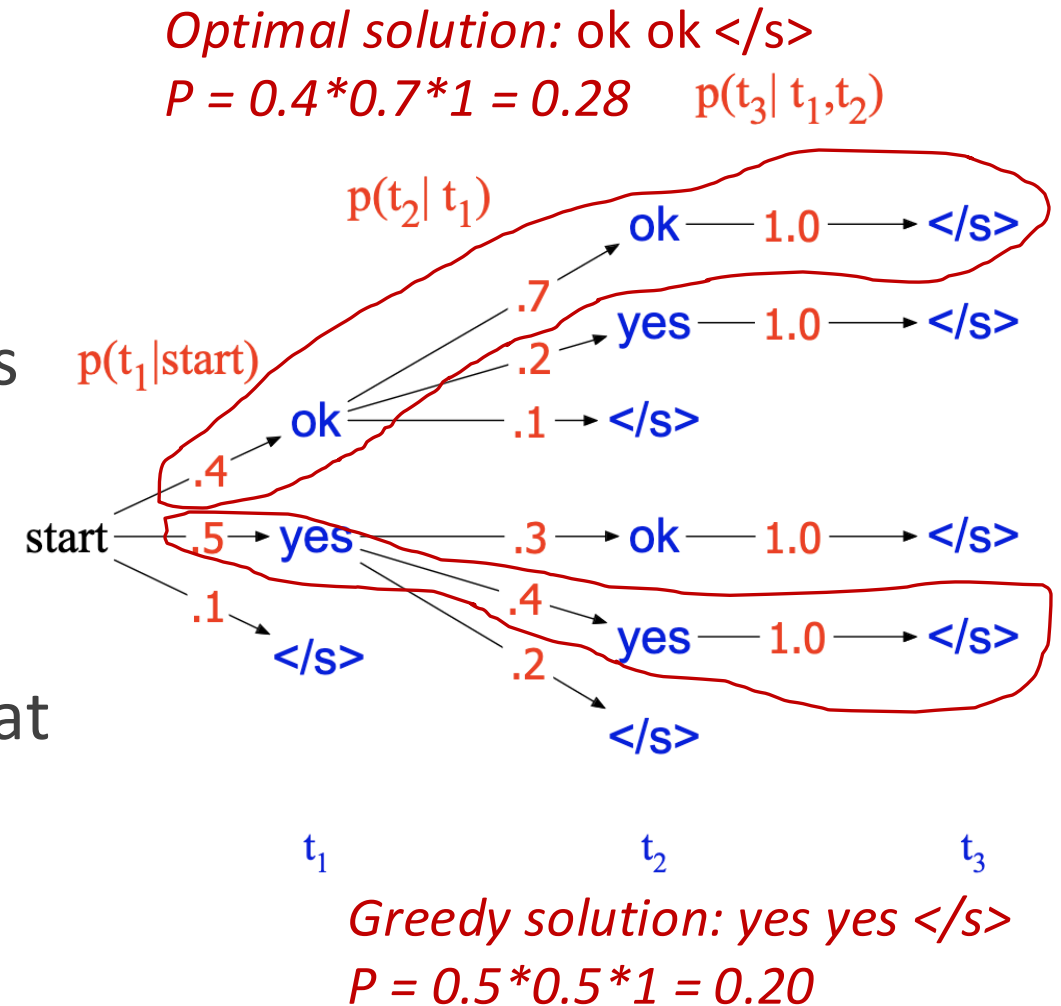


# MT with the encoder-decoder model

Decoding with greedy search is *locally optimal* in the space of solutions

We need to implement search strategies that are optimal in a global way

The only algorithm that guarantees optimal solution is exhaustive search that is infeasible ( $V^T$  possibilities)



# Beam Search

Decoding is achieved keeping the best  $k$  tokens at each step

1. At each step  $t_i$  in the path till the END token:
  - a) For each of the  $k$  best hypotheses so far (i.e. the  $k$  best generated sequences)
    1. Score the entire vocabulary starting from the current hypothesis by  $P(y_i | x, y_{<i})$
  - b) Prune the  $k * V$  scores maintaining the top  $k$  new hypotheses
2. Remove completed paths and set  $k = k - 1$
3. Restart from 1 till  $k == 0$

**function** BEAMDECODE( $c, beam\_width$ ) **returns** best paths

```
 $y_0, h_0 \leftarrow 0$ 
 $path \leftarrow ()$ 
 $complete\_paths \leftarrow ()$ 
 $state \leftarrow (c, y_0, h_0, path)$  ;initial state
 $frontier \leftarrow \langle state \rangle$  ;initial frontier

while  $frontier$  contains incomplete paths and  $beamwidth > 0$ 
     $extended\_frontier \leftarrow \langle \rangle$ 
    for each  $state \in frontier$  do
         $y \leftarrow \text{DECODE}(state)$ 
        for each word  $i \in \text{Vocabulary}$  do
             $successor \leftarrow \text{NEWSTATE}(state, i, y_i)$ 
             $extended\_frontier \leftarrow \text{ADDTOBEAM}(successor, extended\_frontier,$ 
                 $beam\_width)$ 

        for each  $state$  in  $extended\_frontier$  do
            if  $state$  is complete do
                 $complete\_paths \leftarrow \text{APPEND}(complete\_paths, state)$ 
                 $extended\_frontier \leftarrow \text{REMOVE}(extended\_frontier, state)$ 
                 $beam\_width \leftarrow beam\_width - 1$ 
             $frontier \leftarrow extended\_frontier$ 

    return  $completed\_paths$ 
```

**function** NEWSTATE( $state, word, word\_prob$ ) **returns** new state

**function** ADDTOBEAM( $state, frontier, width$ ) **returns** updated frontier

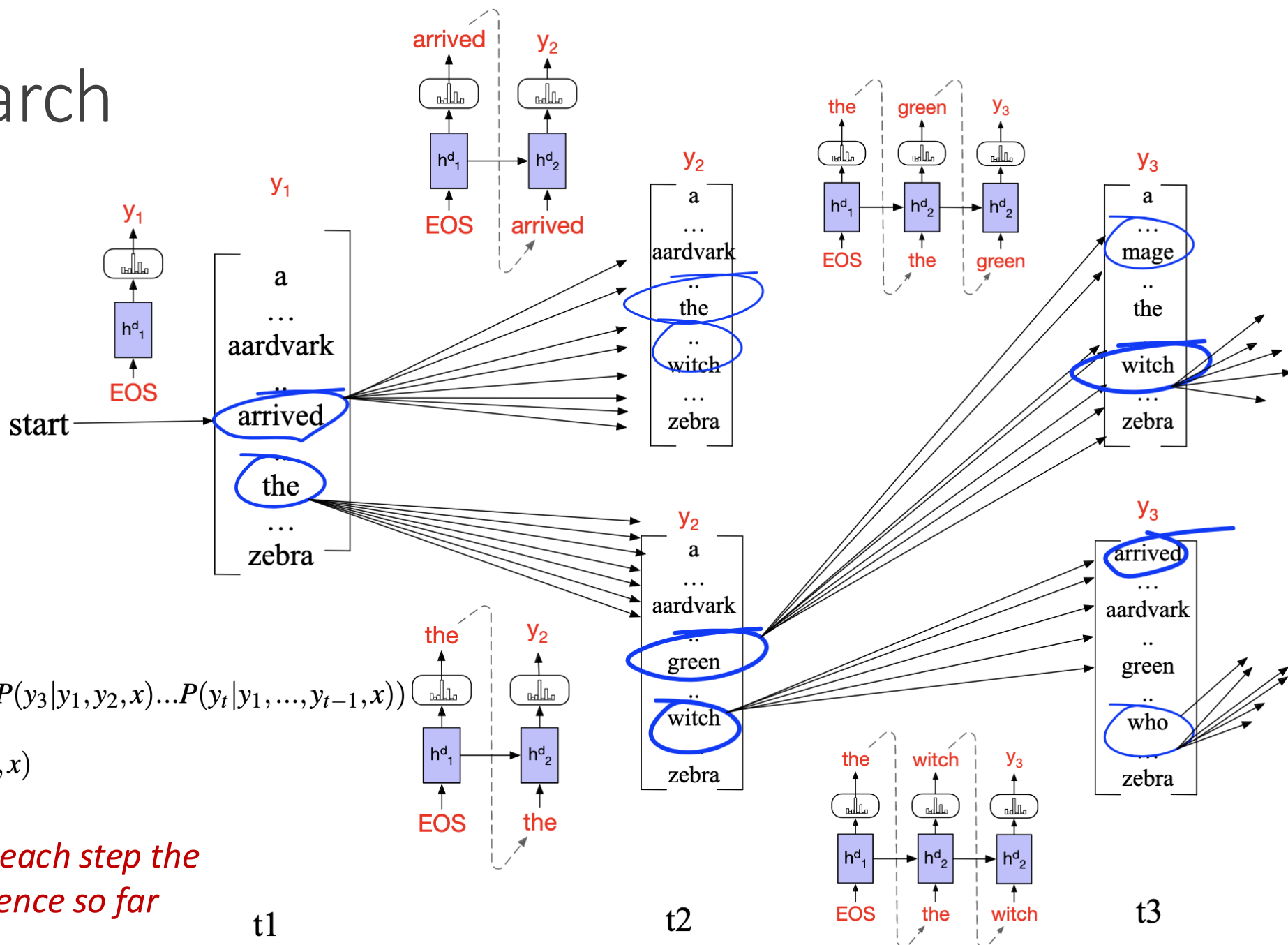
```
if  $\text{LENGTH}(frontier) < width$  then
     $frontier \leftarrow \text{INSERT}(state, frontier)$ 
else if  $\text{SCORE}(state) > \text{SCORE}(\text{WORSTOF}(frontier))$ 
     $frontier \leftarrow \text{REMOVE}(\text{WORSTOF}(frontier))$ 
     $frontier \leftarrow \text{INSERT}(state, frontier)$ 
return  $frontier$ 
```

# Beam Search

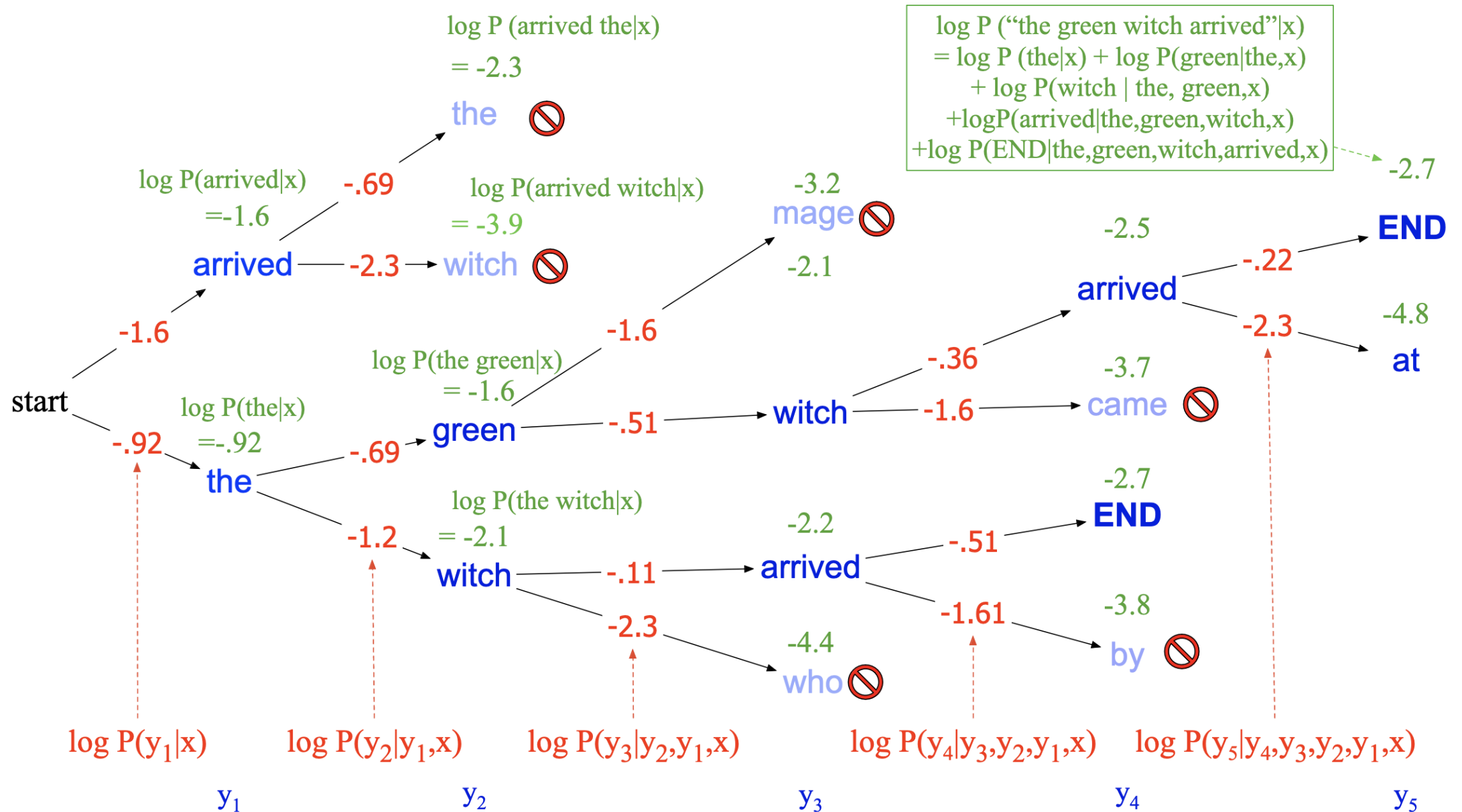
$k = 2$

$$\begin{aligned}
 \text{score}(y) &= \log P(y|x) \\
 &= \log (P(y_1|x)P(y_2|y_1,x)P(y_3|y_1,y_2,x)\dots P(y_t|y_1,\dots,y_{t-1},x)) \\
 &= \sum_{i=1}^t \log P(y_i|y_1,\dots,y_{i-1},x)
 \end{aligned}$$

*We simply need to add at each step the logprob of the prefix sequence so far*



# Beam Search





# Beam Search

The score has to be normalized for sequences with different length

- The additive nature of the score penalizes long sequences

$$score(y) = -\log P(y|x) = \frac{1}{T} \sum_{i=1}^t -\log P(y_i|y_1, \dots, y_{i-1}, x)$$

# Minimum Bayes Risk (MBR) decoding

We use explicitly an evaluation metric  $util$  to measure the goodness of fit of the translation  $c$  with some set  $\mathcal{Y}$  of candidate translation:

$$\hat{y} = \operatorname{argmax}_{y \in \mathcal{Y}} \sum_{c \in \mathcal{Y}} util(y, c)$$

- Metrics should be used to evaluate against perfect human translation
- We do not know it, and approximate the space of all possible translation with  $\mathcal{Y}$

Machine  
Translation,  
RAG, and  
Evaluation

Retrieval Augmented  
Generation

# Question Answering

Question Answering (QA) systems are designed to fill human information needs. Since a lot of information is present in text form they perform similarly to search engines.

Question answering systems often focus on *factoid questions*, questions of fact or reasoning that can be answered with simple facts expressed in short or medium-length texts:

- Where is the Louvre Museum located?
- Where does the energy in a nuclear explosion come from?
- How to get a script I in latex?

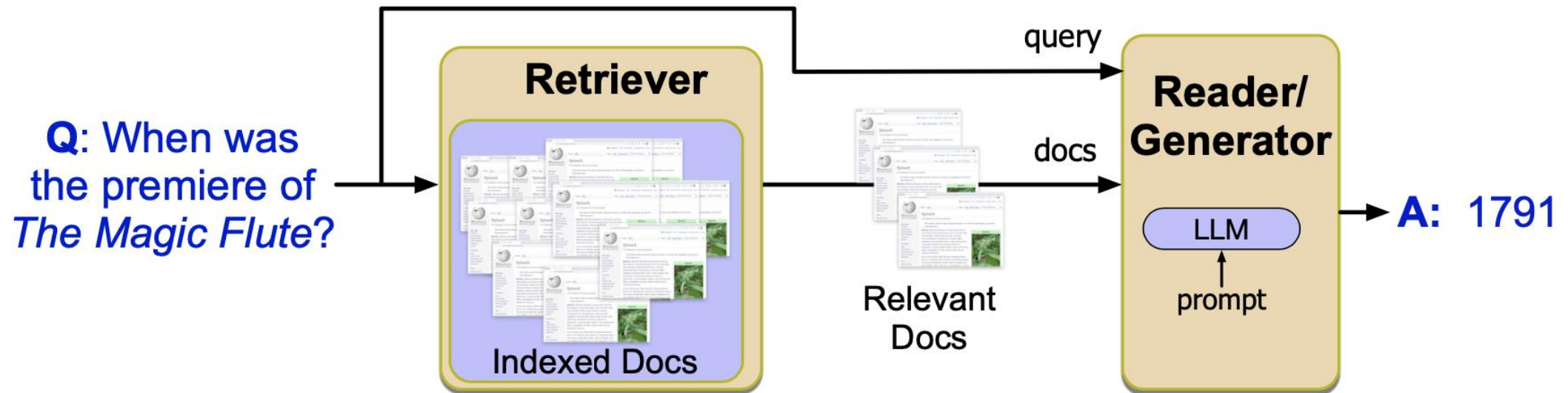
# Question Answering

LLMs can be prompted for QA:

*Q: Where is the Louvre Museum located? A:*

But LLMs can hallucinate and they can be not aligned

# Retrieval Augmented Generation (RAG)



How to make the Retriever?

Vector database, Knowledge base, hybrid ...

Machine  
Translation,  
RAG, and  
Evaluation

Evaluation

# Evaluation

Evaluation is needed in many generative tasks to assess the goodness of the sentence provided by the model.

A corpus (at least one) of reference sentences can be either at disposal or not.

Different approaches in the two cases.



# Evaluation by human raters

In general, no reference corpus is available in this case.

Two properties to be assessed:

- *Adequacy* or *Faithfulness* (the meaning of the text)
- *Fluency* (the grammatical structure)

Human raters (in general crowdworkers) use either scales to score both properties or they can be provided with multiple translations that have to be ranked.

# Automatic evaluation

Automatic evaluation compares a reference sentence  $r$  and a hypothesis  $h$  using either character overlap or embeddings similarity

These metrics are used in all the generative tasks where a reference sentence is available

# character F-score (chrF)

Averages the percentage of common 1-grams, 2-grams, ...,  $k$ -grams in  $r$  and in  $h$ .

**chrP** percentage of character 1-grams, 2-grams, ...,  $k$ -grams in the hypothesis that occur in the reference, averaged.

**chrR** percentage of character 1-grams, 2-grams,...,  $k$ -grams in the reference that occur in the hypothesis, averaged.

$$\text{chrF}\beta = (1 + \beta^2) \frac{\text{chrP} \cdot \text{chrR}}{\beta^2 \cdot \text{chrP} + \text{chrR}}$$

# BiLingual Evaluation Understudy (BLEU)

Given a corpus of candidate and reference sentences  $H = \{h_i\}$  and  $R = \{r_j\}$  BLEU is a function of the n-gram word precision  $p_n(H, R)$  over all the reference sentences combined with a brevity penalty BP computed over the corpus as a whole.

$$p_n(H, R) = \frac{\sum_i \sum_{s \in D_n^{(h_i)}} \min(C(s, h_i), \max_{r_j} C(s, r_j))}{\sum_{s \in D_n^{(h_i)}} C(s, h_i)}$$

$$\text{BLEU}(H, R) = \text{BP} \cdot \exp \left( \frac{1}{N} \sum_{n=1}^N \log p_n(H, R) \right)$$

# Recall-Oriented Understudy for Gisting Evaluation (ROUGE)

A set of metrics:

- ROUGE-N: N-gram overlap (typically ROUGE-1 and ROUGE-2) computed in terms of
- Precision:  $\# \text{ of common N-grams} / \# \text{ of N-grams in } h$
- Recall:  $\# \text{ of common N-grams} / \# \text{ of N-grams in } r$
- $F1 = 2 P * R / (P + R)$

# Recall-Oriented Understudy for Gisting Evaluation (ROUGE)

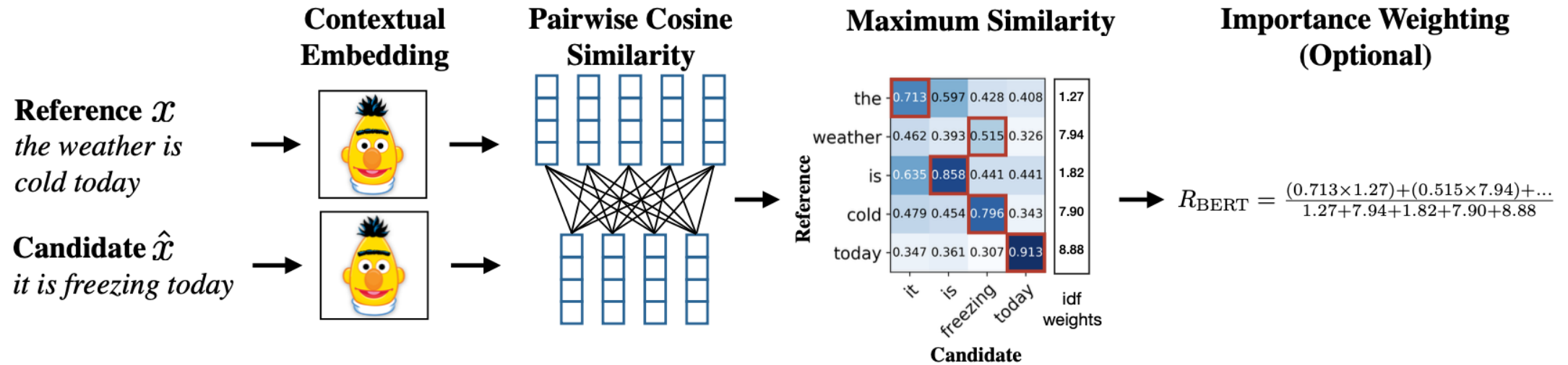
- ROUGE-L: Longest Common Subsequence ( $LCS$ ) between  $h$  and  $r$
- Precision:  $LCS(h, r) / \# \text{ of 1-grams in } h$
- Recall:  $LCS(h, r) / \# \text{ of 1-grams in } r$
- $F1 = 2 P * R / (P + R)$

# BERT score

N-grams overlap metrics do not perform well in case of synonyms or paraphrases.

BERT score intuition: *compare token embeddings instead!!*

# BERT score



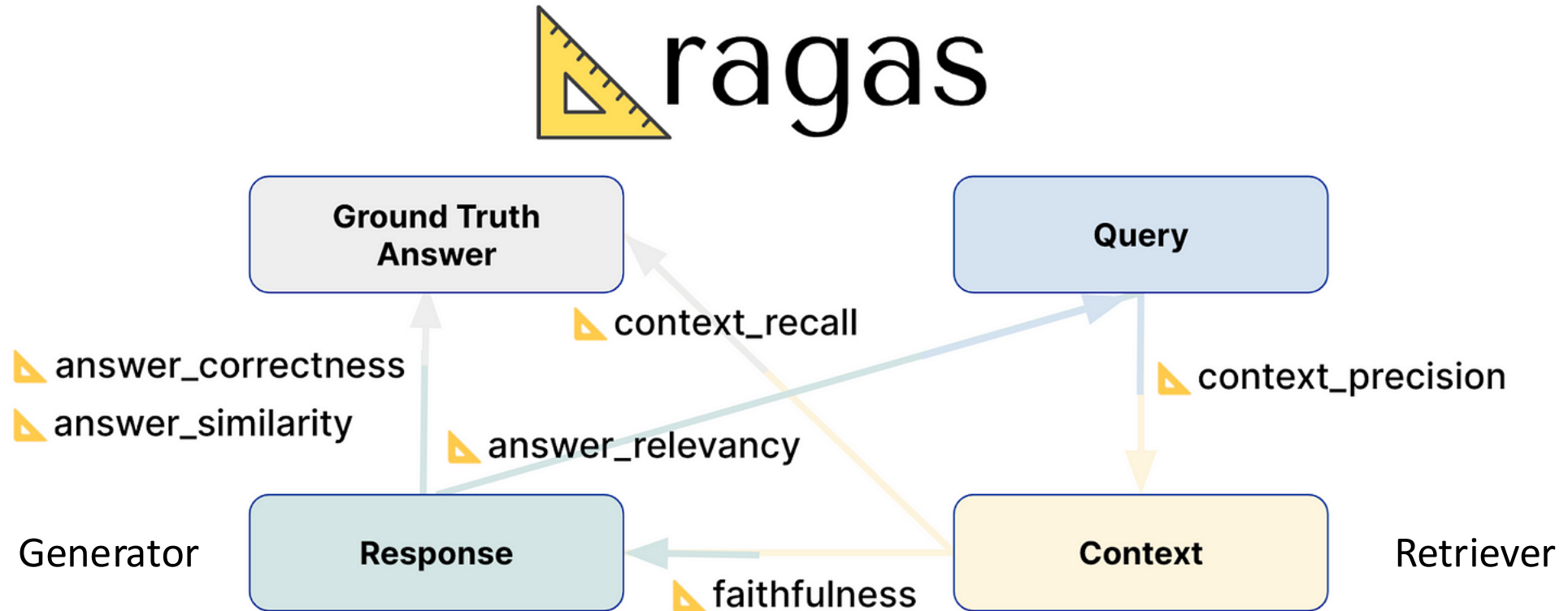
$$R_{\text{BERT}} = \frac{1}{|\mathcal{X}|} \sum_{x_i \in \mathcal{X}} \max_{\tilde{x}_j \in \tilde{\mathcal{X}}} x_i \cdot \tilde{x}_j$$

$$P_{\text{BERT}} = \frac{1}{|\tilde{\mathcal{X}}|} \sum_{\tilde{x}_j \in \tilde{\mathcal{X}}} \max_{x_i \in \mathcal{X}} x_i \cdot \tilde{x}_j$$



# RAGAS

Framework designed to evaluate answers in RAG systems



# RAGAS

## Faithfulness

assesses the factual accuracy of the generated answers by checking if the statements made in the answers are supported by the provided context. It is computed by analyzing the validity of each statement derived from the generated answer against the context.

$$\text{Faithfulness} = \# \text{ of valid statements} / \# \text{ of statements}$$

# RAGAS

## Answer Relevancy

measures how relevant and directly related the generated answer is to the posed question. It evaluates the similarity between probable questions that could elicit the same answer and the actual question asked.

Answer relevancy = mean cosine similarity between the user input a number of answers reverse engineered from the response

# RAGAS

## Context Recall

measures how many relevant documents or pieces of information (the Ground Truth – GT) were successfully retrieved. It focuses on not missing important results.

Context recall =  $\frac{|\text{GT claims attributed to the context}|}{|\text{GT claims}|}$

# RAGAS

## Context Precision

Context Precision is a metric that measures the proportion of relevant chunks in the retrieved contexts. It is calculated as the mean of the precision@ $k$  for each chunk in the context.

Precision@ $k$  is the ratio of the number of relevant chunks at rank  $k$  to the total number of chunks at rank  $k$  (that is in the top  $k$  results)

# RAGAS

## Context Precision

$$\text{Context Precision@K} = \frac{\sum_{k=1}^K ( \text{Precision@k} \times v_k )}{\text{Total number of relevant items in the top } K \text{ results}}$$
$$\text{Precision@k} = \frac{\text{true positives@k}}{(\text{true positives@k} + \text{false positives@k})}$$

$v_k$  is the relevance indicator at rank  $k$