

# Interpretability of Generative Information Retrieval

Anja Reusch – Postdoctoral Researcher – NLP Lab, Technion



**TECHNION**

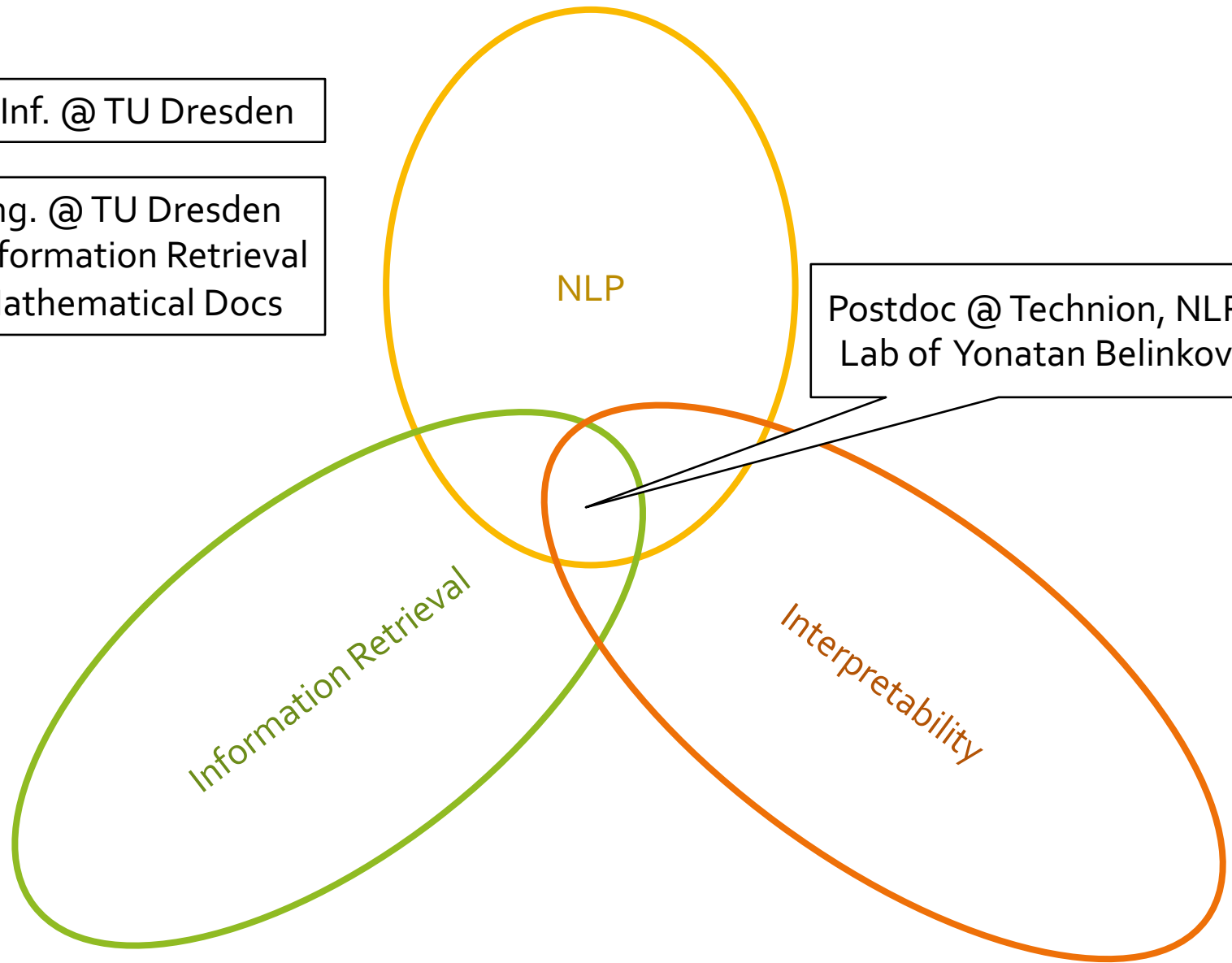


# About Me

Dipl. Inf. @ TU Dresden

Dr.-Ing. @ TU Dresden  
on Information Retrieval  
for Mathematical Docs

Postdoc @ Technion, NLP  
Lab of Yonatan Belinkov



# Interpretability in Retrieval

# Interpretability in Retrieval



I want to know  
about the latest  
developments in  
robotics.

# Interpretability in Retrieval



I want to know  
about the latest  
developments in  
robotics.



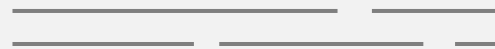
SEARCH



*Top Search Results*



Science News - Robotics



Top Developments in ...



# Interpretability in Retrieval

SEARCH



*Top Search Results*



Science News - Robotics



Top Developments in ...



I want to know about the latest developments in robotics.

CHAT

I want to know about the latest developments in robotics.

*Chat:*

**Latest Developments in Robotics**

Here are some of the key developments from this year:

<b>Trend</b>	<b>Description</b>	<b>Impact</b>
--------------	--------------------	---------------

# Interpretability in Retrieval

## CHAT

I want to know about the latest developments in robotics.



LLM - Large  
Transformer based  
Language Model

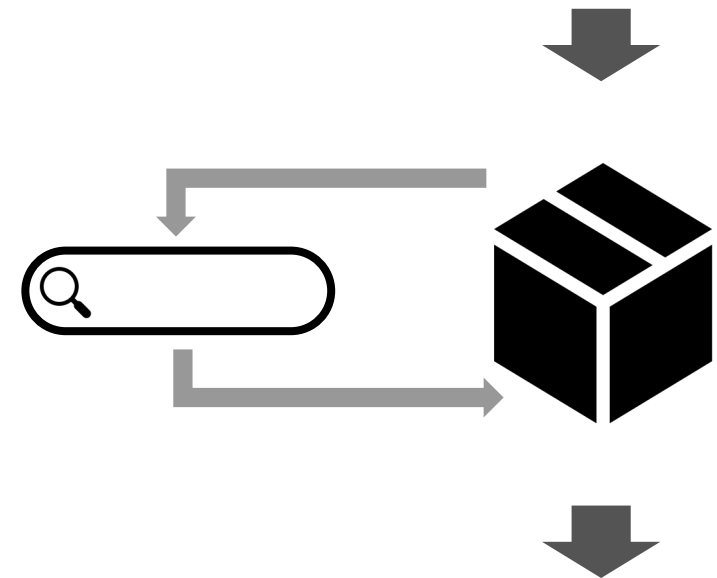


**Latest Developments in Robotics**  
Here are some of the key...

# Interpretability in Retrieval

## CHAT

I want to know about the latest developments in robotics.



**Latest Developments in Robotics**  
Here are some of the key...

# Interpretability in Retrieval

## SEARCH

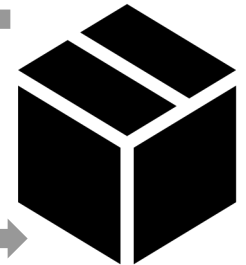
robotics news



Science News - Robotics

## CHAT

I want to know about the latest developments in robotics.



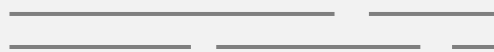
Latest Developments in Robotics  
Here are some of the key...

# Interpretability in Retrieval

SEARCH



Science News - Robotics



## Interpretability of Information Retrieval

- How does the model determine **relevance** of an article?
- How to **improve** this ranking?

Information  
Retrieval is  
everywhere!



**Product  
Recommendation**



**Music  
Recommendation**



**Multi-Media  
Search**

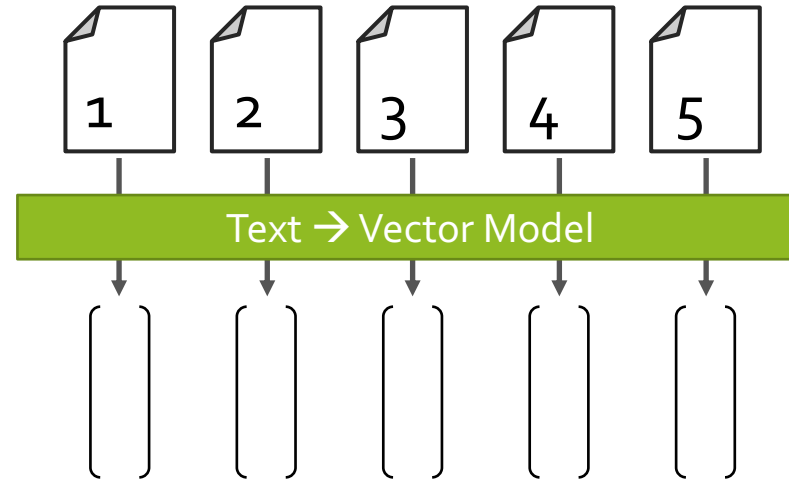
“Why was a product/ song/ album/ video recommended?”

# Interpretability of Generative Information Retrieval

# Interpretability of Generative Information Retrieval

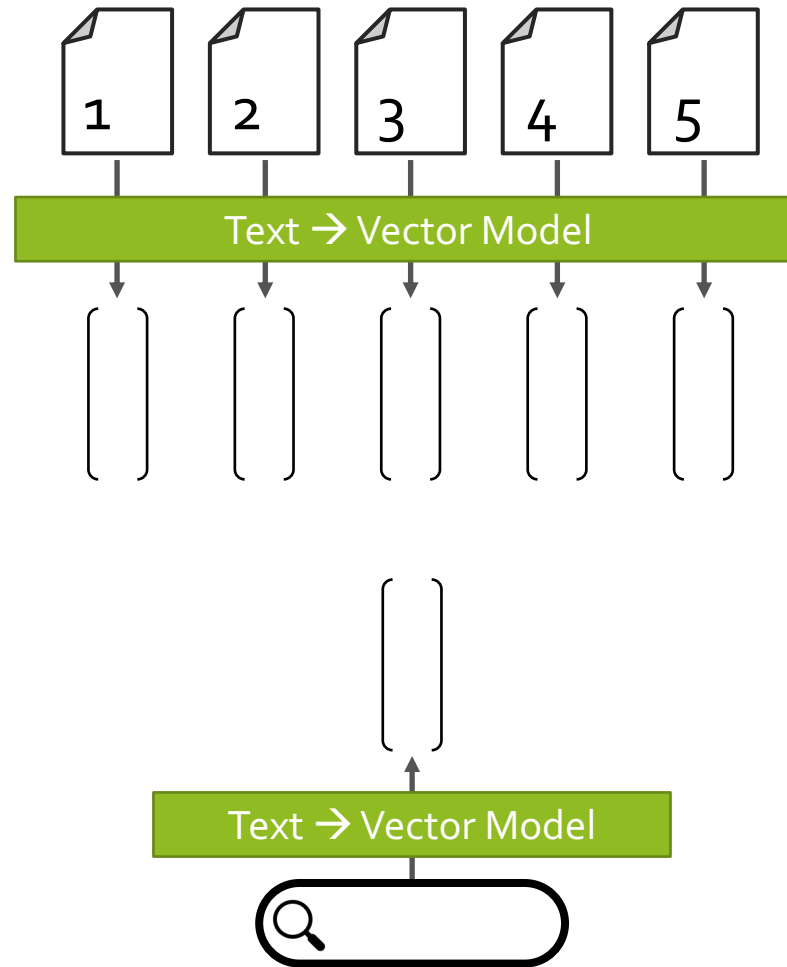
# Current Retrieval Systems

## Dense Retrieval



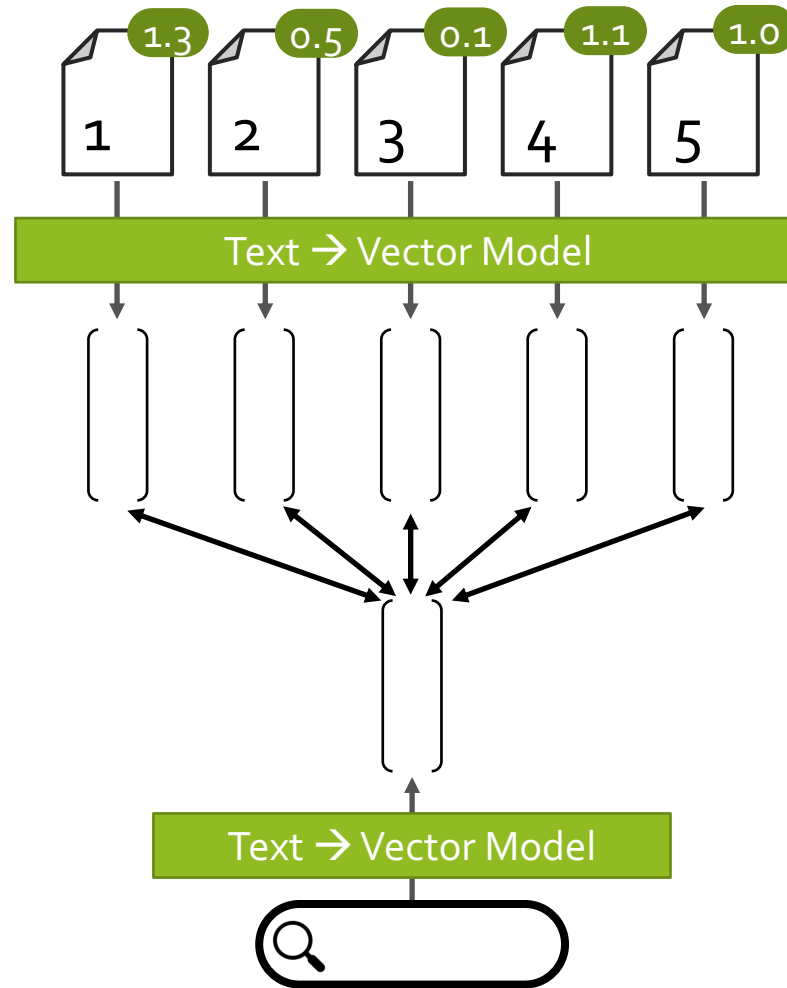
# Current Retrieval Systems

## Dense Retrieval



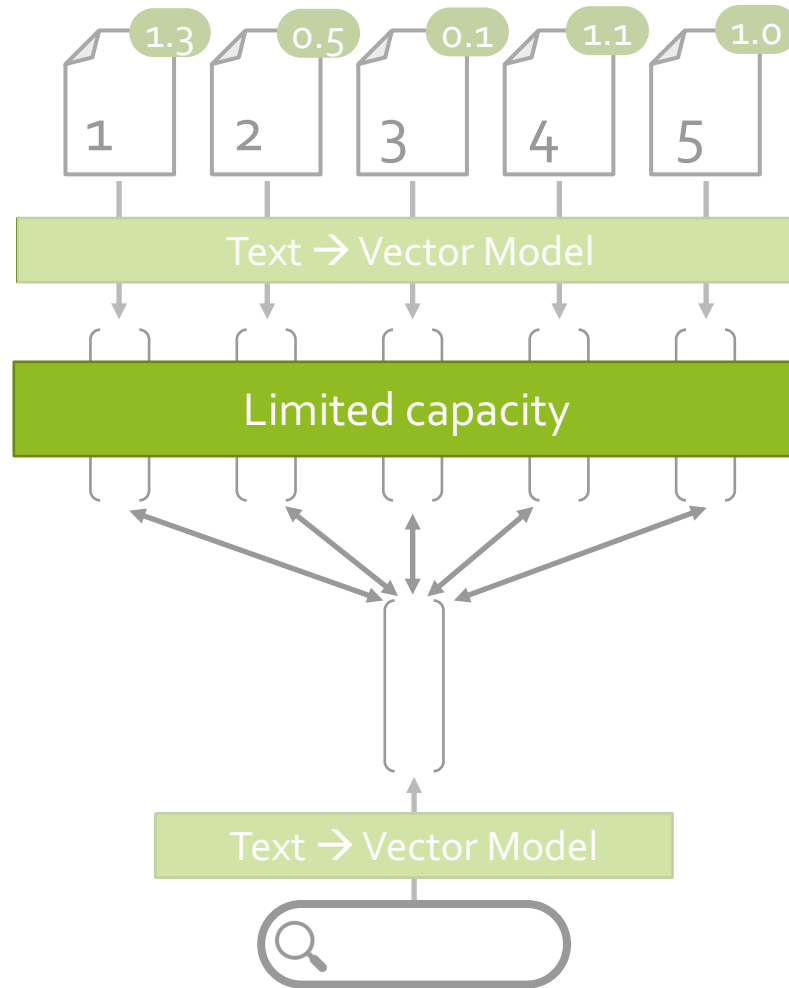
# Current Retrieval Systems

## Dense Retrieval



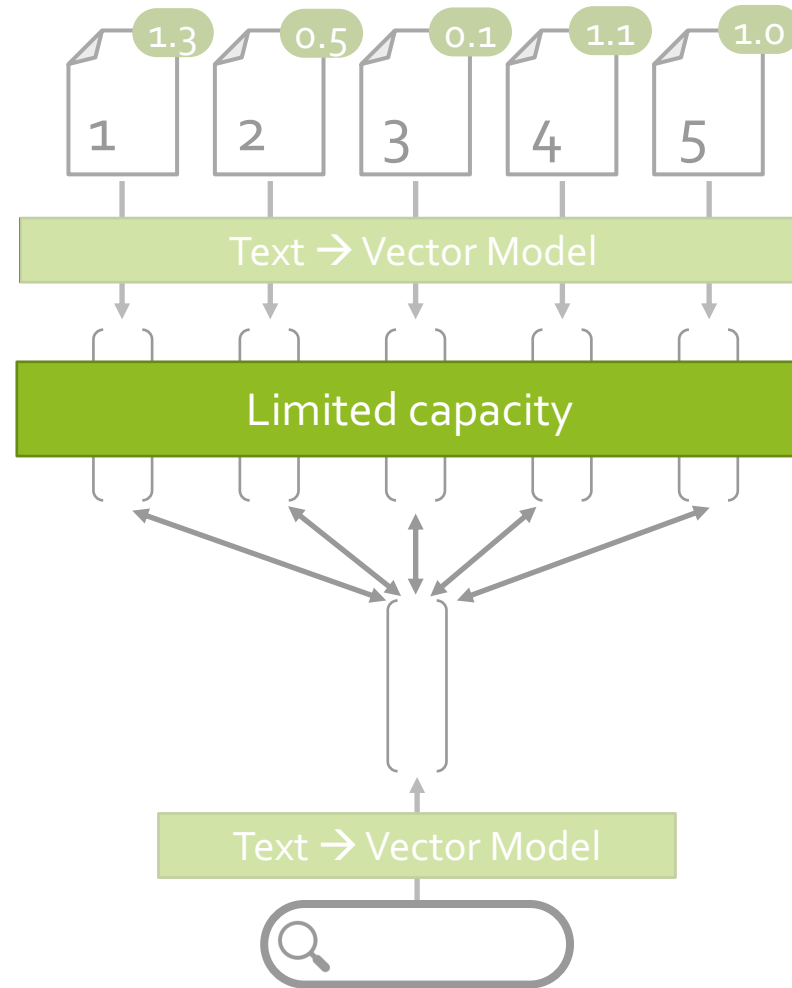
# Current Retrieval Systems

## Dense Retrieval

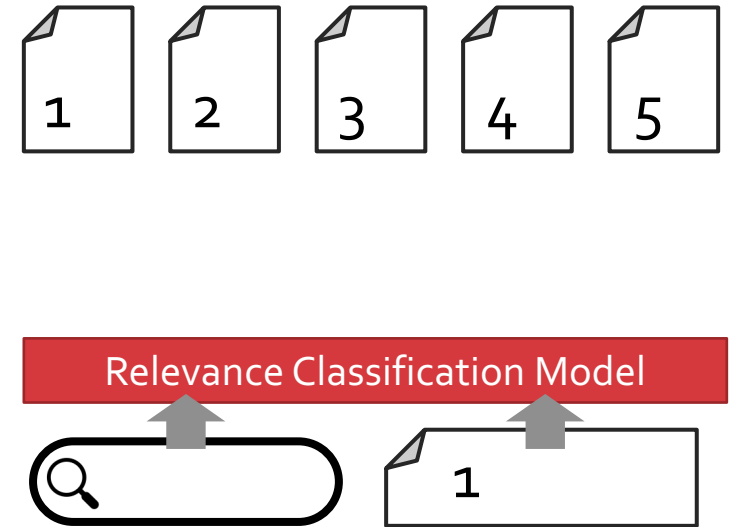


# Current Retrieval Systems

## Dense Retrieval

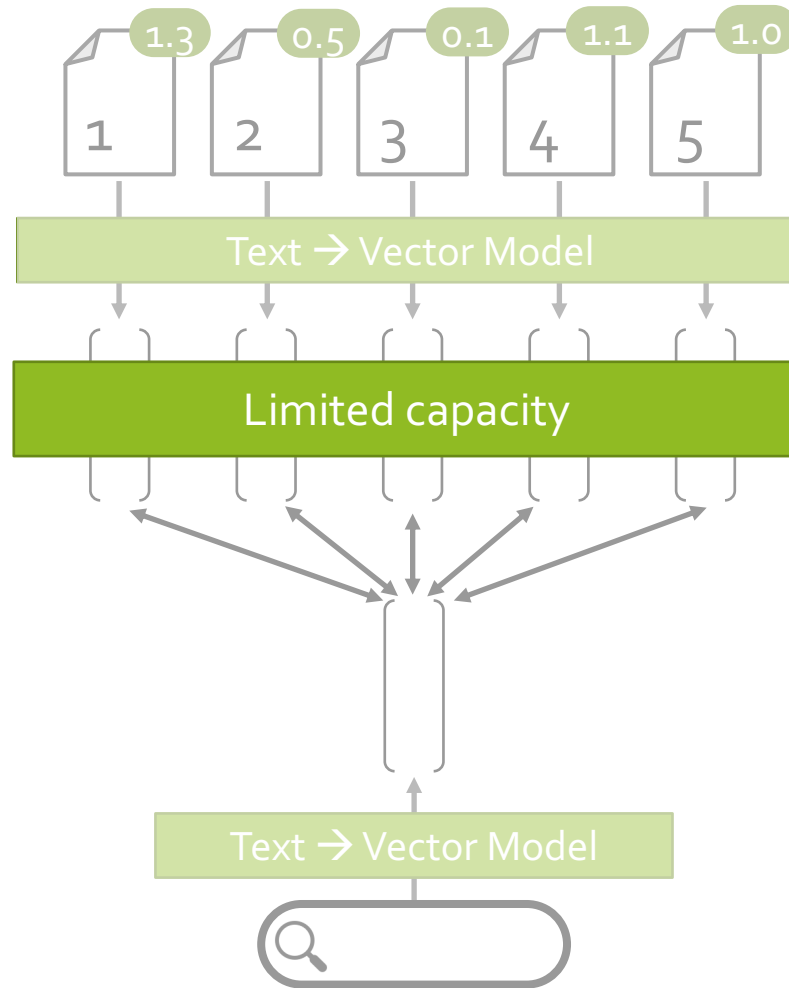


## Cross-Encoders

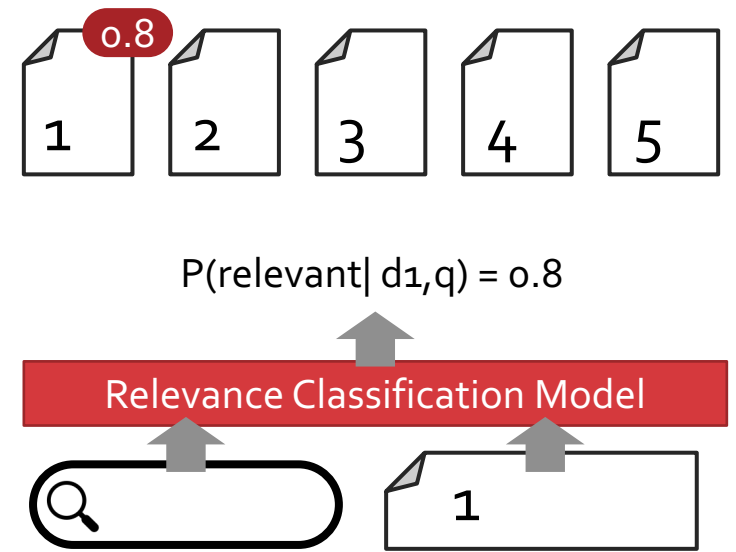


# Current Retrieval Systems

## Dense Retrieval

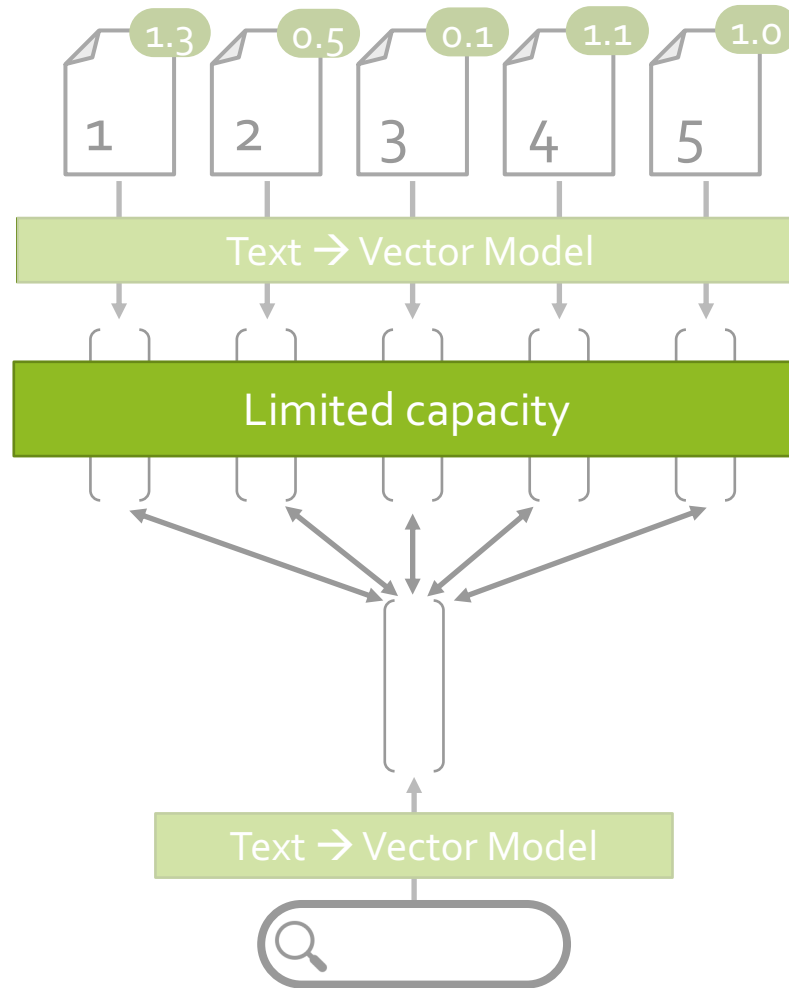


## Cross-Encoders

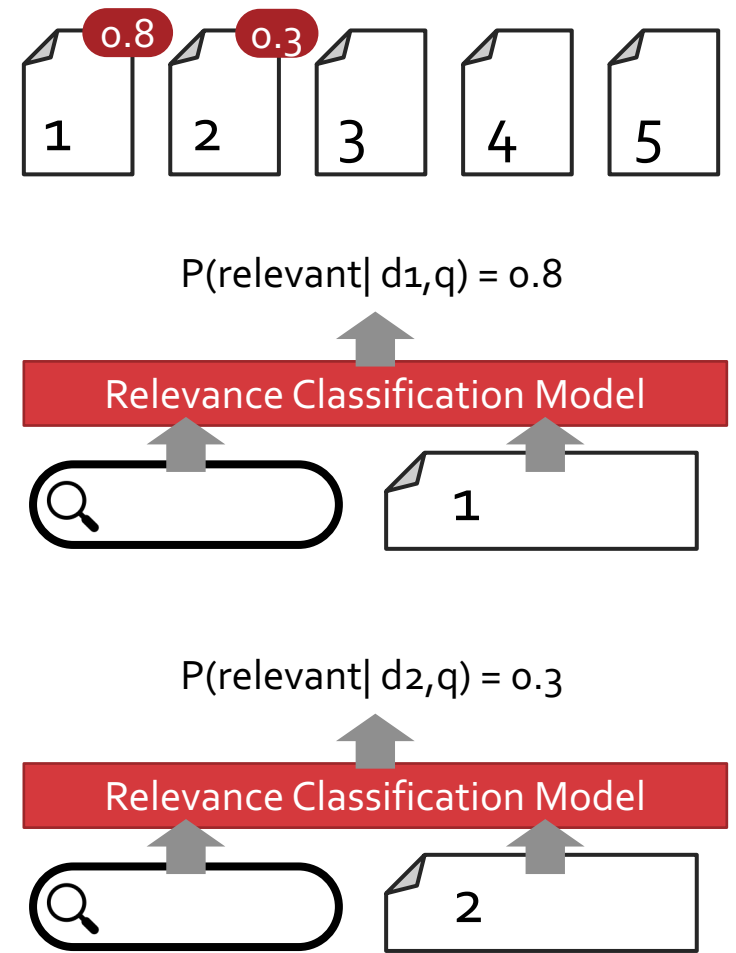


# Current Retrieval Systems

## Dense Retrieval

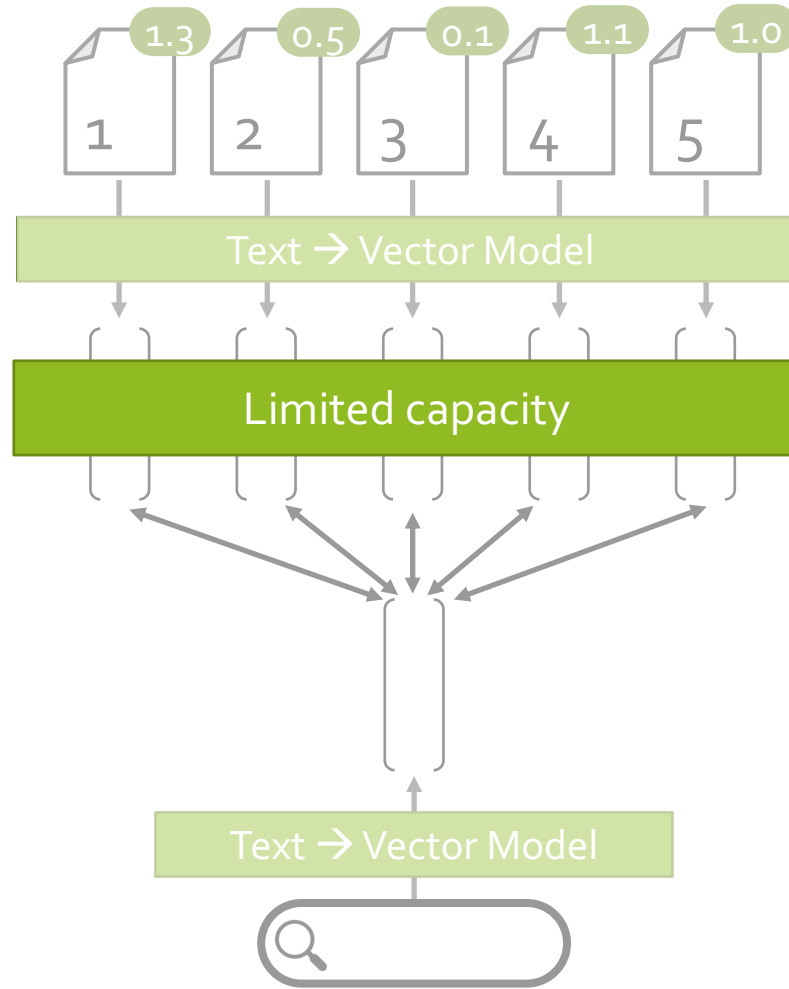


## Cross-Encoders

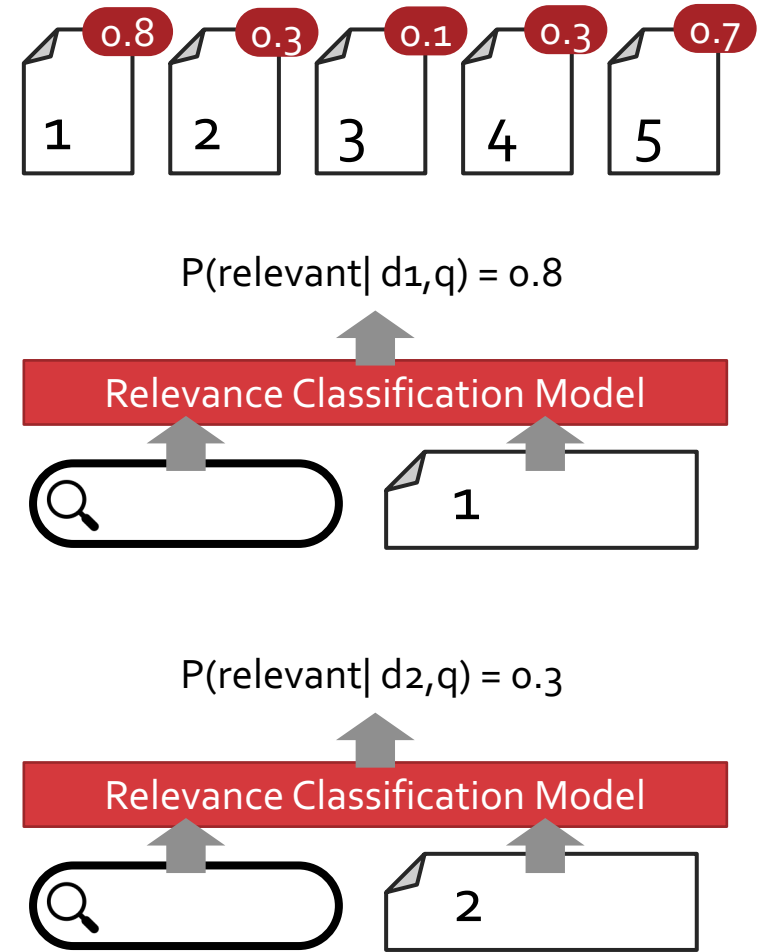


# Current Retrieval Systems

## Dense Retrieval

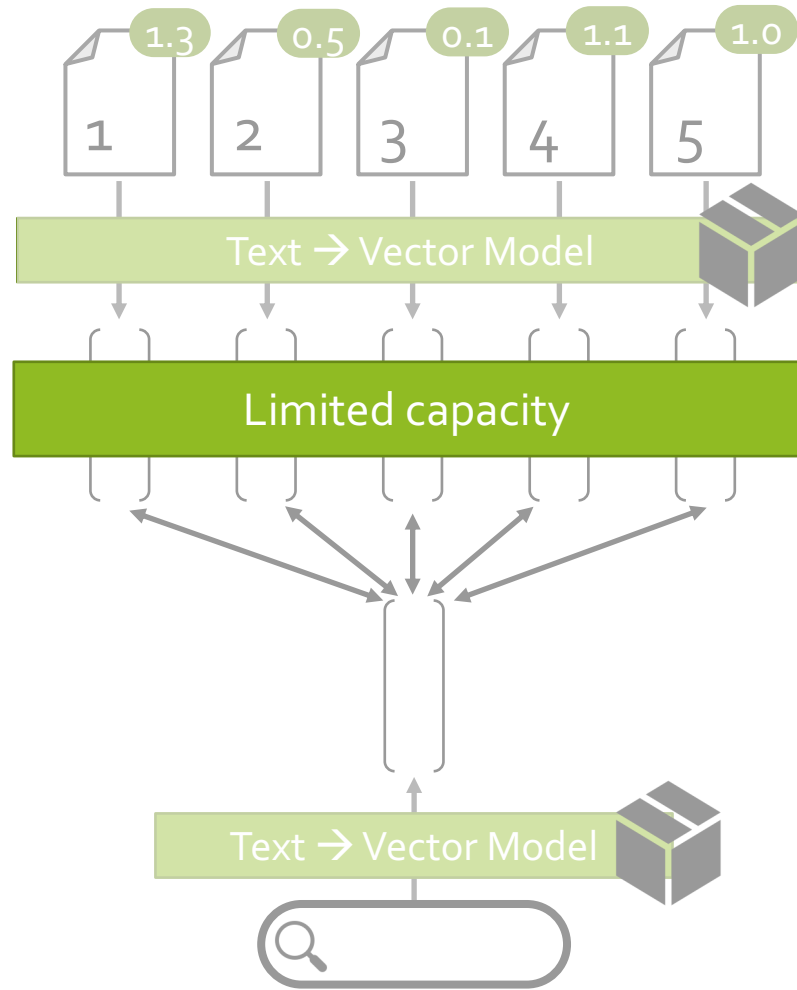


## Cross-Encoders

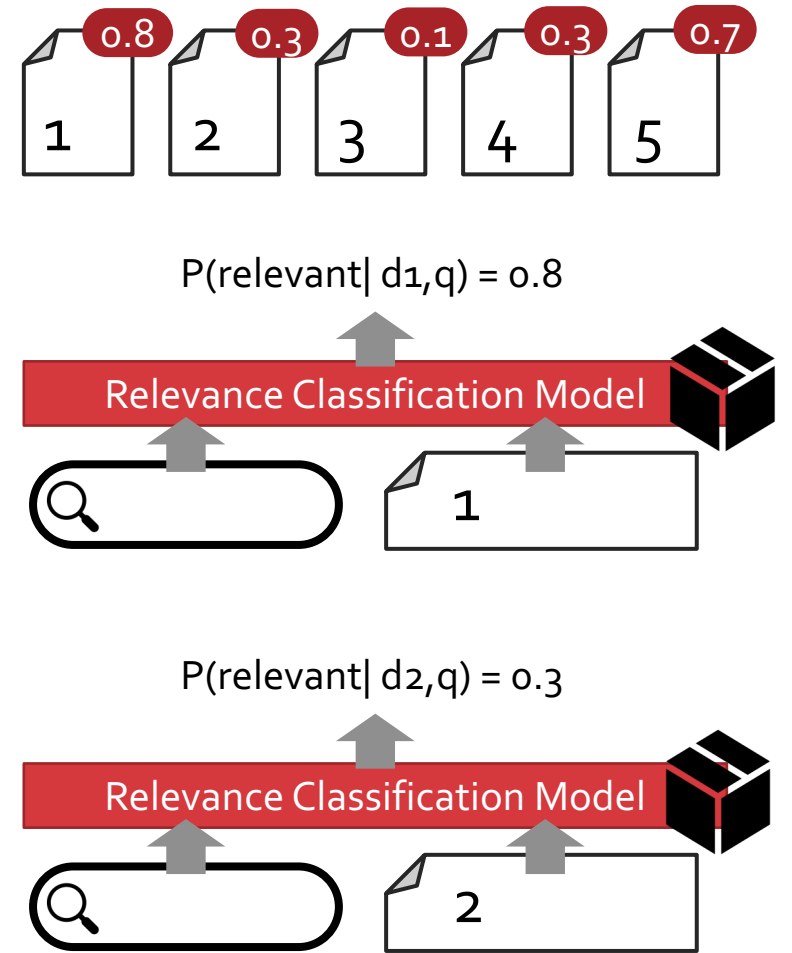


# Current Retrieval Systems

## Dense Retrieval

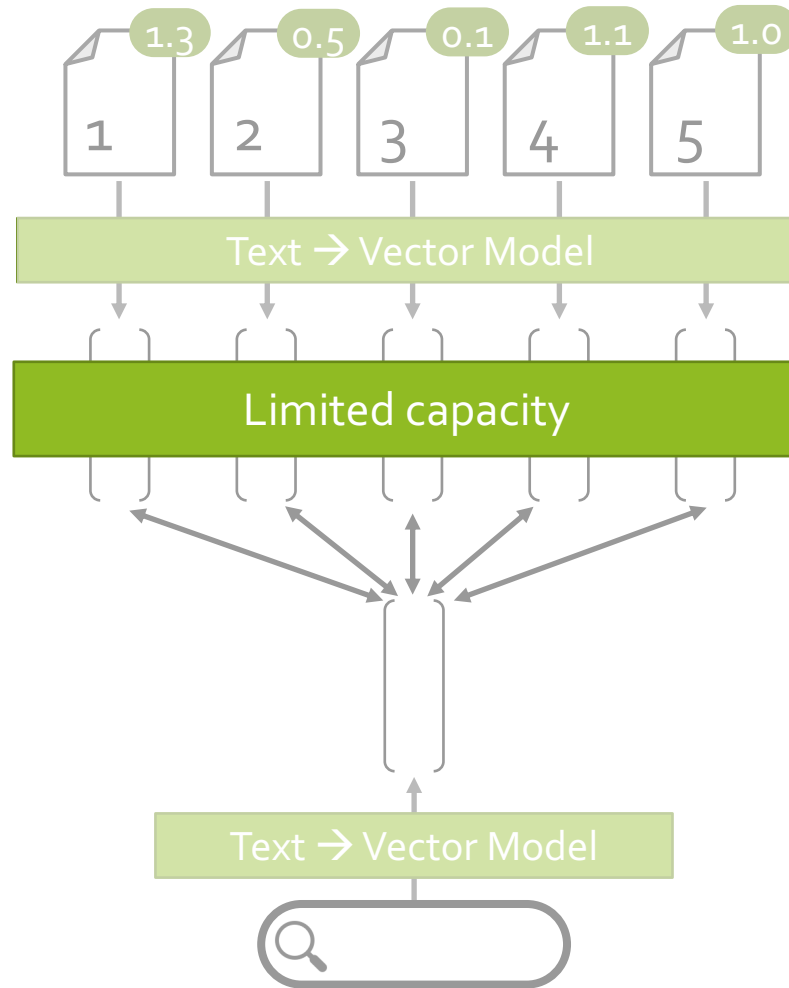


## Cross-Encoders

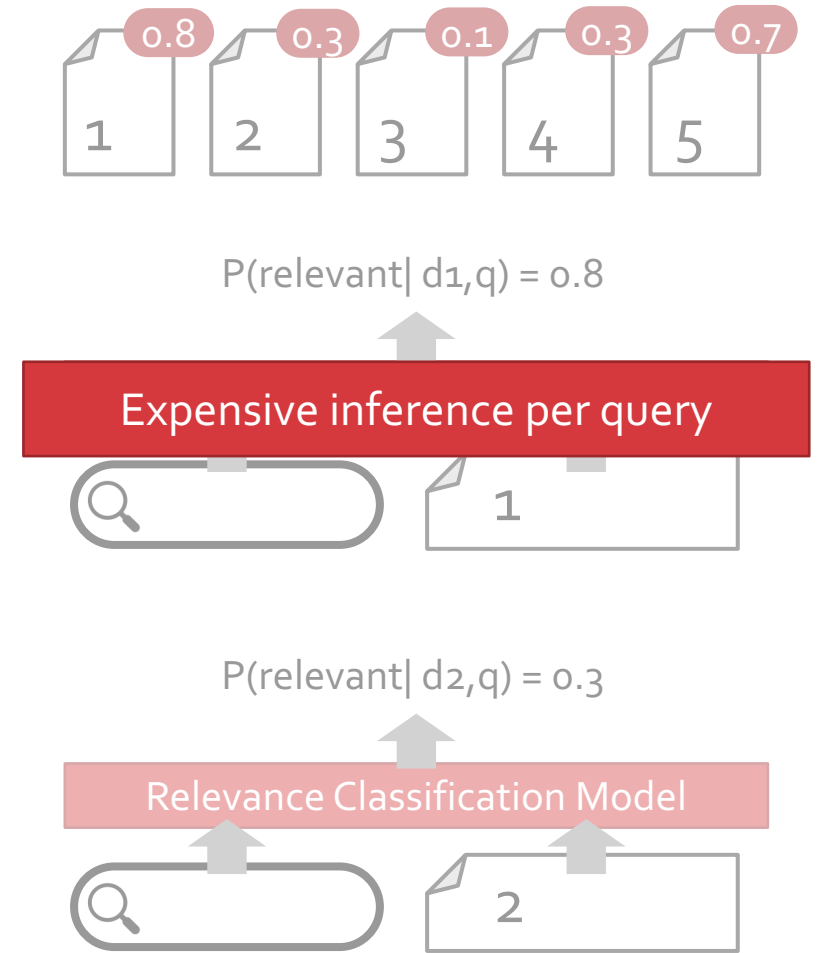


# Current Retrieval Systems

## Dense Retrieval

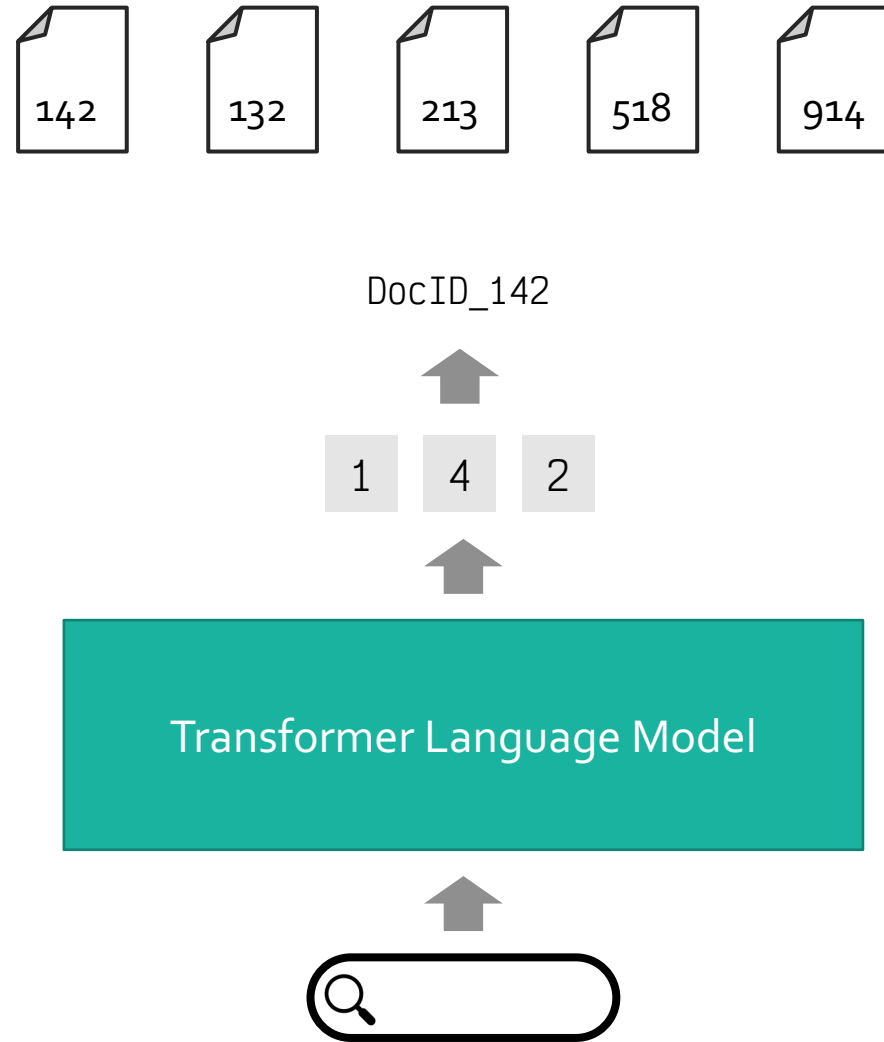


## Cross-Encoders

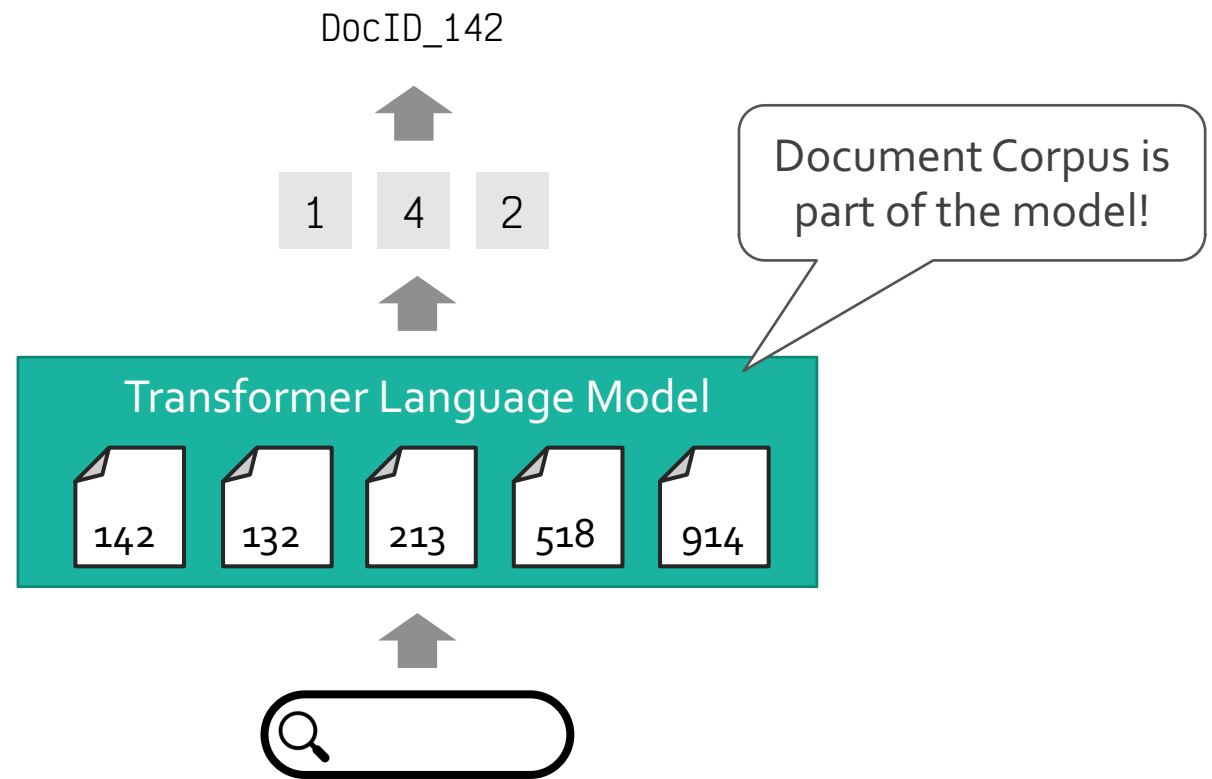


# Generative Information Retrieval

# Generative Information Retrieval



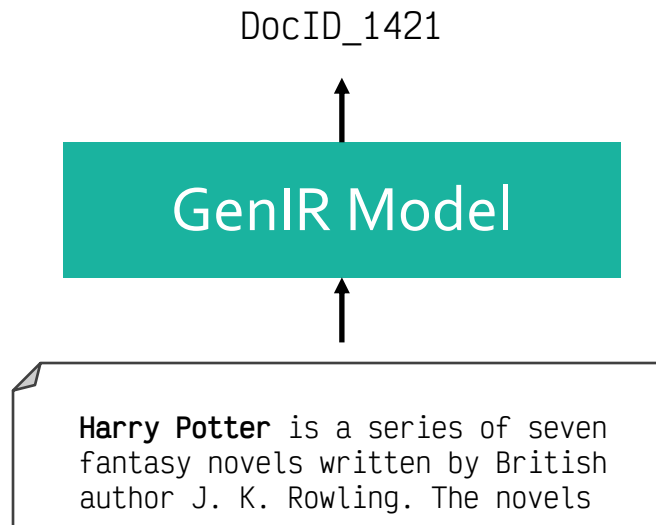
# Generative Information Retrieval



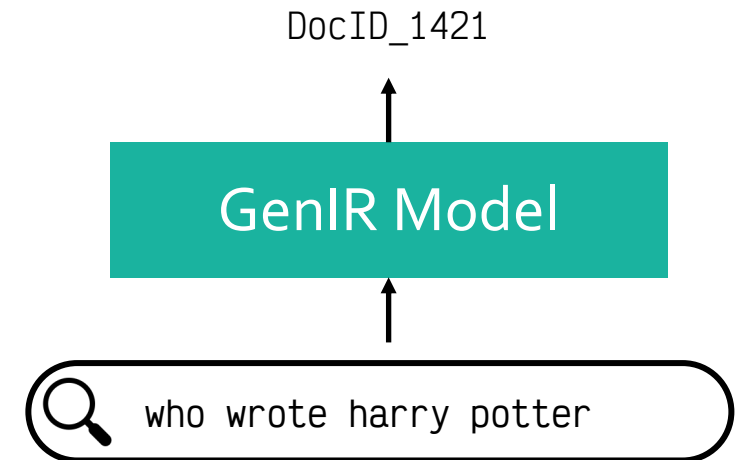
# GenIR Training

Goal: Associate the content of a document with its ID

*Indexing Phase*



*Retrieval Phase*

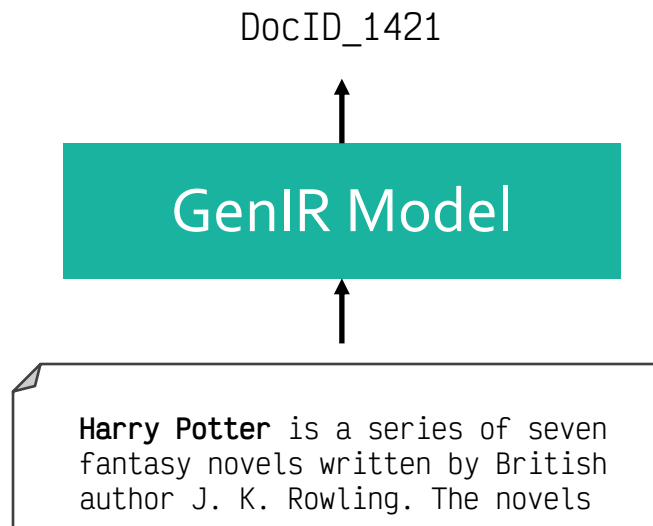


# GenIR Training

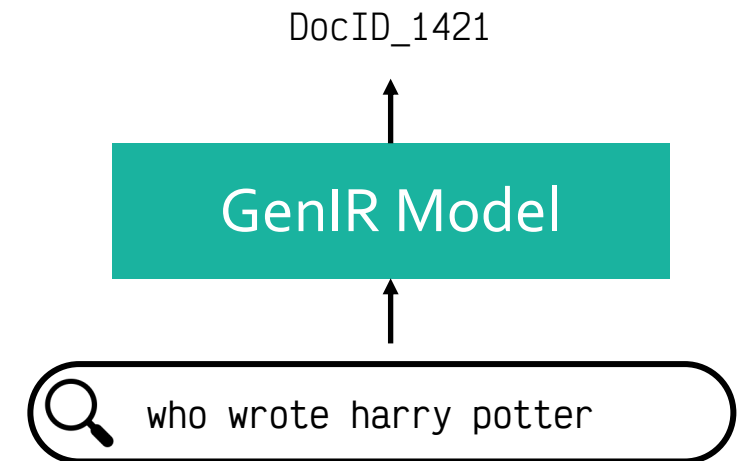
Goal: Associate the content of a document with its ID

Problem: Corpus is "hard-coded" as part of the model.

*Indexing Phase*

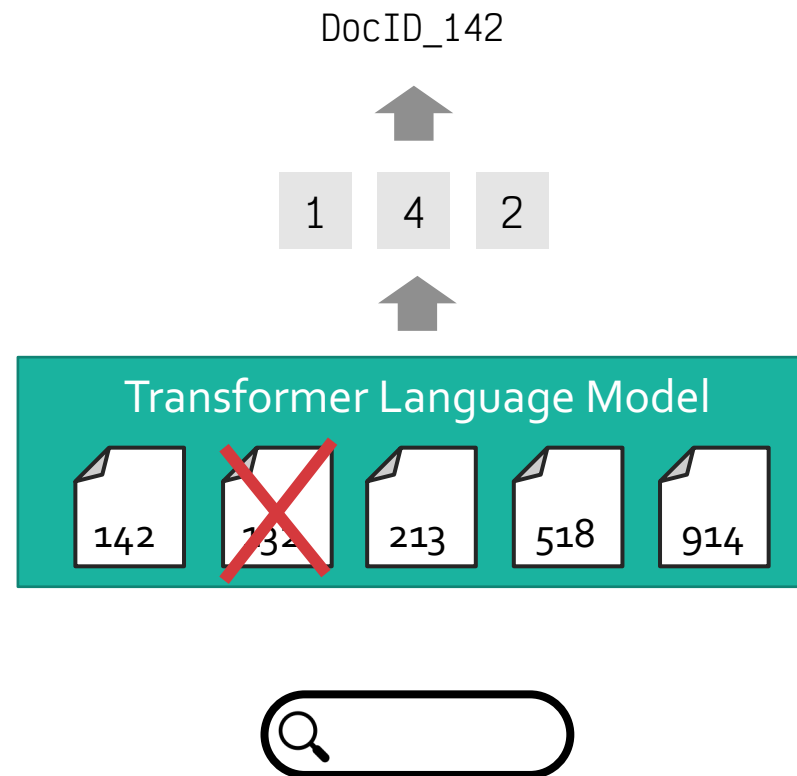


*Retrieval Phase*

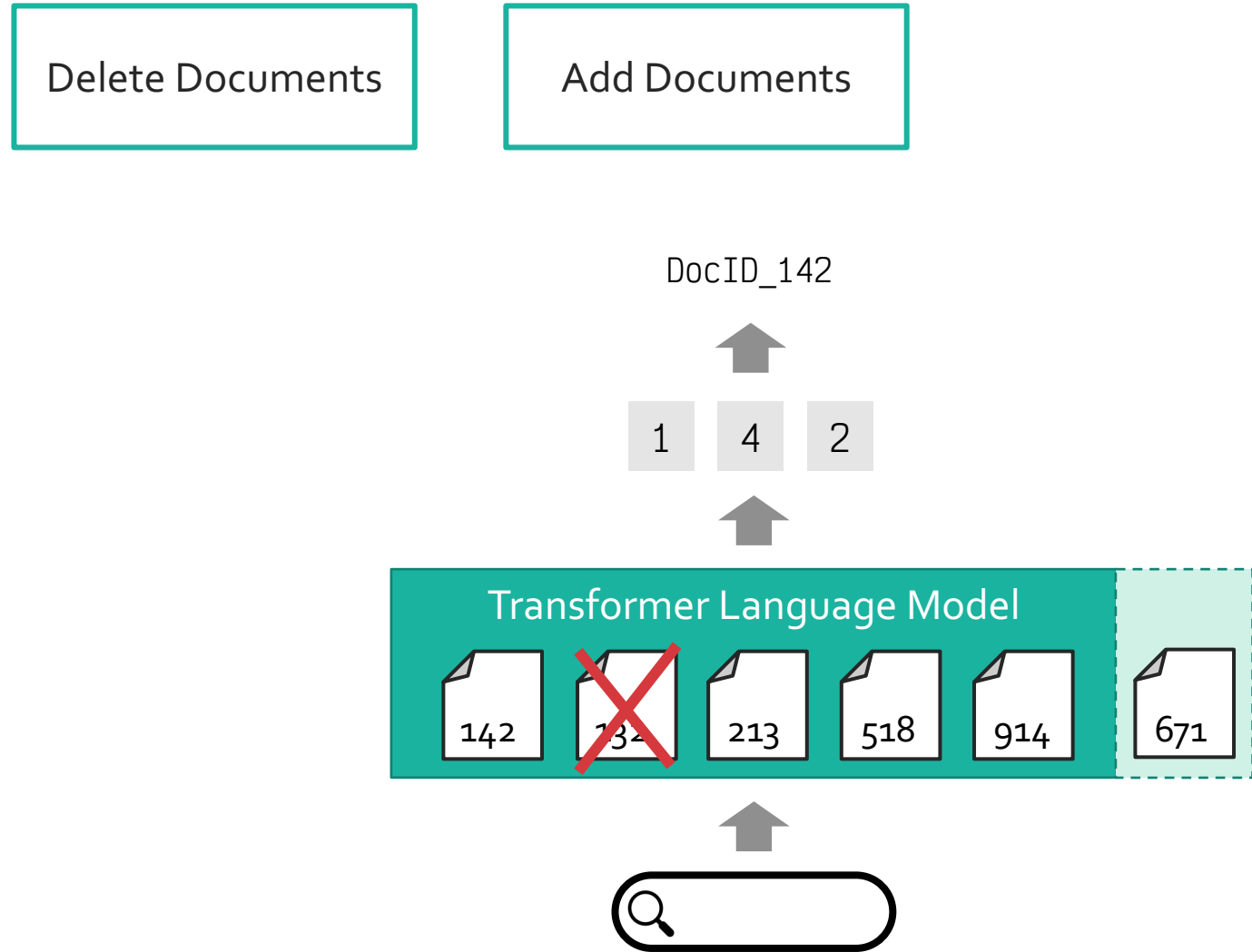


# The Problem

Delete Documents



# The Problem



# The Problem

Delete Documents

Add Documents

Change Documents

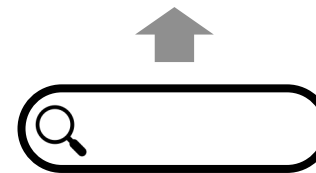
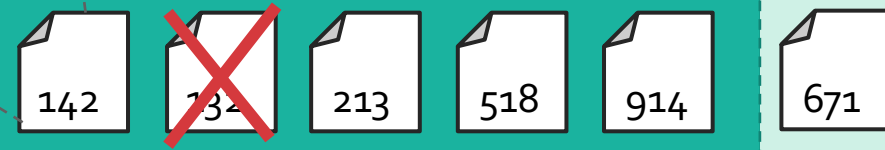
The state-of-the-art for robotic movement is based on ~~diffusion~~.

*VLAs*

DocID\_142



Transformer Language Model



# My Research

Delete Documents

Add Documents

Change Documents

Goal: Edit the corpus without re-training the model

## Challenges

- Find out how these models work.
- Develop editing method.

# My Research

Delete Documents

Add Documents

Change Documents

Goal: Edit the corpus without re-training the model

## Challenges

- Find out how these models work.
- Develop editing method.

# Reverse-Engineering the Retrieval Process in GenIR Models

Anja Reusch & Yonatan Belinkov – SIGIR 2025



*Funding:*



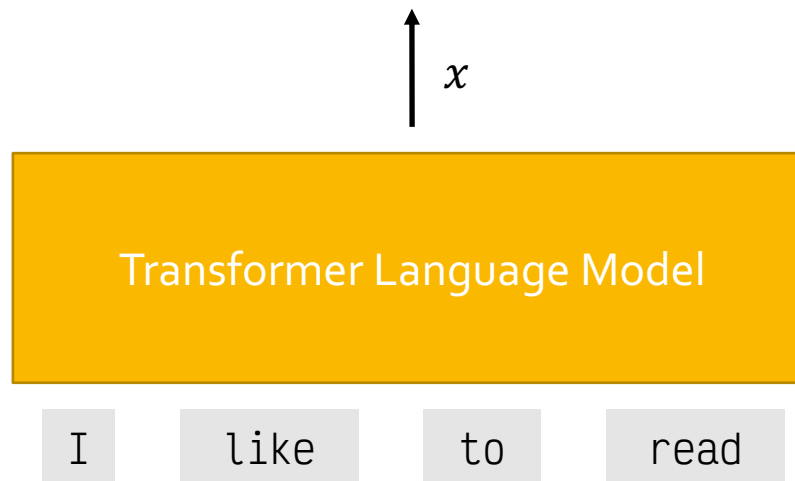
Funded by  
the European Union



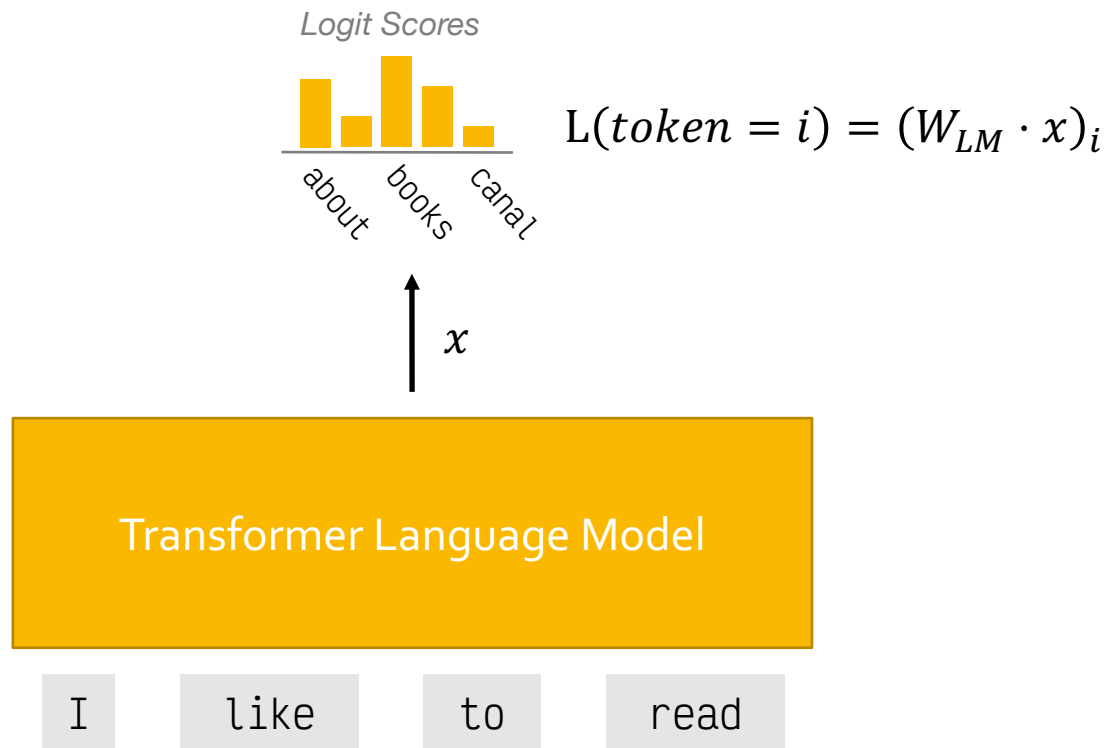
European Research Council  
Established by the European Commission



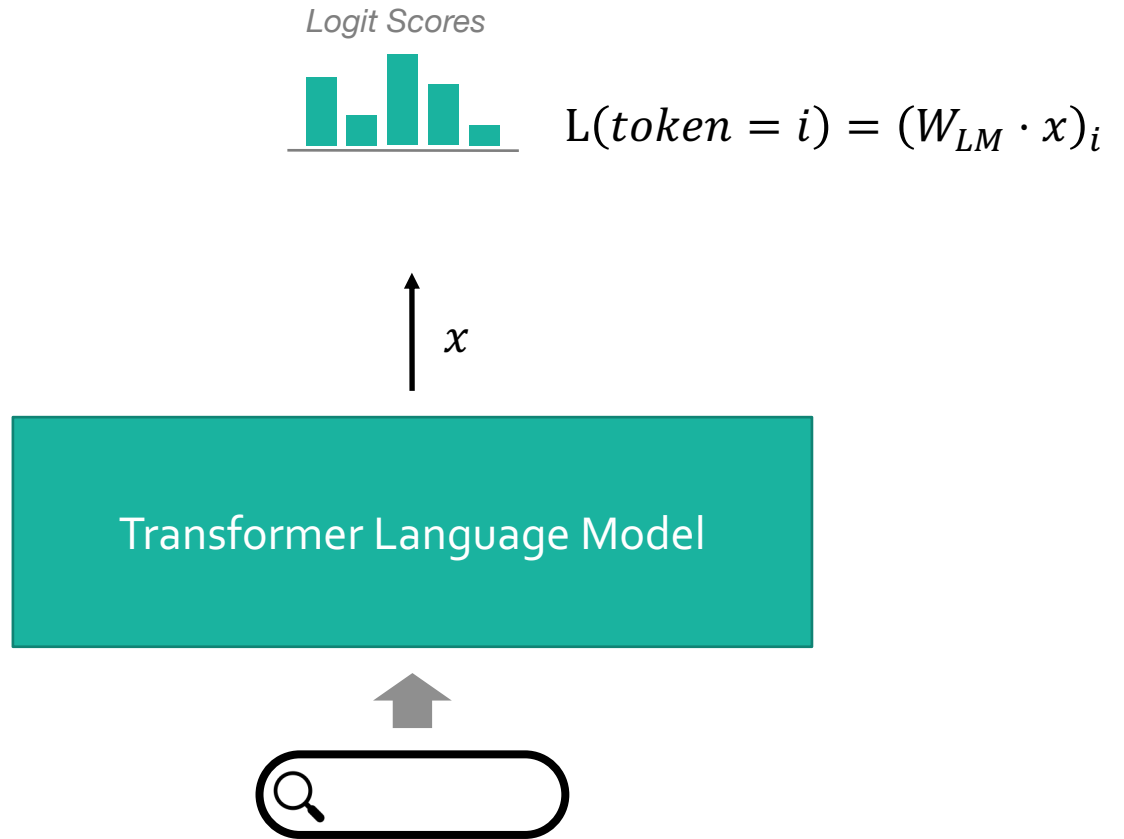
# Generative Models



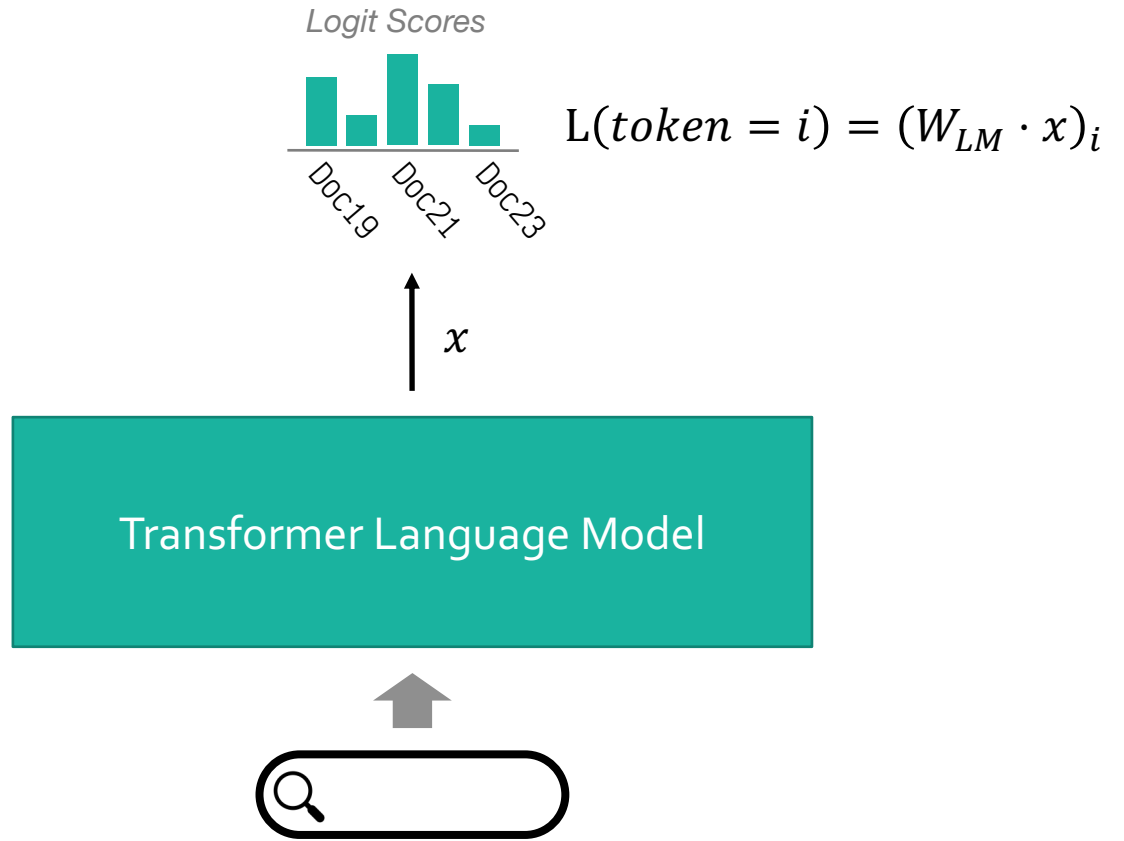
# Generative Models



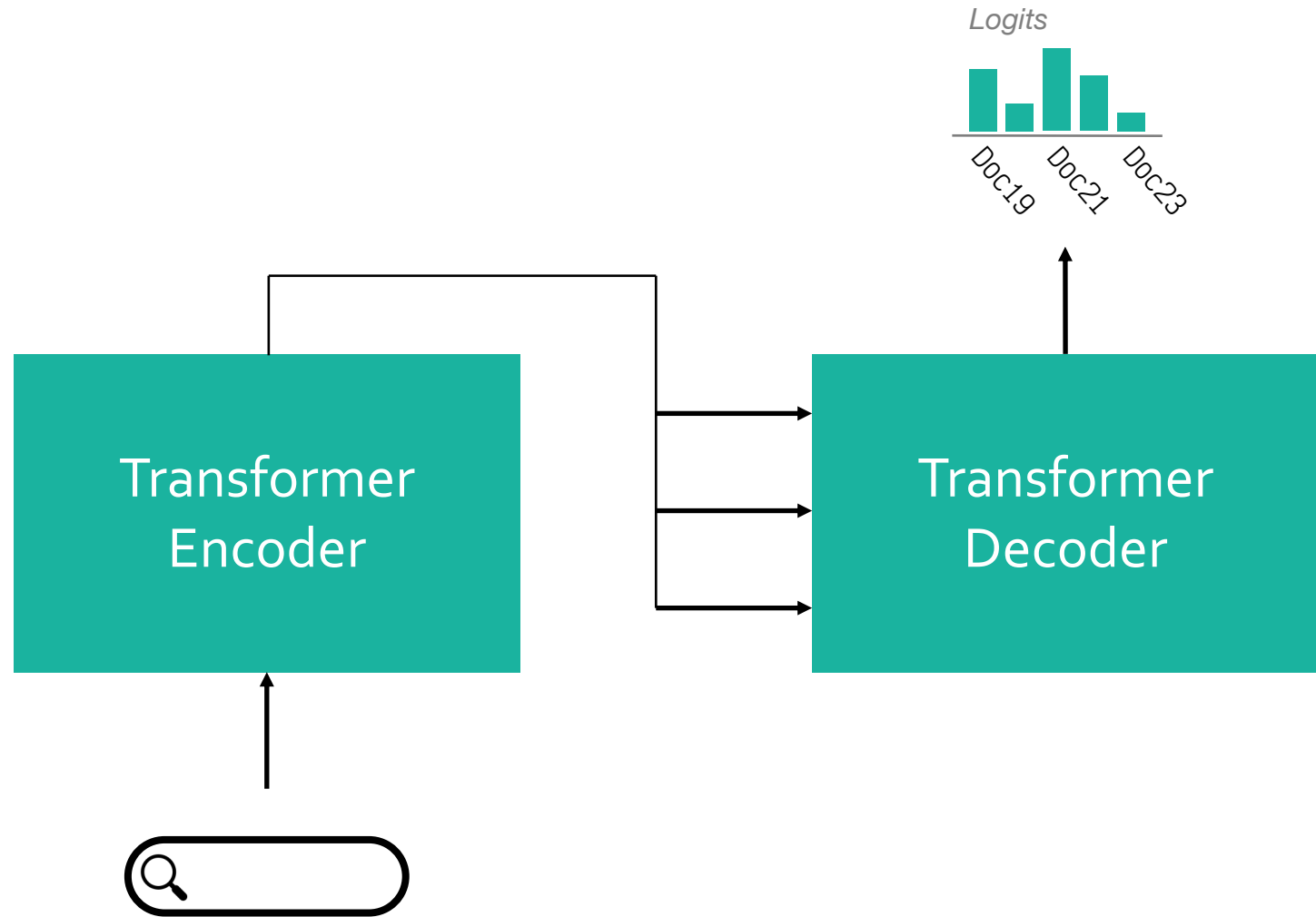
# Generative Information Retrieval



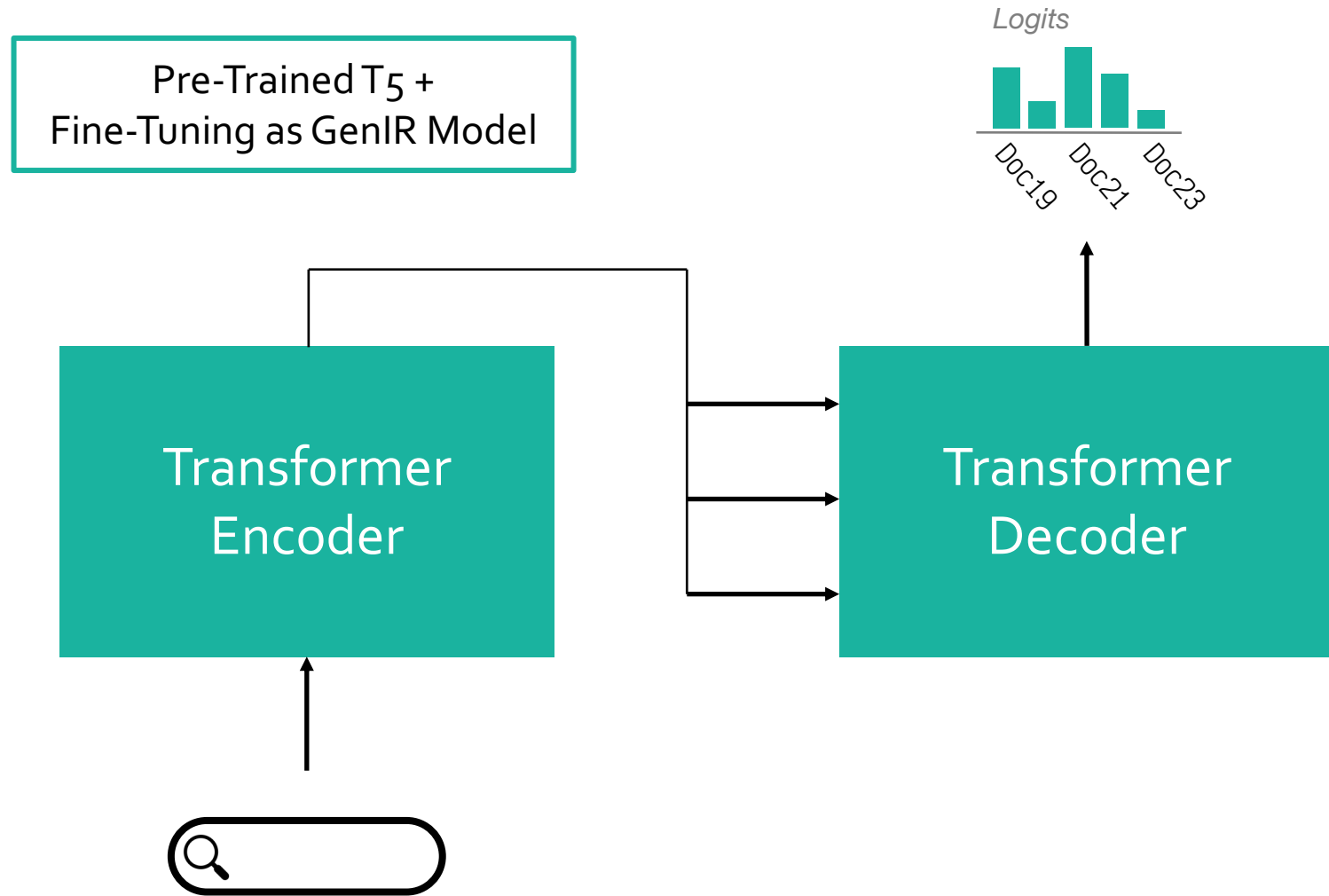
# Generative Information Retrieval



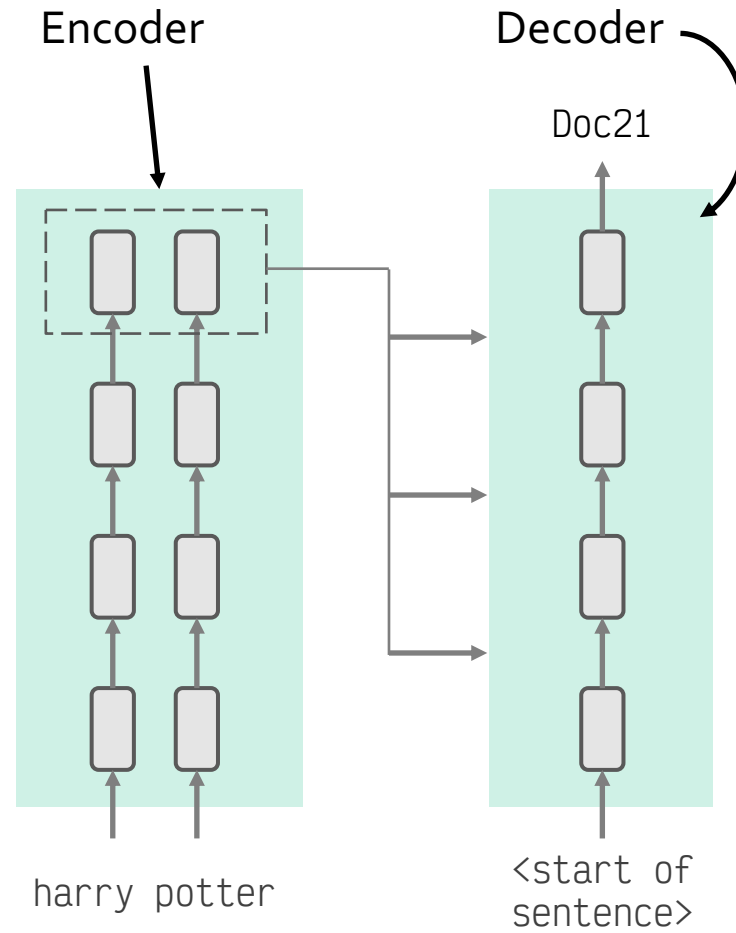
# Generative Encoder-Decoder



# Generative Encoder-Decoder

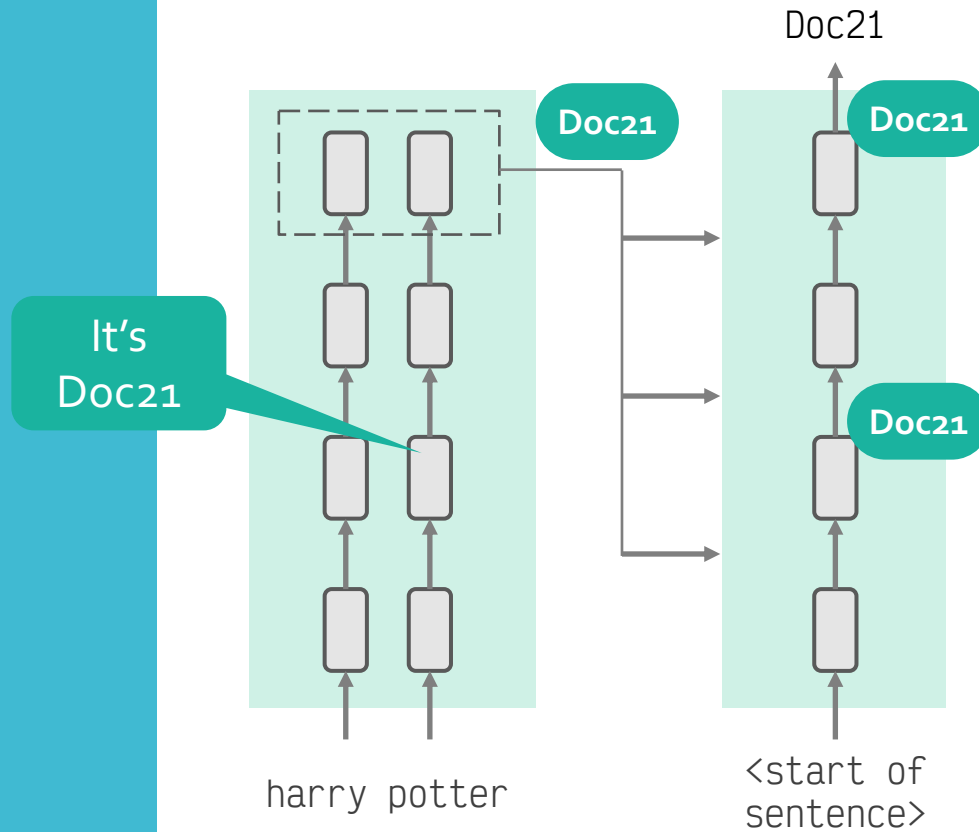


Which part identifies the document?



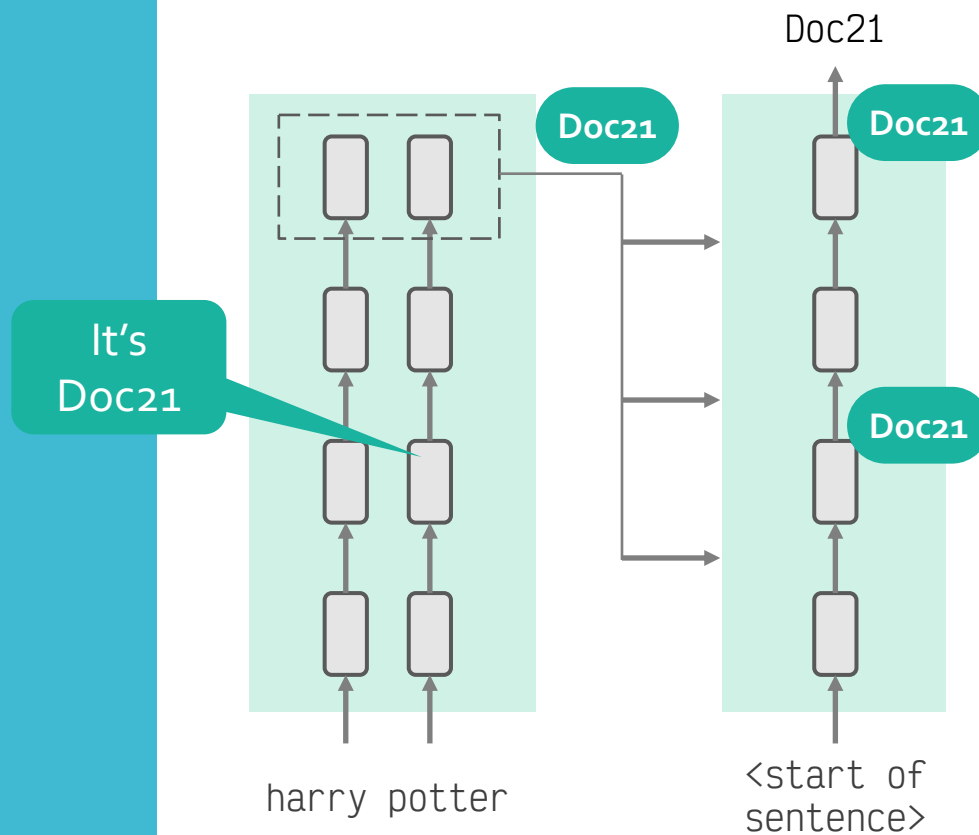
# Two Hypotheses

Encoder identifies the document

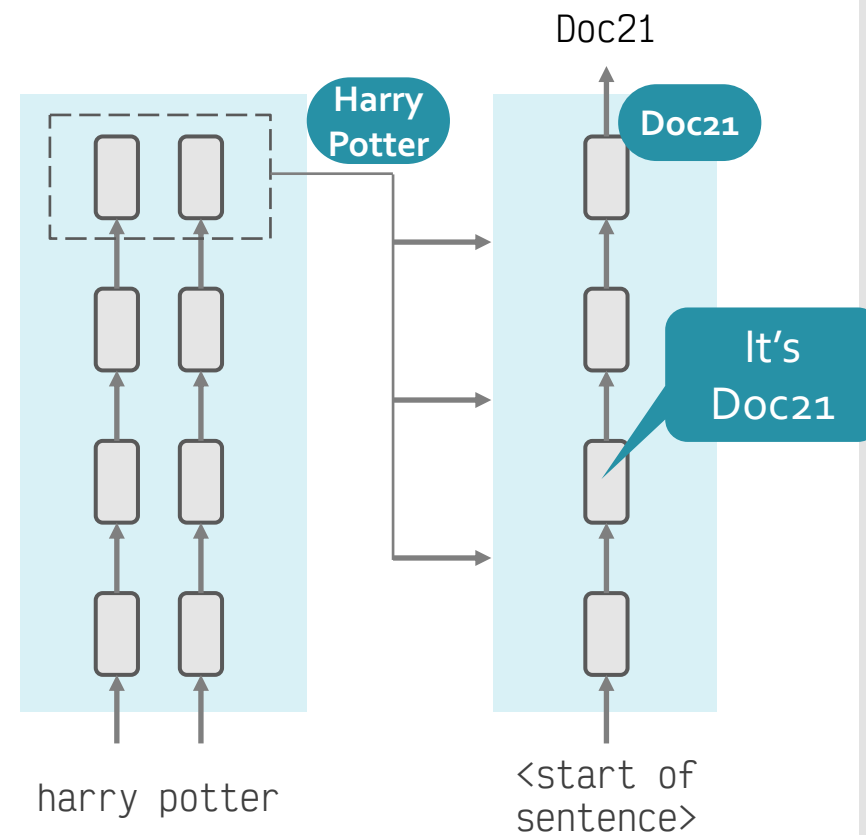


# Two Hypotheses

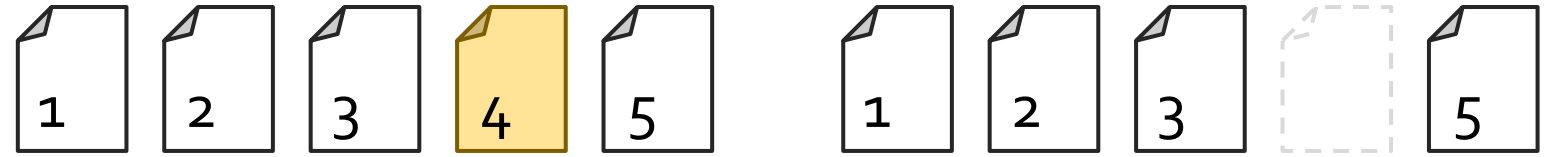
Encoder identifies the document



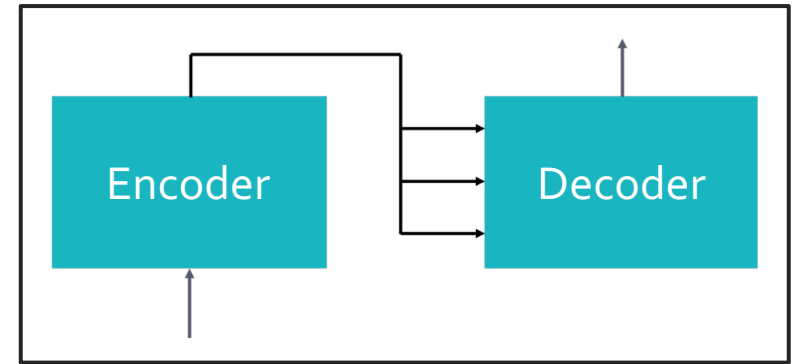
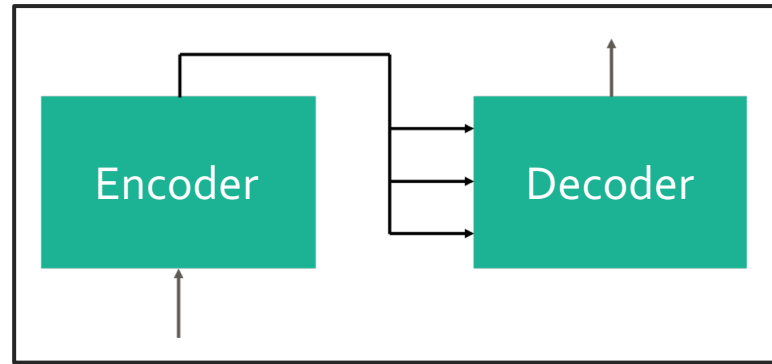
Decoder identifies the document



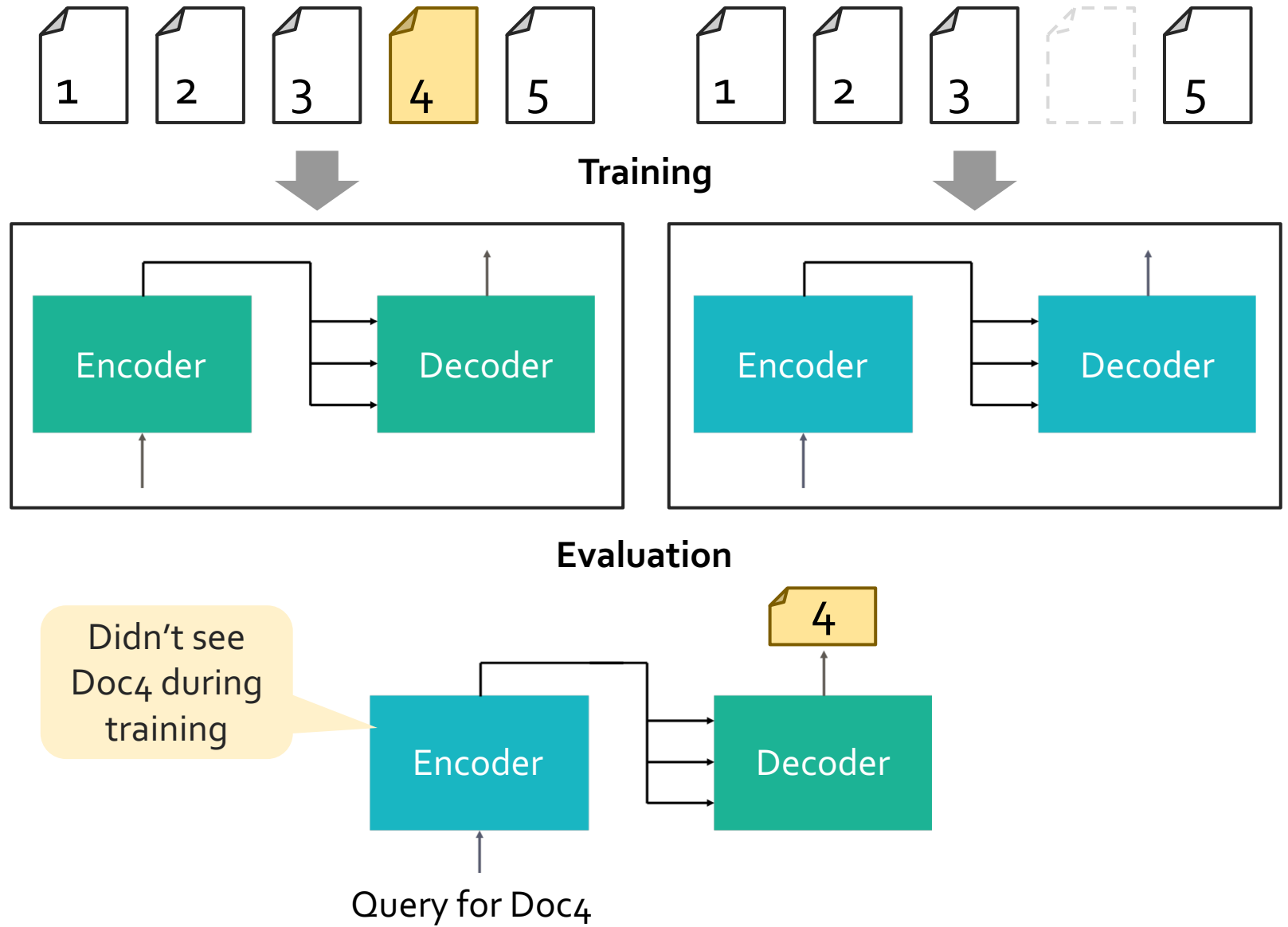
# Looking at the Encoder



Training

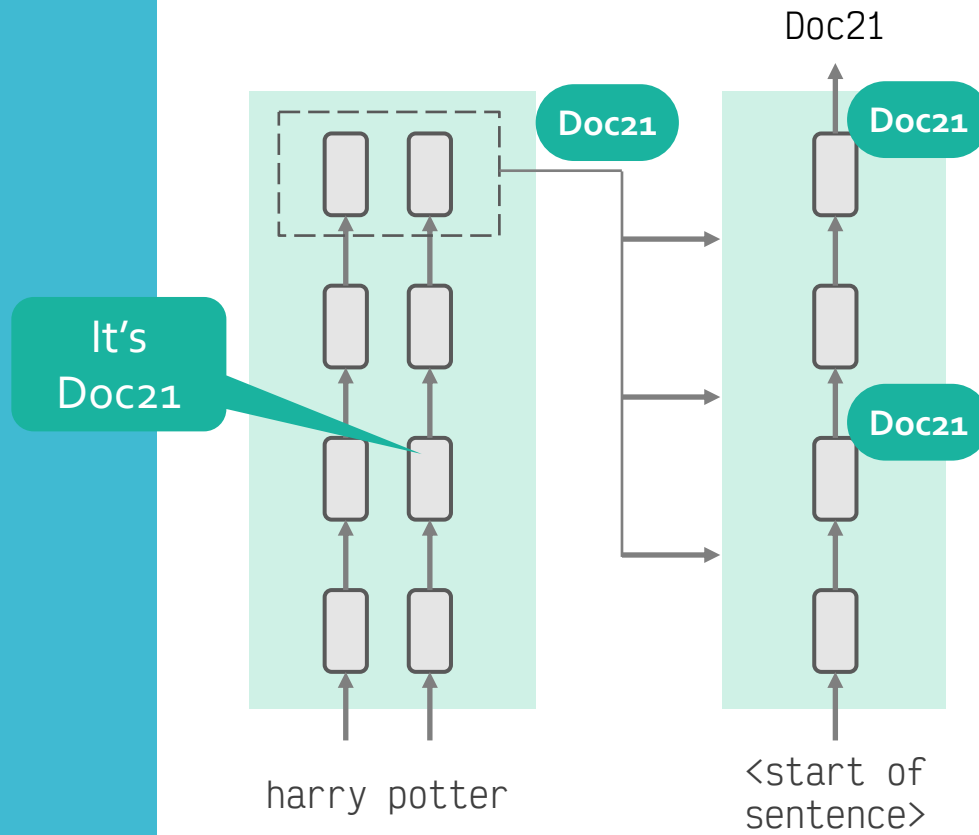


# Looking at the Encoder



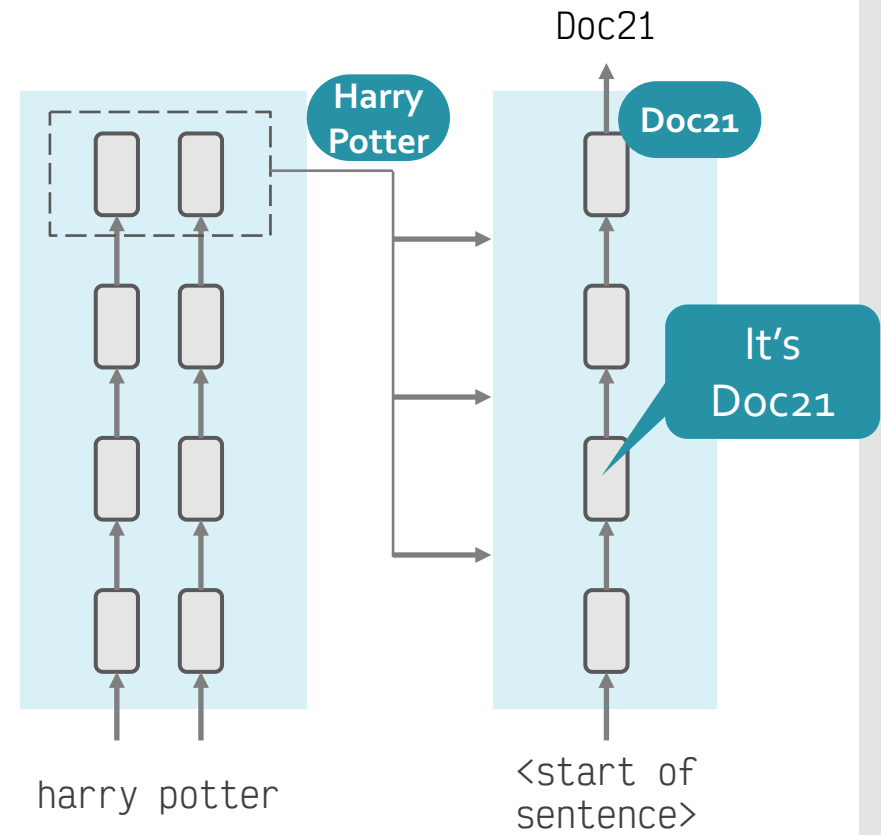
# Two Hypotheses

Encoder identifies the document



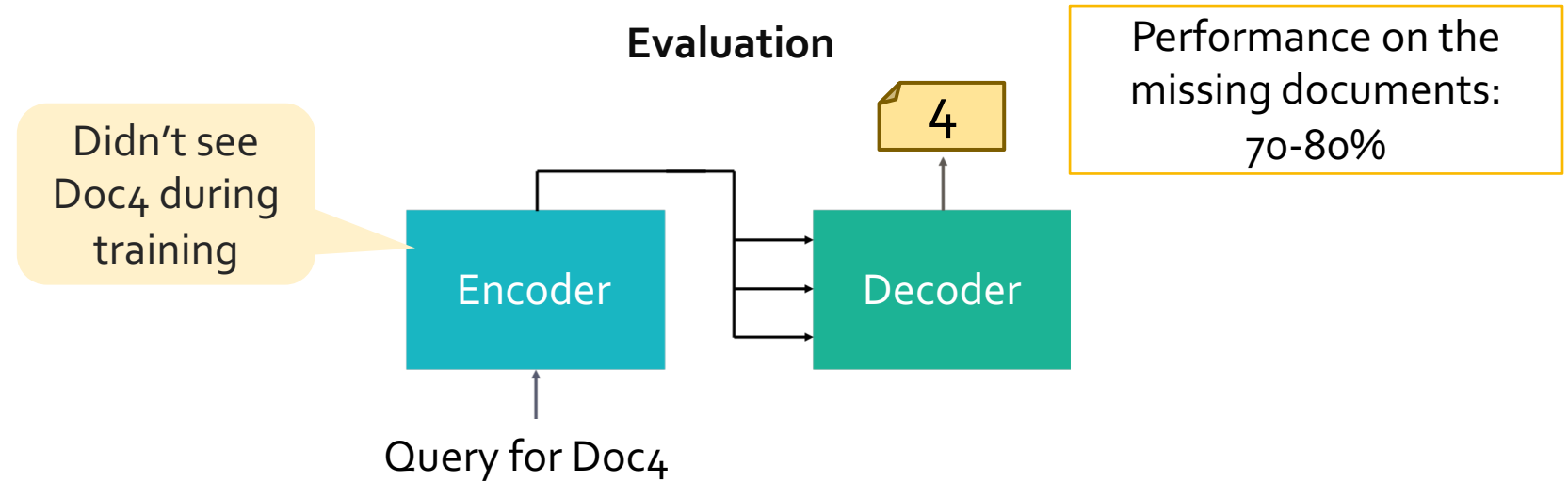
Expectation: 0% Accuracy

Decoder identifies the document



Expectation: ~High Accuracy

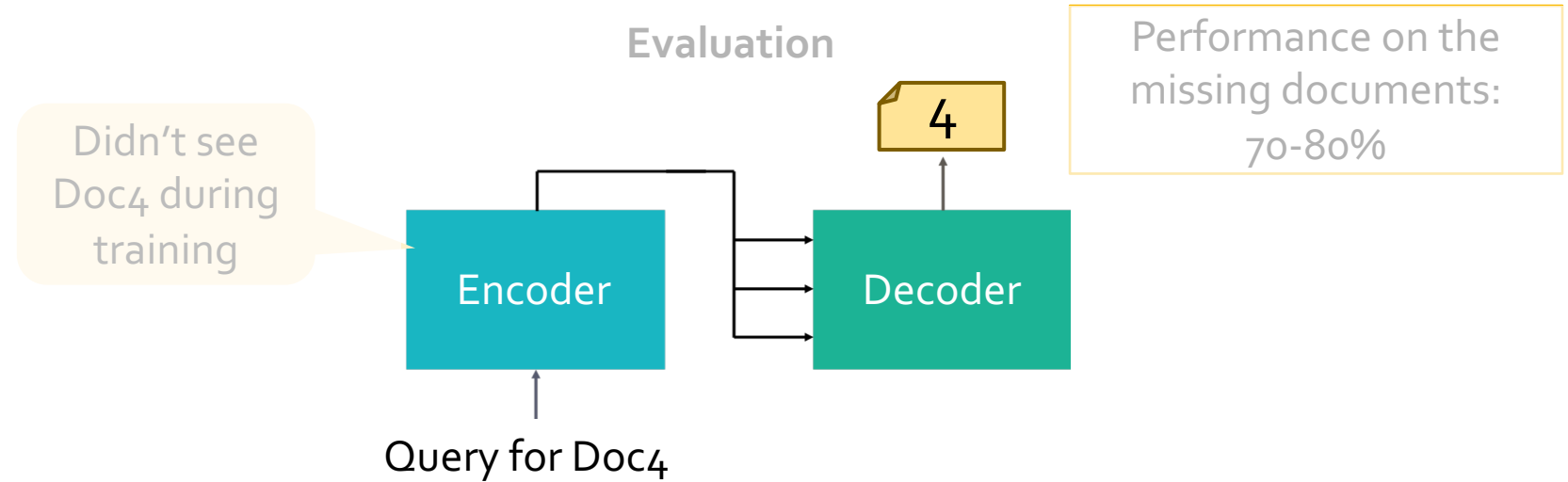
# Looking at the Encoder



## Setup

- Base Model: Trained on 10k documents
- 4x "Missing Docs Model": Trained on same 9990 documents
- Hits@1 (Correct doc at Position 1): 70-80%

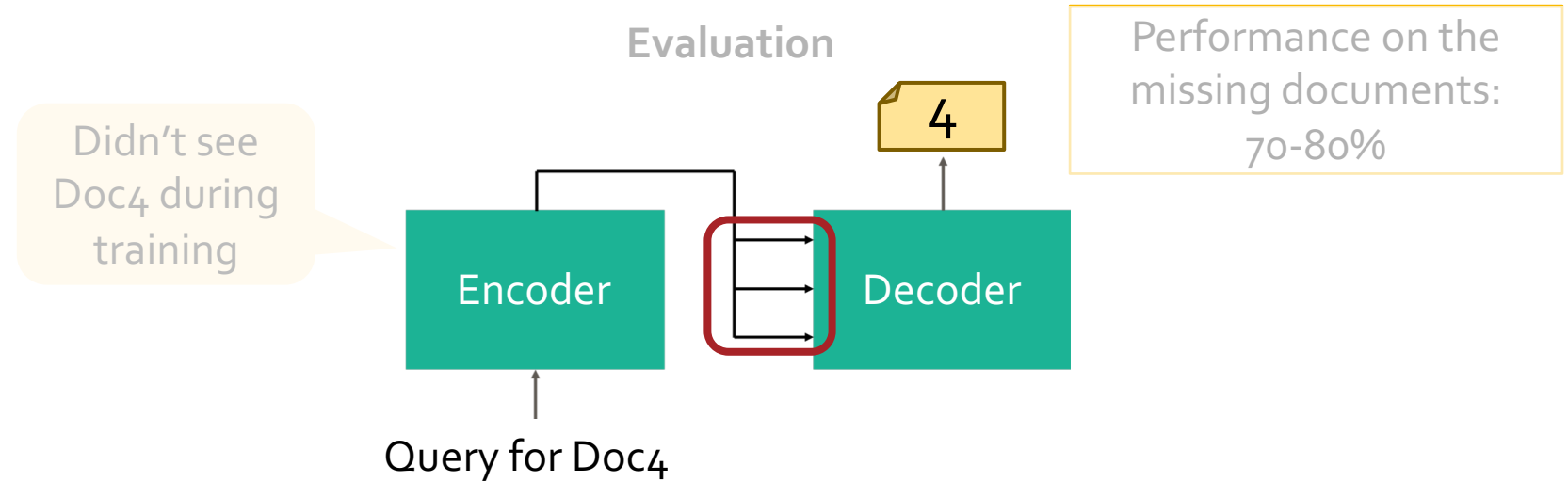
# Looking at the Encoder



**Conclusion:** The decoder identifies the document.

- Base Model: Trained on 10k documents
- 4x "Missing Docs Model": Trained on same 9990 documents
- Hits@<sub>1</sub> (Correct doc at Position 1): 70-80%

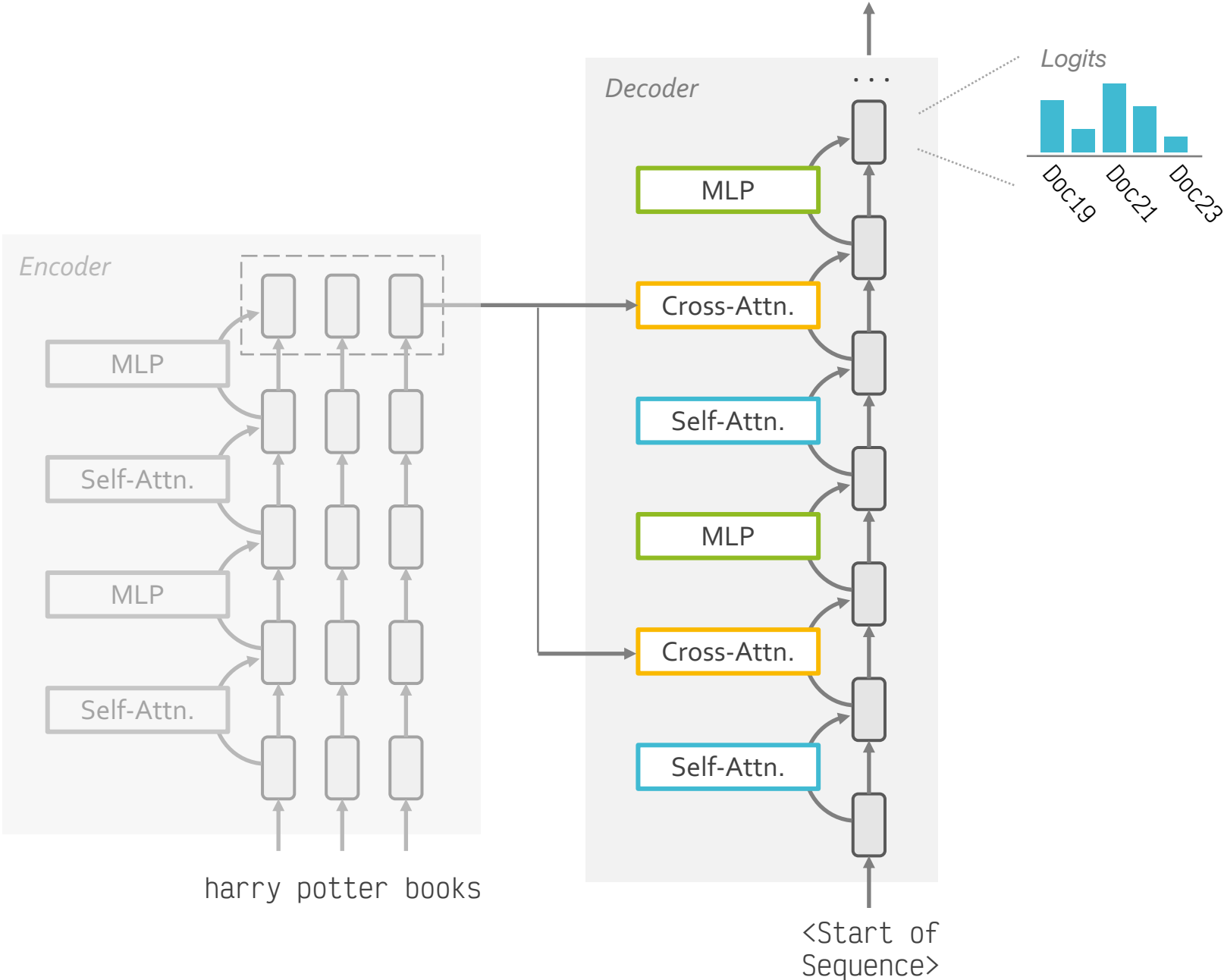
# Looking at the Encoder



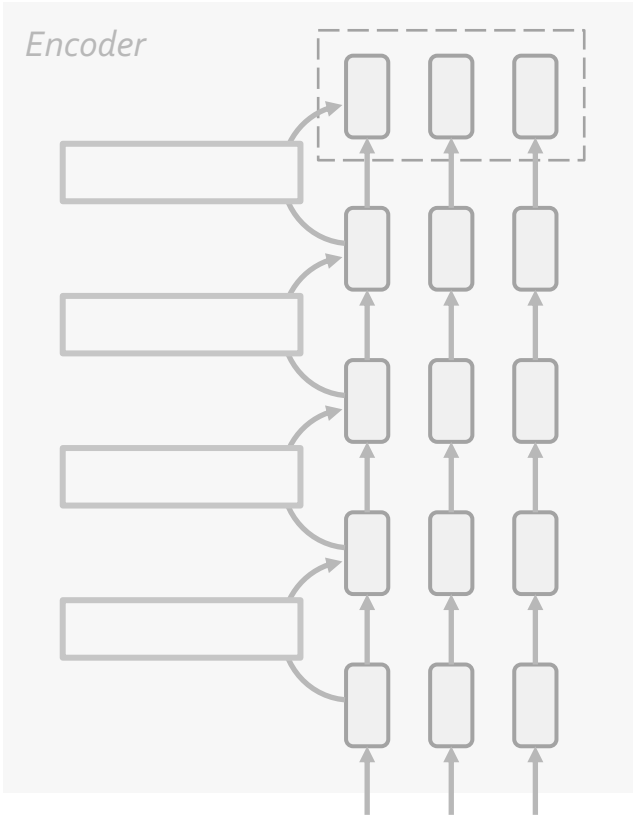
## Setup

- Base Model: Trained on 10k documents
- 4x "Missing Docs Model": Trained on same 9990 documents
- Hits@<sub>1</sub> (Correct doc at Position 1): 70-80%

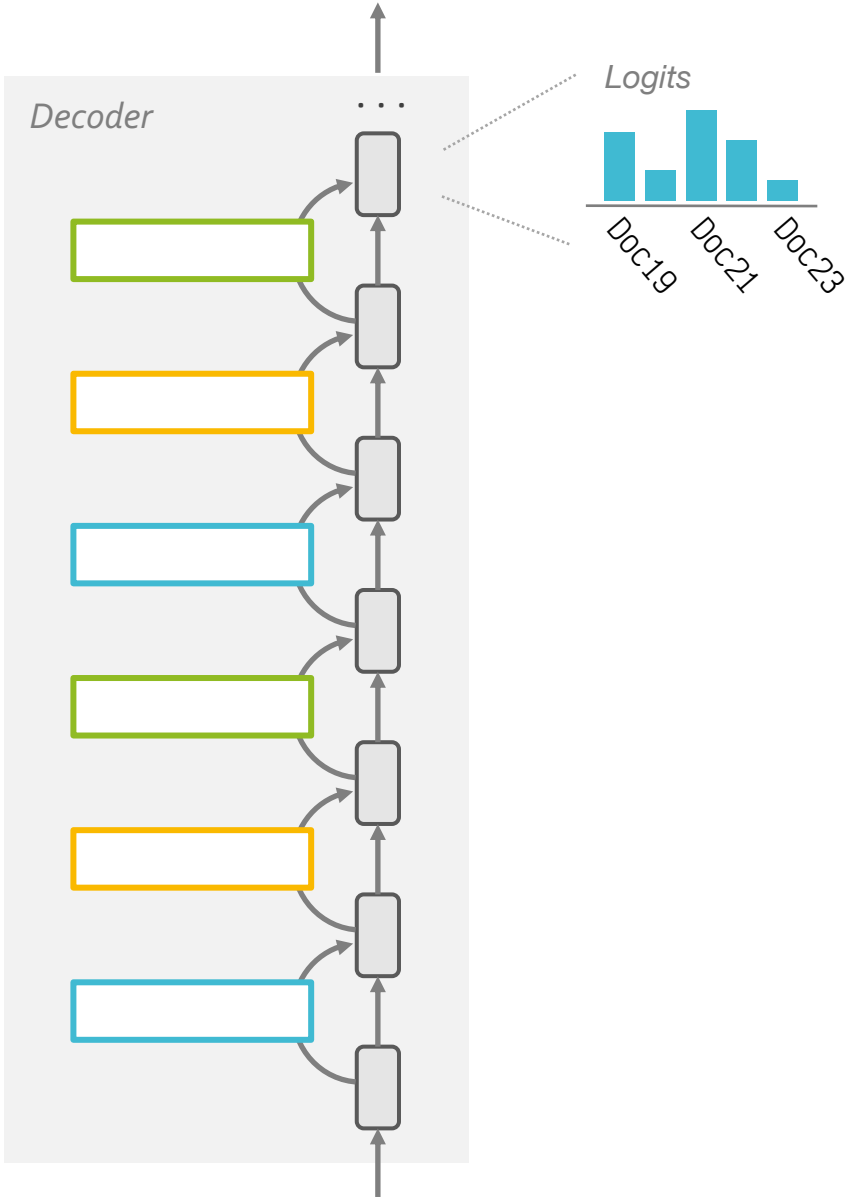
# How does a GenIR Model look like?



# How does a GenIR Model look like?

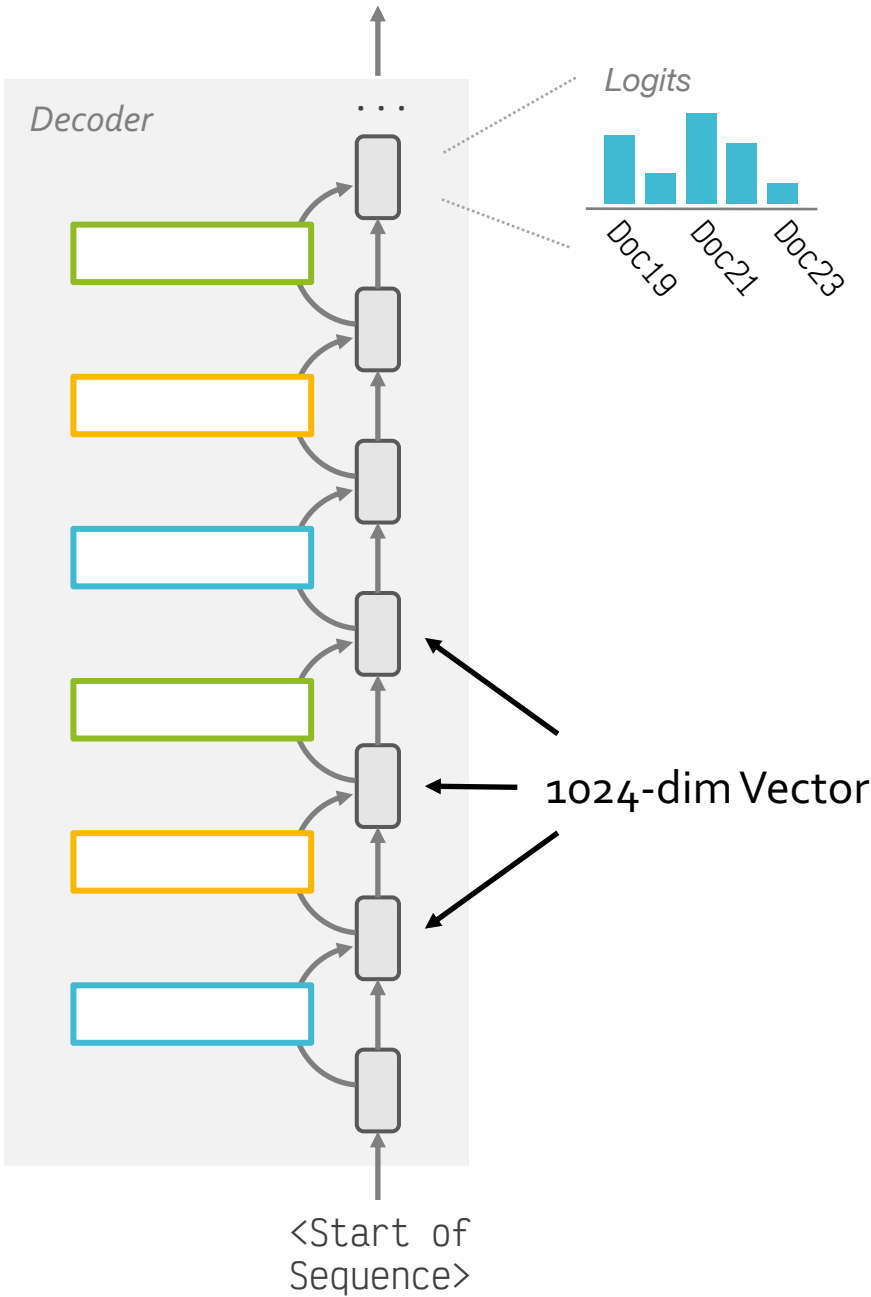
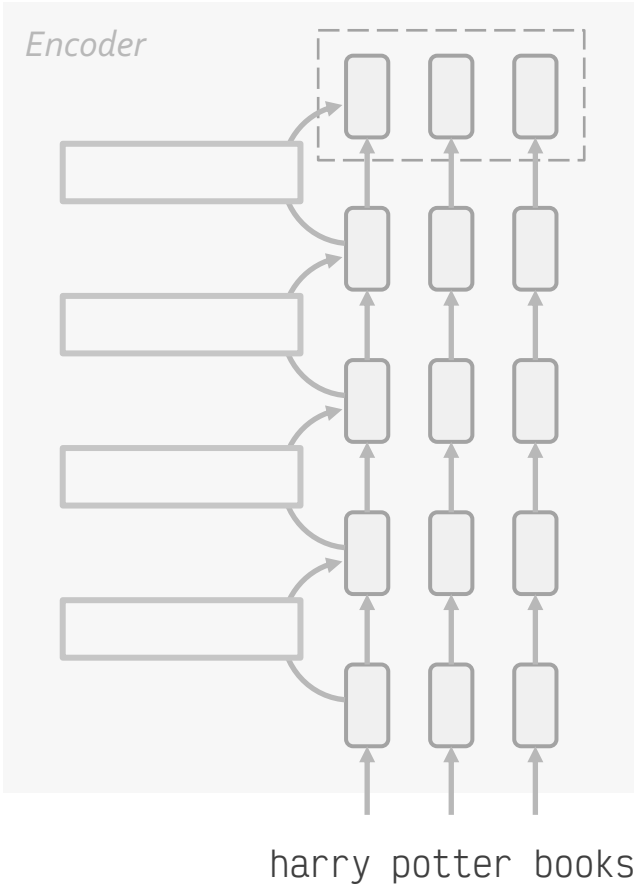


harry potter books

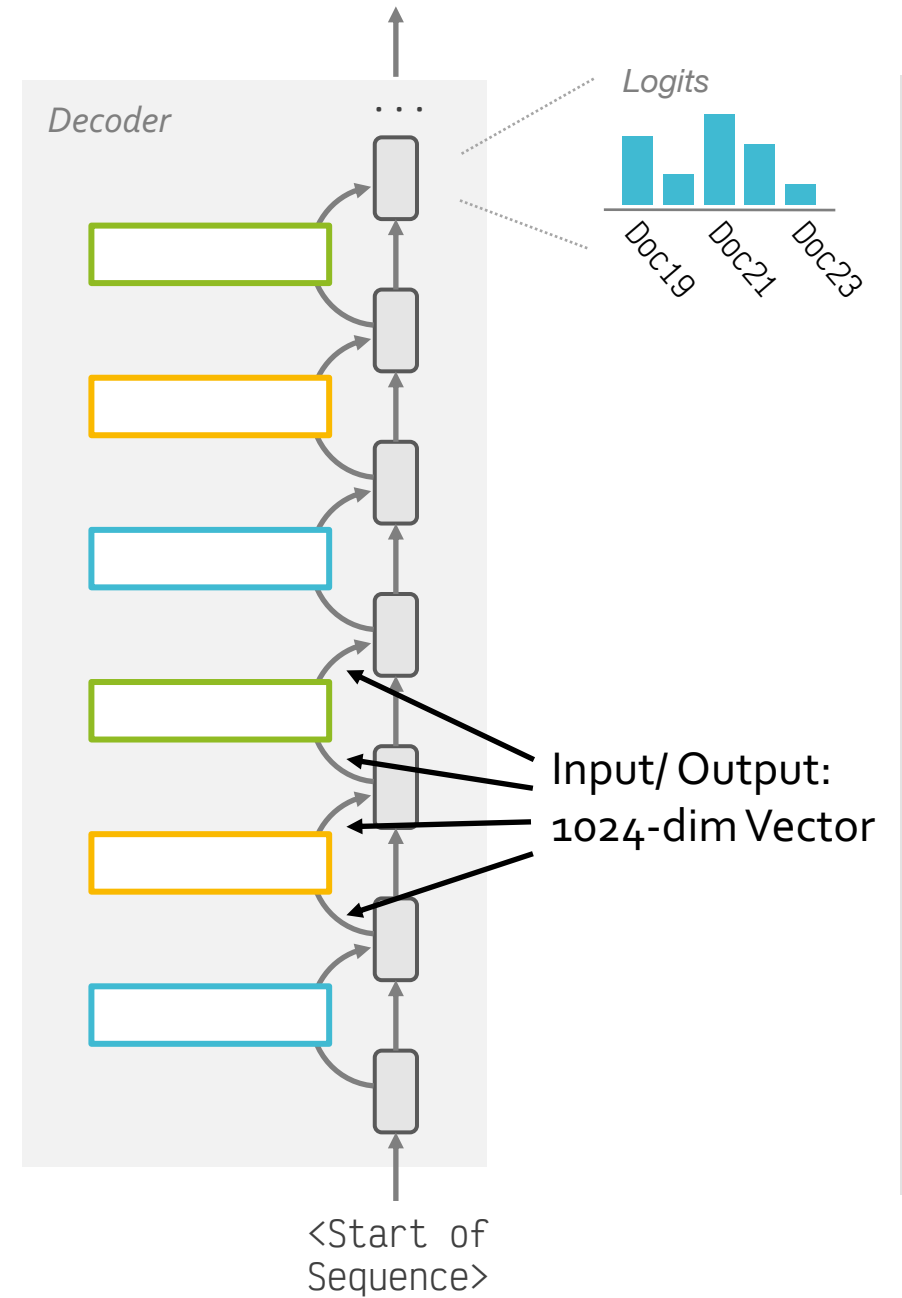
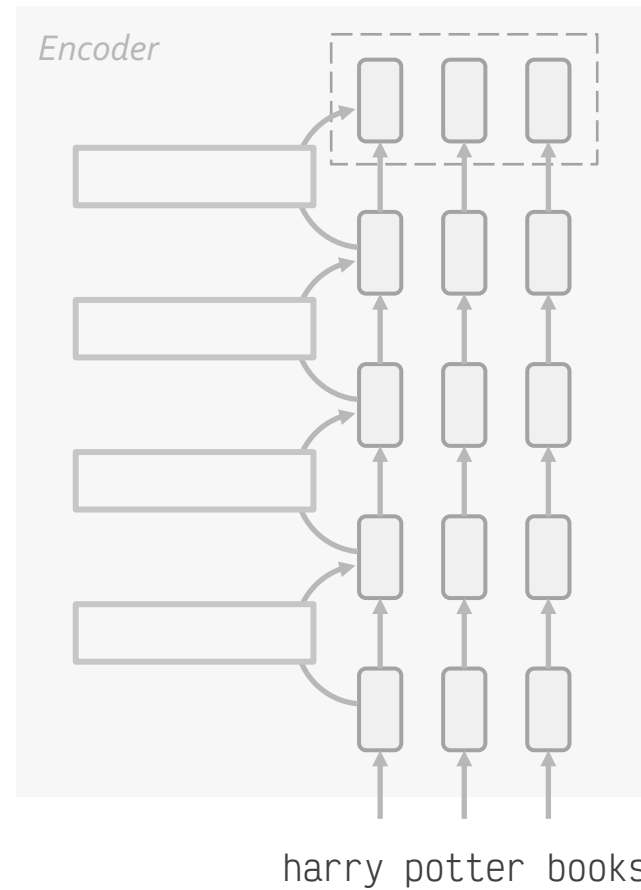


<Start of Sequence>

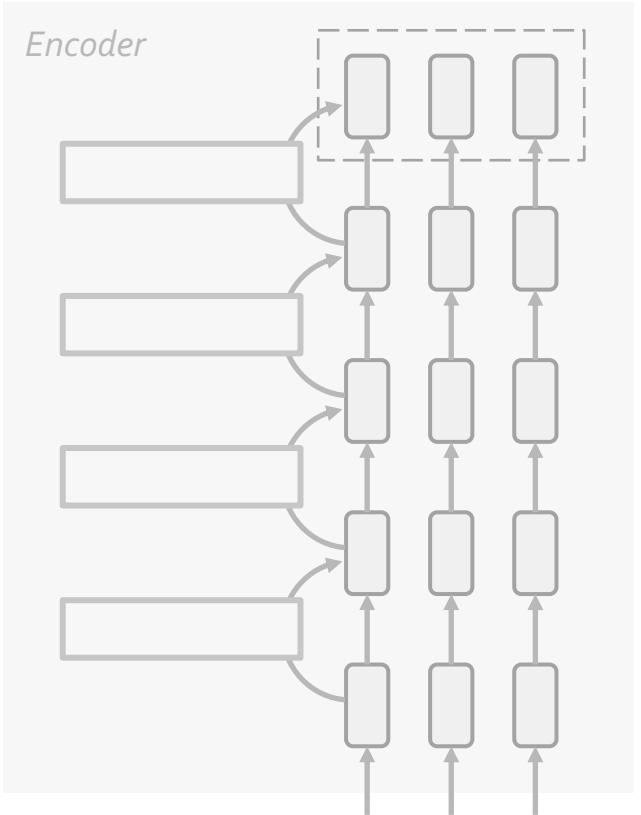
# How does a GenIR Model look like?



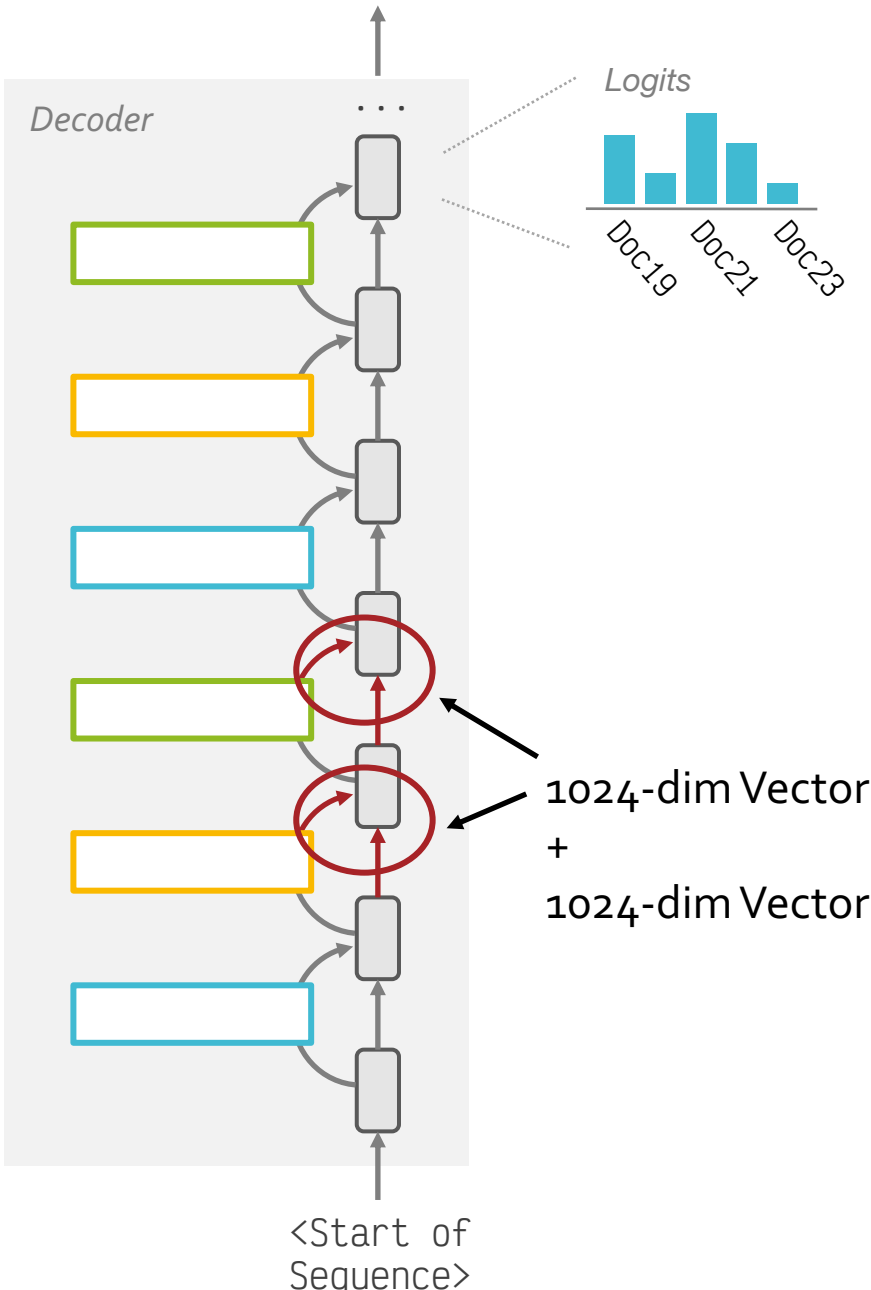
# How does a GenIR Model look like?



# How does a GenIR Model look like?



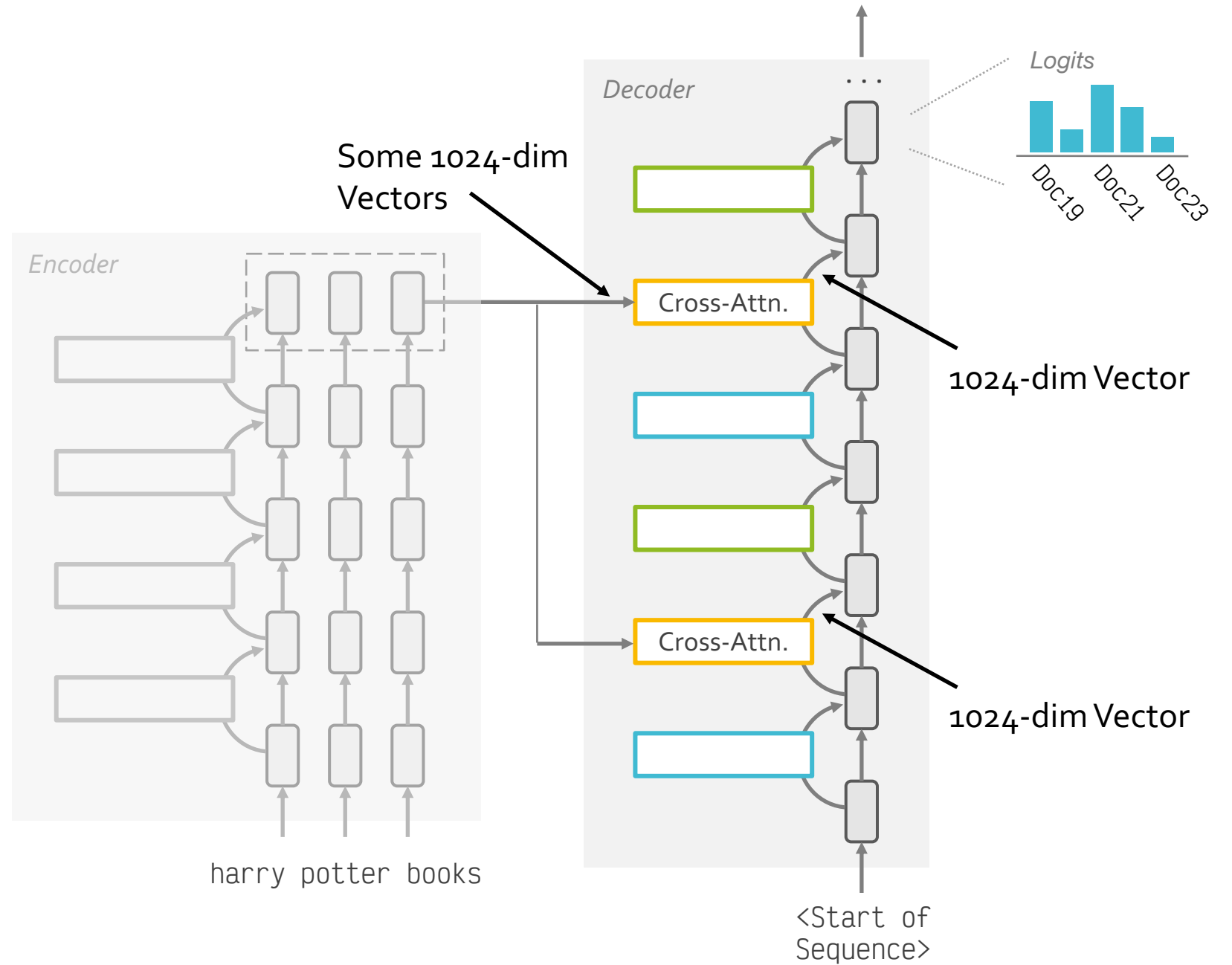
harry potter books



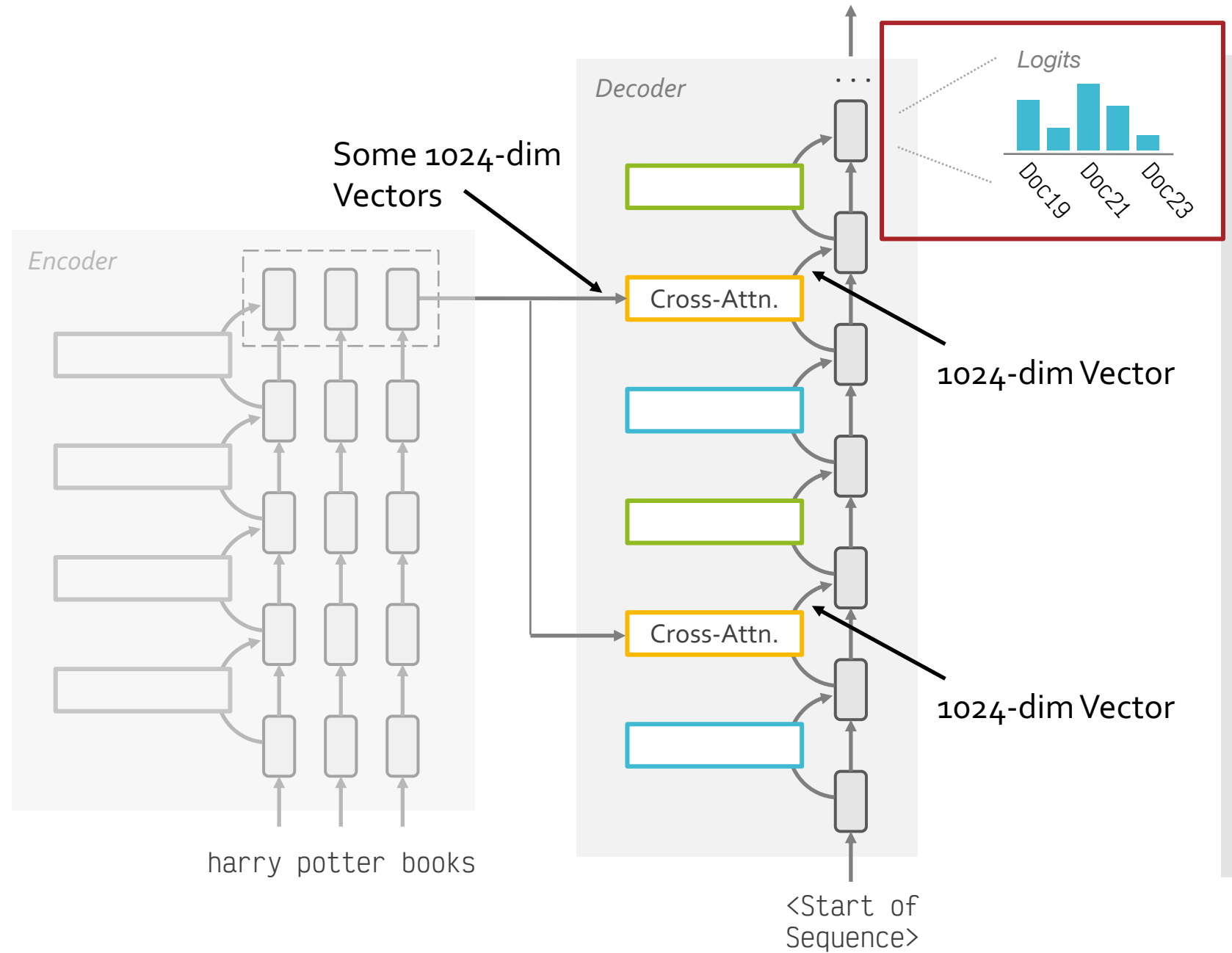
<Start of Sequence>

1024-dim Vector + 1024-dim Vector

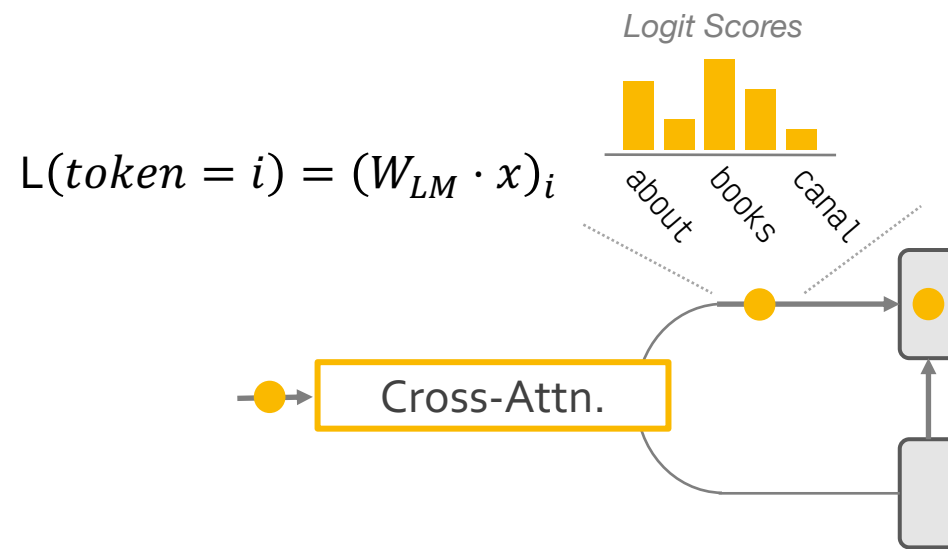
# How does a GenIR Model look like?



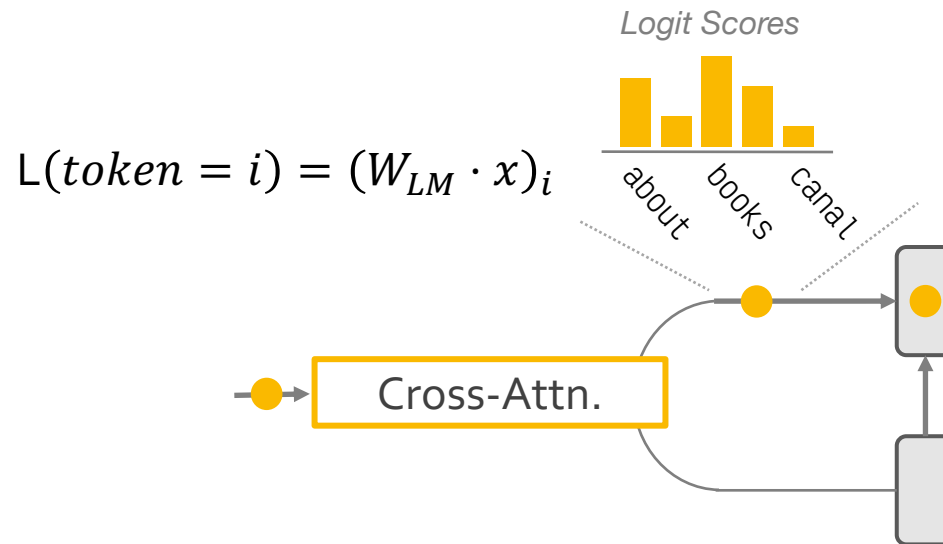
# How does a GenIR Model look like?



# What does Cross- Attention write?

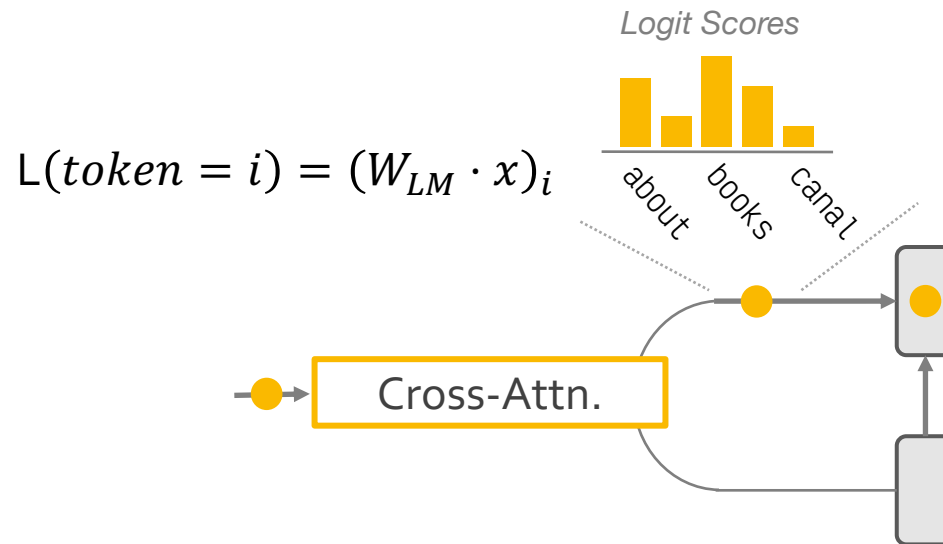


# What does Cross-Attention write?



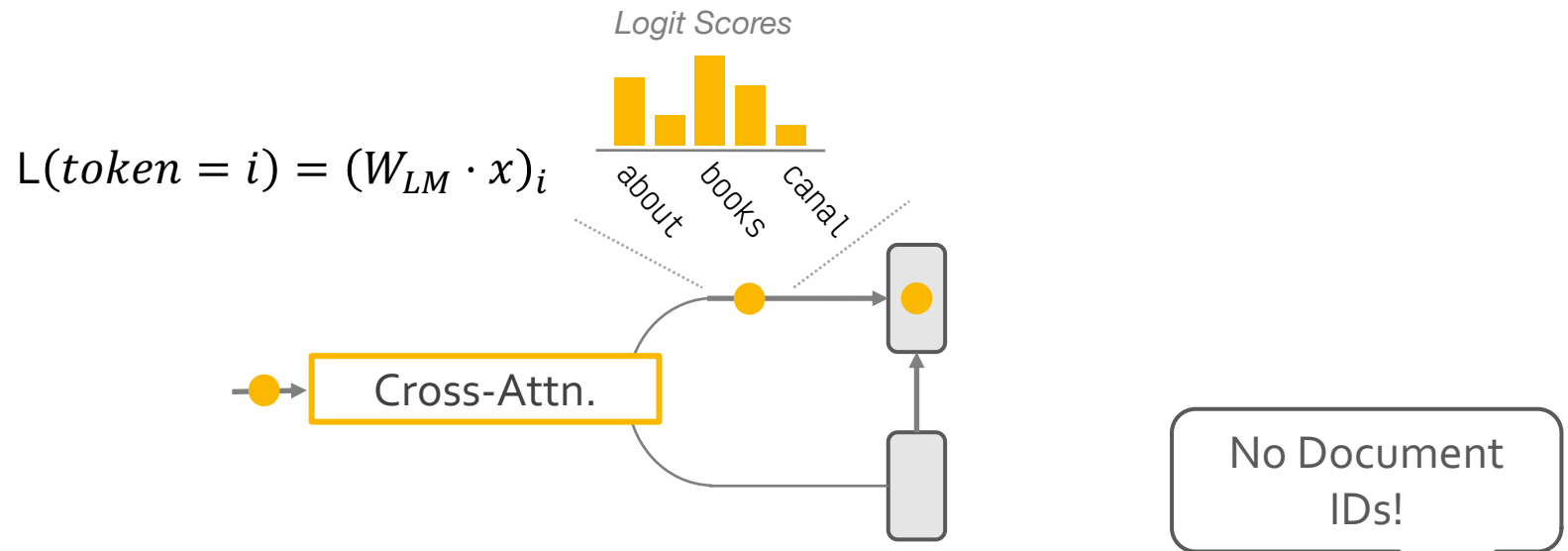
QUERY	LAYER	HEAD	TOP 5 TOKENS
who wrote the harry potter books	Layer 14	Head 2	
	Layer 16	Head 1	
who won the football championship in 2002	Layer 16	Head 1	
	Layer 16	Head 13	

# What does Cross-Attention write?



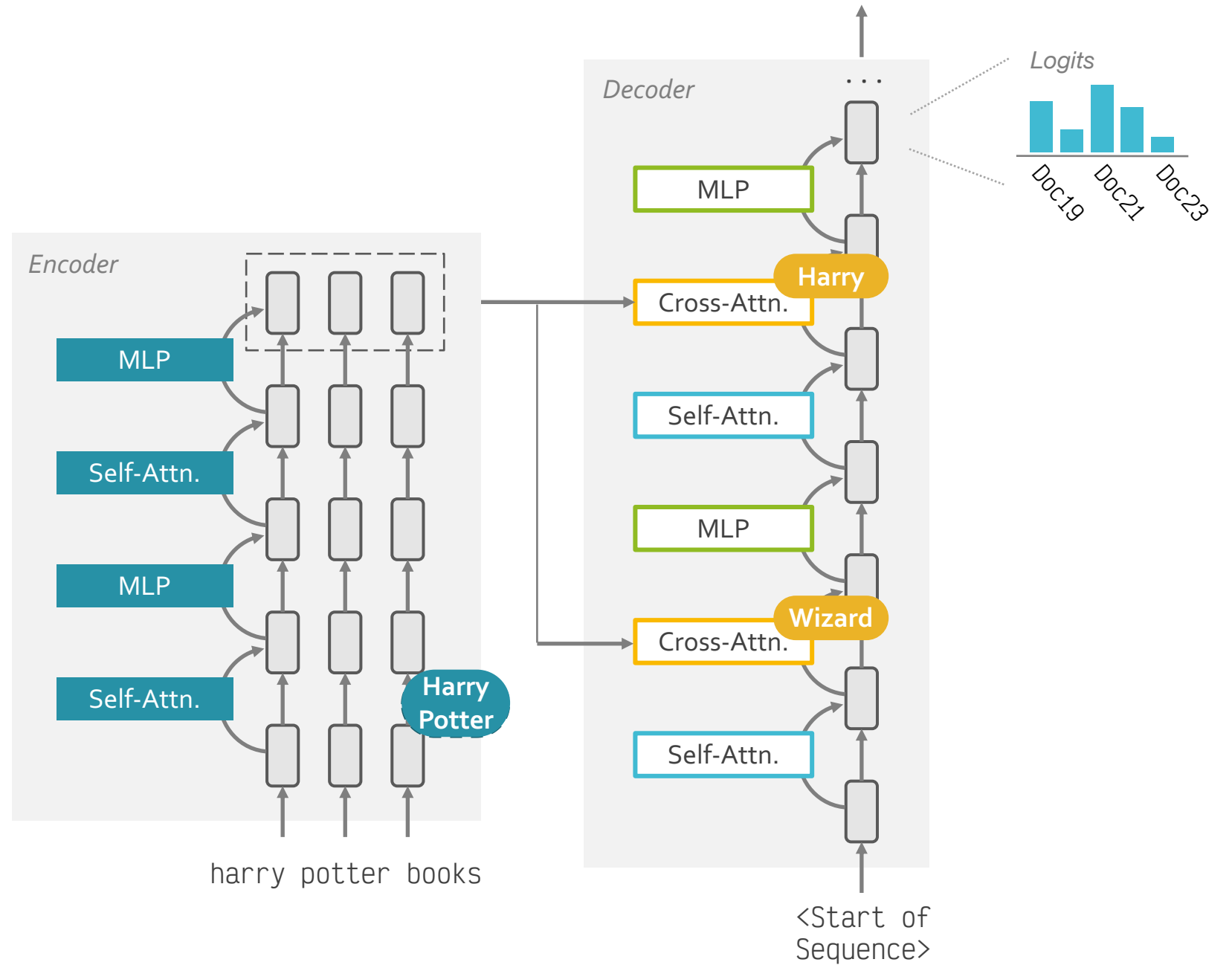
QUERY	LAYER	HEAD	TOP 5 TOKENS
who wrote the harry potter books	Layer 14	Head 2	about, written, about, tailored, privire
	Layer 16	Head 1	books, ouvrage, books, authors, book
who won the football championship in 2002	Layer 16	Head 1	year, YEAR, Year, year, jahr
	Layer 16	Head 13	football, Football, fotbal, soccer, NFL

# What does Cross-Attention write?

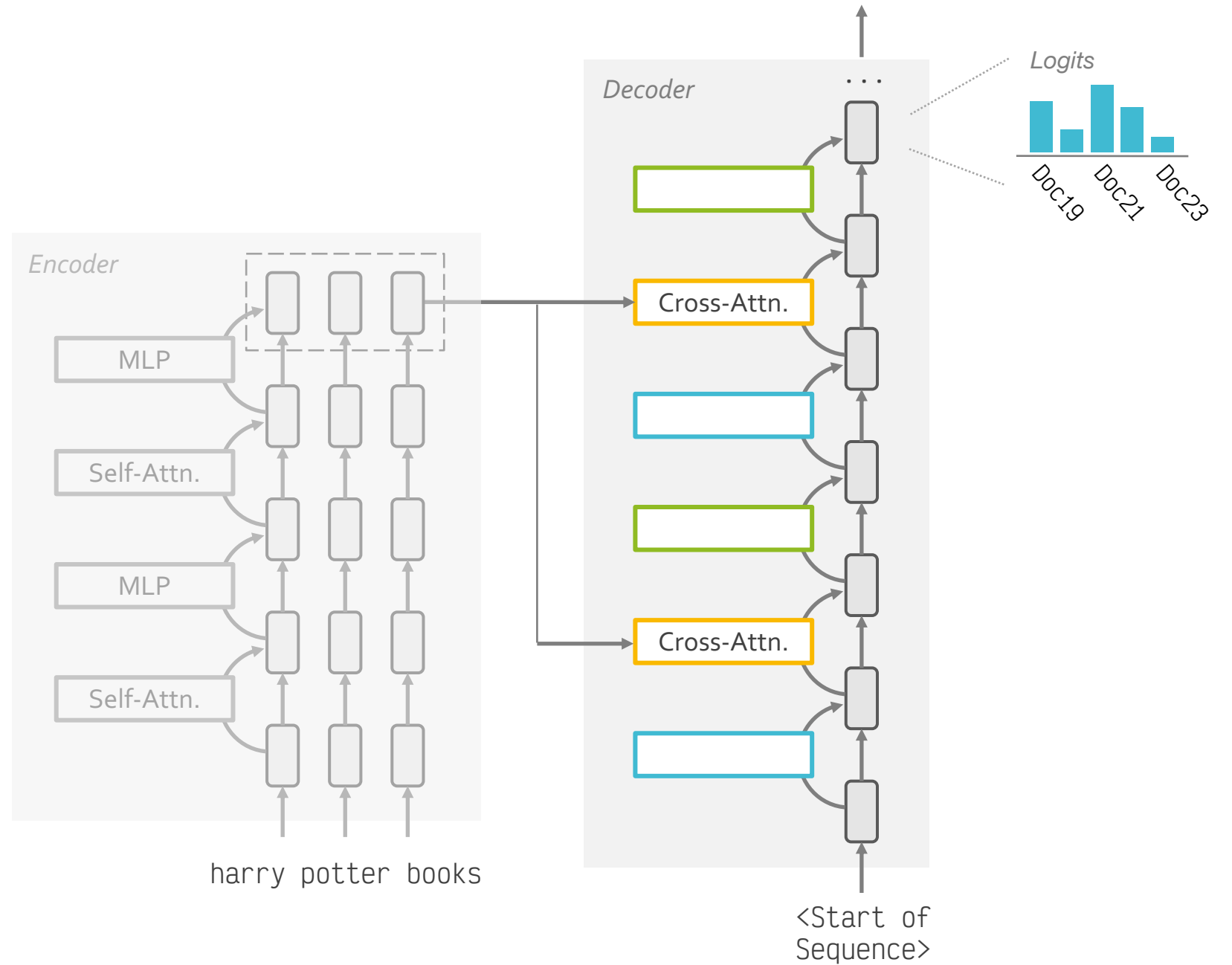


QUERY	LAYER	HEAD	TOP 5 TOKENS
who wrote the harry potter books	Layer 14	Head 2	about, written, about, tailored, privire
	Layer 16	Head 1	books, ouvrage, books, authors, book
who won the football championship in 2002	Layer 16	Head 1	year, YEAR, Year, year, jahr
	Layer 16	Head 13	football, Football, fotbal, soccer, NFL

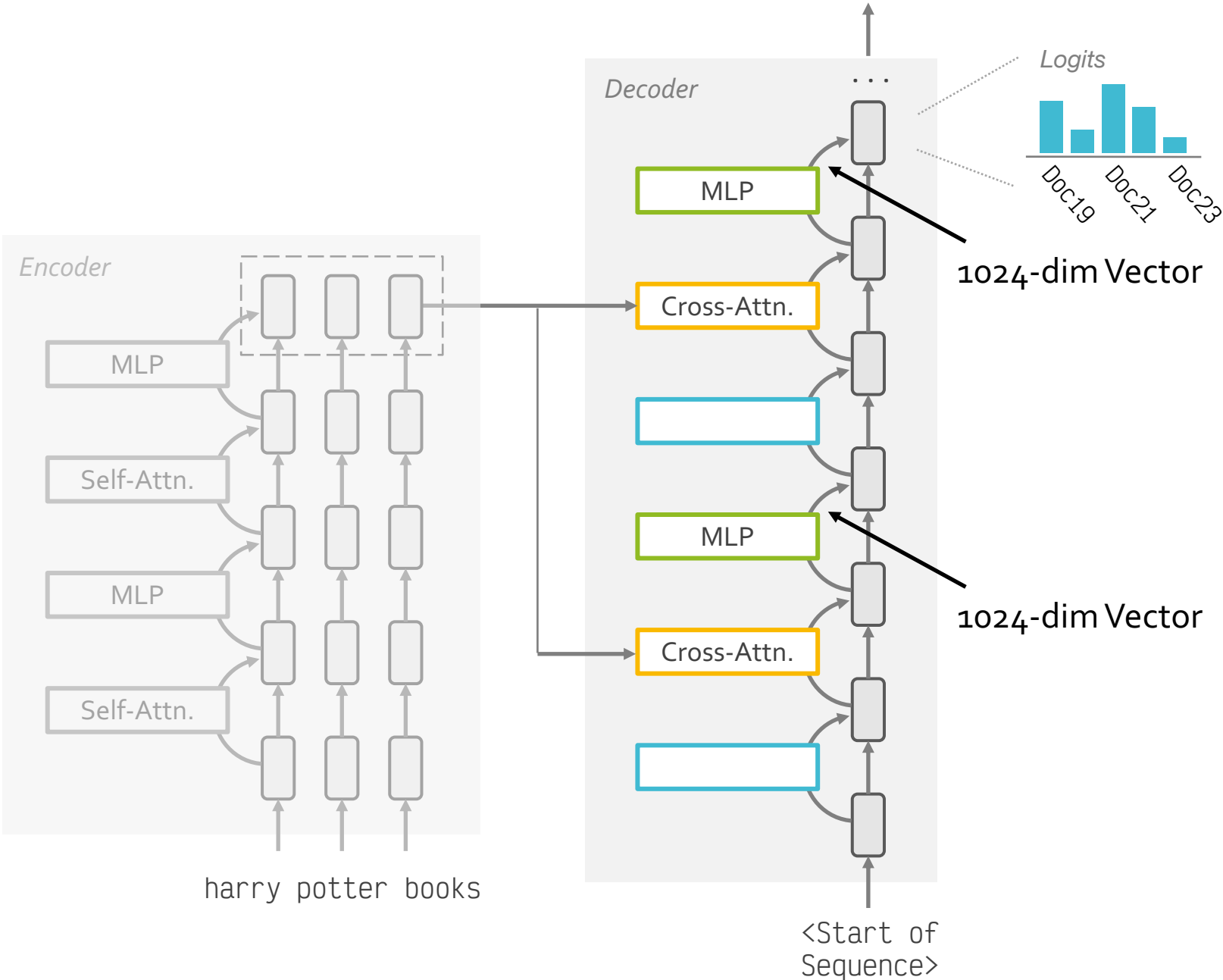
# The Process so far



# The GenIR Model

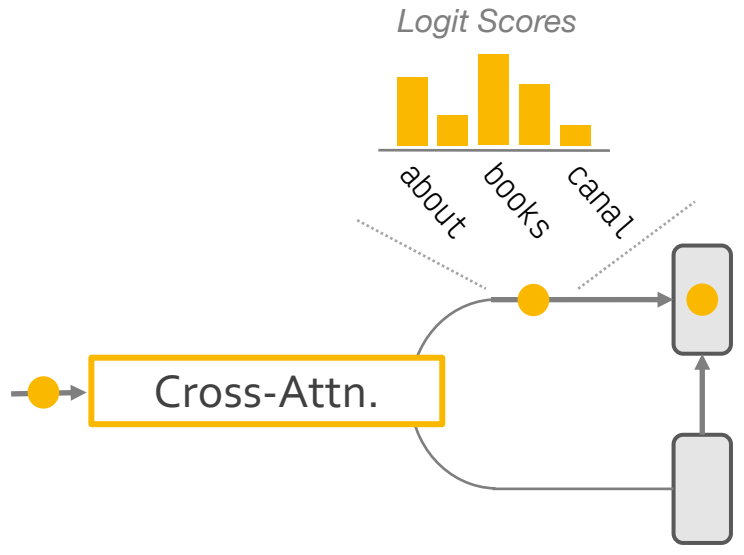


# The GenIR Model

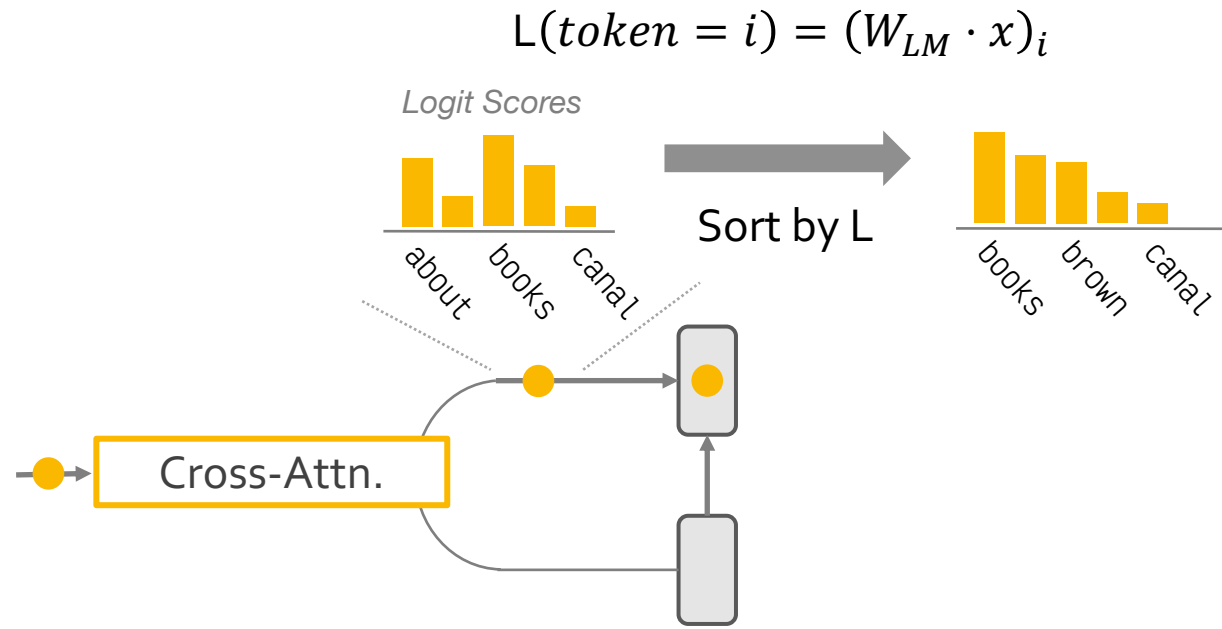


# Component Output

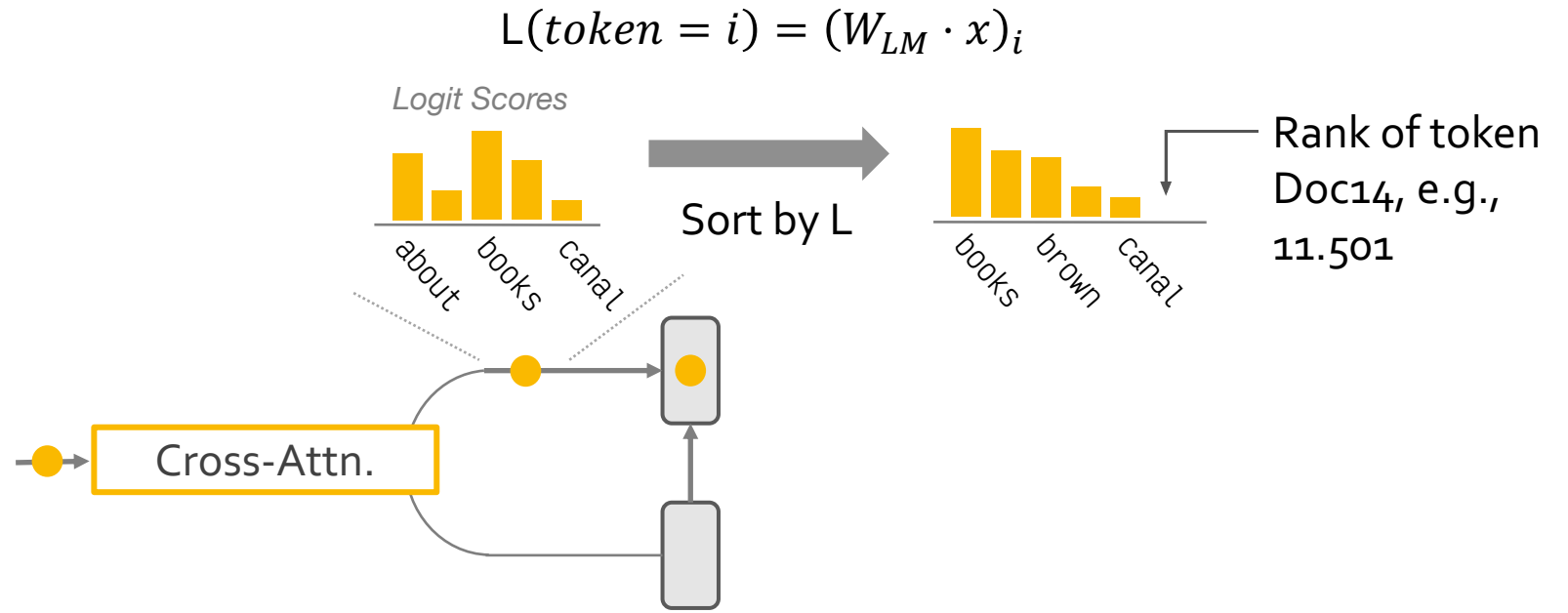
$$L(\text{token} = i) = (W_{LM} \cdot x)_i$$



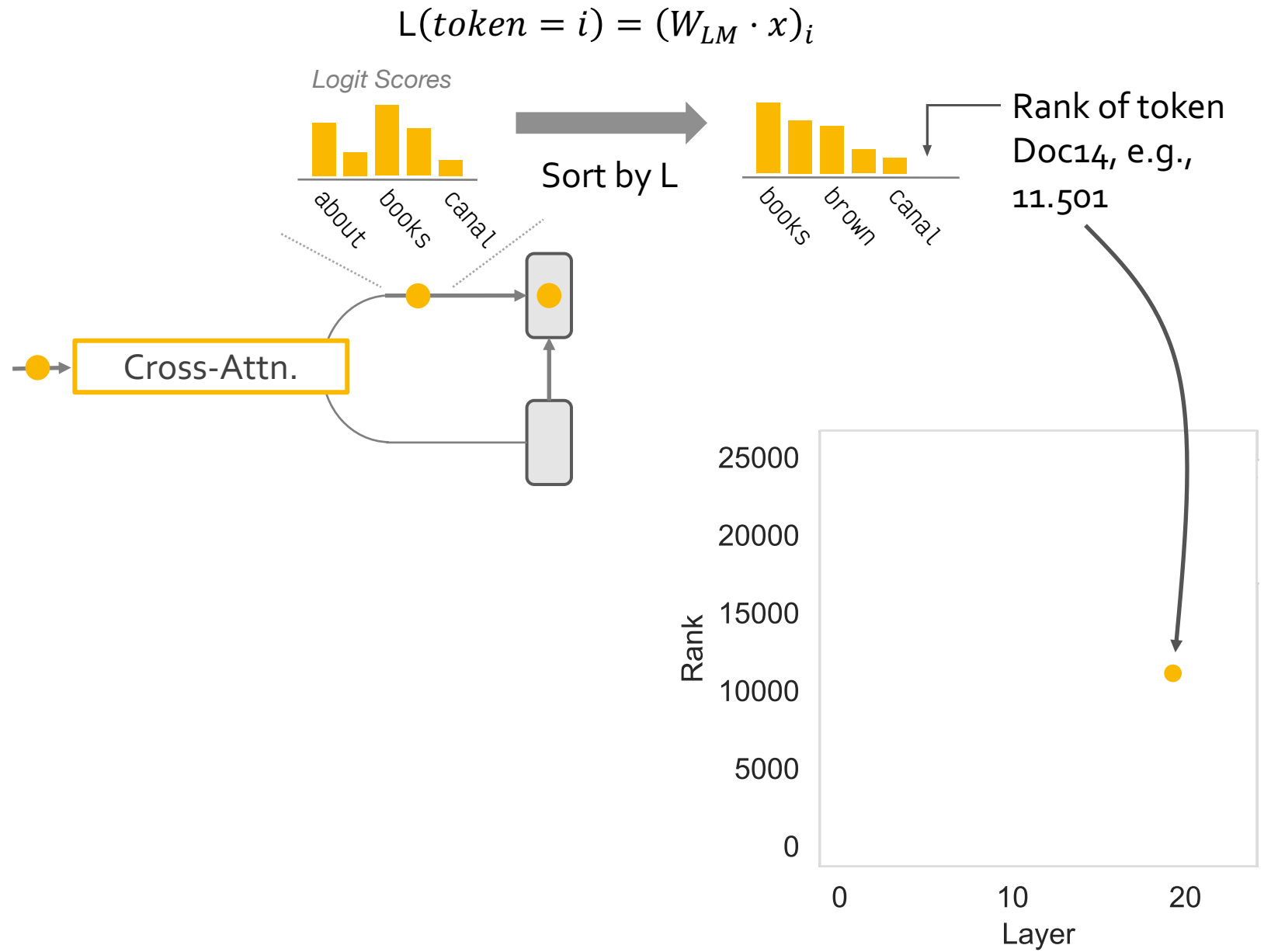
# Component Output



# Component Output

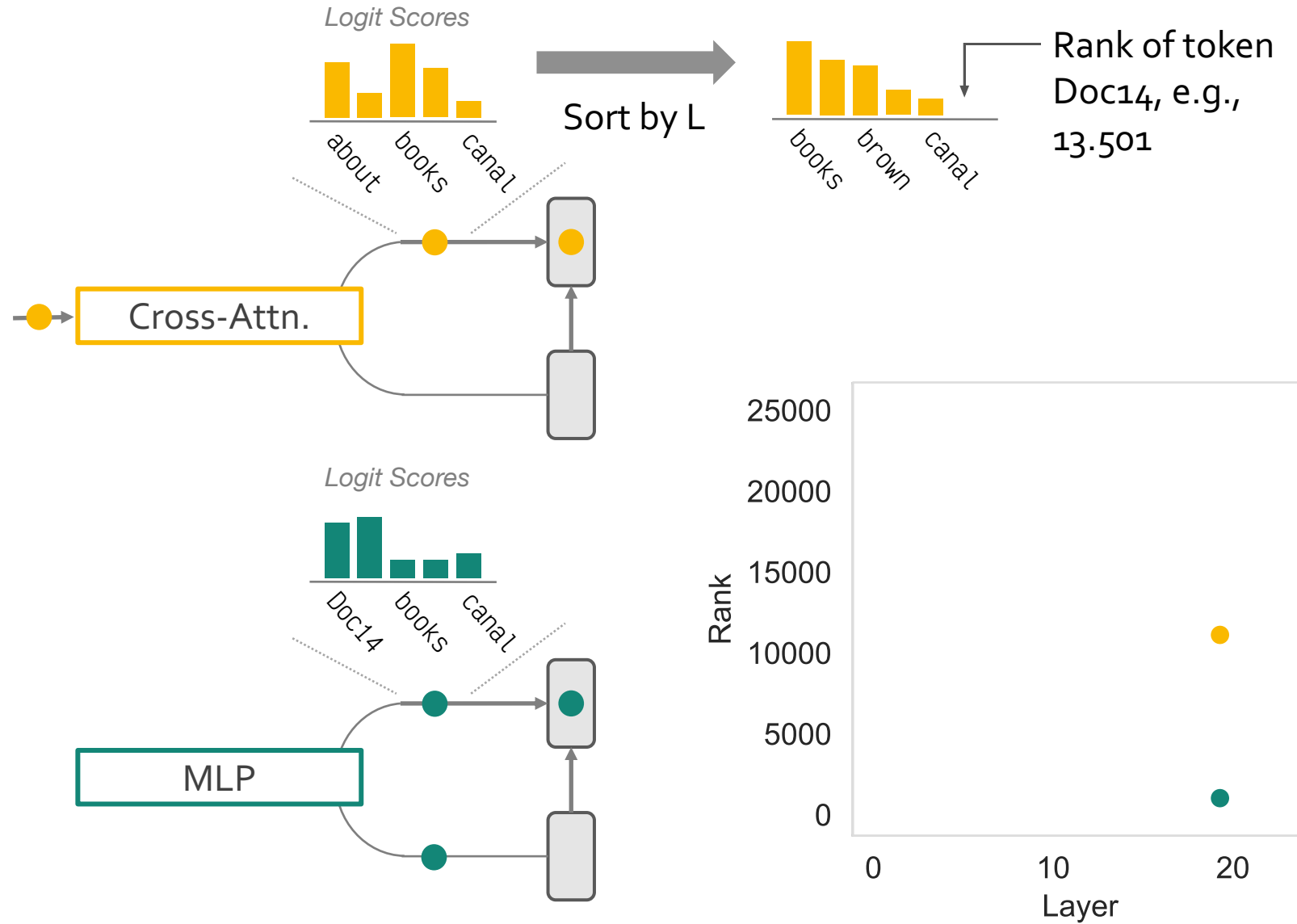


# Component Output



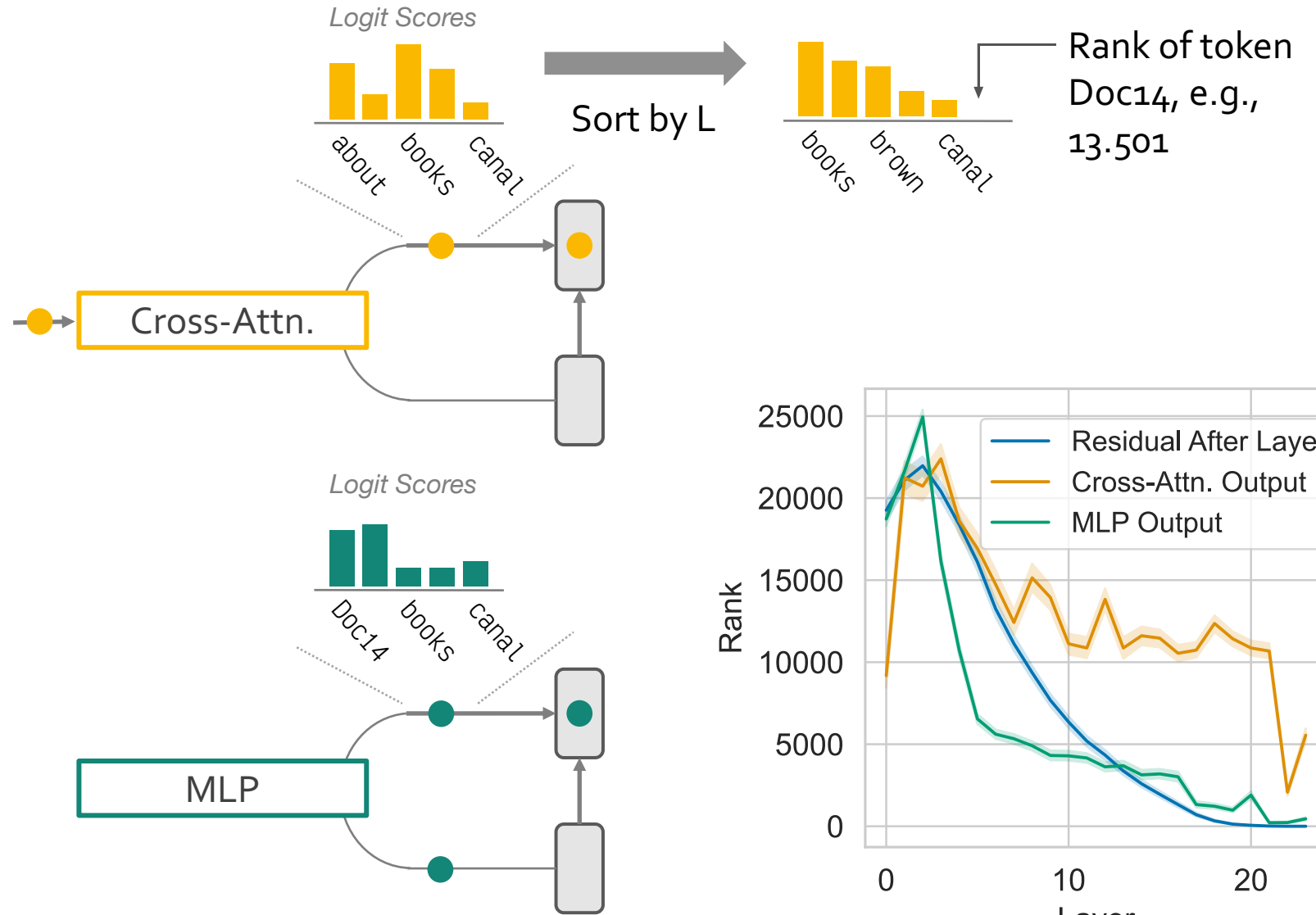
# Component Output

$$L(\text{token} = i) = (W_{LM} \cdot x)_i$$

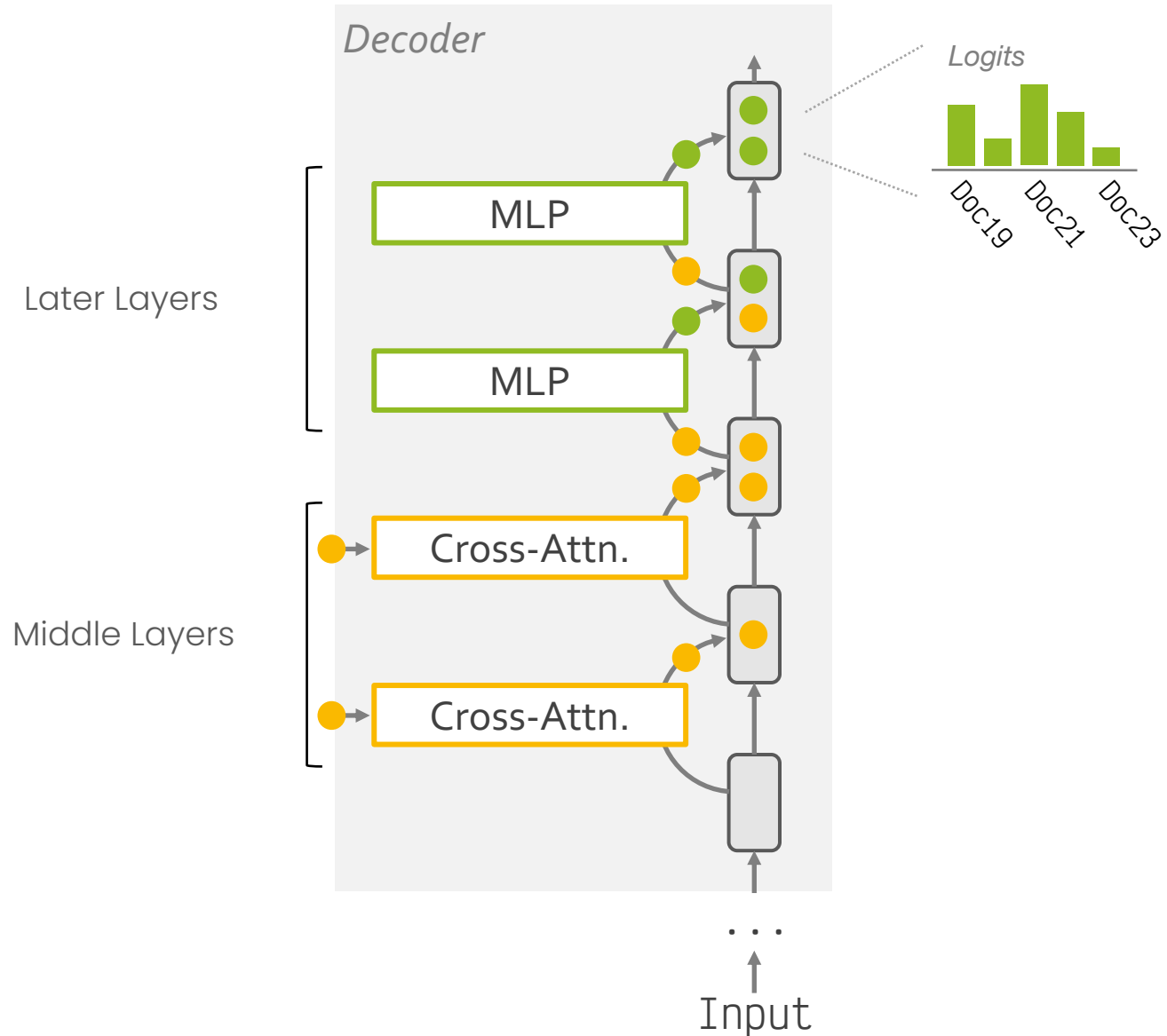


# Component Output

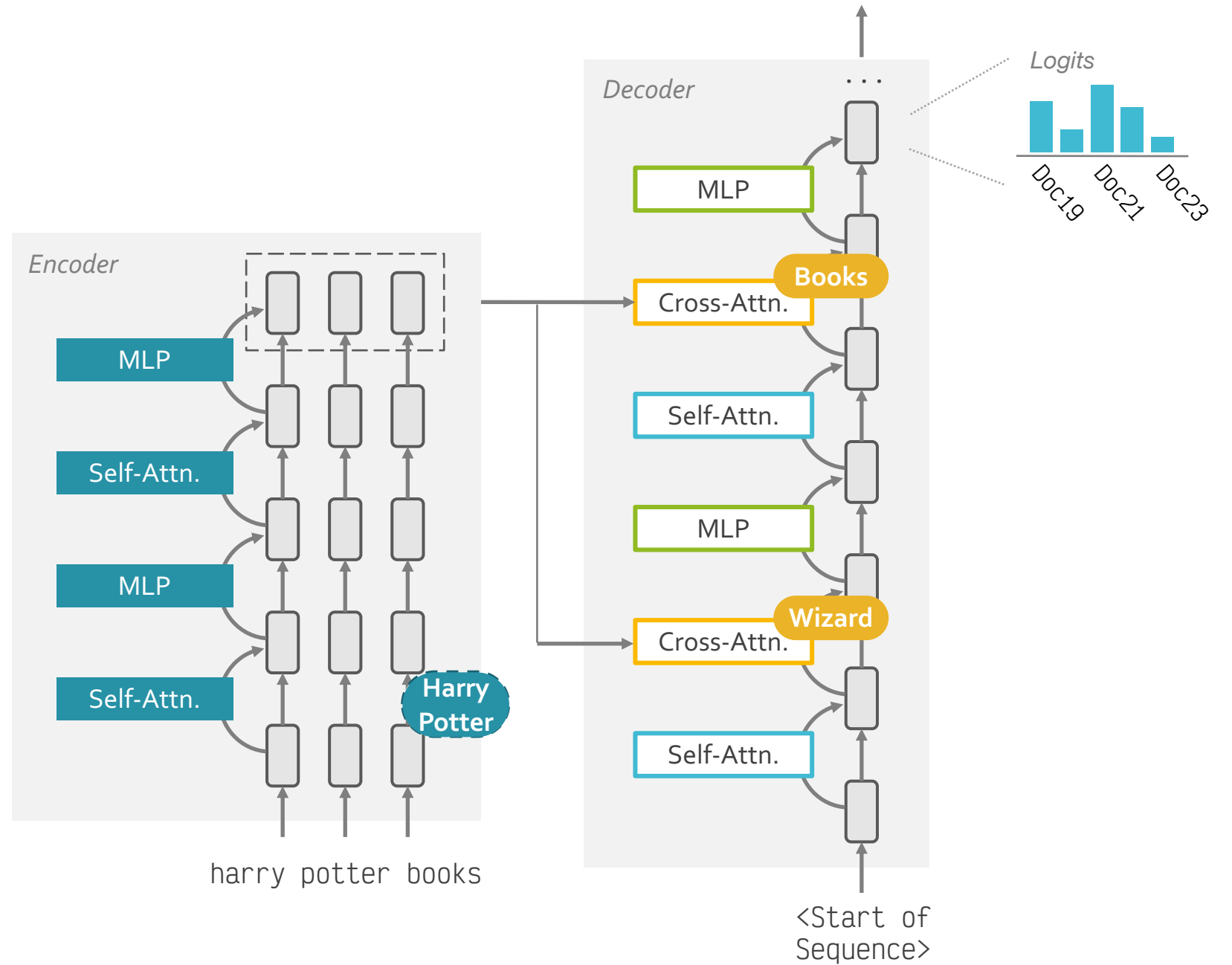
$$L(\text{token} = i) = (W_{LM} \cdot x)_i$$



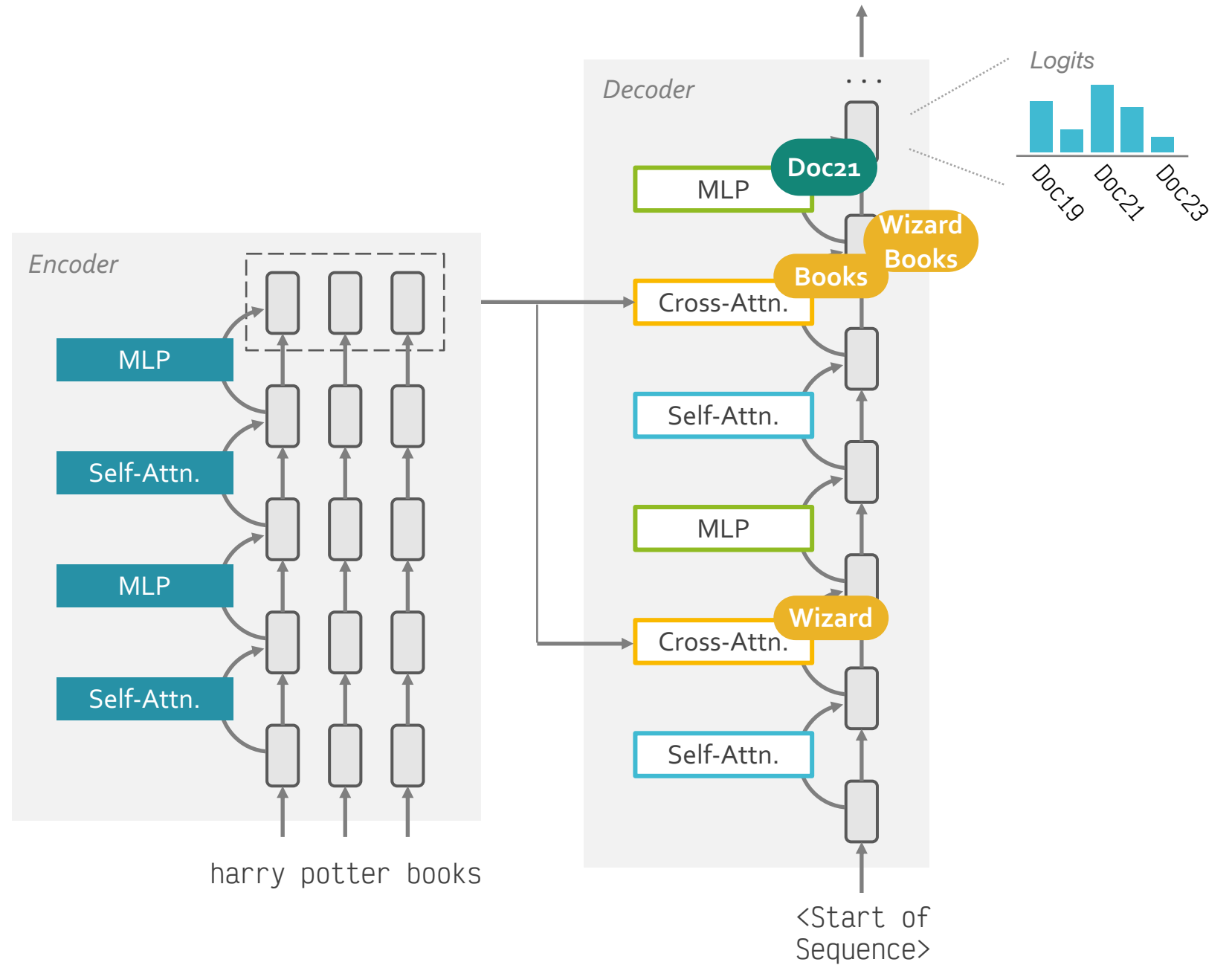
Cross-Attention writes, what MLPs later read.



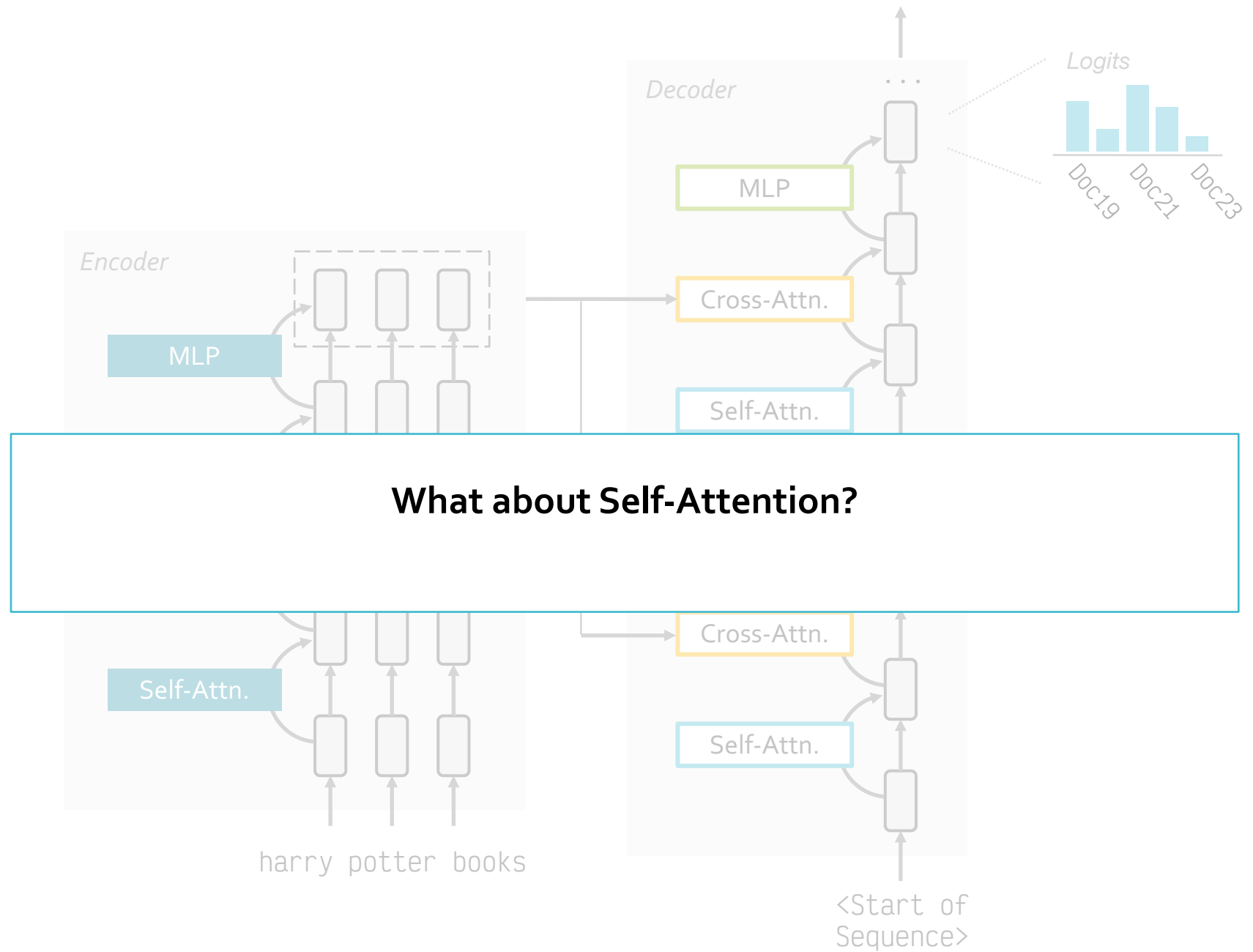
# The whole picture (?)



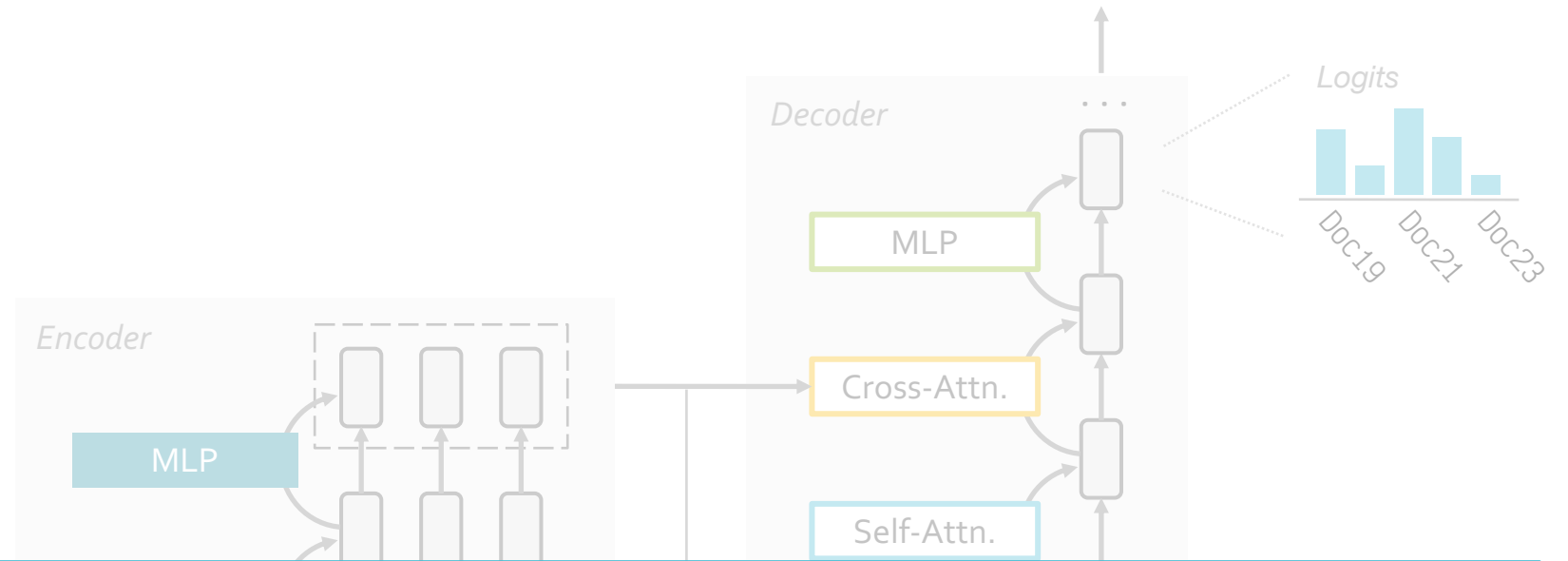
# The whole picture (?)



# The whole picture (?)

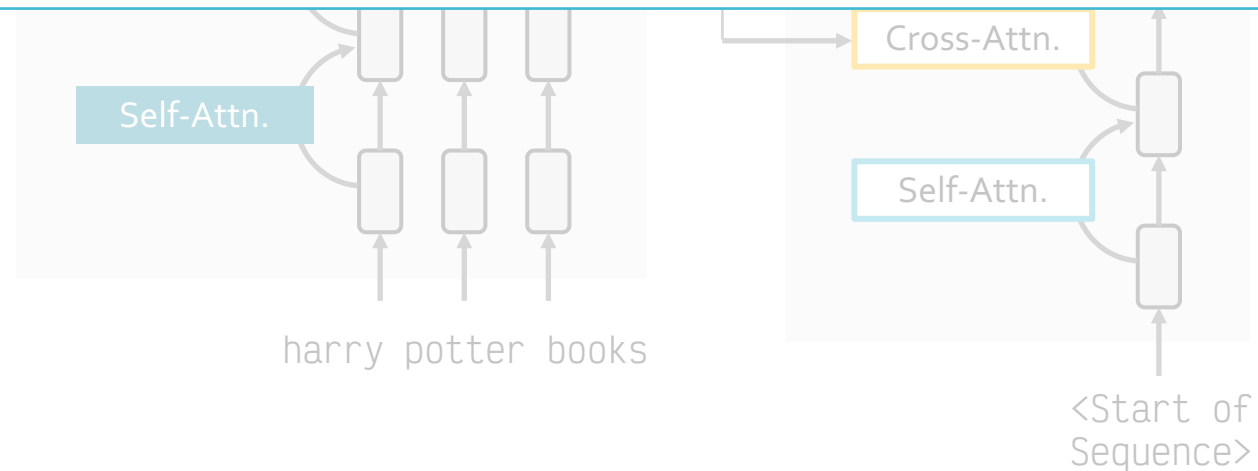


# The whole picture (?)



## What about Self-Attention?

→ We found that Self-Attention is irrelevant and can be removed!

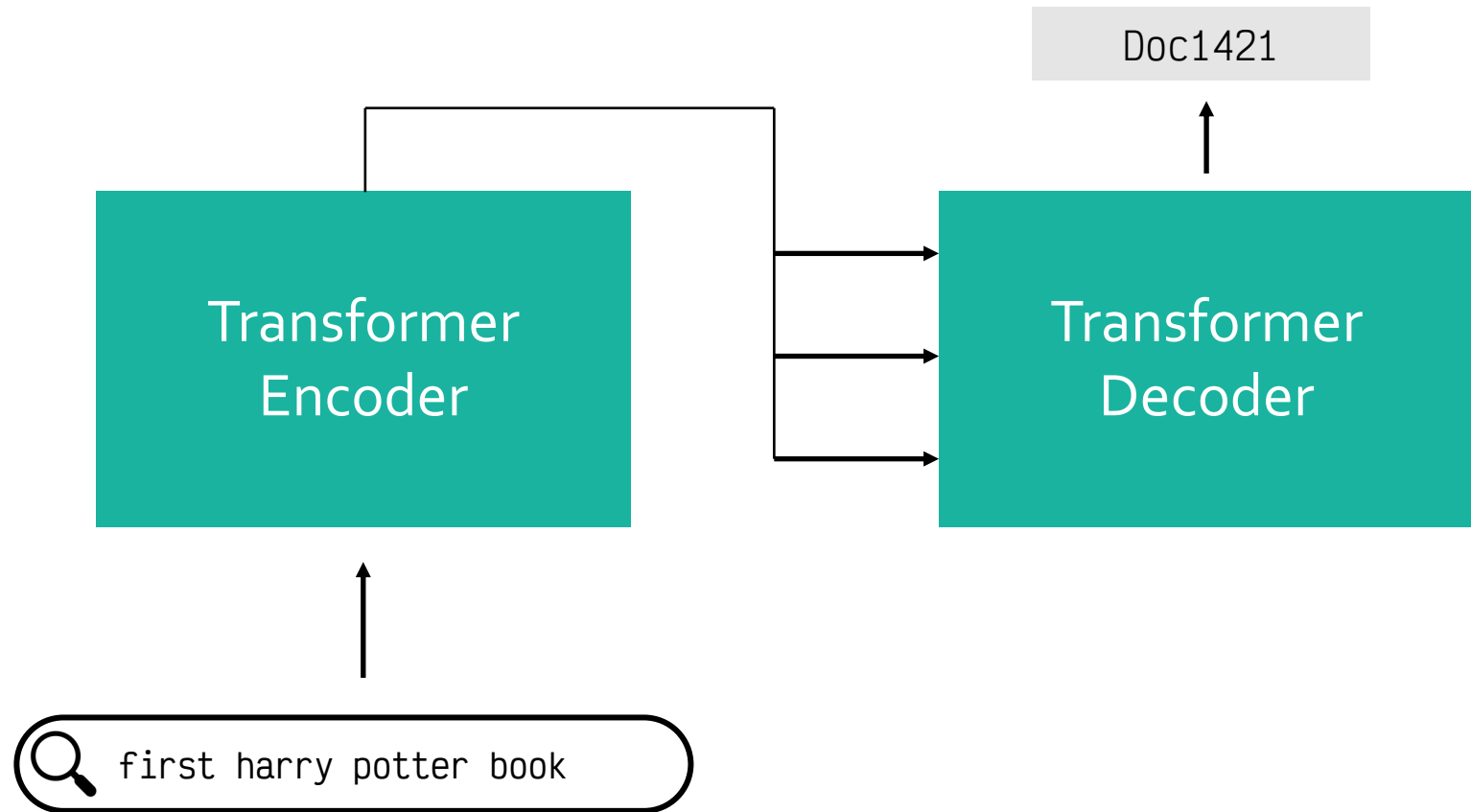




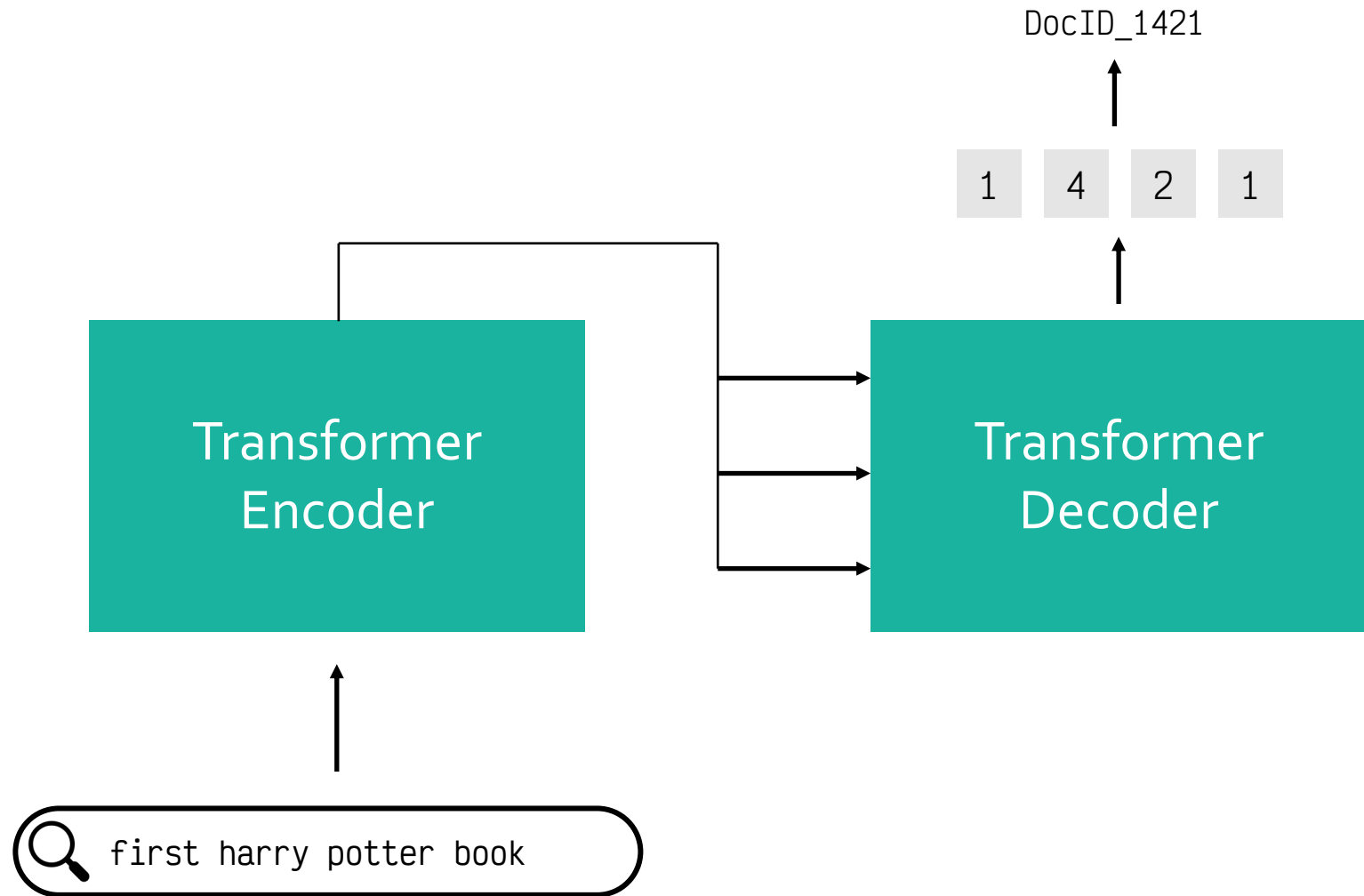
# The Role of Self-Attention in Generative IR



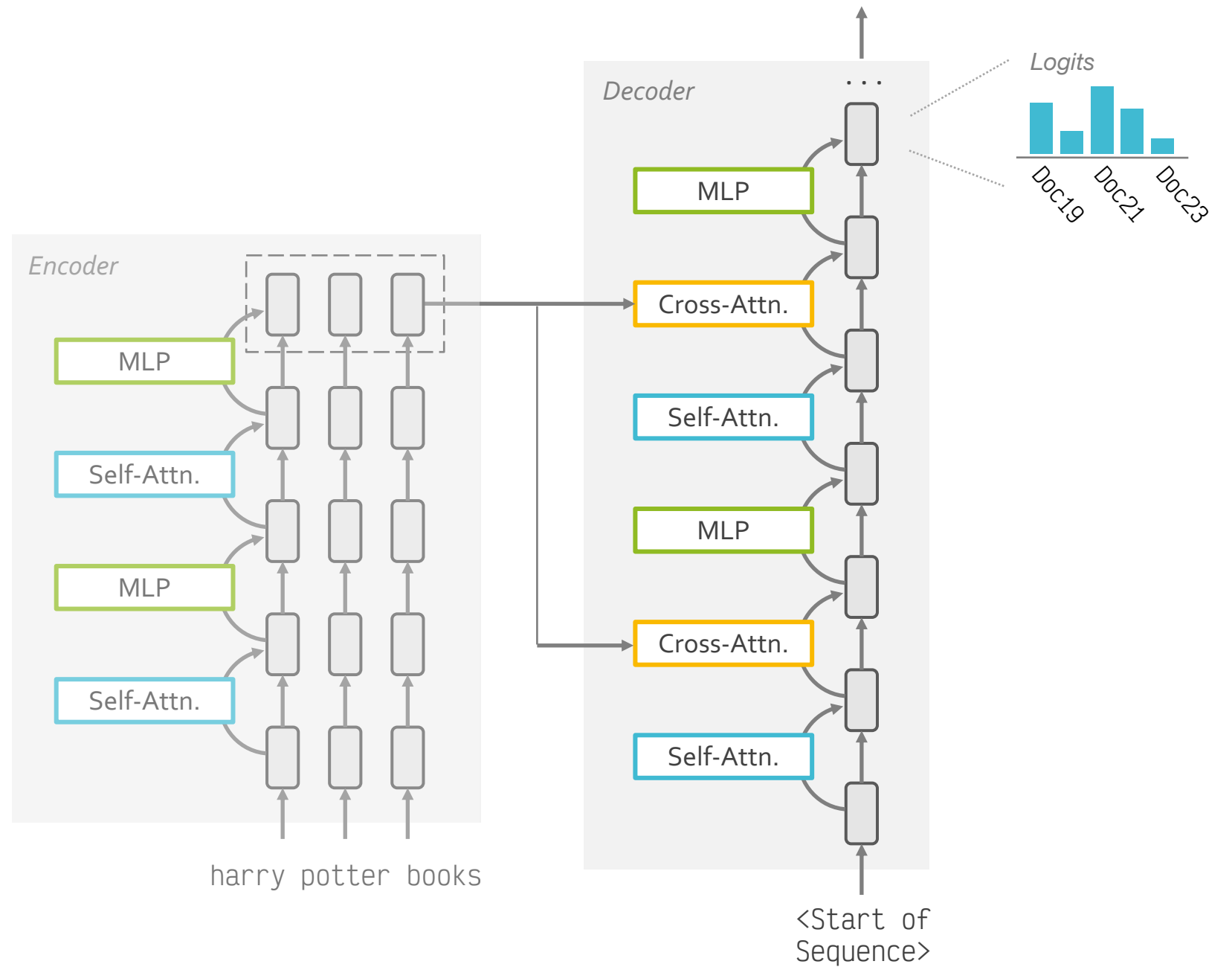
# Single-Token vs. Multi-Token



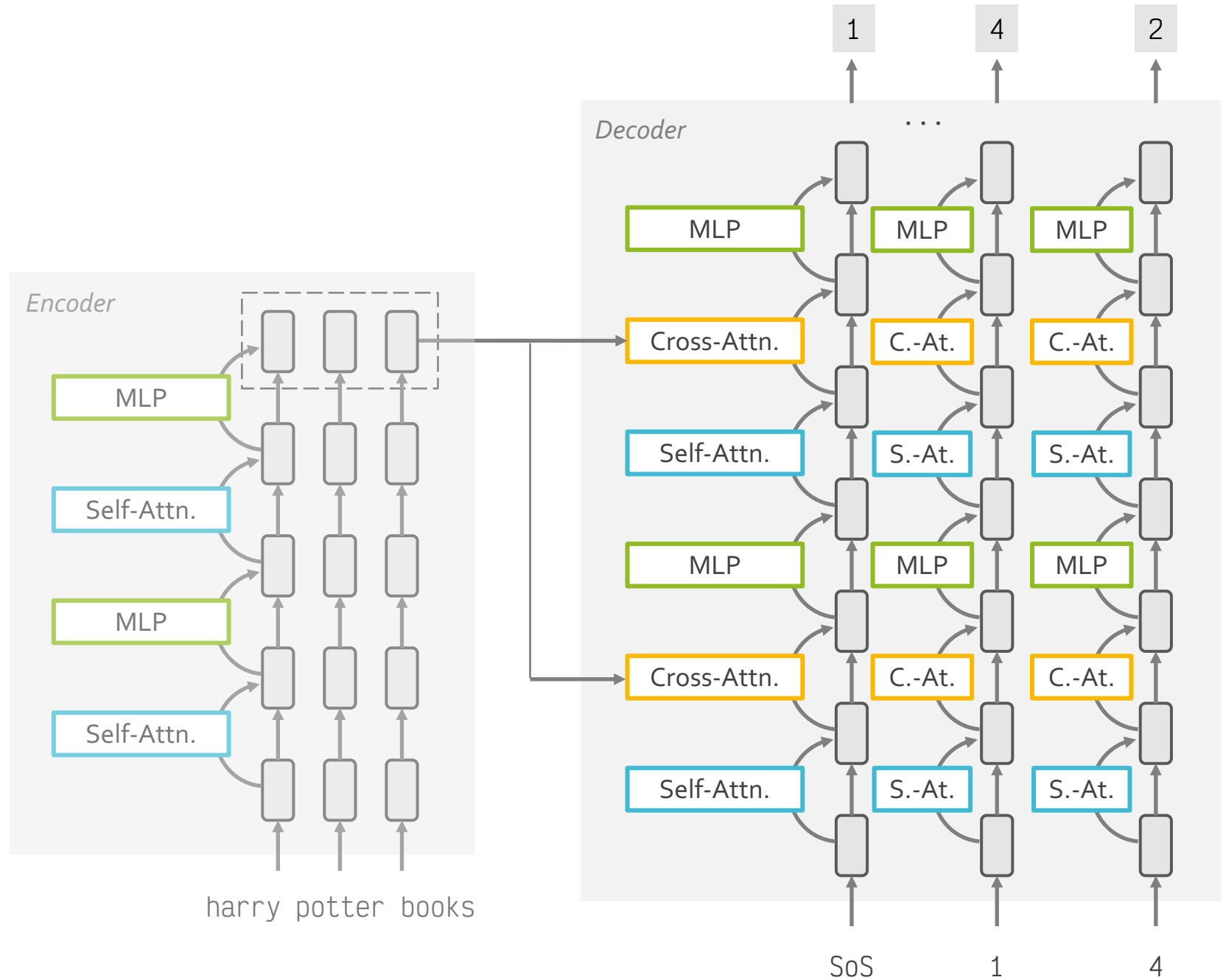
# Single-Token vs. Multi-Token



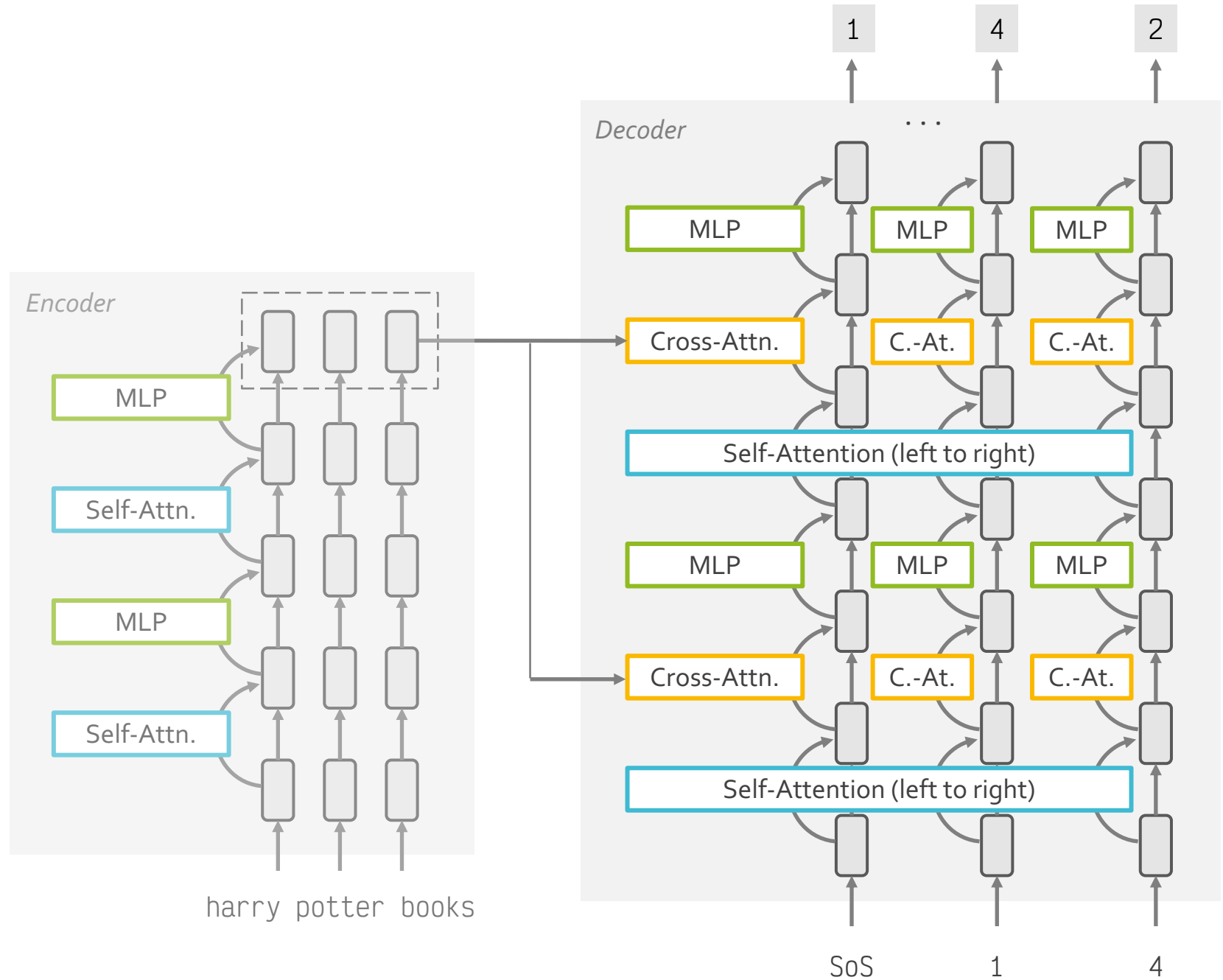
# Single-Token Prediction



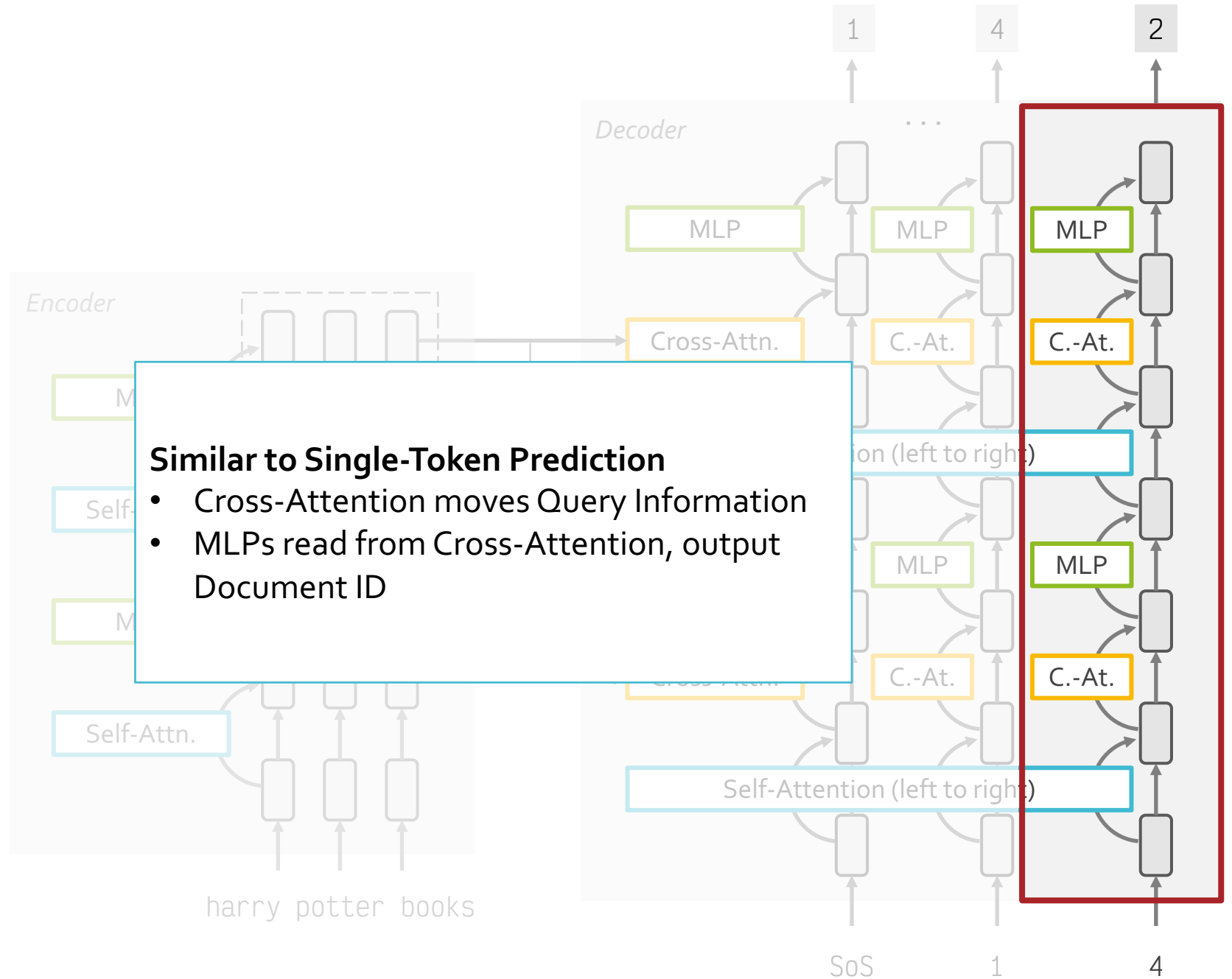
# Multi-Token Prediction



# Multi-Token Prediction



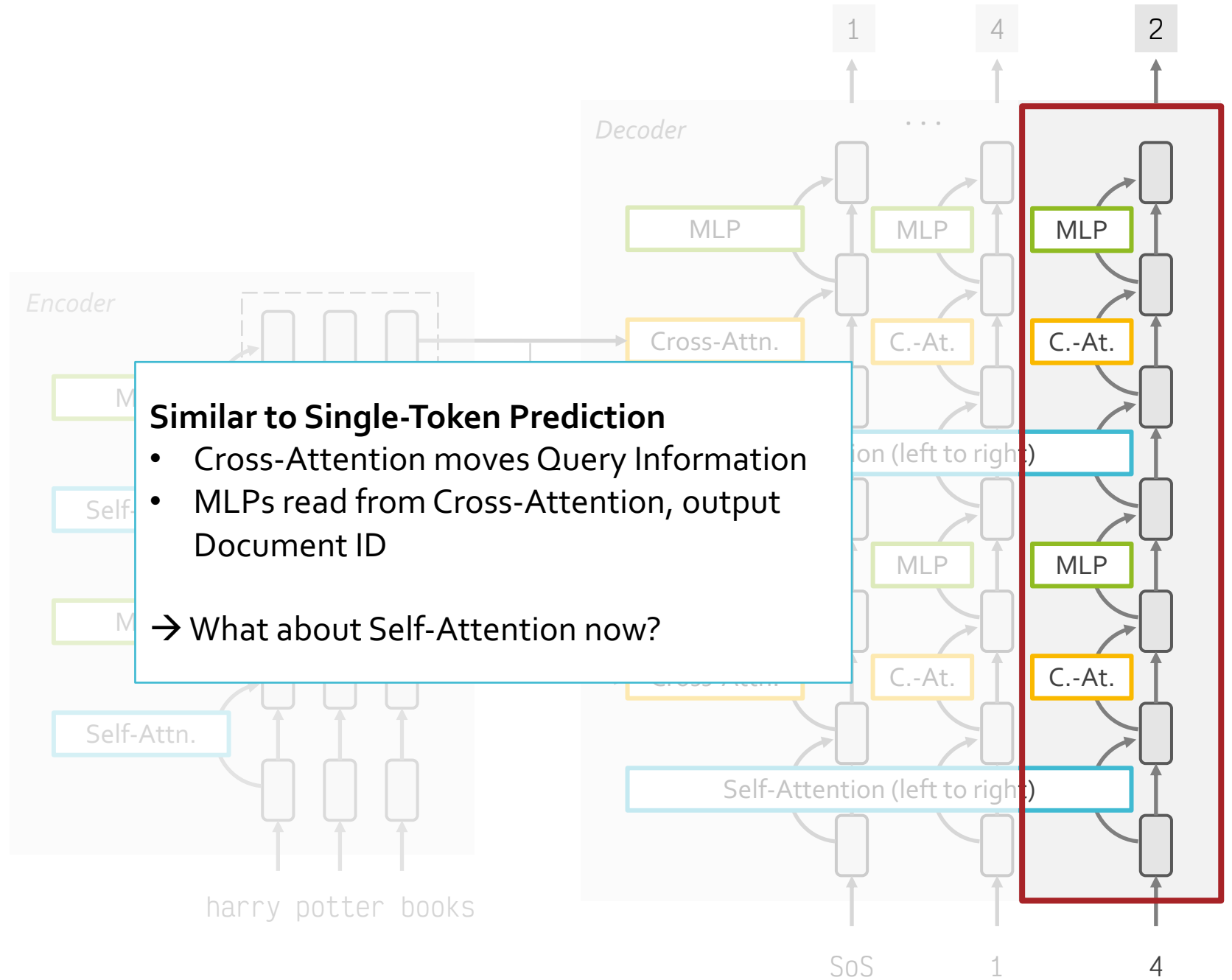
# Multi-Token Prediction



## Similar to Single-Token Prediction

- Cross-Attention moves Query Information
- MLPs read from Cross-Attention, output Document ID

# Multi-Token Prediction

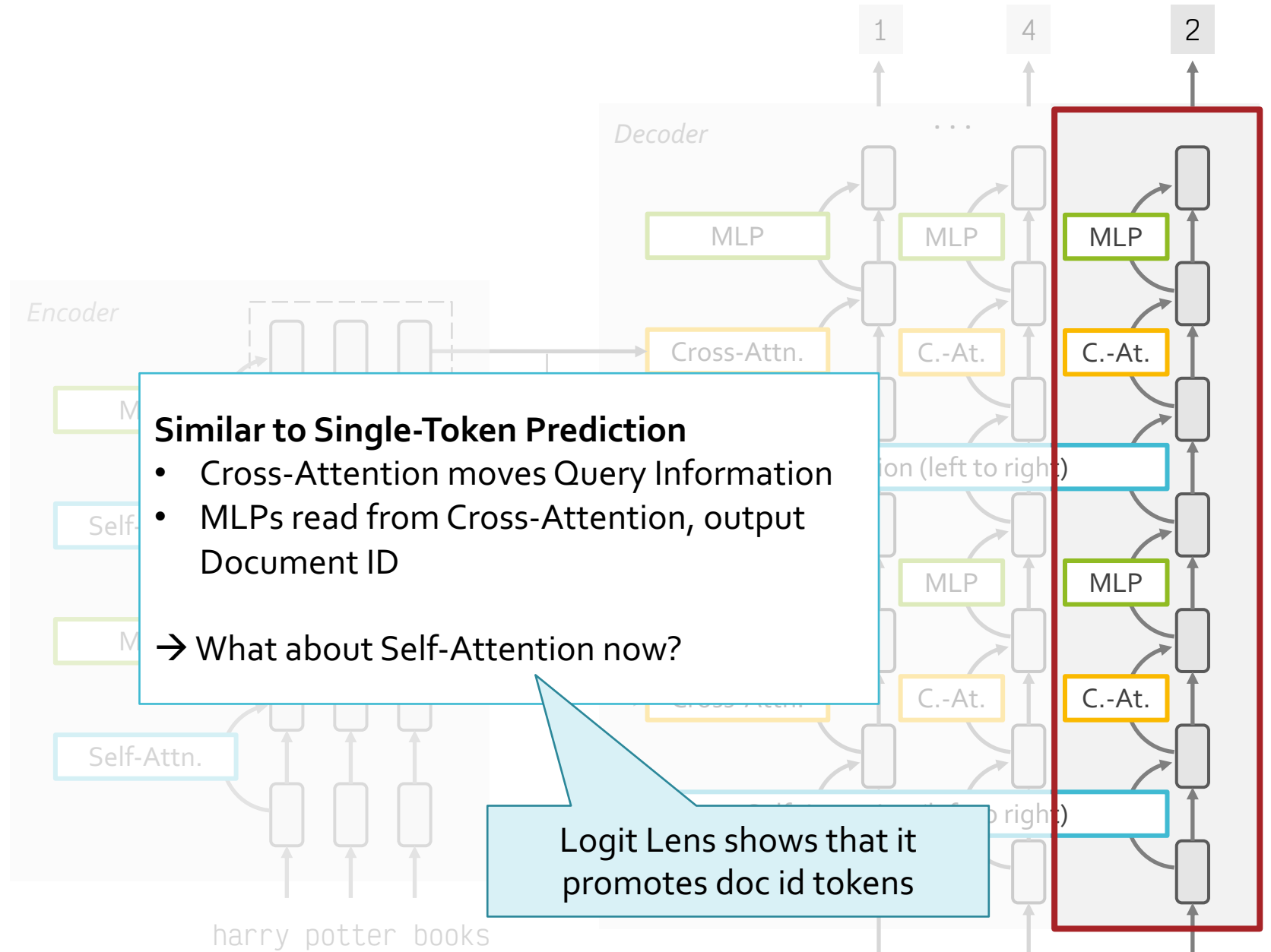


## Similar to Single-Token Prediction

- Cross-Attention moves Query Information
- MLPs read from Cross-Attention, output Document ID

→ What about Self-Attention now?

# Multi-Token Prediction



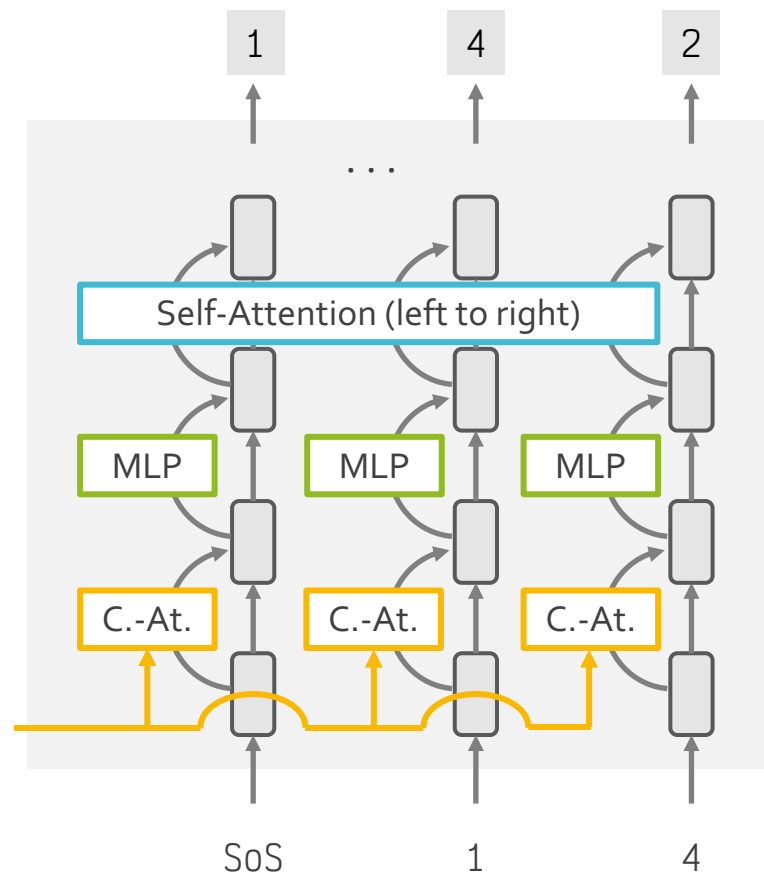
## Similar to Single-Token Prediction

- Cross-Attention moves Query Information
- MLPs read from Cross-Attention, output Document ID

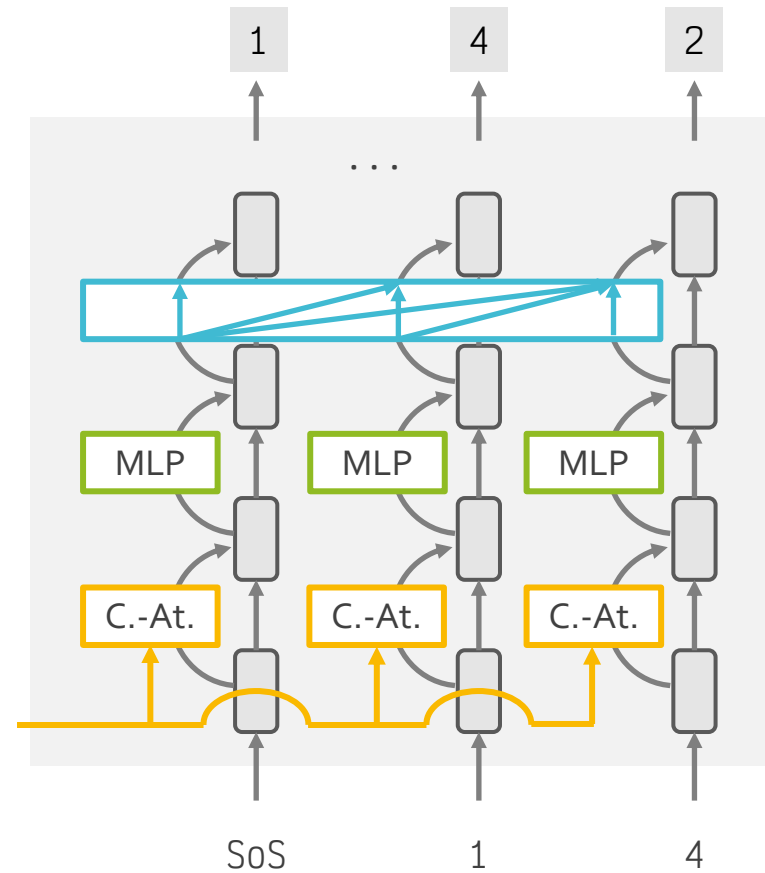
→ What about Self-Attention now?

Logit Lens shows that it promotes doc id tokens

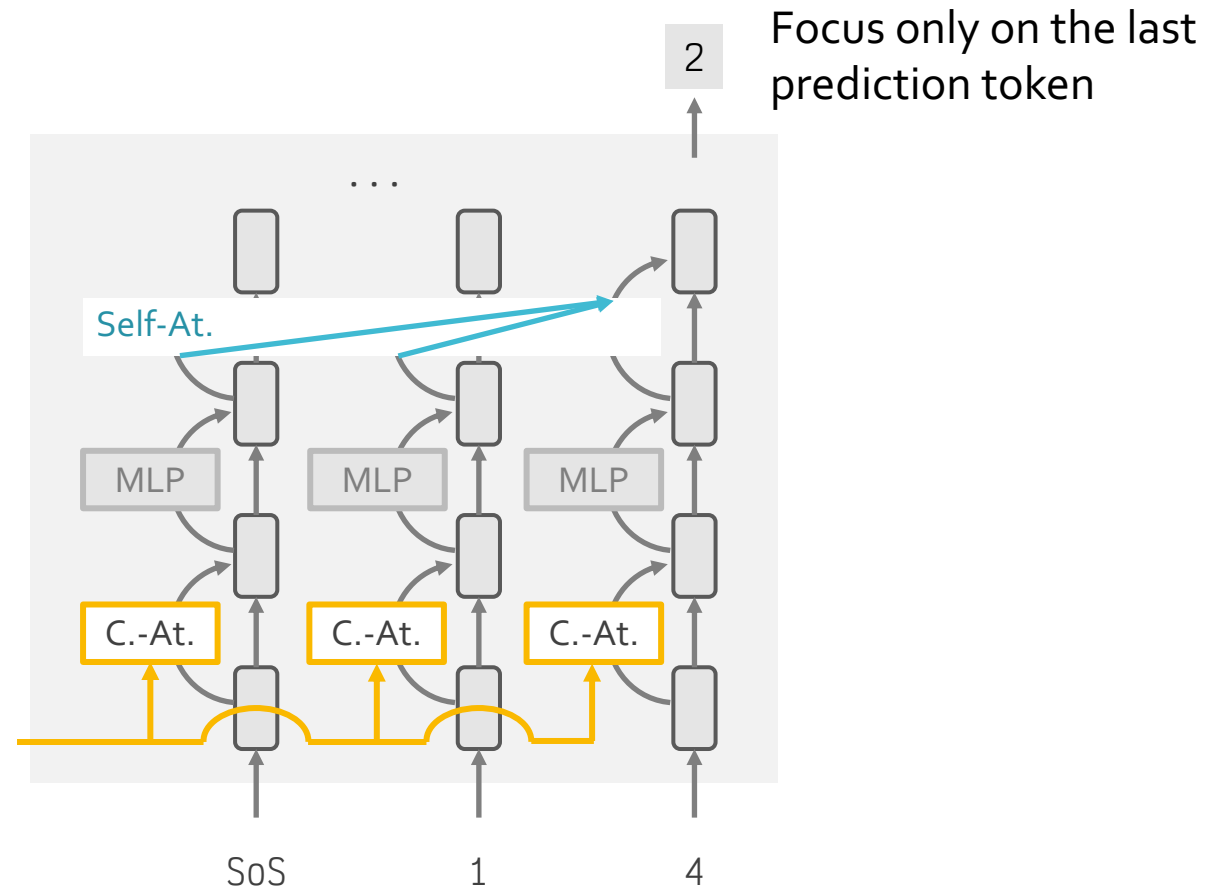
# Information Movement through Self-Attention



# Information Movement through Self-Attention

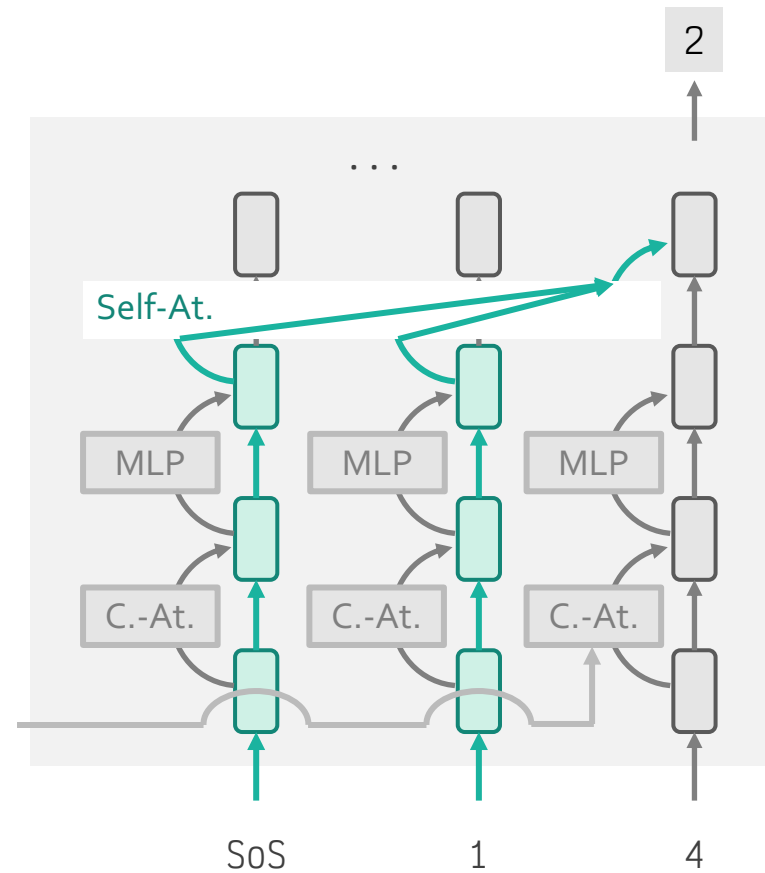


# Information Movement through Self-Attention



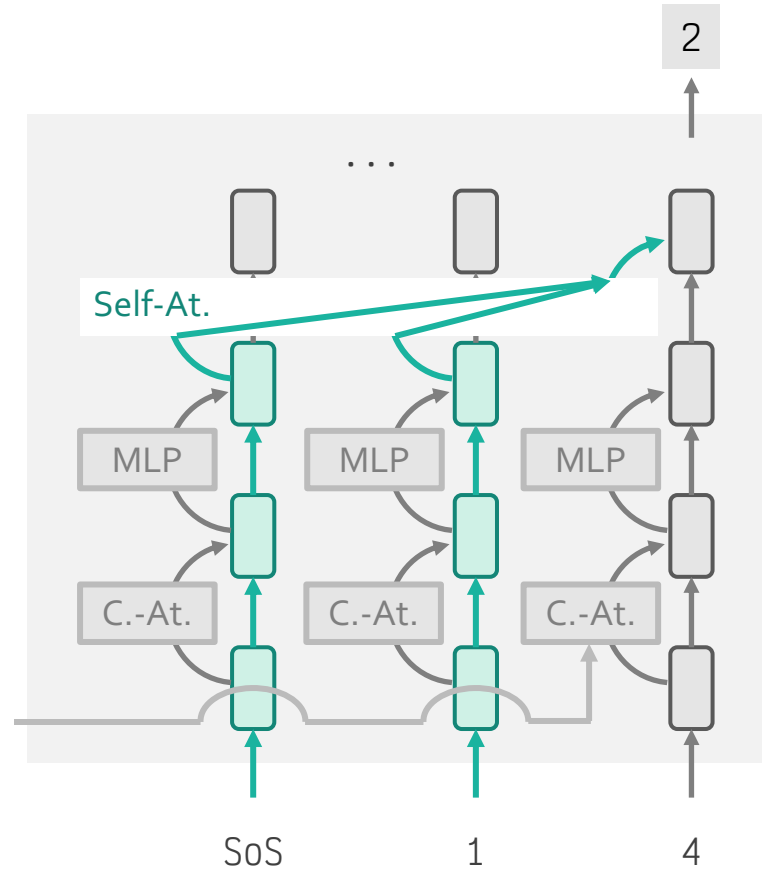
# Information Movement through Self-Attention

## Previous Token Information

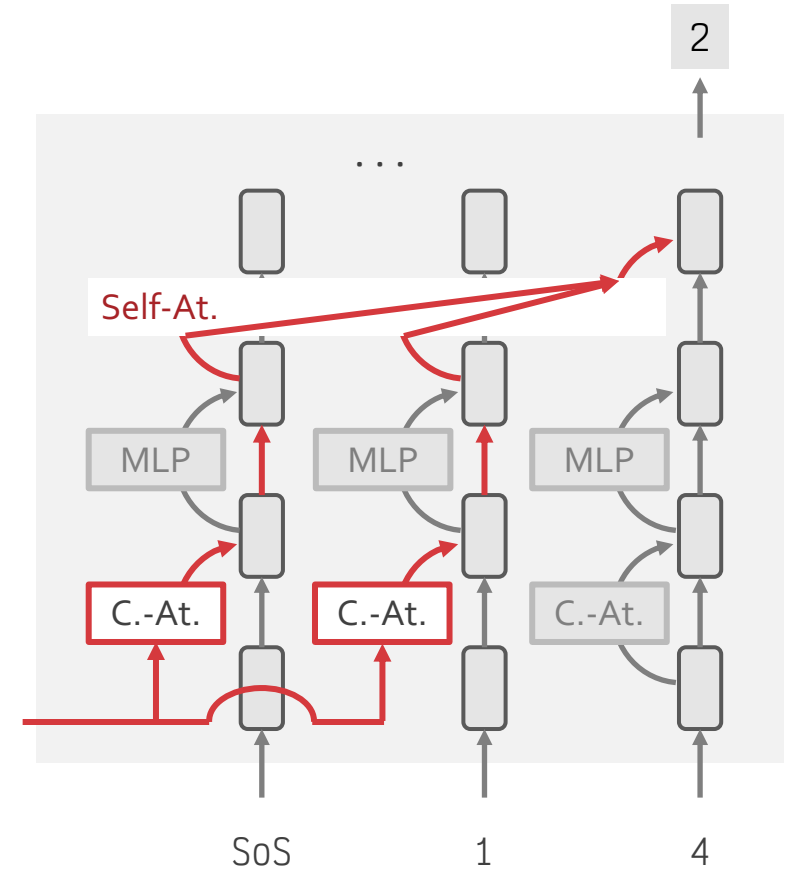


# Information Movement through Self-Attention

## Previous Token Information

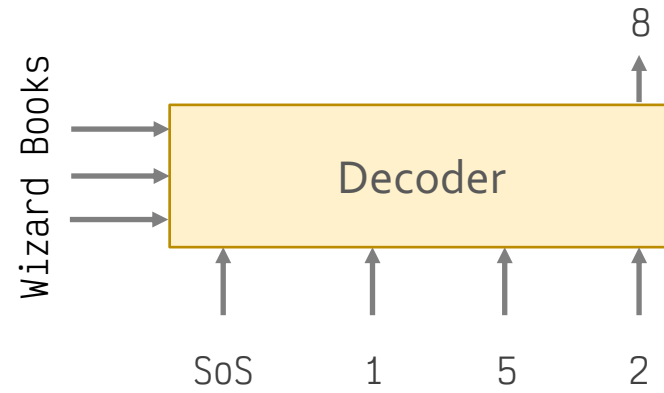


## Query Information

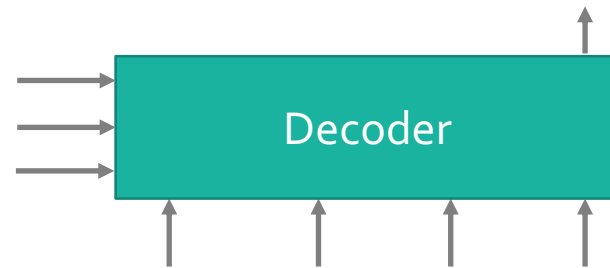


# Two Patching Experiments

Previous Token Information

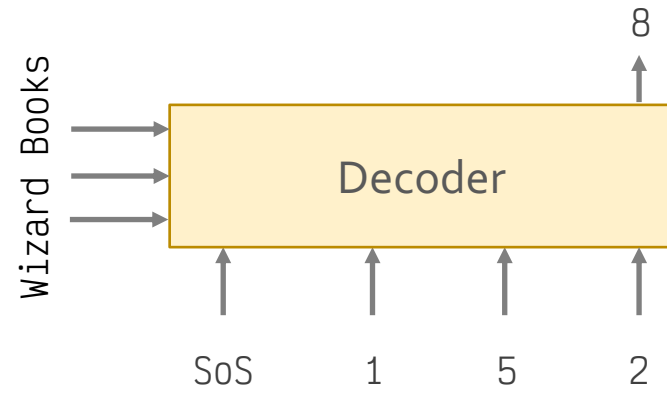


Query Information

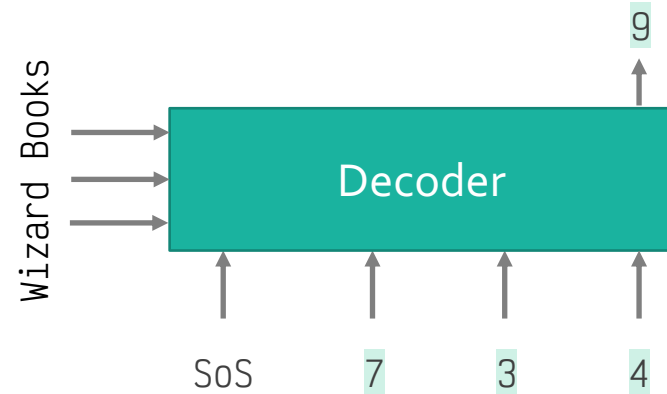


# Two Patching Experiments

## Previous Token Information

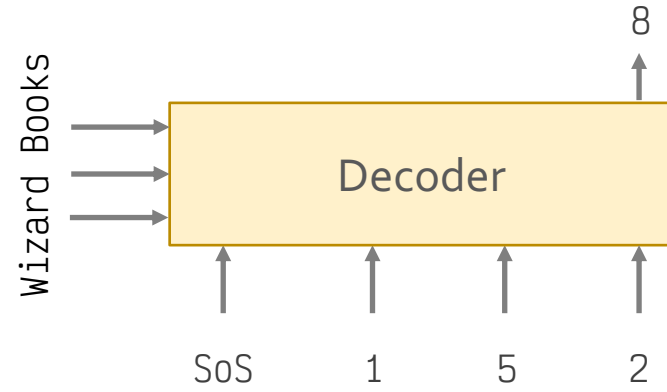


## Query Information

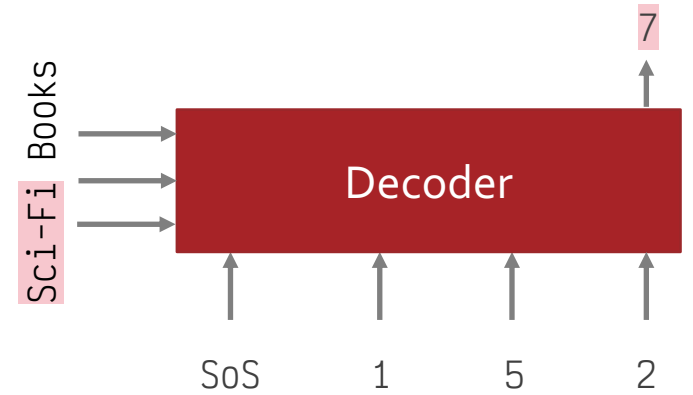
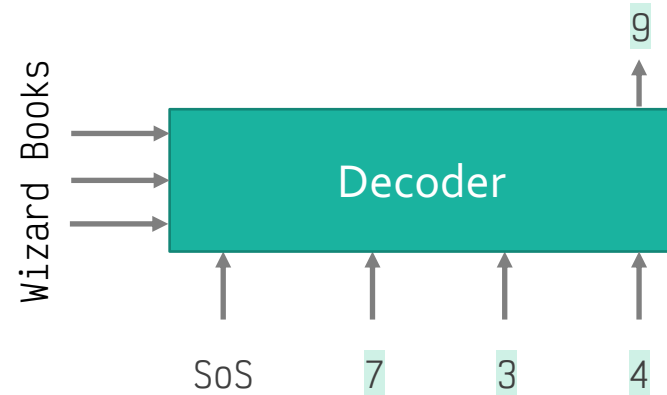
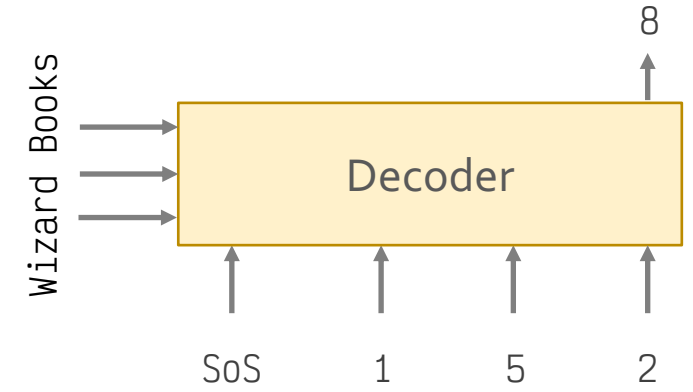


# Two Patching Experiments

### Previous Token Information

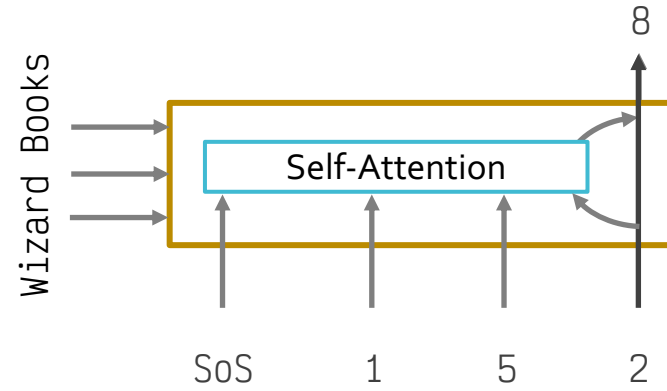


### Query Information

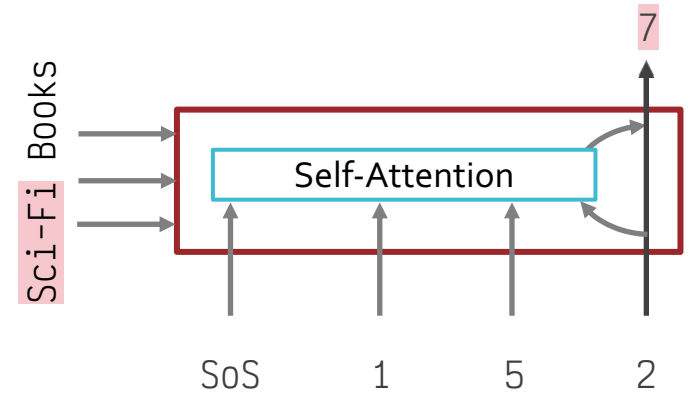
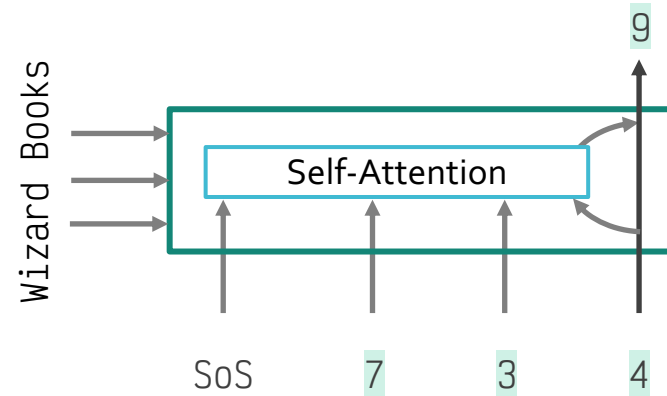
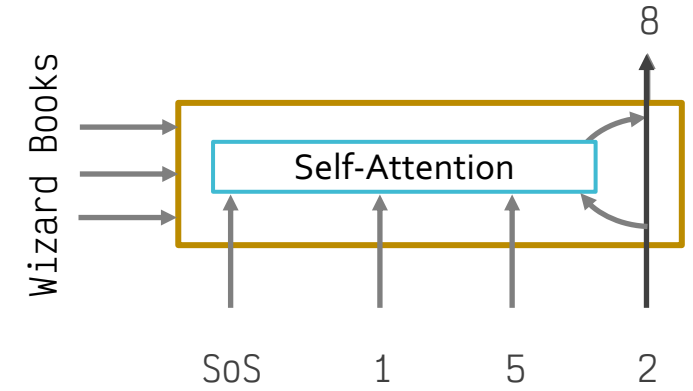


# Two Patching Experiments

### Previous Token Information

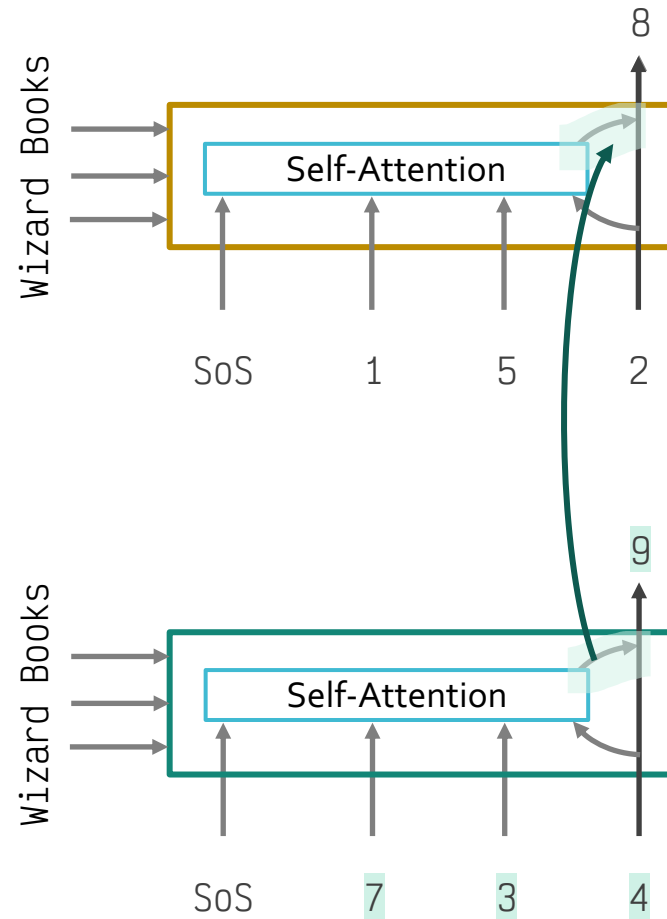


### Query Information

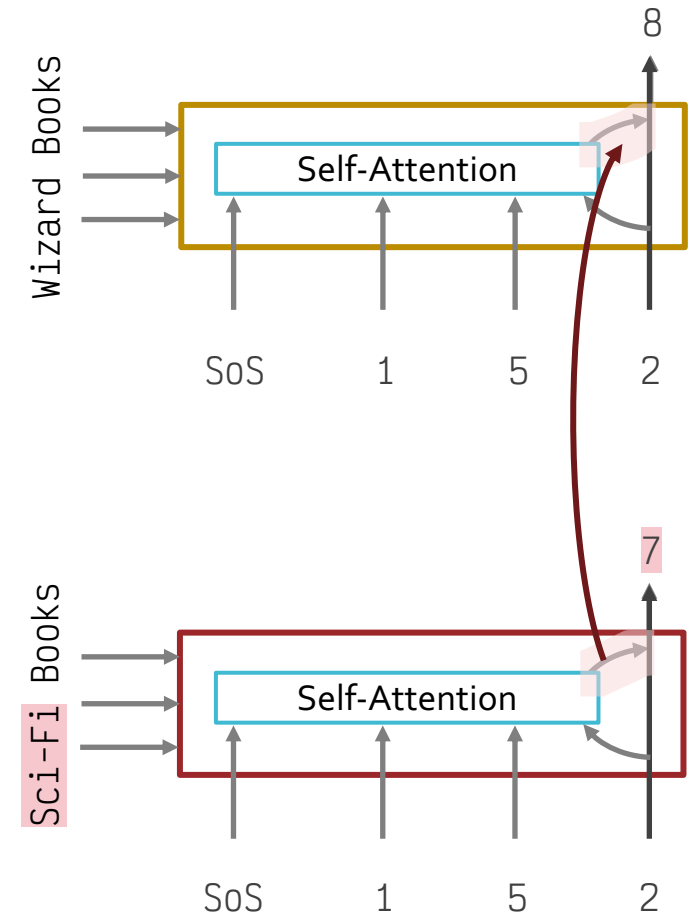


# Two Patching Experiments

### Previous Token Information

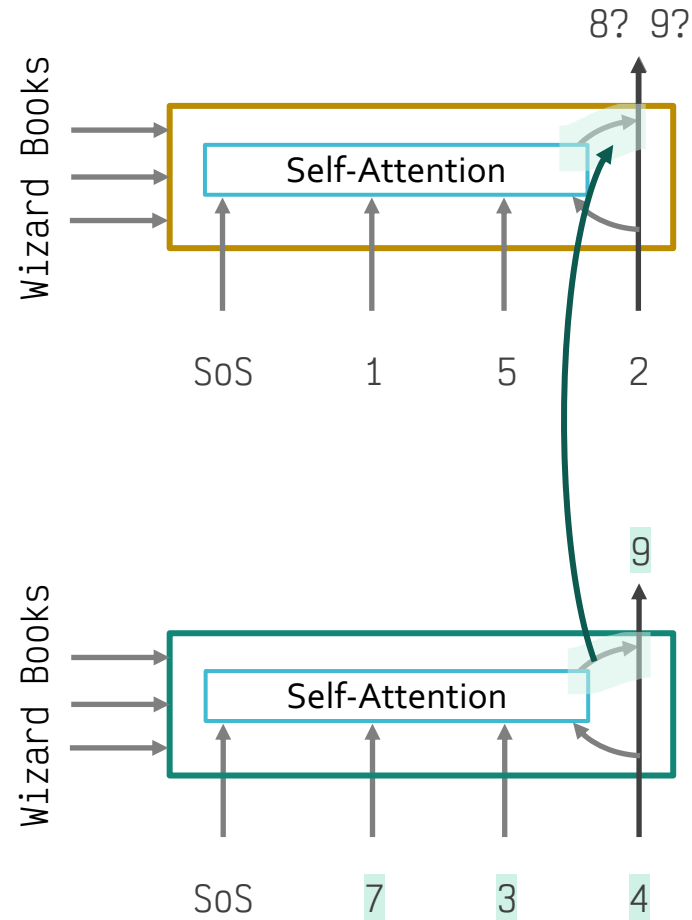


### Query Information

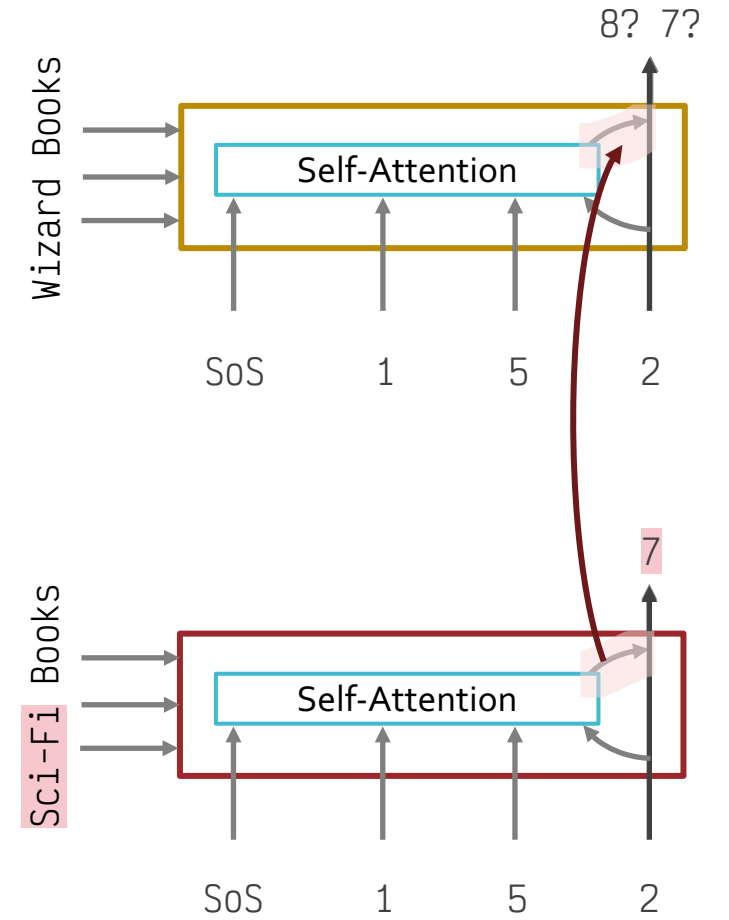


# Two Patching Experiments

## Previous Token Information

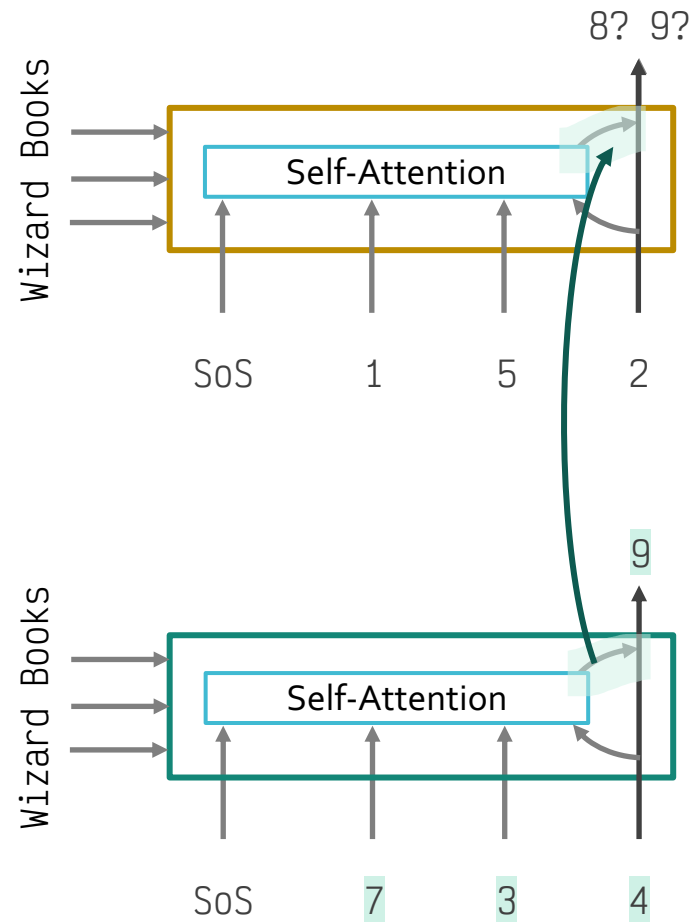


## Query Information

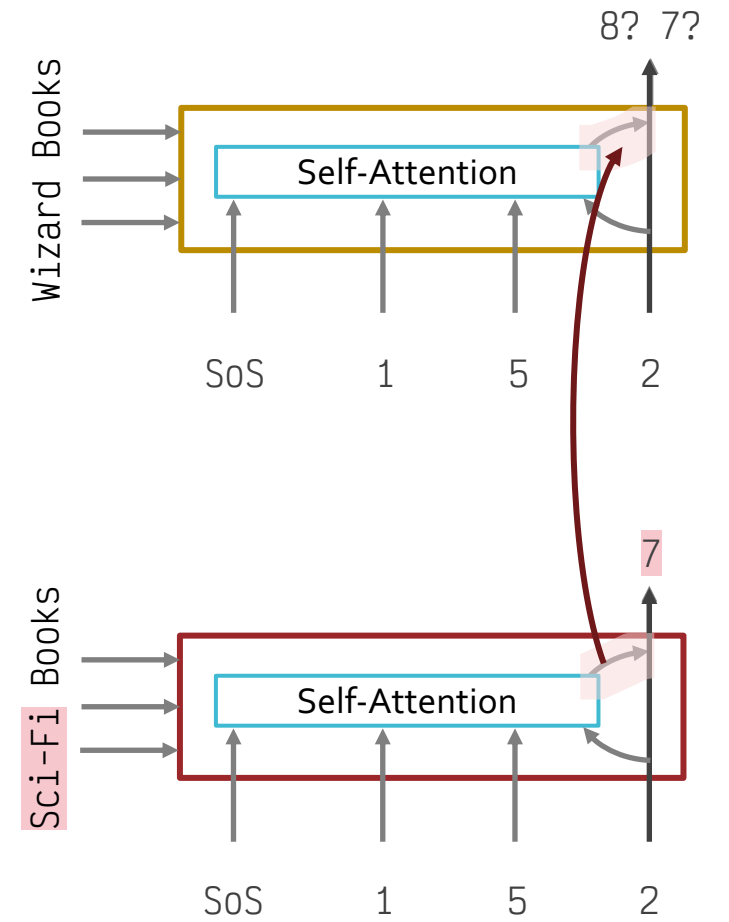


# Two Patching Experiments

### Previous Token Information

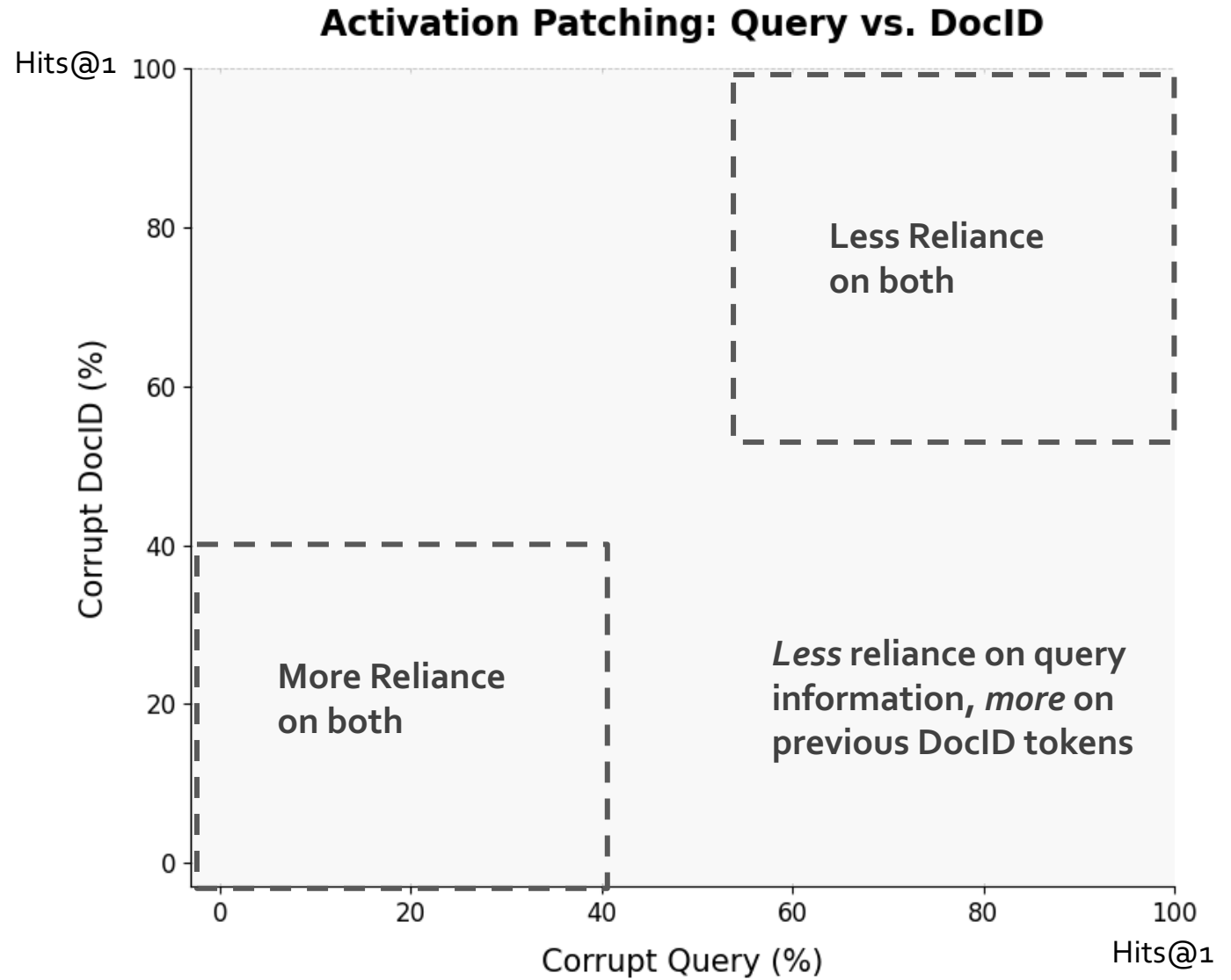


### Query Information



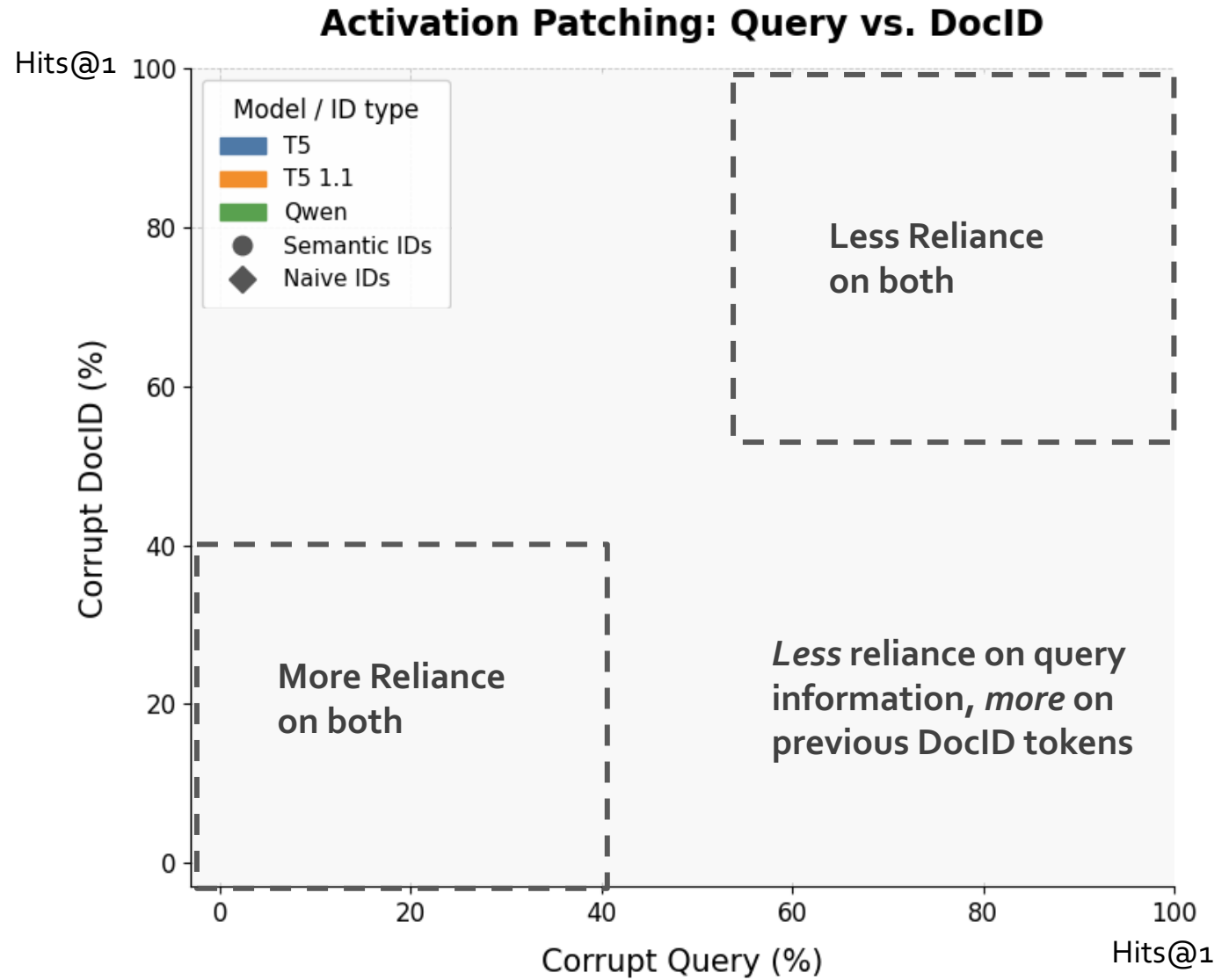
Worse Performance = More Reliance

# Two Patching Experiments



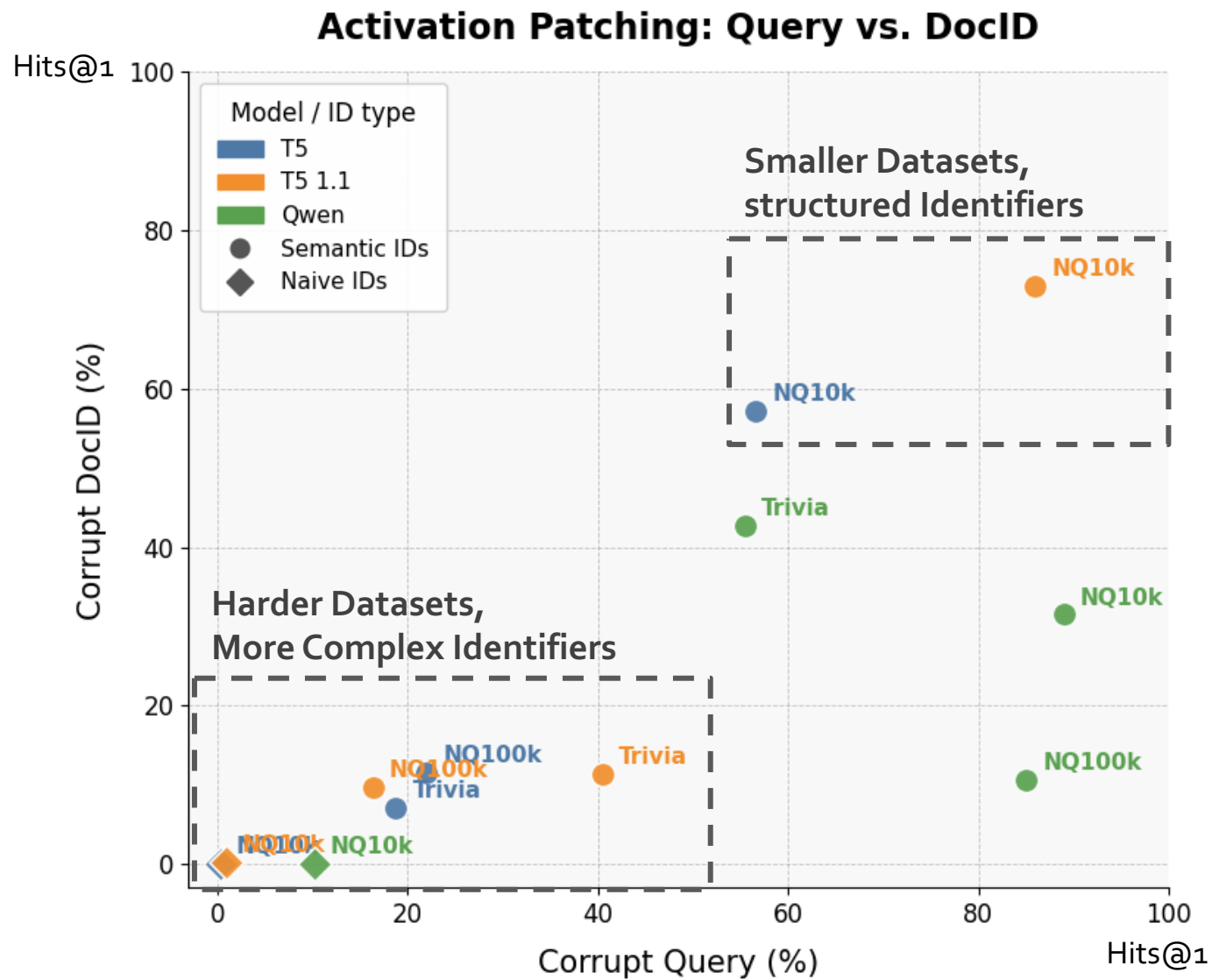
Worse Performance = More Reliance

# Two Patching Experiments



Worse Performance = More Reliance

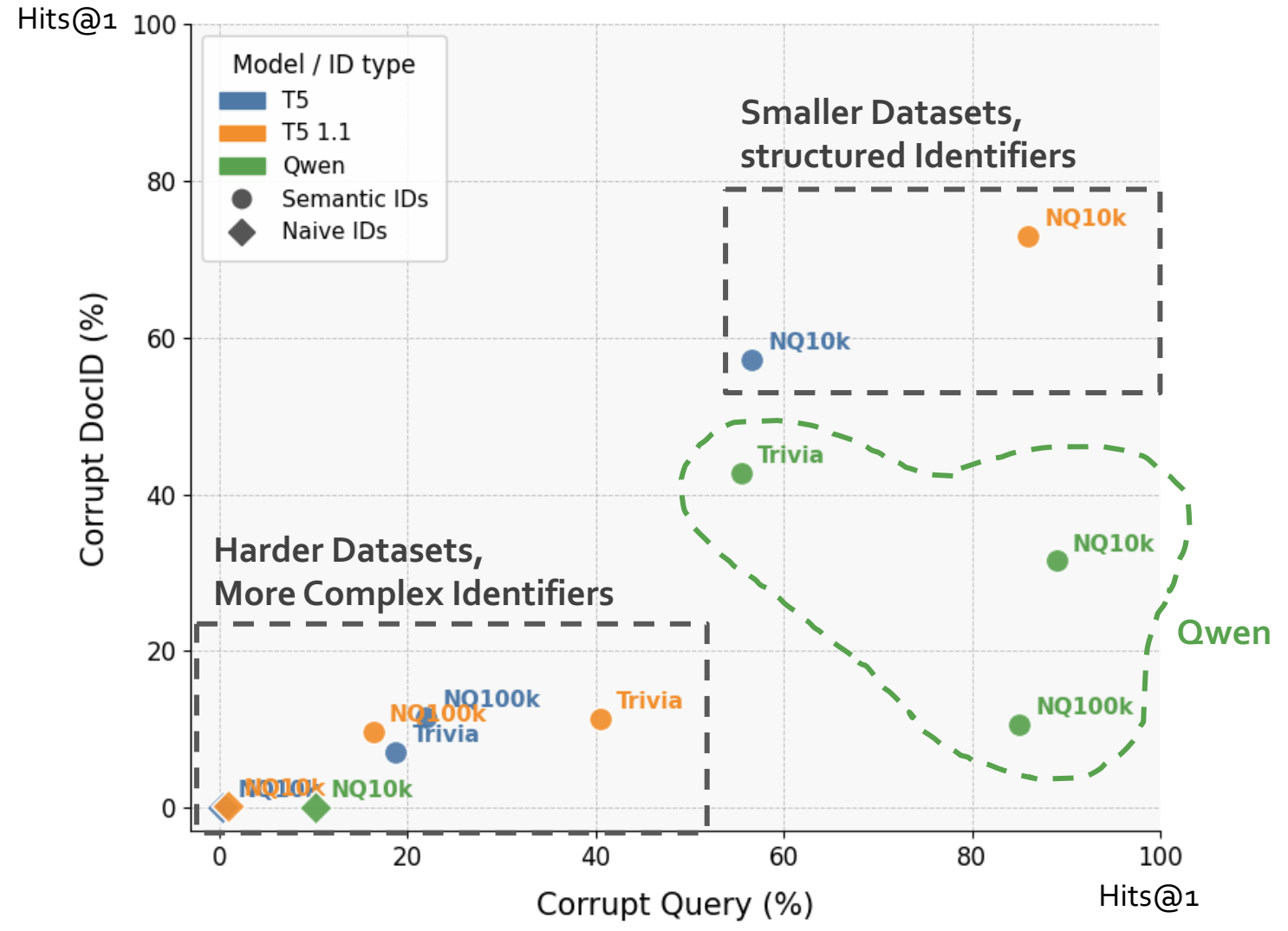
# Two Patching Experiments



Worse Performance = More Reliance

# Two Patching Experiments

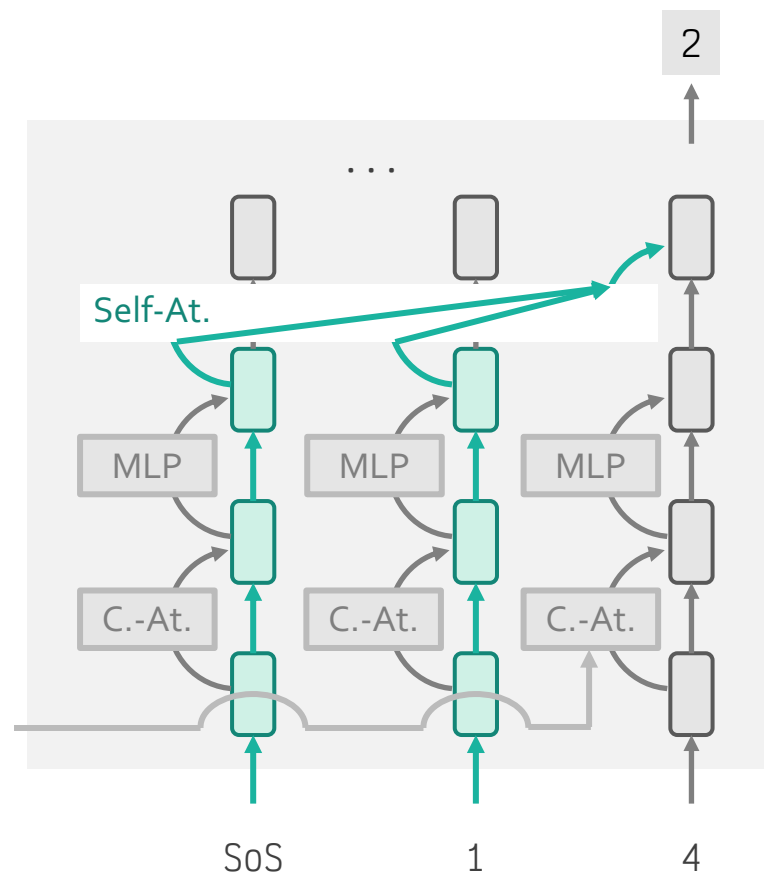
## Activation Patching: Query vs. DocID



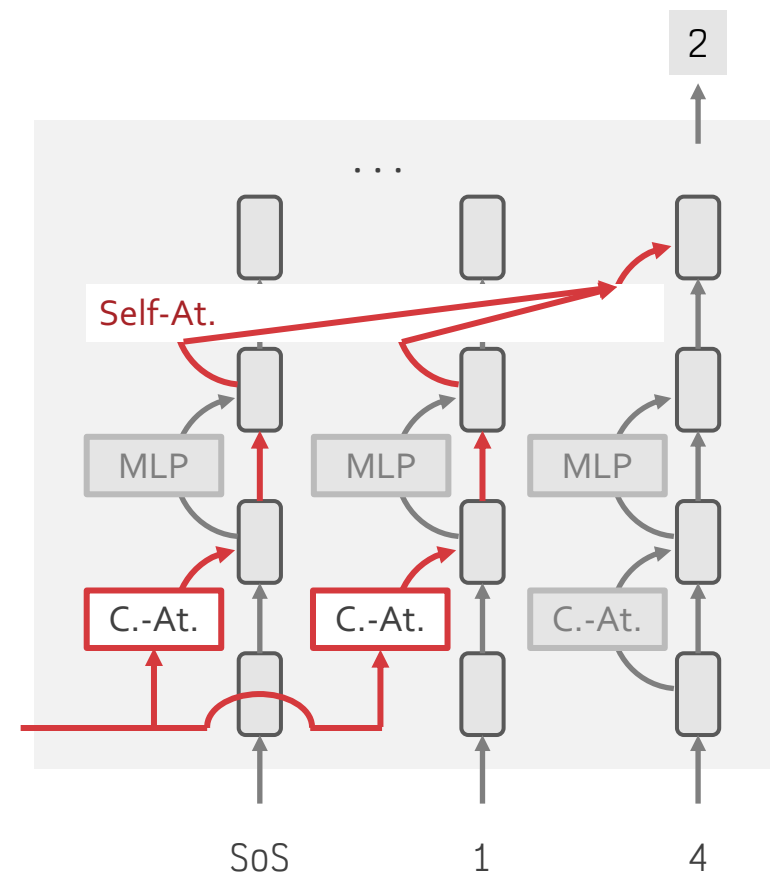
Worse Performance = More Reliance

# Two Possibilities

## Previous Token Information

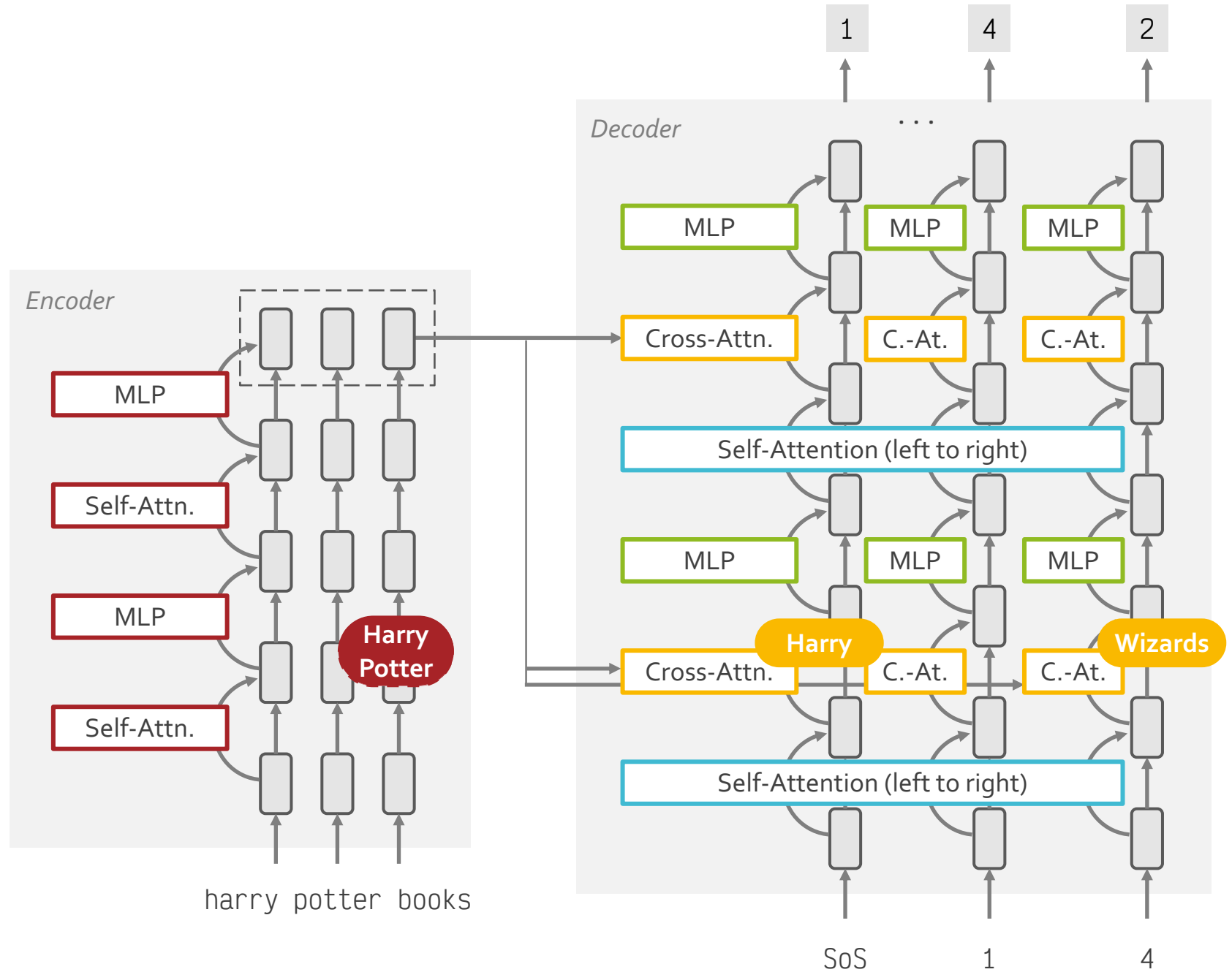


## Query Information

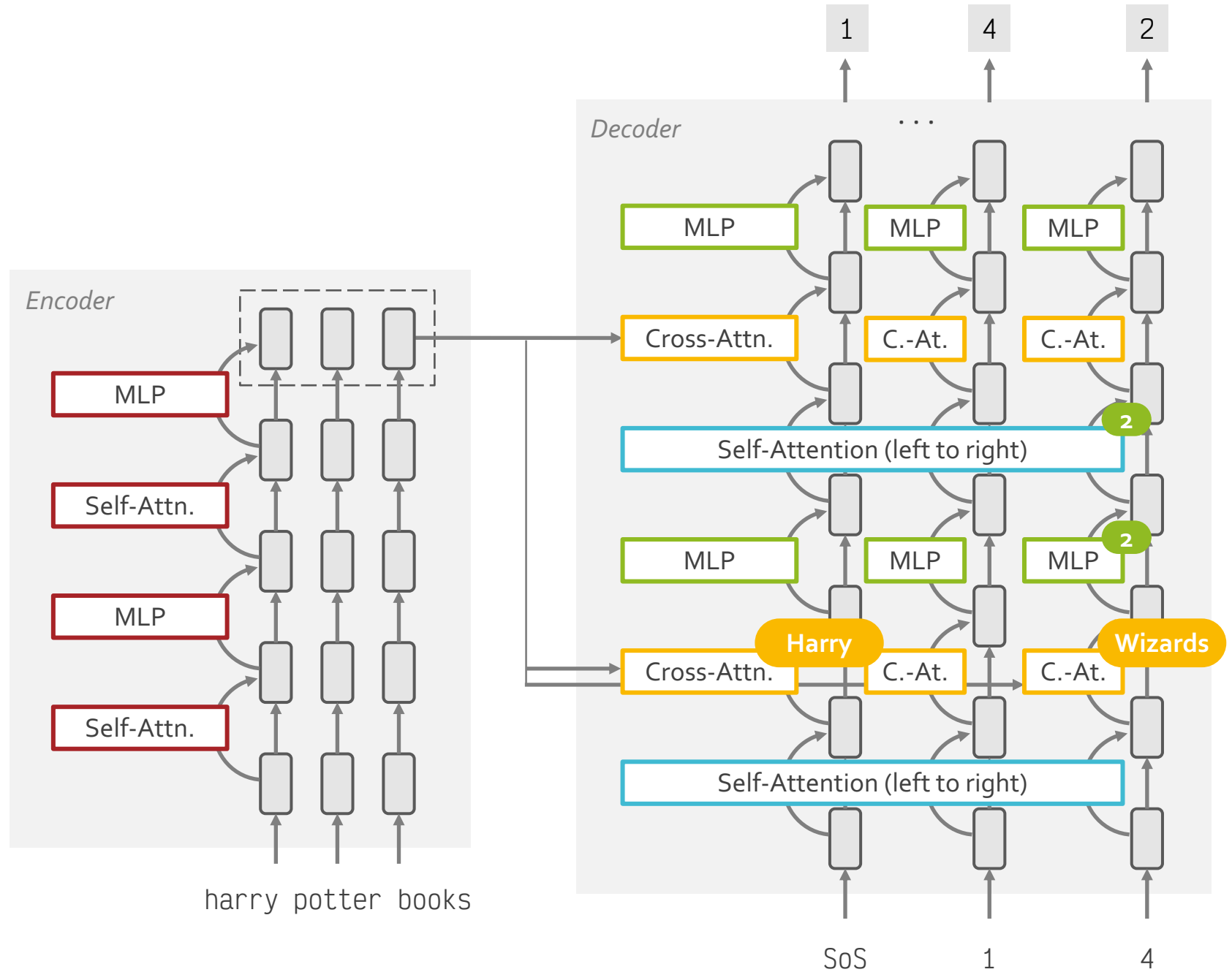


Self-Attention moves both information!

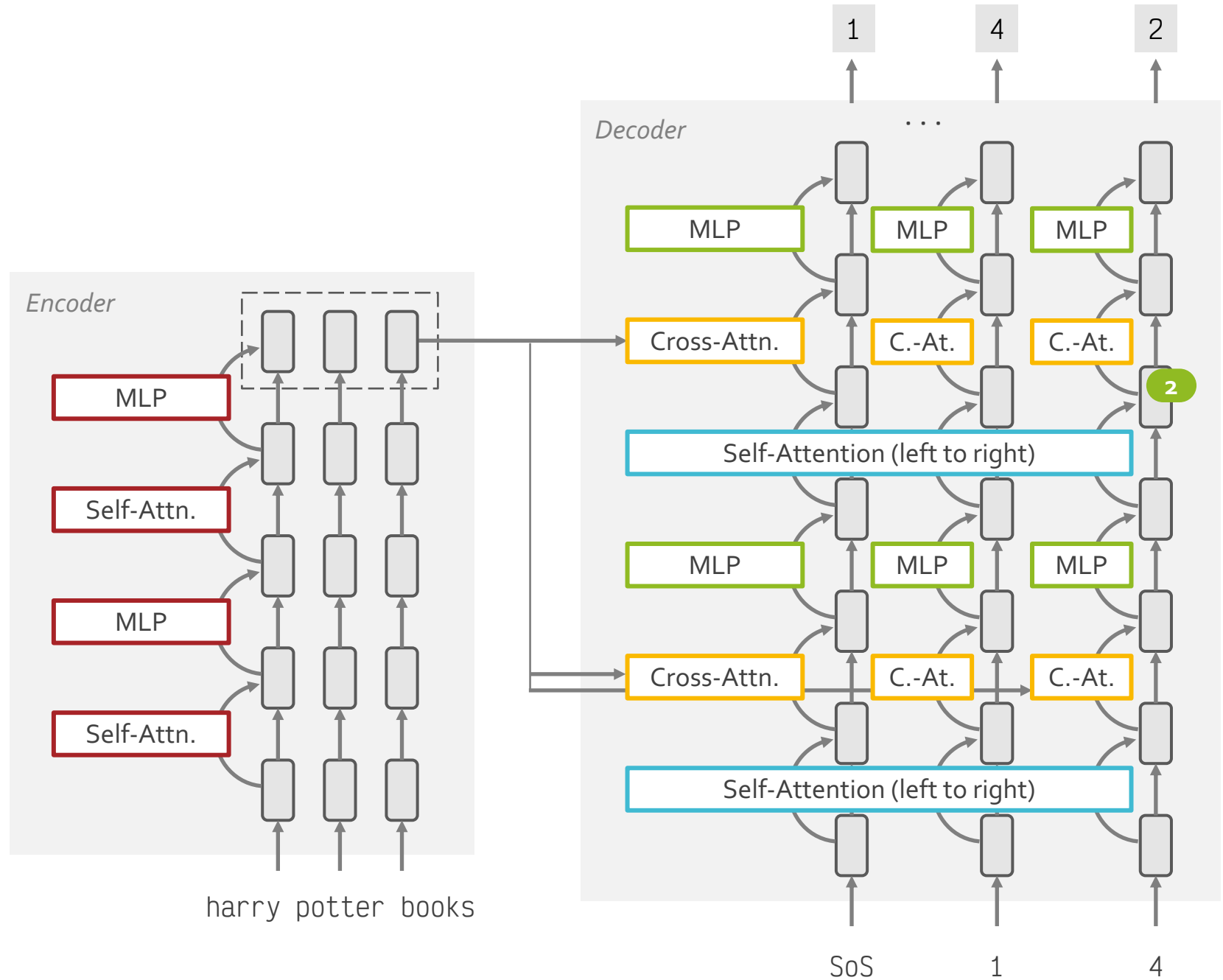
# Full Picture



# Full Picture



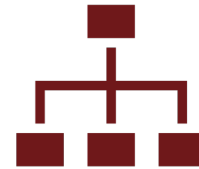
# Full Picture



# Summary



Analysis of Generative  
Retrieval Models



Adaption to Multi-  
Token Models



Faster, better and  
transparent retrieval

# Interpretability of Generative Information Retrieval

Anja Reusch – Postdoctoral Researcher – NLP Lab, Technion



**TECHNION**

