

# Using adaptive knowledge graphs in neural dialogue generation

Patrik Purgai

# Chatbots

## Domain specific

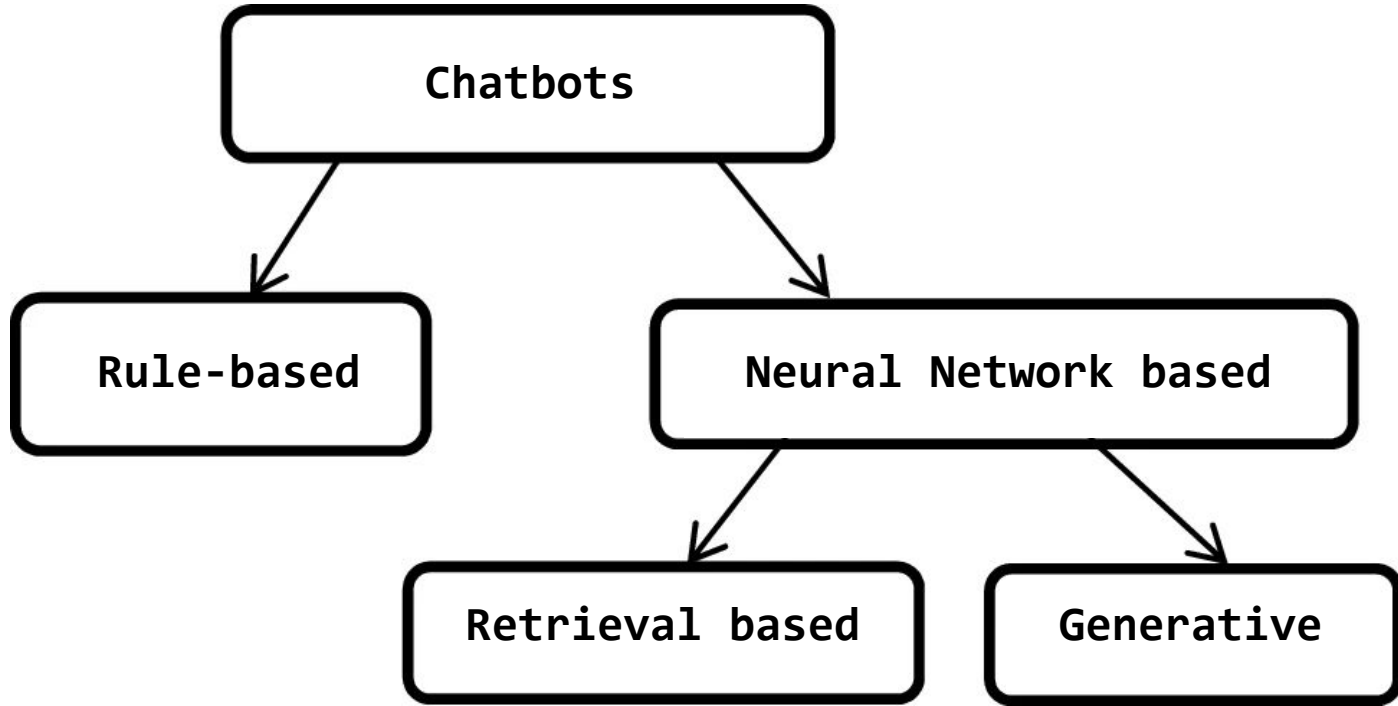
- Rasa
- Luis
- Dialogflow



## Open domain

- Blender
- Cleverbot





# Rule-based

- **Pattern matching - ELIZA**
- **AIML (Artificial Intelligence Markup Language) - Mitsuku**
- **Chatscript - SUZETTE**

# AIML

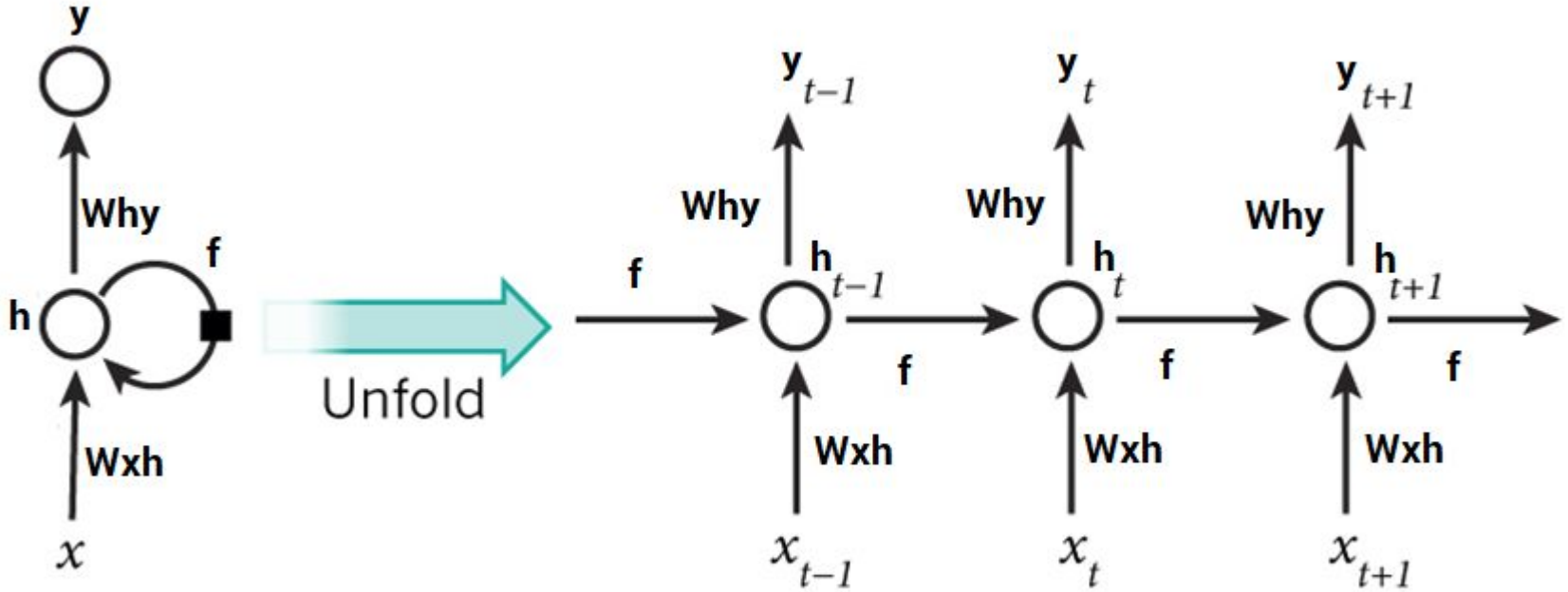
```
<category>  
  <pattern>MY NAME IS *</pattern>  
  <template>Hi there, my name is Andrew.</template>  
</category>
```

```
<category>  
  <pattern>GUESS A NUMBER</pattern>  
  <template>  
    <random>  
      <li>1</li>  
      <li>6</li>  
      <li>227</li>  
    </random>  
  </template>  
</category>
```

# Neural Network based

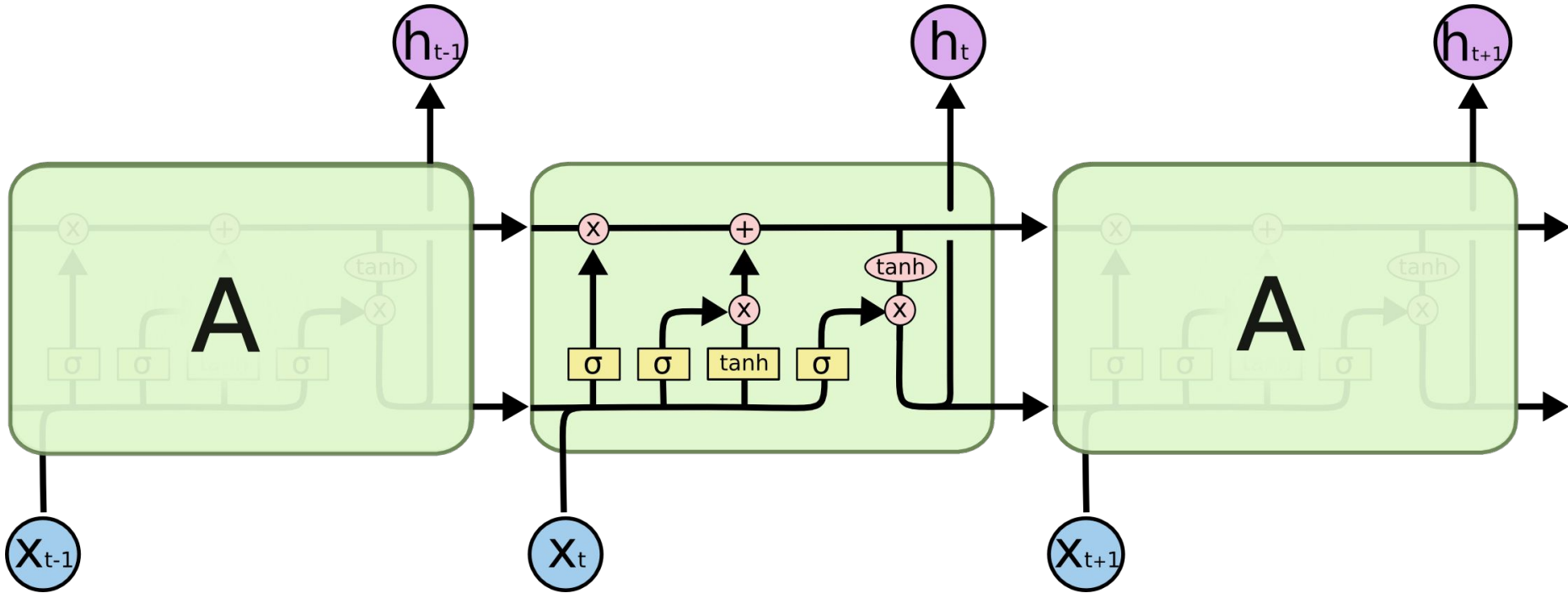
- RNN (LSTM)
- Transformer
- BERT
- GPT-2

# RNN



# LSTM

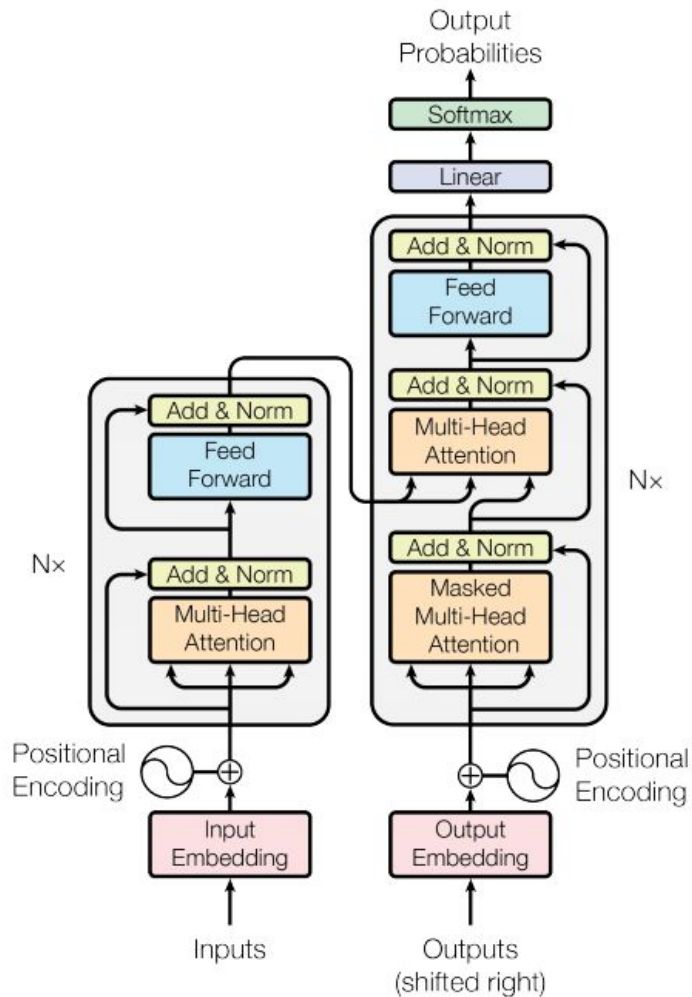
Hochreiter & Schmidhuber (1997)

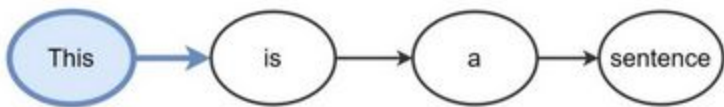




# Transformer

Vaswani et al. (2017)

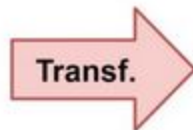
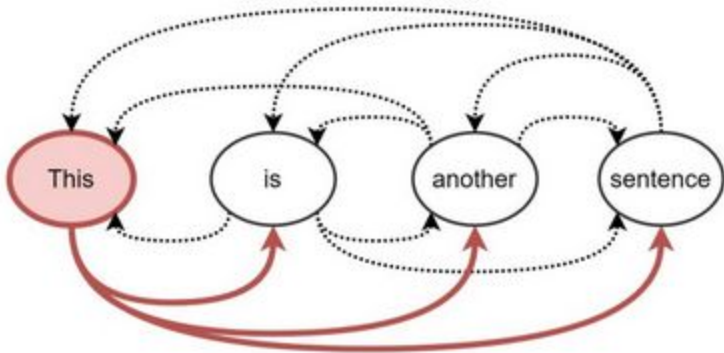




Translation?

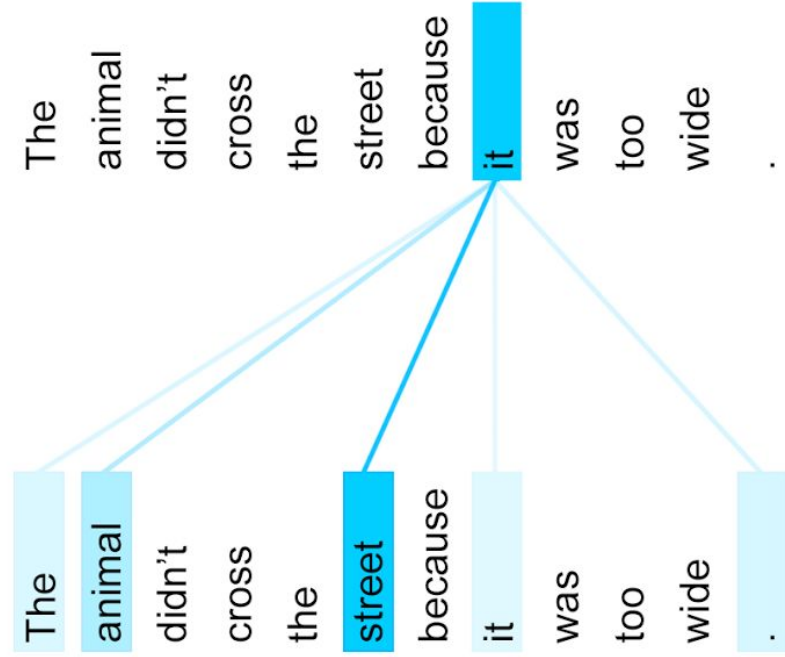
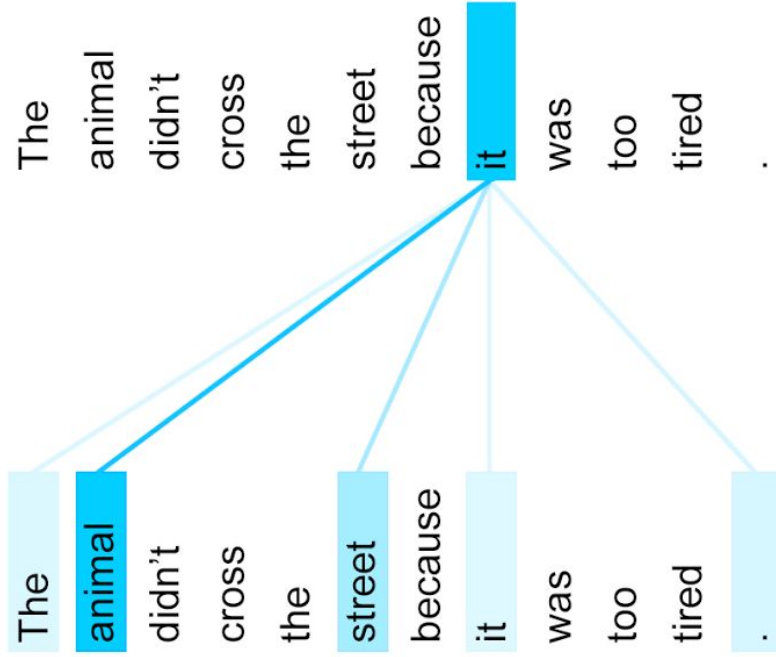
Sentiment?

Next word?

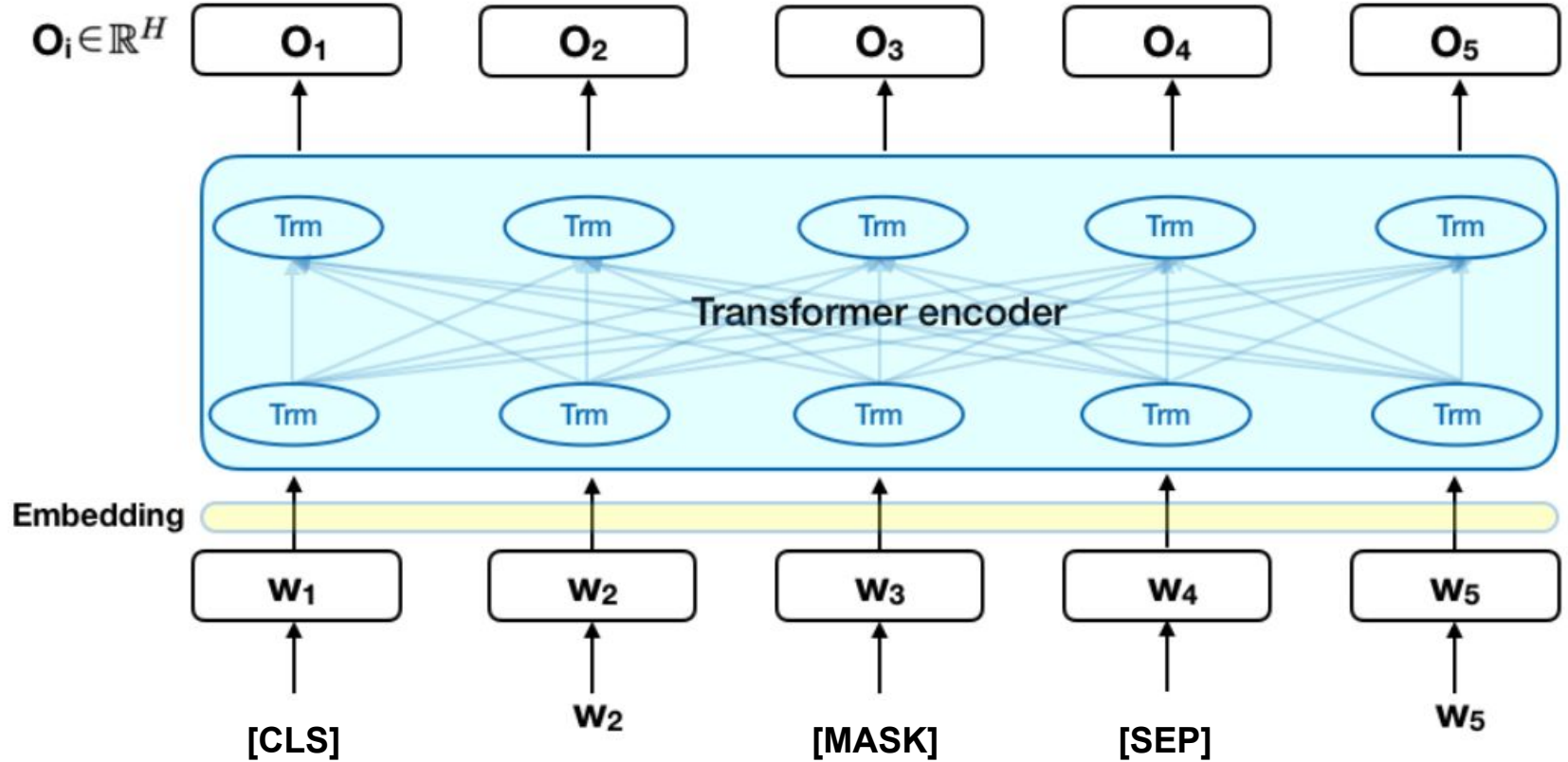


Part-of-speech tags?

# Attention Mechanism

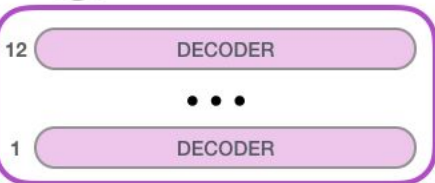


# BERT



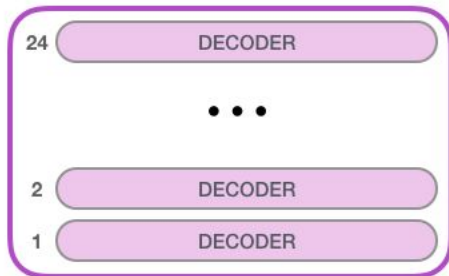
# Radford et al. (2019)

 GPT-2  
SMALL



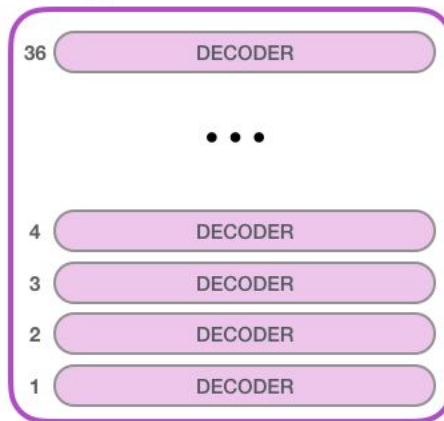
Model Dimensionality: 768

 GPT-2  
MEDIUM



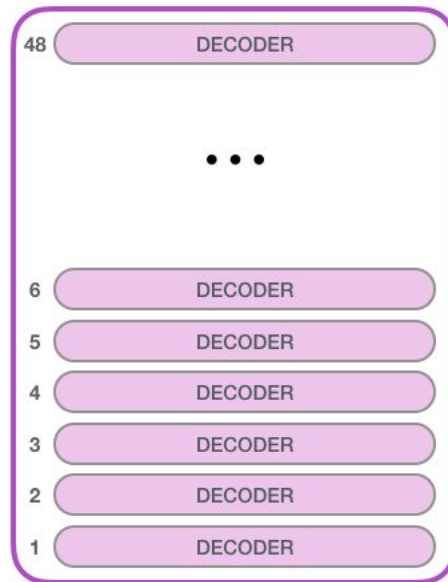
Model Dimensionality: 1024

 GPT-2  
LARGE



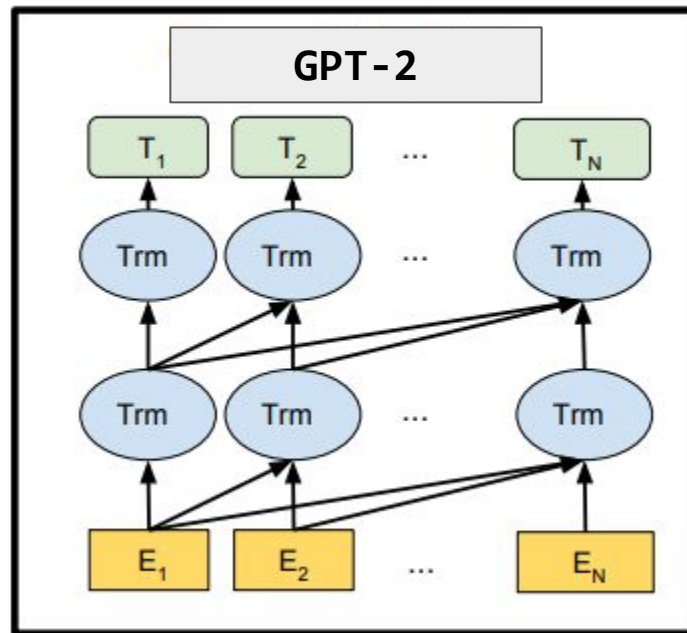
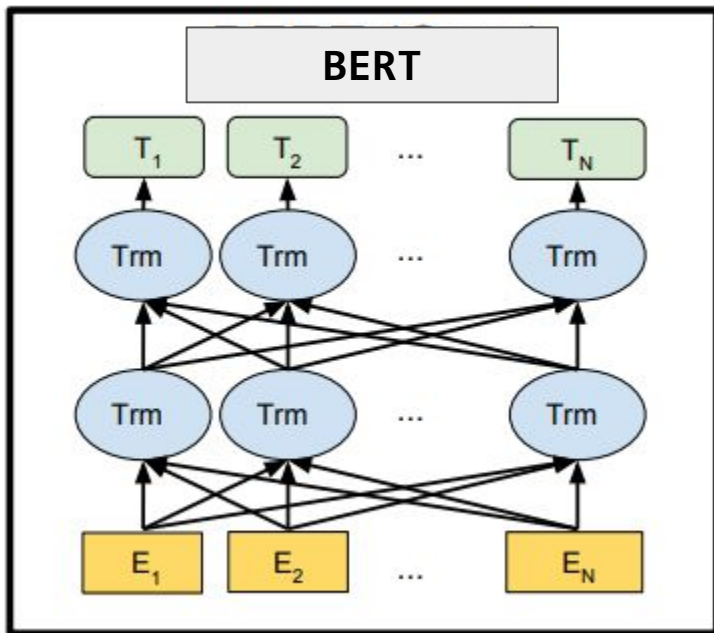
Model Dimensionality: 1280

 GPT-2  
EXTRA  
LARGE



Model Dimensionality: 1600

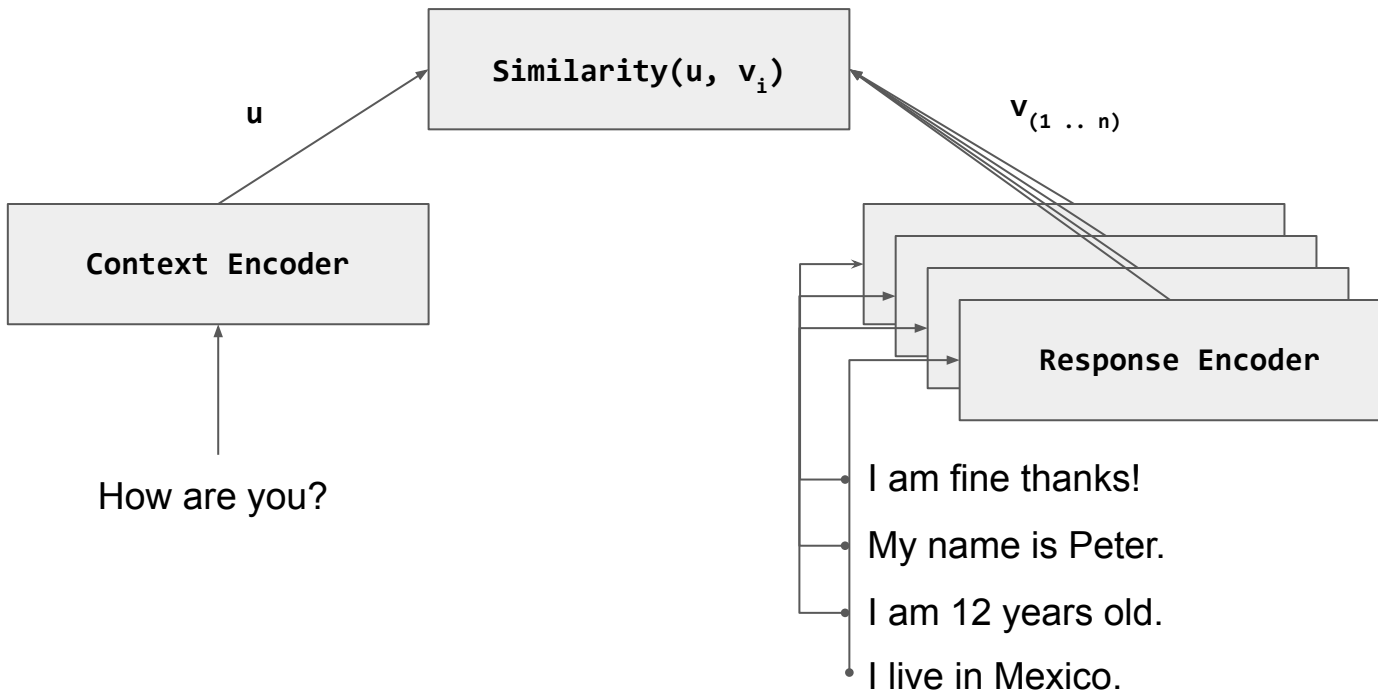
# BERT vs GPT-2



# Retrieval-based

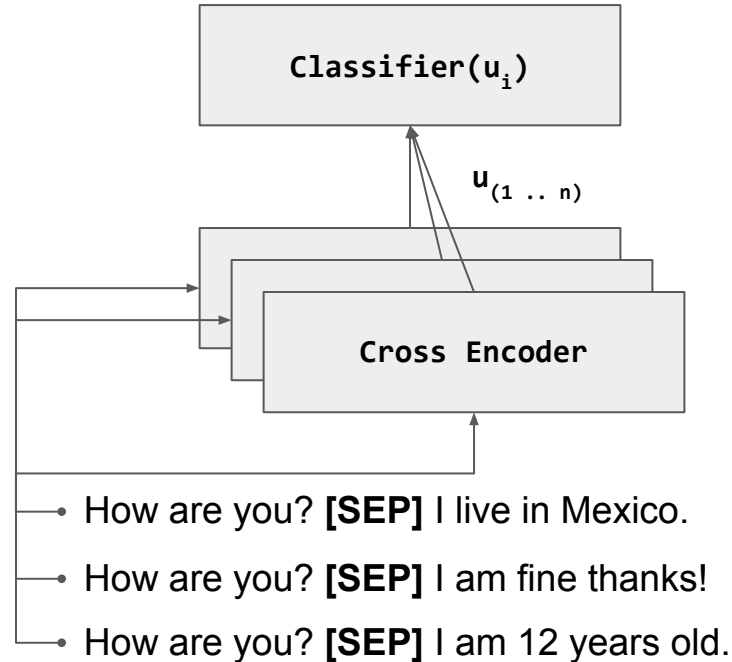
- **Speaker-Aware BERT** - Gu et al. (2020)
- **Poly-encoders** - Humeau et al. (2019)

# Bi-Encoder

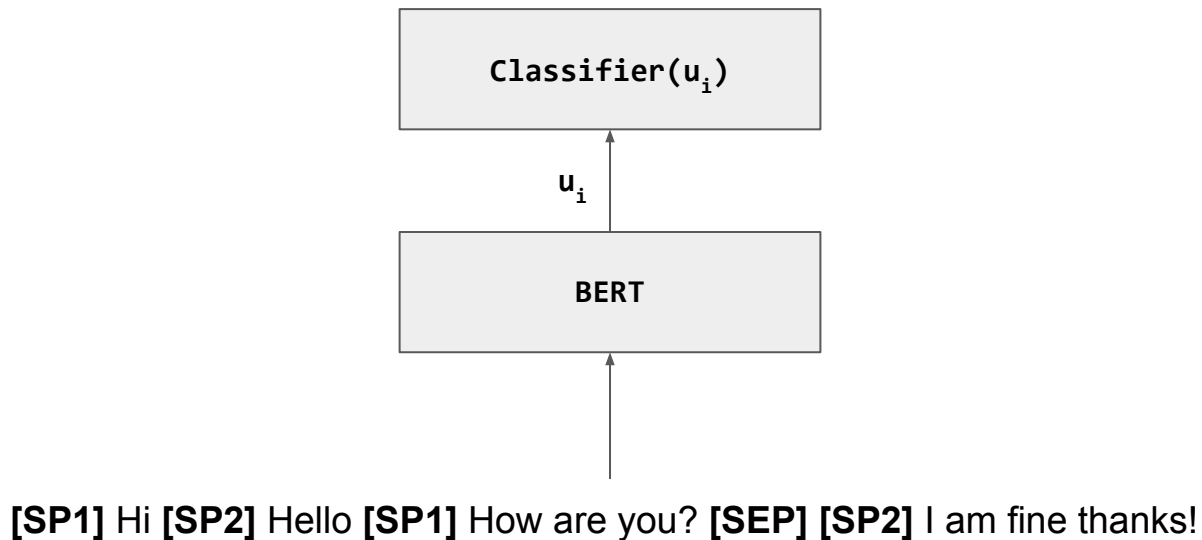




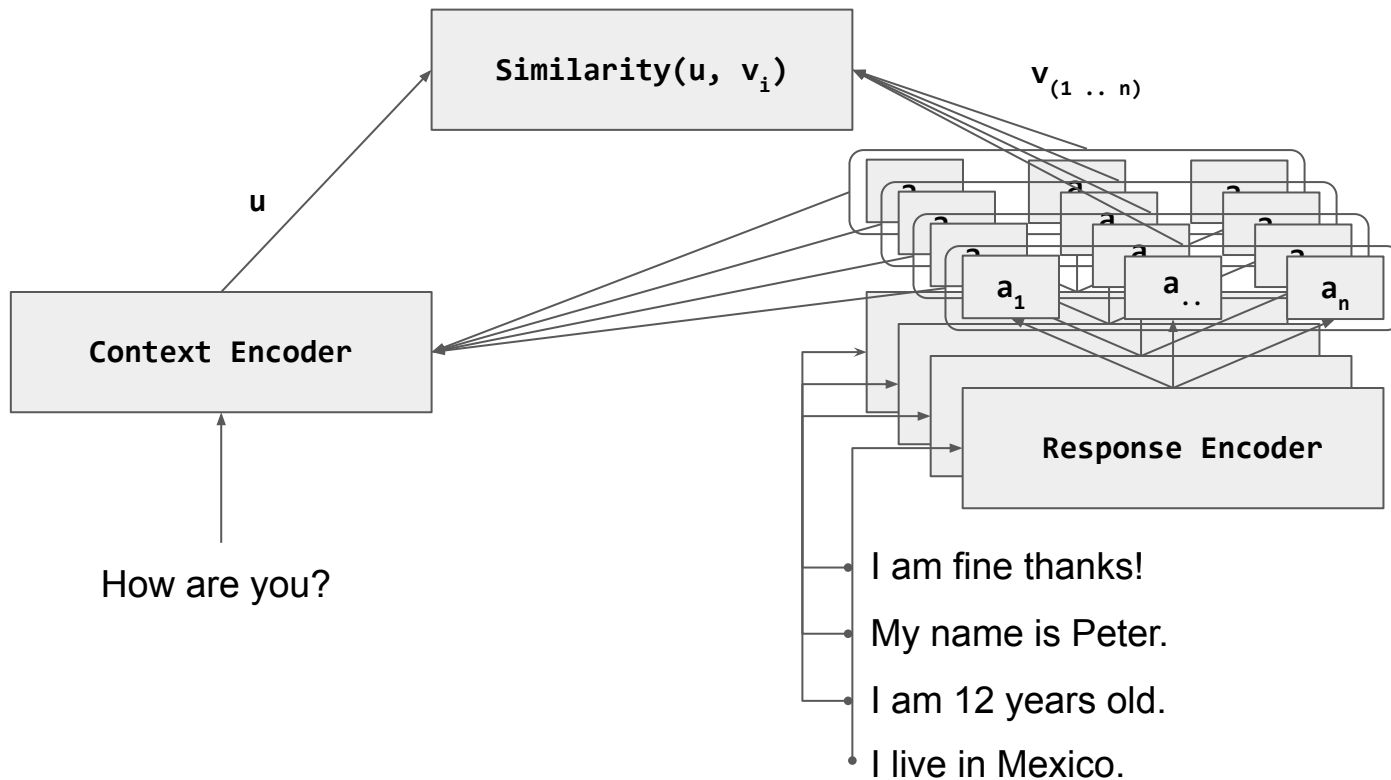
# Cross - Encoder



# Speaker Aware BERT



# Poly-Encoder



# Neural Network based - Generative

- A Neural Conversational Model - Vinyals et al. (2015)
- DialogPT - Zhang et al. (2019)
- Towards a Human-like Open-Domain Chatbot (Meena) - Adiwardana et al. (2020)
- Recipes for building an open-domain chatbot - Roller et al. (2020)

Can you please come **here** ?



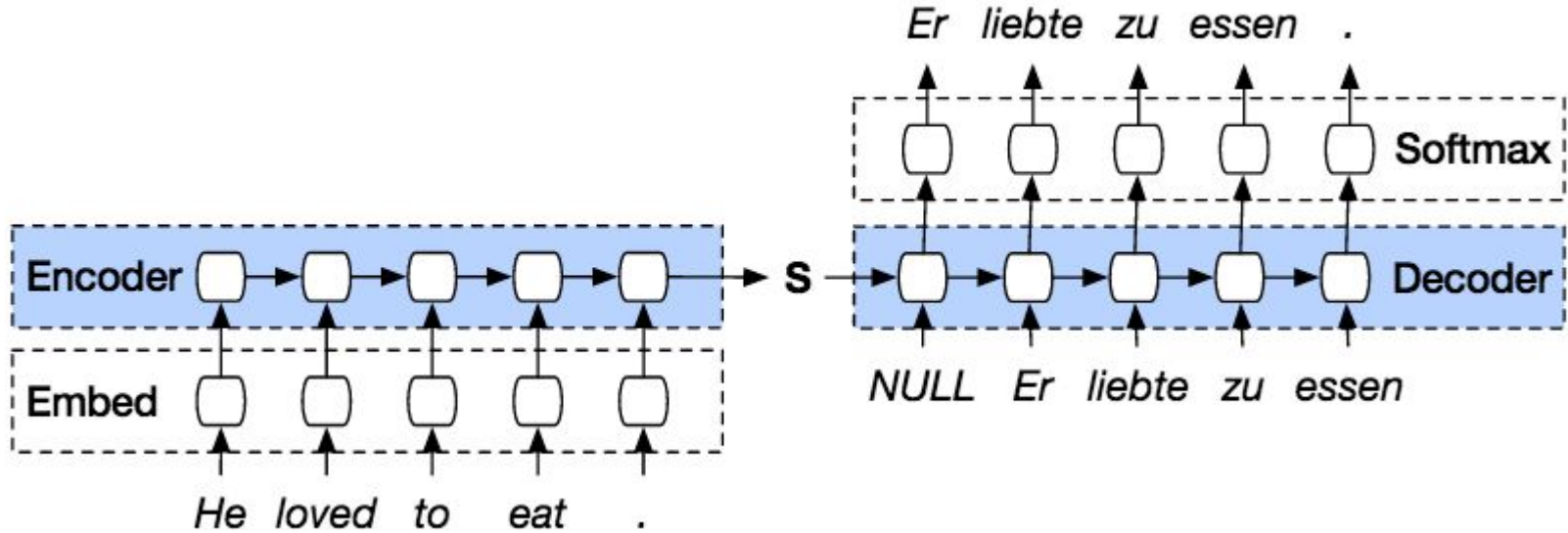
History



Word being predicted

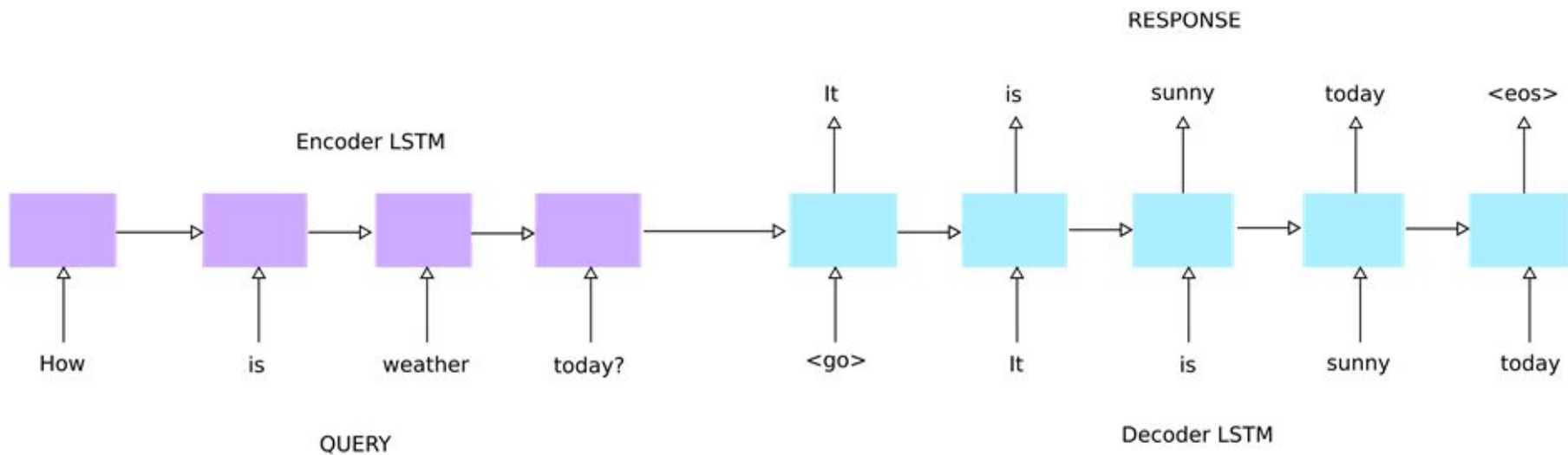
# Seq2Seq

Sutskever et al. (2015)



# Neural Conversational Model

Vinyals et al. (2015)



# DialogPT

Zhang et al. (2019)

- **Reddit conversations**
  - 146.8M training instances
  - 3.87B source tokens
  - 2.14B target tokens
- **Different Models**
  - 117M
  - 345M + Reverse
  - 762M



<User> Does money buy happiness?

<Bot> Depends how much money you spend on it .

<User> What is the best way to buy happiness ?

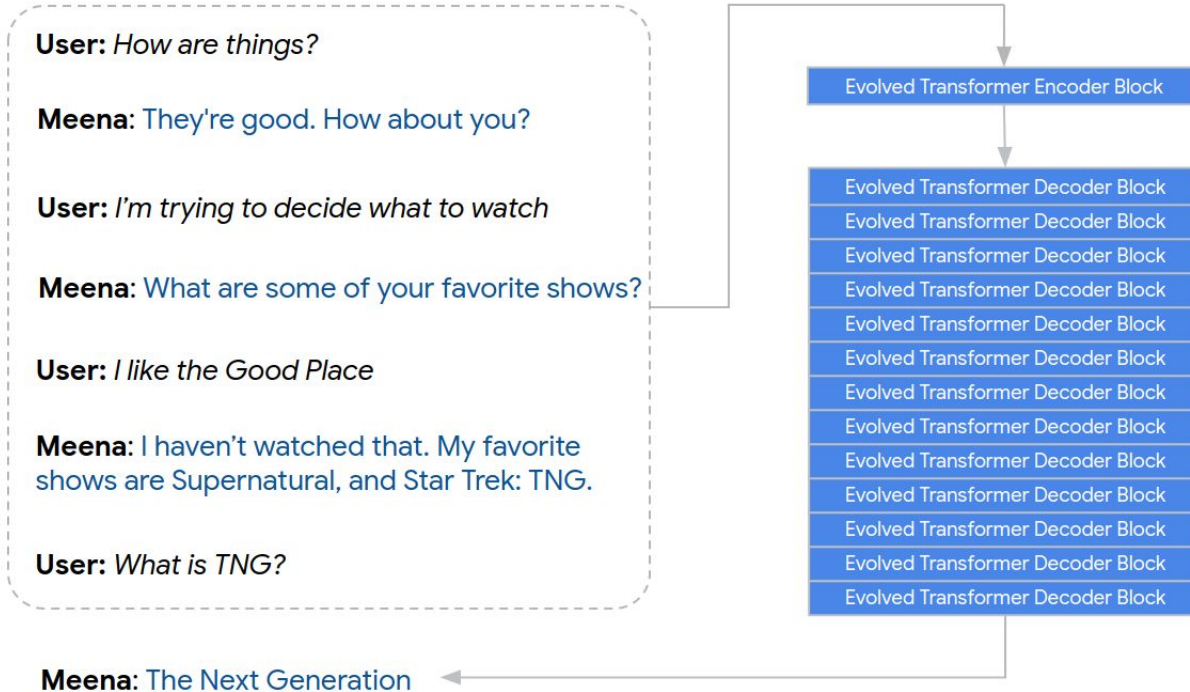
<Bot> You just have to be a millionaire by your early 20s, then you can be happy .

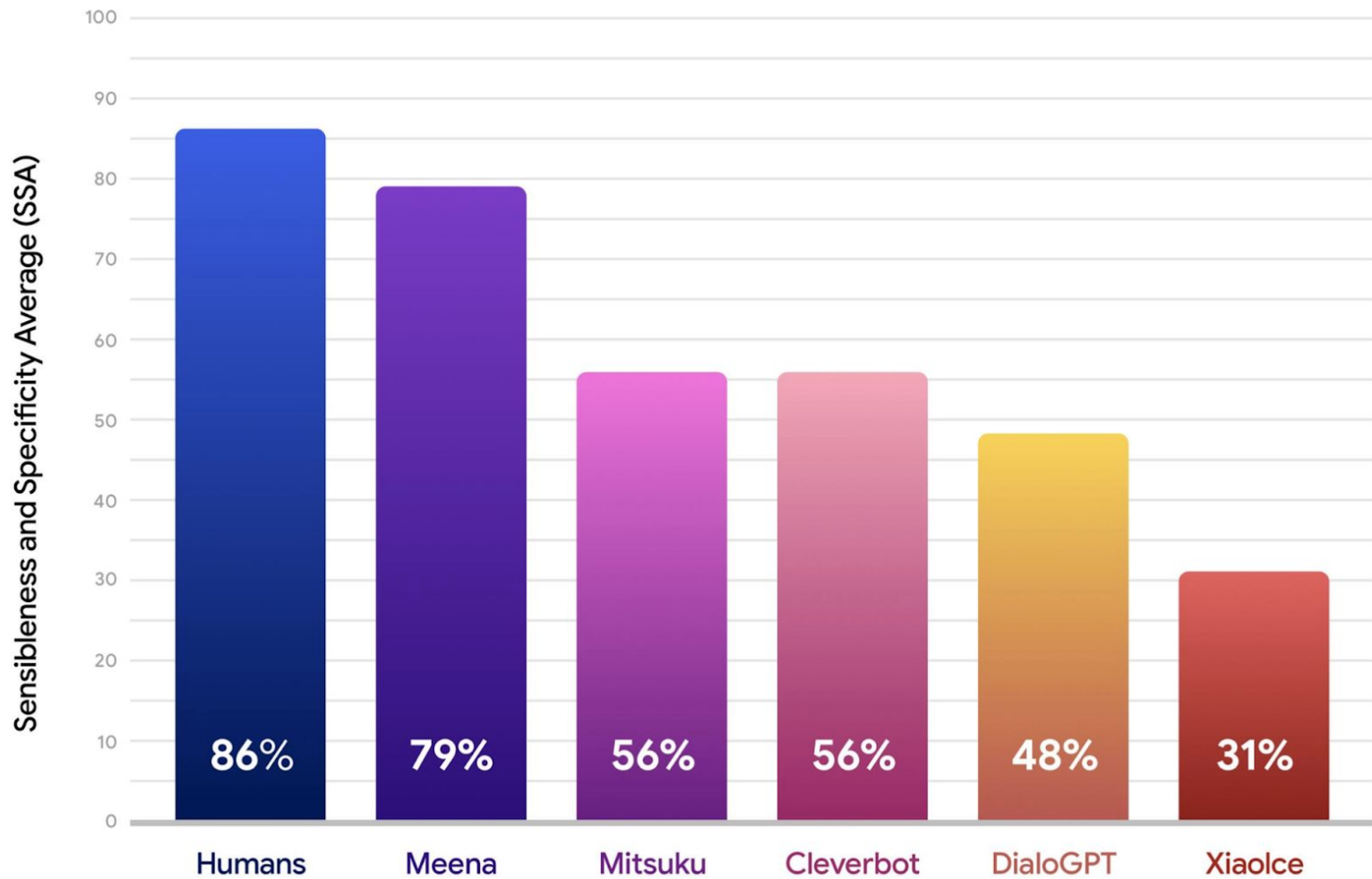
<User> This is so difficult !

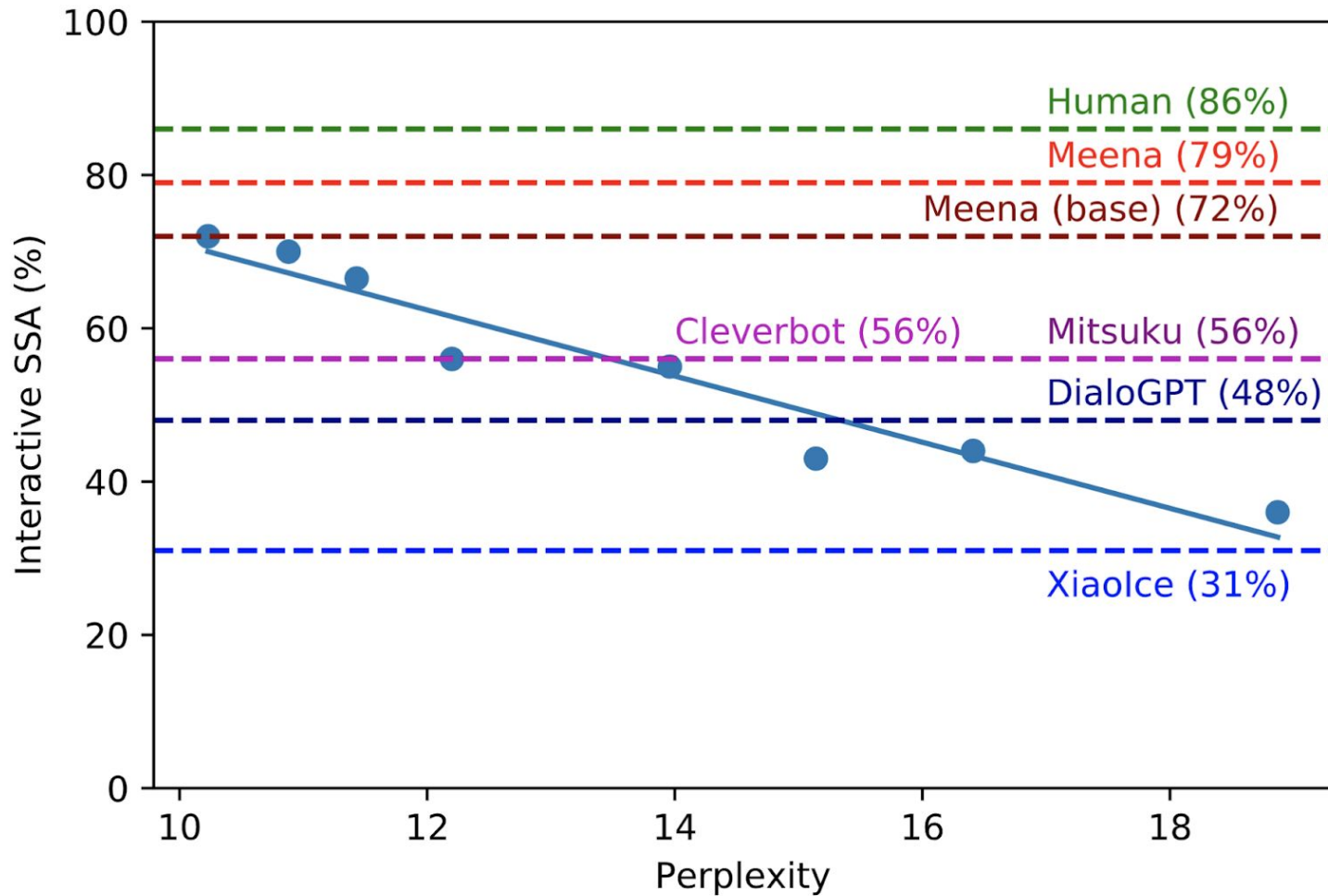
<Bot> You have no idea how hard it is to be a millionaire and happy . There is a reason the rich have a lot of money

# Meena

## Adiwardana et al. (2020)







# Blender

Roller et al. (2020)

- **Blending skills**
  - Engaging use of personality (PersonaChat)
  - Engaging use of knowledge (Wizard of Wikipedia)
  - Display of empathy (Empathetic Dialogues)
  - Ability to blend all three seamlessly (BST)
- **Retrieve and refine**

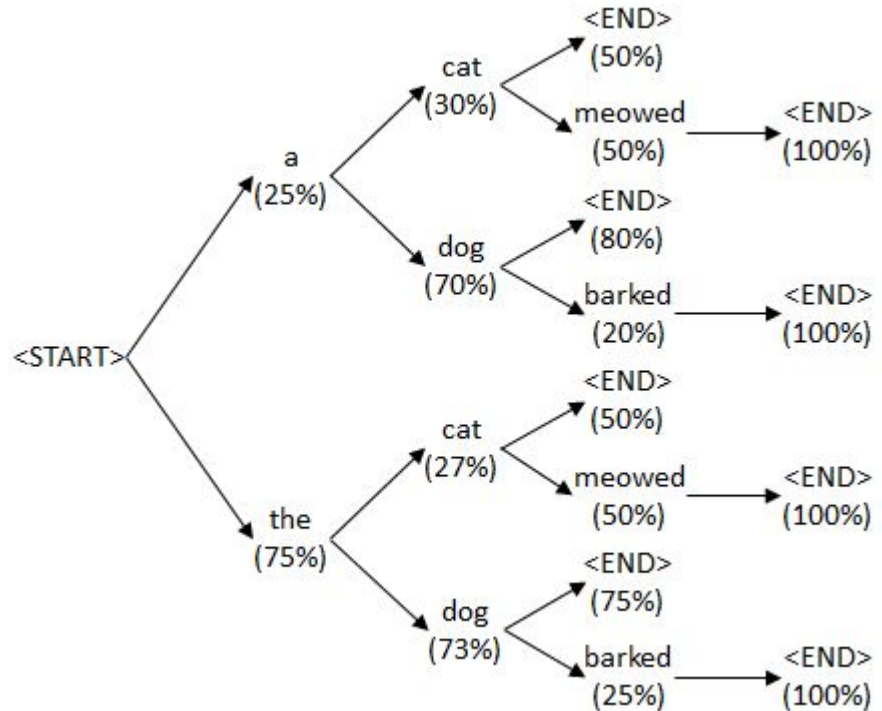
# Generative Decoding

- **Deterministic**
  - Greedy Decoding
  - Beam Search
- **Sampling**
  - Top-k
  - Nucleus
- **Similarity search (Continuous output)**

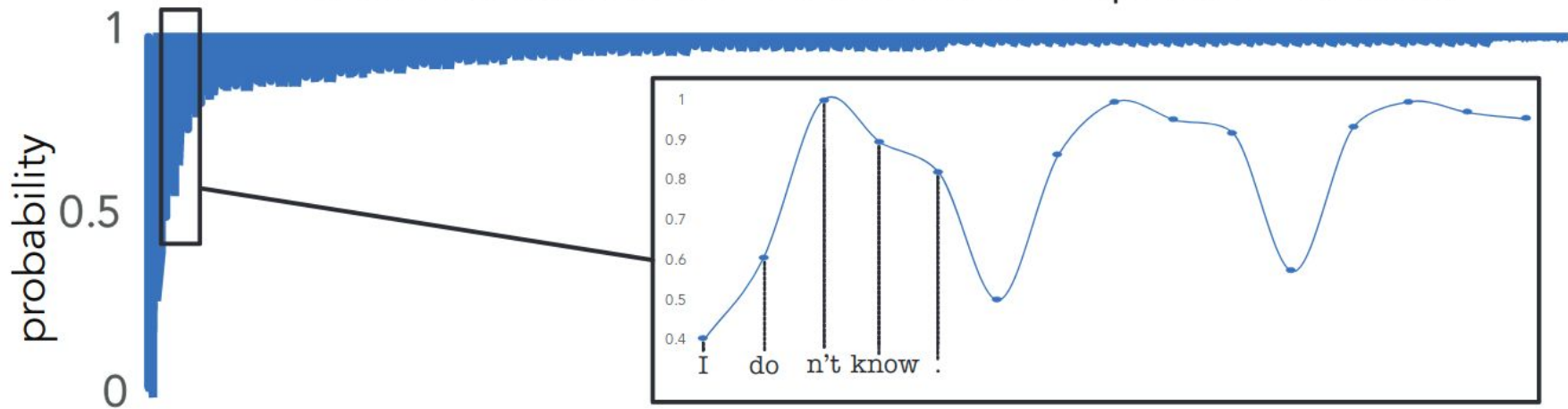
# Beam search

Hyperparameter:  
K (beam width)

Greedy decoding  
K = 1



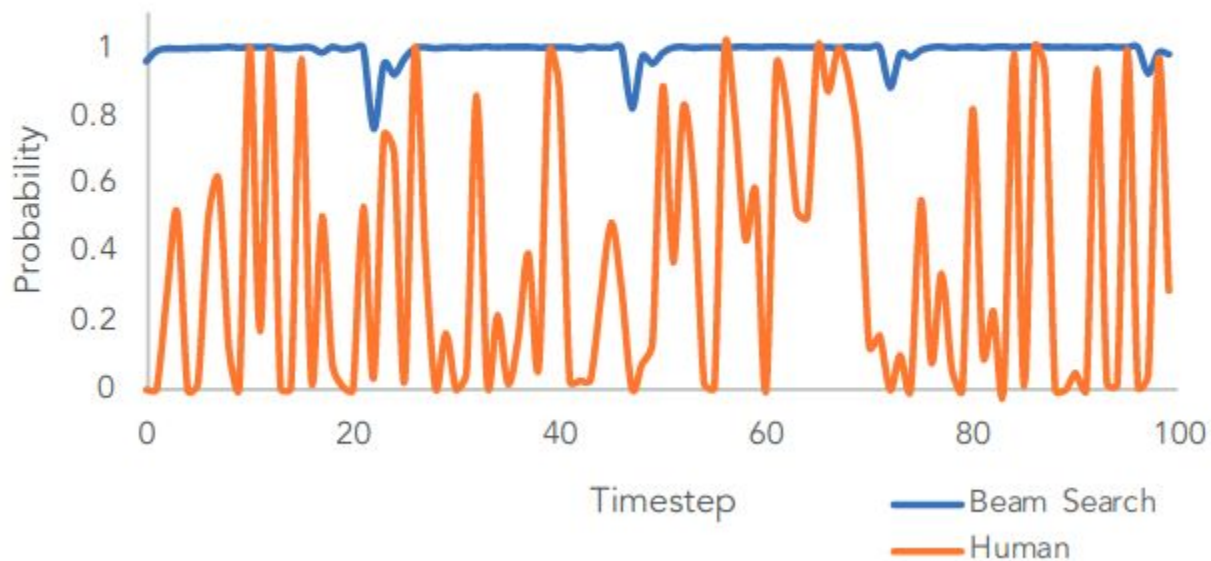
Token Probabilities for "I don't know." Repeated 200 times





# Holtzman et al. (2020)

Beam Search Text is Less Surprising



# Enhancing Beam search

- **Response length**

- Hard coded minimum

- Response length prediction

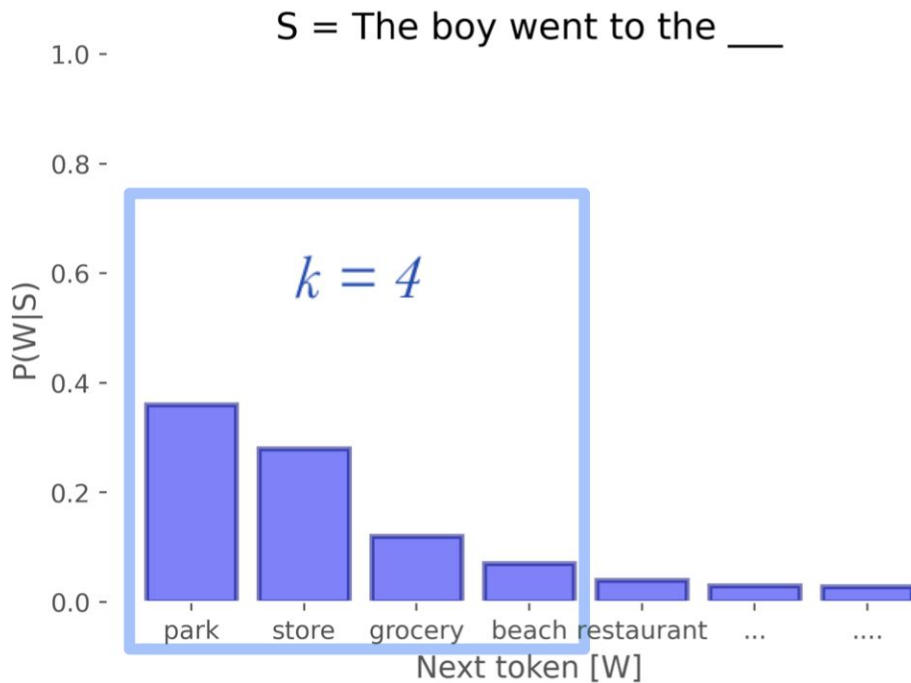
- (Blender) 4 class classifier (<10, <20, <30, >30)

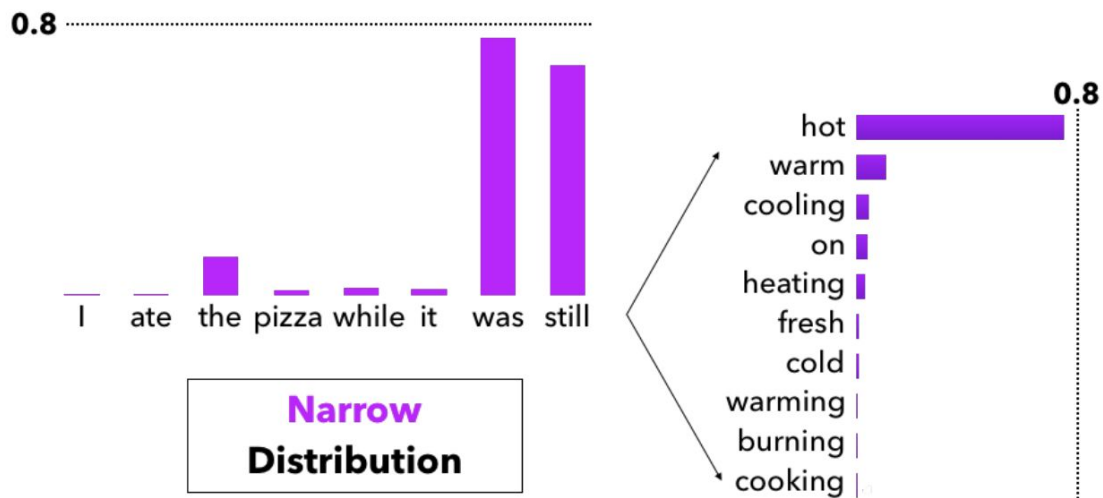
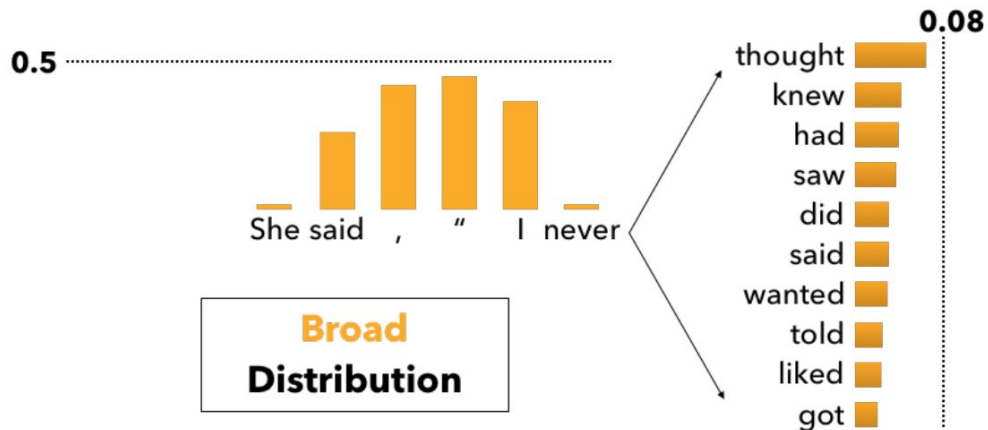
- **Subsequence blocking**

- Not allowing existing n-grams

# Top-k decoding

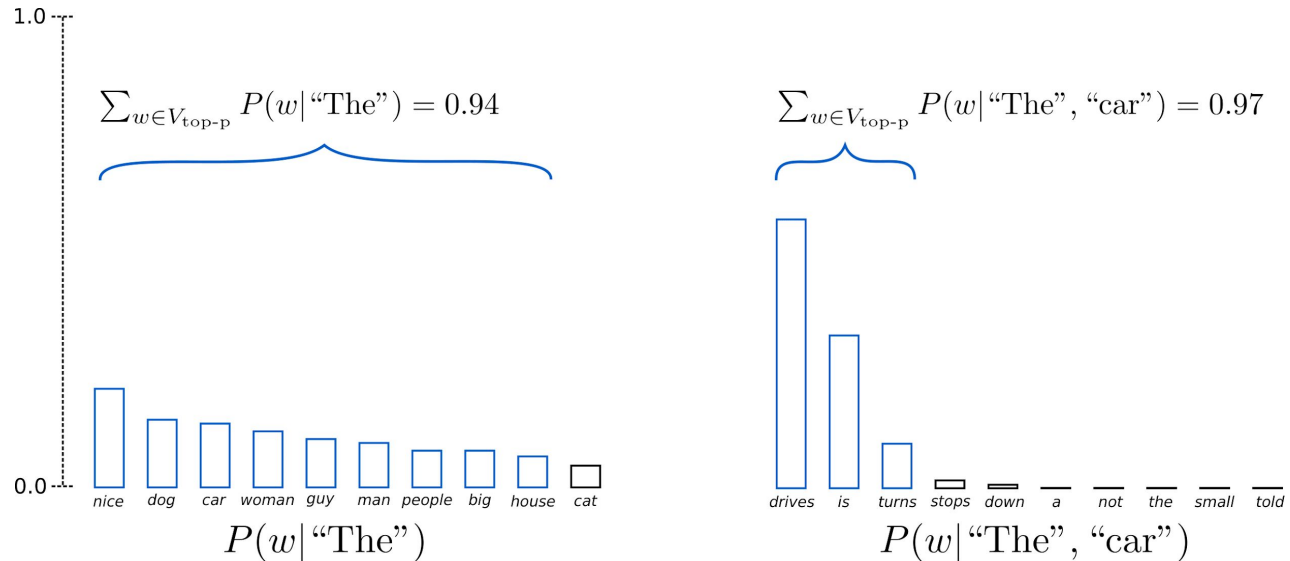
Hyperparameter:  
K (sample size)

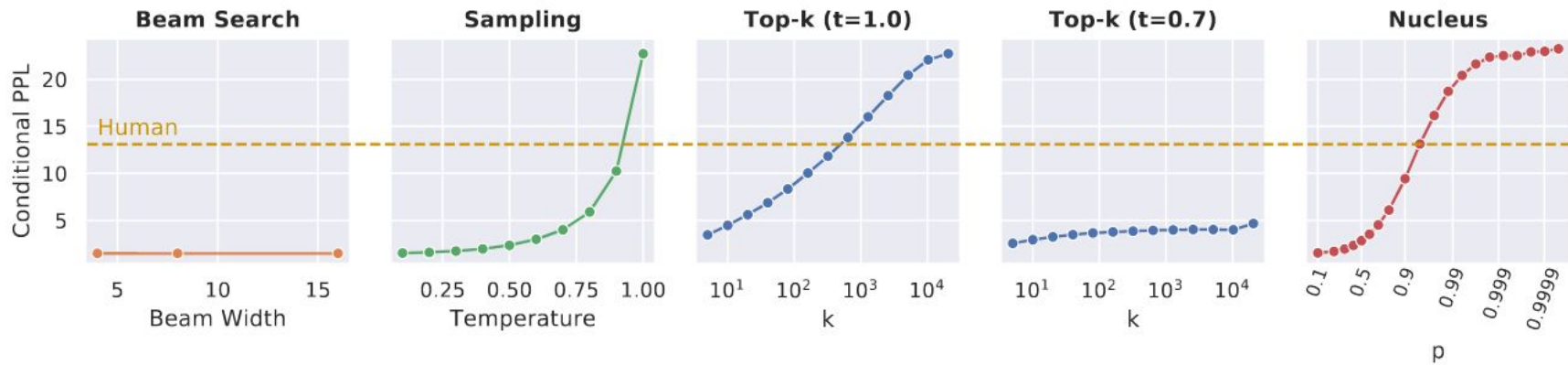




# Nucleus decoding

Hyperparameter:  
**p** (probability mass)





# MMI Decoding

## Forward + Backward language model

1. The forward model generates the output and computes  $P_{\text{forward}}$
2. The backward model computes  $P_{\text{backward}}$

Output is then scored by  $P_{\text{backward}} + P_{\text{forward}}$

# Similarity Search (Continuous output)

## Method:

- Model output is not projected onto the vocabulary
- Objective is to minimize  $\text{Sim}(v_{\text{output}}, v_{\text{target}})$
- $v_{\text{target}}$  is drawn from pre-trained embeddings

## Advantages:

- Faster training
- Interpretable output

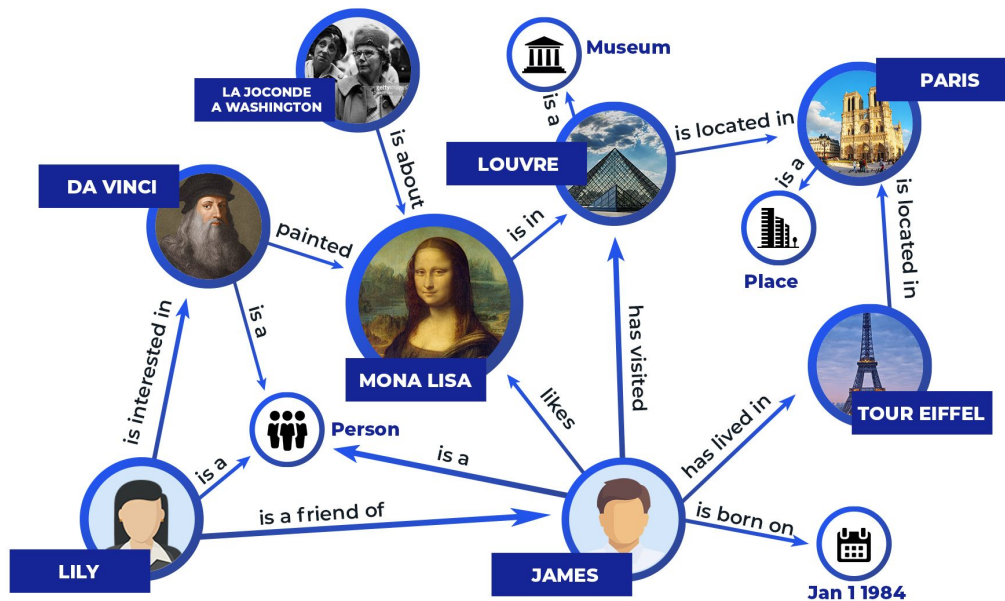


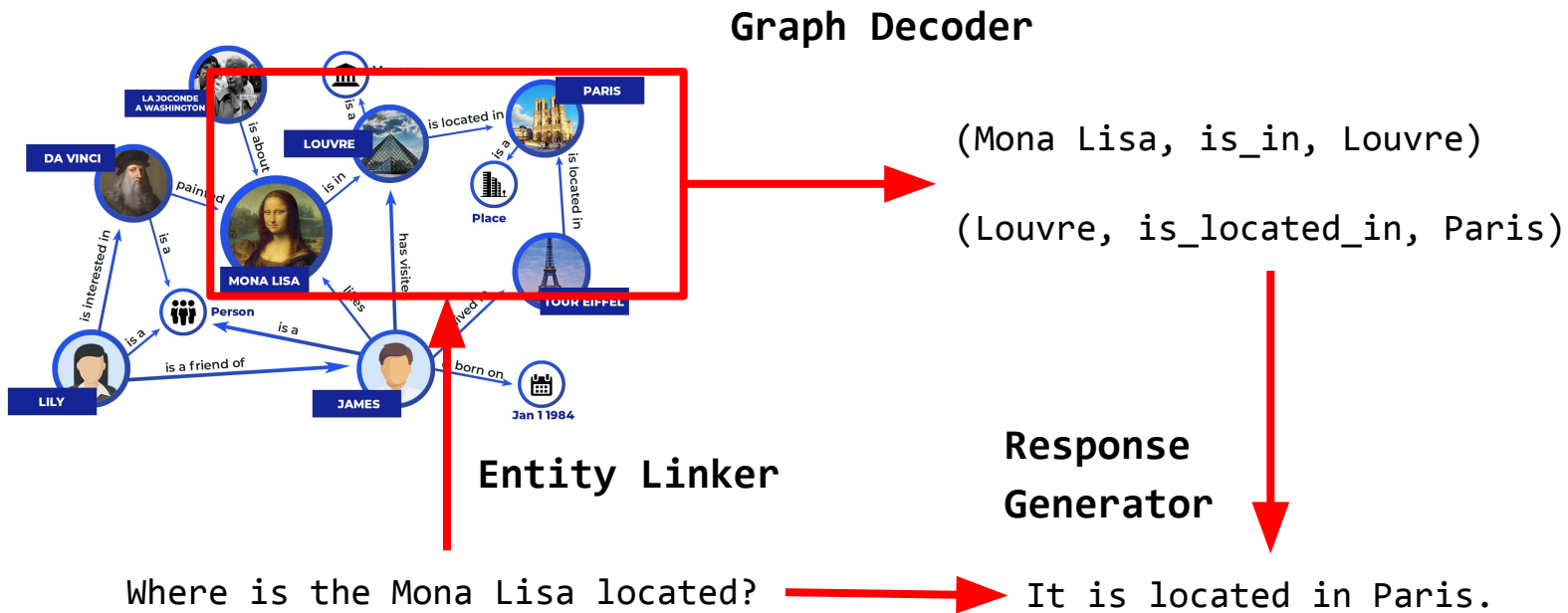
# Issue with end-to-end solutions

- **Bad interpretability**
  - Knowledge is stored as parameter values
  - We can't inspect the cause of a bad answer
- **Information hallucination**

Example: "Cats have three legs"

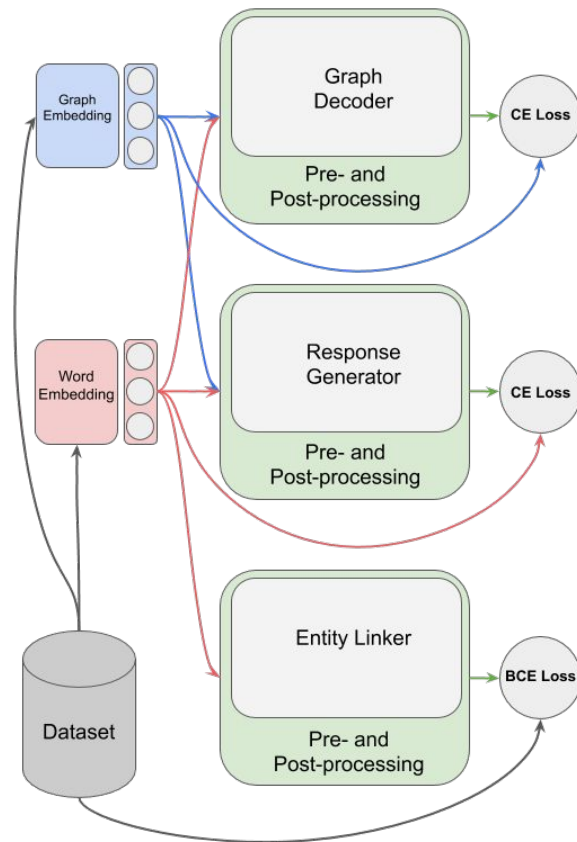
# Knowledge graphs





# Architecture

- **Word Embedding**
  - Fasttext
- **Graph Embedding**
  - TransE
- **Graph Decoder**
  - LSTM
- **Response Generator**
  - Transformer
  - Continuous output
- **Entity Linker**
  - BERT



# Dataset

## OpenDialKG

Moon et al. (2019)

91031	messages
100813	entities
1358	relations
1190658	triplets

```
[
  {
    "Message": "Do you like Iron Man",
    "Metadata": [
      [
        "Iron Man",
        "starred_actors",
        "Robert Downey Jr."
      ]
    ]
  }
]
```

# Word embeddings (Fasttext)



# Graph embeddings

- DeepWalk - Perozzi et al. (2014)
- TransE - Bordes et al. (2013)
- TransR - Lin et al. (2015)

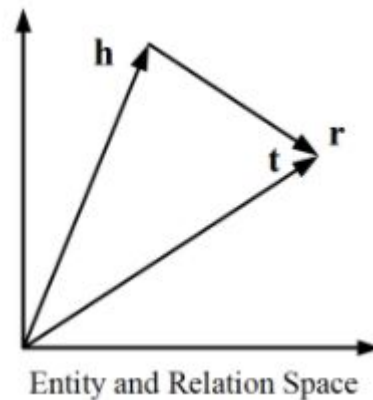
# TransE

Bordes et al. (2013)

## Knowledge triplet:

Mona Lisa (head), is\_in (relation), Louvre (tail)

$$V_{\text{head}} + V_{\text{relation}} = V_{\text{tail}}$$



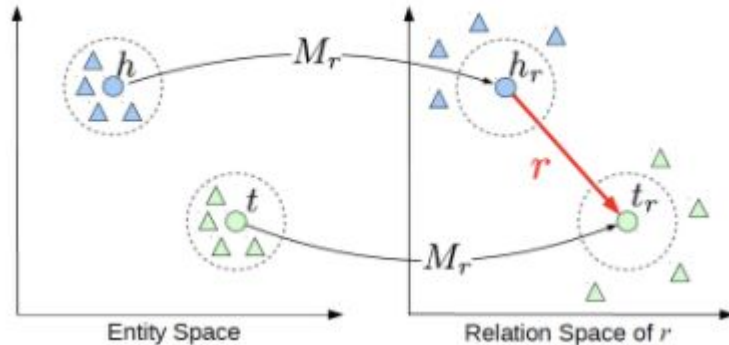


# TransR

Lin et al. (2014)

Each relation has its own transformation matrix:  $M_r$

$$M_r v_h + v_r = M_r v_t$$



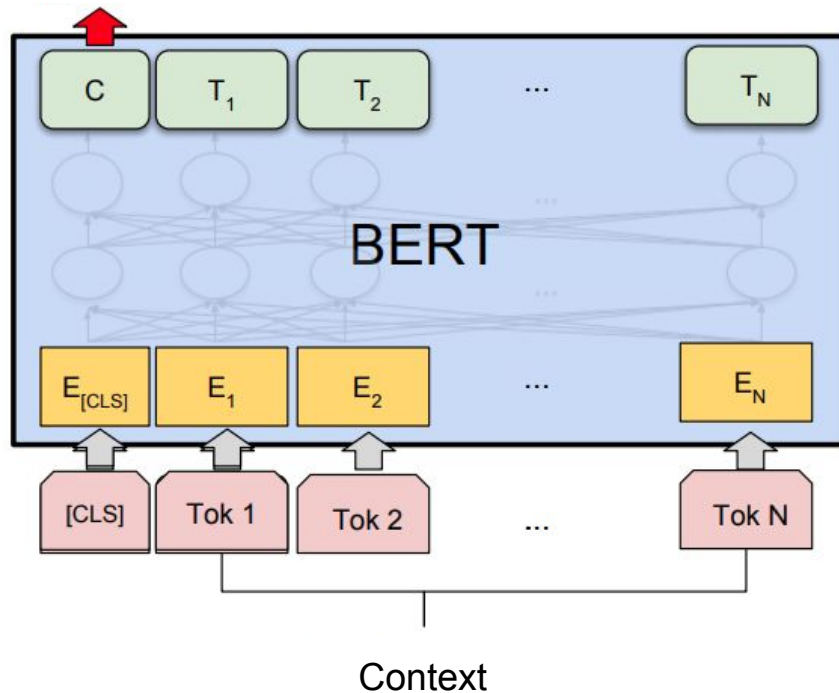
# Entity Linker

Limited information  
in OpenDialKG

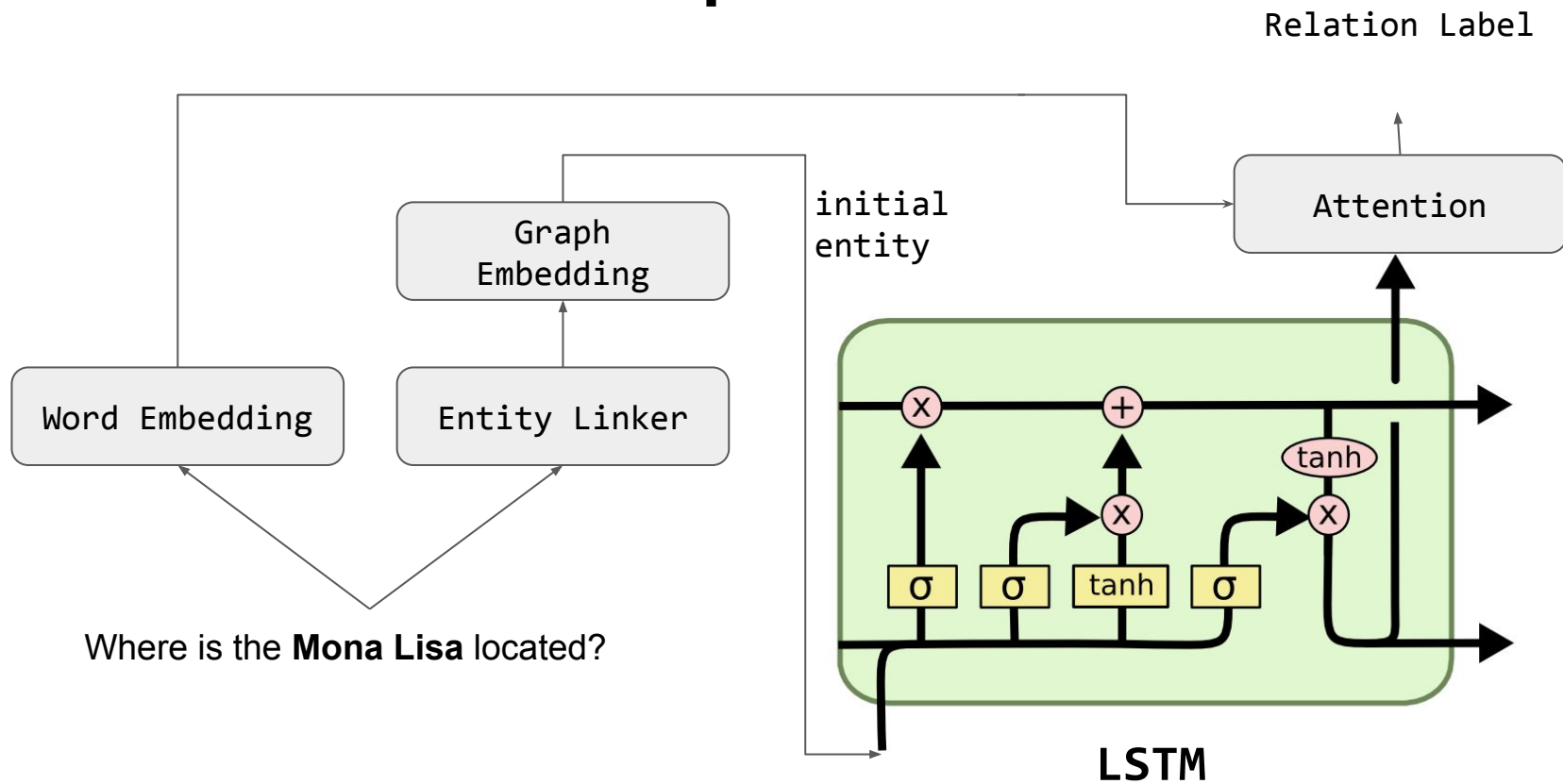


Labels are learned  
with BCE objective

Entity Labels



# Graph Decoder



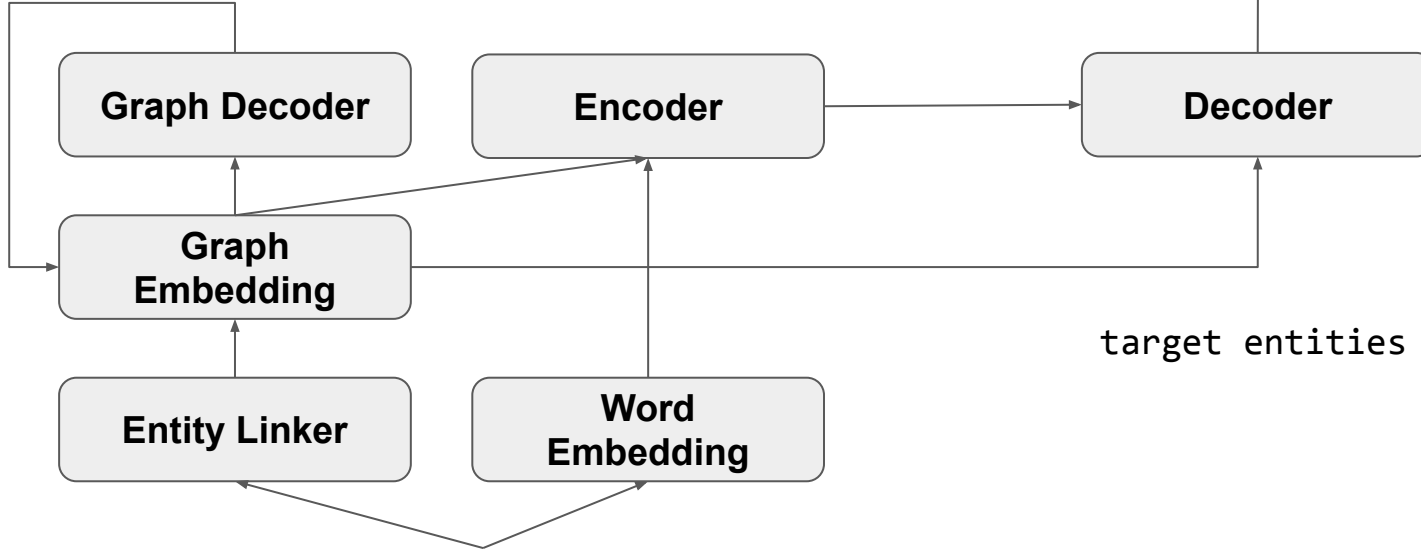
# Response Generator

- **Transformer encoder-decoder architecture**
  - additional cross-attention layer for the entity node representations
- **Continuous outputs with dot-product loss + negative sampling**
  - (Future work) Decoding entity nodes directly

# Response Generator

graph paths for  
target entities

It is located in Paris.



target entities

Where is the **Mona Lisa** located?

# Baselines

- GPT-2 (345M)
- XLNet (110M)
- Vanilla Transformer
- Continuous Transformer

# Automatic Metrics

- **Length**

Generated data: Average length of model outputs

Test data: Delta length compared to gold responses

- **BLEU**

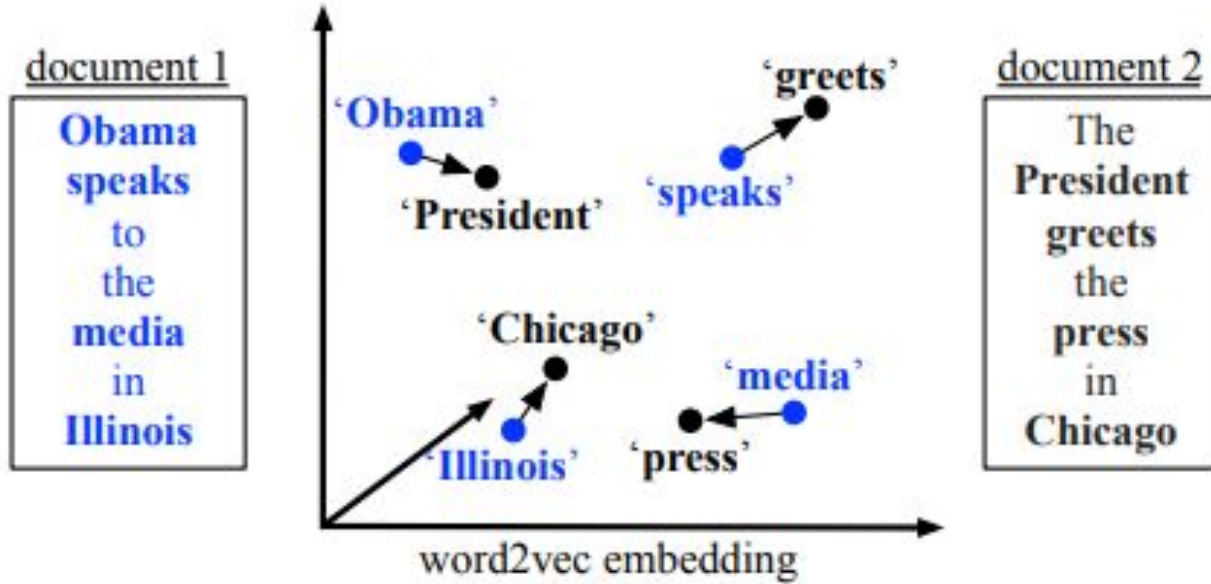
Modified precision score of n-gram overlaps

Common metric in NMT

Low correlation with human judgement

- **WMD**

# WMD





# Results

Model	Decoding	Length		WMD		BLEU	
		Mean	Std	Mean	Std	Mean	Std
GPT-2 (small)	Greedy decoding	10.2	3.5	2.0	0.4	0.1	0.0
	Beam Search $k = 3$	12.1	3.6	1.7	0.5	0.1	0.0
	Beam Search $k = 5$	12.0	3.5	<b>1.6</b>	<b>0.4</b>	0.1	0.0
	Top-k Sampling $k = 100$	14.2	3.8	1.6	0.6	0.1	0.0
	Nucleus decoding $p = 0.1$	12.1	3.9	1.7	0.6	<b>0.2</b>	<b>0.1</b>
XLNet (base)	Greedy decoding	9.4	3.4	2.2	0.3	0.1	0.0
	Beam Search $k = 3$	11.2	3.2	1.9	0.4	0.1	0.0
	Beam Search $k = 5$	12.9	3.3	1.8	0.3	<b>0.2</b>	<b>0.1</b>
	Top-k Sampling $k = 100$	<b>14.3</b>	<b>4.3</b>	1.9	0.6	0.1	0.0
	Nucleus decoding $p = 0.1$	14.1	4.5	1.8	0.5	0.1	0.0
Transformer	Greedy decoding	7.1	2.2	2.3	0.2	0.0	0.0
	Beam Search $k = 3$	8.5	2.0	2.3	0.4	0.0	0.0
	Beam Search $k = 5$	9.4	2.8	2.2	0.4	0.0	0.0
	Top-k Sampling $k = 100$	9.4	4.1	2.3	0.4	0.1	0.0
	Nucleus decoding $p = 0.1$	9.6	4.7	2.3	0.3	0.1	0.0
Continuous Transformer	Similarity Search	9.8	3.3	2.1	0.3	0.0	0.0
<b>KG-Enhanced Transformer</b>	Similarity Search	12.8	3.4	<b>1.6</b>	<b>0.4</b>	<b>0.2</b>	<b>0.1</b>

**Generated dialogues**

Input	Method	Model Output	Target
I like the movie "The Martian", could you recommend a movie similar to it?	GPT-2 (small) Beam Search $k = 5$	You might enjoy Cloverfield with Drew Barrymore or The Martian. They are both Sci-Fi films.	Do you like The Martian because it stars Matt Damon?
	GPT-2 (small) Nucleus decoding $p = 0.1$	Sure, do you prefer Sci-Fi or Adventure movies?	
	XLNet (base) Beam Search $k = 5$	It is starring Michael Fassbender, who also starred in X-Men: First Class and X-Man 2. Do you like movies that he stars in?	
	XLNet (base) Nucleus decoding $p = 0.1$	I do! I love Sci Fi movies as well. Have you seen either of those?	
	Transformer Nucleus decoding $p = 0.1$	This is a good one	
	Continuous Transformer	Have you seen The Running Man?	
	KG-Enhanced Transformer	It is written by James Patterson, starring Matt Damon and also stars Tracie Thoms. Do you like them?	

# KG-enhanced transformer

- [Human] Hi
- [Model] Hello, have you heard of the movie The Hobbit: Desolation of Smaug?
- [Human] No, can you tell me more about it?
- [Model] It is a fantasy adventure film starring James Nesbitt. Do you like this actor or would you like to know other movies he has starred in?
- [Human] List some of his movies please.
- [Model] Sure! He also starred in The Way and The Amateurs. Have you seen either of those?
- [Human] What is "The Way" about?
- [Model] It's a teen film, but it was directed by Marc Forster who also wrote the movie The Valley of Elah which stars Tom Hanks.
- [Human] Oh I see. Can you recommend me a Sci-Fi movie?
- [Model] Sure! White Television and Undertow are good ones as well. Have you seen either of those?

# GPT-2

- **[Human]** hi
- **[Model]** greetings, i'm gerald.
- **[Human]** how are you ?
- **[Model]** good, thanks. getting ready for another hike, i love hiking!
- **[Human]** oh really? where are you going?
- **[Model]** pikes peak. i am curently exploring famous mountains when i have time.
- **[Human]** that is cool are you going alone?
- **[Model]** my brother is in the us, we go hiking with him every month.
- **[Human]** how old are you by the way?
- **[Model]** a little under 46 years old
- **[Human]** are you in a relationship?
- **[Model]** yes i am married, 3 boys, i am pregnant again. i am too busy to focus

# Conclusion

- Lack of automatic metrics for checking factual accuracy
  - Lack of datasets for this purpose
  - Unsatisfactory results from the proposed architecture
- 
- + I paved the way for future work
  - + I created a framework for training chatbots from pre-trained transformers:

<https://github.com/bme-chatbots/dialogue-generation>

## Future work

- Implementing pre-trained transformer as the response generator
- Improving the OpenDialKG dataset
- Finding better automatic metric

# References

- [arxiv.org/abs/1706.03762](https://arxiv.org/abs/1706.03762) ( Transformer )
- [jalammar.github.io/illustrated-gpt2/](https://jalammar.github.io/illustrated-gpt2/) ( GPT-2 )
- [yashueth.blog/2019/10/08/](https://yashueth.blog/2019/10/08/) ( Knowledge graph )