# Multi-Hop Reasoning for Question Answering with Hyperbolic Representations

**Simon Welz**[1]    **Lucie Flek**[1,2]    **Akbar Karimi**[1,2]

[1]Bonn-Aachen International Center for Information Technology, University of Bonn, Germany
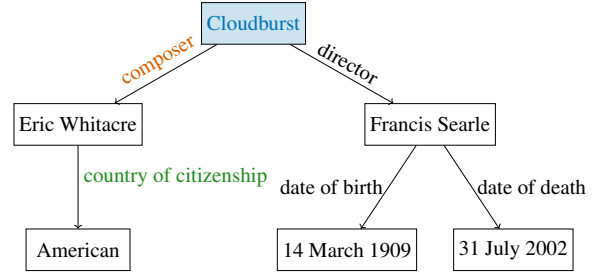[2]Lamarr Institute for ML and AI, Germany
`ak@bit.uni-bonn.de`

## Abstract

Hyperbolic representations are effective in modeling knowledge graph data which is prevalently used to facilitate multi-hop reasoning. However, a rigorous and detailed comparison of the two spaces for this task is lacking. In this paper, through a simple integration of hyperbolic representations with an encoder-decoder model, we perform a controlled and comprehensive set of experiments to compare the capacity of hyperbolic space versus Euclidean space in multi-hop reasoning. Our results show that the former consistently outperforms the latter across a diverse set of datasets. In addition, through an ablation study, we show that a learnable curvature initialized with the delta hyperbolicity of the utilized data yields superior results to random initializations. Furthermore, our findings suggest that hyperbolic representations can be significantly more advantageous when the datasets exhibit a more hierarchical structure.

## 1 Introduction

Multi-hop reasoning is a complex task that requires models to integrate information across multiple pieces of evidence to arrive at accurate conclusions. For instance, to answer *Which country is the composer of the song Cloudburst from?* using a simple knowledge graph like Figure 1, the model has to first find the answer to the first relation (composer of Cloudburst), and then look for the answer to the second relation (Whitacare's country of citizenship). This inherently involves traversing hierarchical relationships, making it particularly challenging for traditional language models that rely on Euclidean representations. While Euclidean representations are commonly used and can capture hierarchical structures to some extent (Nguyen et al., 2023; Misra et al., 2023; Zhang et al., 2024a), recent research has shown that hyperbolic representations are more effective in handling such data



**Question:** Which country is the composer of the song Cloudburst from?

**Path:** Cloudburst → composer → Eric Whitacre → country of citizenship → *American*

Figure 1: A knowledge graph with entities as the nodes and the relations as the edges, illustrating a 2-hop question answering process.

due to their superior ability to model hierarchical and relational information (Nickel and Kiela, 2017; Dhingra et al., 2018; Balažević et al., 2019; Chami et al., 2020; Xu et al., 2022).

In multi-hop reasoning tasks, hierarchical reasoning often manifests in navigating knowledge graphs or layered question answering frameworks, where the complexity increases with each additional hop. Given the hierarchical nature of multi-hop reasoning tasks, hyperbolic geometry presents a compelling alternative to Euclidean space for embedding representations. Specifically, hyperbolic space provides a larger capacity and more expressive way of encoding tree-like and graph-like structures (Chami et al., 2019), which are prevalent in knowledge graphs used for reasoning.

Although many studies successfully implement hyperbolic architectures and report performance gains (Zhou et al., 2021; Xiao et al., 2022; Wang et al., 2024), they often fail to disentangle the effects of geometric properties from the introduction of additional trainable parameters and the resulting model architecture.

A controlled comparison incorporating equiva-

lent Euclidean architectures with comparable parametric complexity would be necessary to isolate the geometric contribution to model performance. In addition, models in the literature often require significant architectural changes, which can increase model complexity and computational costs. In contrast, our approach focuses on incorporating hyperbolic geometry into existing language model architectures with minimal changes.

In this paper, **we address the lack of a carefully controlled comparative study for evaluating the differences between hyperbolic and Euclidean spaces in multi-hop reasoning**. As such, we incorporate hyperbolic representations into an encoder-decoder via the addition of a single layer and exponential mapping operation in the Poincaré ball model of hyperbolic space. We conduct a comprehensive set of experiments on multiple datasets, demonstrating that **adding a hyperbolic layer to increase the learning capability of a language model consistently outperforms its Euclidean counterpart**. In addition, we perform an ablation study to evaluate the impact of initialization of the curvature. Our results indicate that initializing the curvature using the $\delta$-hyperbolicity of the dataset leads to superior performance compared to random initialization. Furthermore, we show that the **performance gain from hyperbolic representations is more pronounced for datasets with more hierarchical structures** (defined based on the number of out-going relations for the nodes). Our comprehensive ablation studies deepen our understanding of geometric learning advantages in the context of language models and underscore the importance of aligning the geometric properties of the model with the inherent structure of the data.

## 2 Related Work

**Knowledge-based multi-hop reasoning.** Multi-hop reasoning requires traversing relational paths to synthesize new knowledge, making it essential for question answering. In the early approaches, path-based methods were utilized, in which reasoning was conducted using predefined rules or relational paths within the KB (Lao et al., 2011). Although these approaches were interpretable, they were frequently constrained by the availability and completeness of the KB. Neural-based reasoning models, including embedding-based methods (Bordes et al., 2013; Wang et al., 2014; Yang et al., 2015; Sun et al.), introduced vectorized representations of

entities and relations, enabling reasoning through learned relational patterns. More recent work has integrated graph-based neural architectures, such as Graph Convolutional Networks (Schlichtkrull et al., 2018) and Graph Attention Networks (Veličković et al., 2018), to propagate information across multi-hop relational structures in KBs. Reinforcement learning has been effectively applied to multi-hop reasoning over knowledge bases, enabling models to navigate complex relational paths and infer missing information. In this context, an RL agent is trained to traverse a KB by selecting a sequence of relations and entities, forming a reasoning path that leads to the desired answer (Ma et al., 2024; Wan et al., 2021; Lin et al., 2018; Zhu et al., 2022).

**Hyperbolic multi-hop reasoning.** Recent studies have proposed frameworks that leverage hyperbolic geometry to enhance multi-hop reasoning capabilities (Zhou et al., 2021; Xiao et al., 2022; Wang et al., 2024). Hyperbolic knowledge graph embeddings have demonstrated significant potential for multi-hop reasoning to model the hierarchical relationships in knowledge graphs (Chami et al., 2020; Balaževic et al., 2019; Montella et al., 2021; Kolyvakis et al., 2019). More recently, hyperbolic graph neural networks (HGNNs) have emerged as a promising direction for improving multi-hop reasoning (Liu et al., 2024, 2019; Chami et al., 2019). These models extend traditional GNNs by incorporating hyperbolic message passing, allowing for better hierarchical aggregation of multi-hop dependencies.

While these works highlight the potential of fully hyperbolic architectures, they often make significant architectural changes and fail to disentangle the advantages of hyperbolic geometry from those modifications in the model. By simply adding a hyperbolic layer to an existing language model, we efficiently integrate it without sacrificing the scalability or flexibility of these models. Furthermore, we utilize an identical model with the same number of trainable parameters to compare the two geometries, resulting in a deeper understanding of the characteristics of the two spaces.

## 3 Background

### 3.1 Multi-Hop Reasoning

Multi-hop reasoning can be defined as a process by which conclusions or answers are derived by sequentially combining information from multiple pieces of evidence. In contrast to single-hop

reasoning, which relies on direct connections between a query and its answer, multi-hop reasoning involves traversing intermediate steps or relationships to reach the outcome. This capability is essential for tasks that require analyzing interconnected data or reasoning through layered information. The process of multi-hop reasoning necessitates the retrieval of relevant pieces of information and their coherent integration, often involving the navigation of hierarchical relationships, temporal sequences, or contextual dependencies within data. A common approach to facilitate multi-hop reasoning is through the utilization of knowledge graphs, which represent entities and their relationships as a network. In a knowledge graph, reasoning involves following edges between entities to combine information across multiple nodes, thereby enabling complex inference over interconnected facts.

## 3.2 Poincaré Ball Model

As previous work has demonstrated the effectiveness of the Poincaré ball model (Nickel and Kiela, 2017; Ganea et al., 2018; Chami et al., 2020; Khrulkov et al., 2020; Chen et al., 2024) in capturing hierarchical relationships, we adopt it to enhance our ability to model such structures efficiently. The Poincaré ball provides a hyperbolic space where points are confined within the unit ball, enabling the representation of complex hierarchal structures with increasing precision near the boundary. Similarly to the approach in Nickel and Kiela (2017), we define the Poincaré ball model as $\mathbb{B}_c^n = \{x \in \mathbb{R}^n : c\|x\|^2 < 1, c \geq 0\}$. This space is equipped with a conformal factor given by: $\lambda_x^c = \frac{2}{1-c\|x\|^2}$ where the hyperparameter $c$ determines the curvature of the space, with larger values of $c$ corresponding to spaces of higher negative curvature.

**Möbius addition.** For a pair $x, y \in \mathbb{B}_c^n$, the Möbius addition is defined as follows:

$$x \oplus_c y := \frac{(1 + 2c\langle x,y \rangle + c\|y\|^2)x + (1 - c\|x\|^2)y}{1 + 2c\langle x,y \rangle + c^2\|x\|^2\|y\|^2}$$

where $\langle x, y \rangle$ denotes the Euclidean inner product, and $\|\cdot\|$ represents the Euclidean norm.

**Distance.** The induced distance function in this model is expressed as:

$$d_c(x, y) := \frac{2}{\sqrt{c}} \arctan(\sqrt{c}\| - x \oplus_c y\|)$$

which captures the geodesic distance between points $x$ and $y$ within the Poincaré ball.

To transition between Euclidean and hyperbolic spaces, we utilize the exponential and logarithmic mappings:

**Exponential mapping** maps a Euclidean vector $v \in T_o\mathbb{D}$ (the tangent space at the origin) to a point $y \in \mathbb{D}_c^n$ on the Poincaré ball:

$$\exp_0^c(v) = \frac{\tanh(\|v\| \cdot \sqrt{c})}{\|v\| \cdot \sqrt{c}} \cdot v$$

**Logarithmic mapping** maps a point $y \in \mathbb{B}_c^n$ back to the tangent space at the origin $T_o\mathbb{B}_c^n$:

$$\log_0^c(y) = \frac{\tanh^{-1}(\|y\| \cdot \sqrt{c})}{\|y\| \cdot \sqrt{c}} \cdot y$$

**Poincaré linear layer.** The Poincaré linear layer, adapted from Ryohei et al. (2021); van Spengler et al. (2023), extends the concept of a Euclidean linear layer into hyperbolic space, enabling models to effectively process hierarchical data. For an input $x \in \mathbb{D}_c^n$ the layer computes a hyperbolic transformation parameterized by weights $Z = \{z_k \in \mathbb{R}^n\}$ and biases $r = \{r_k \in \mathbb{R}\}_{k=1}^m$. The transformed output for each class $k$ is obtained through the following formulation of hyperbolic multinomial logistic regression:

$$v_k(x) =$$
$$\frac{2}{\sqrt{c}}\|z_k\| \sinh^{-1}\left(\lambda_x^c \langle \sqrt{c}x, \frac{z_k}{\|z_k\|} \rangle \cosh(2\sqrt{c}r_k)\right.$$
$$\left. - (\lambda_x^c - 1)\sinh(2\sqrt{c}r_k)\right) \quad (1)$$

where $\langle .,. \rangle$ represents the Euclidean inner product. The final output of the Poincaré linear layer is computed as:

$$y = \frac{w}{1 + \sqrt{1 + c\|w\|}},$$

where $w = (\frac{1}{\sqrt{c}}\sinh(\sqrt{c}v_k(x)))_{k=1}^m$.

## 3.3 Delta Hyperbolicity

Delta hyperbolicity ($\delta$-hyperbolicity) quantifies the extent to which a space is similar to a tree. This property makes it particularly relevant for analyzing and optimizing the curvature of multi-hop reasoning datasets.

**Gromov product.** The Gromov product is defined for points $x, y, w$ in a metric space $(X, d)$ as:

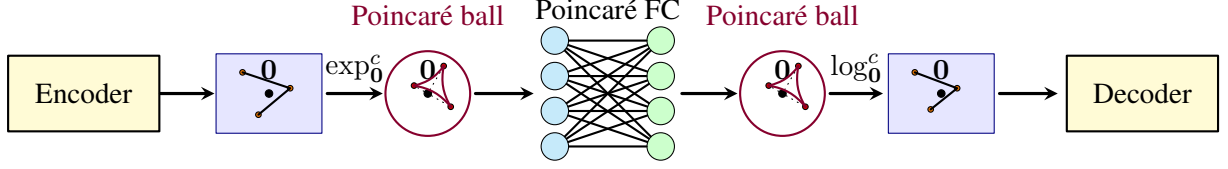$$(x, y)_w = \frac{1}{2}[d(x, w) + d(y, w) - d(x, y)].$$

Figure 2: Model architecture for our approach. After the T5 encoder, we project the resulting Euclidean embeddings onto the Poincaré ball. These hyperbolic embeddings are then refined through a trainable Poincaré layer. For compatibility with the T5 decoder, we project the hyperbolic embeddings back to the Euclidean space. While the T5 encoder and decoder parameters remain frozen throughout the training, the input contains trainable soft prompts.

A metric space is $\delta$-hyperbolic if, for any four points $w, x, y, z \in X$, the inequality:

$$(x, z)_w \geq \min\{(x, y)_w, (y, z)_w\} - \delta$$

is satisfied. Smaller $\delta$ values indicate a closer resemblance to a tree-like structure. In our work, we adopt the approach outlined in Khrulkov et al. (2020); Sawhney et al. (2024); Ermolov et al. (2022), using $\delta$-hyperbolicity as a scale-invariant measure to assess the hyperbolic nature of the dataset. Specifically, we estimate the hyperbolicity constant $\delta(X)$, which represents the smallest possible $\delta$ satisfying the four-point condition for all quadruples of points in $X$. To account for variations in the scale of the dataset, we compute the relative hyperbolicity as

$$\delta_{rel}(X) = \frac{2\delta(X)}{\text{diam}(X)},$$

where $\text{diam}(X)$ represents the diameter of the dataset, defined as the maximum pairwise distance between points. Since $\delta_{rel}(X)$ is normalized by the diameter of the dataset, it remains invariant under uniform rescaling of distances, ensuring comparability across datasets of different scales. By construction $\delta_{rel}(X) \in [0, 1]$, with values closer to zero indicating a strong resemblance to hyperbolic spaces. Using the estimated $\delta_{rel}(X)$, we compute the curvature $c(X)$ of the embedding space following the formula provided by Khrulkov et al. (2020):

$$c(X) = \left(\frac{0.144}{\delta_{rel}(X)}\right)^2 \quad (2)$$

This calculation enables us to determine the curvature hyperparameter $c$.

# 4 Method

Our approach builds on the PaTH method Misra et al. (2023), a two-step framework that fine-tunes the T5 model using soft prompts, added as trainable parameters to the input embeddings.

## 4.1 PaTH Method Overview

The PaTH method involves two primary stages of **knowledge integration** and **soft prompt tuning**. First, the T5 model is fine-tuned on the knowledge graph using the triples of the entity-relation-entity form $(e_1, r_1, e_2)$, enabling the model to internalize the foundational entity-relation structures. For each dataset, we only use the subgraph of triples relevant to our 2-hop questions (e.g., the paths connecting entities in each question) similar to Misra et al. (2023). In soft prompt tuning, two distinct soft prompts, called **parsing prompt** and **hopping prompt**, are trained to facilitate question parsing and reasoning tasks. The parsing prompt parses a question into an incomplete sequence $(e_1, r_1, r_2, ..., r_n)$, which serves as the input for the hopping prompt in the reasoning step. The hopping prompt is trained using uniform random walks over the knowledge graph. The walks from the dev and test sets are excluded. Given an incomplete sequence representing the starting entity and intermediate relations $(e_1, r_1, r_2, ..., r_n)$, the model is tasked with predicting the complete sequence, including the intermediate entities and relations $(e_1, r_1, e_2, r_2, ..., r_{n-1}, e_n)$. This enables the model to infer reasoning paths by referring to the incomplete path.

## 4.2 Incorporating Hyperbolic Representations

Figure 2 illustrates our simple integration of a hyperbolic layer into the T5 model.

The Euclidean embeddings generated by the T5 encoder are mapped onto the Poincaré ball using exponential mapping. Then they are processed through the hyperbolic layer, specifically designed for operations within the Poincaré space to preserve their geometric properties. After the transformation, the embeddings are mapped back to the Euclidean space using logarithmic mapping. This step enables compatibility with the T5 decoder for effective downstream processing. For the hyper-

| Dataset | Nodes | Edges | Relations |
|---|---|---|---|
| 2WikiMultiHopQA | 97,298 | 95,116 | 29 |
| MetaQA | 31,374 | 58,974 | 9 |
| MLPQ | 51,402 | 53,327 | 72 |
| PQ | 1,056 | 1,211 | 13 |

Table 1: Knowledge graph statistics of the datasets

| Dataset | Train | Dev | Test |
|---|---|---|---|
| 2WikiMultiHopQA | 72,760 | 8,085 | 6,768 |
| MetaQA | 47,108 | 5,951 | 5,942 |
| MLPQ | 57,283 | 7,160 | 7,161 |
| PQ | 1,698 | 210 | 191 |

Table 2: Number of questions in train/dev/test splits.

bolic operations and Poincaré layer, we use the open-sourced implementation given by van Spengler et al. (2023).

## 5   Experimental Setup

For all experiments, we used the T5-Large model (770M parameters) (Raffel et al., 2020). This model was fine-tuned using checkpoints adapted through the prefix LM objective (Liu et al., 2018) over 100,000 steps. We adopt the hyperparameters presented in Misra et al. (2023), with a modification to the batch size, reducing it to 64 to accommodate hardware limitations. This adjustment applies to both the knowledge integration and prompt tuning processes. The optimizer is AdaFactor (Shazeer and Stern, 2018) and for the additional hyperbolic layer, we use the same learning rate of 0.001 used to fine-tune the T5 model.

The curvature $c$ is initialized using Formula 2 in Section 3.3, which is based on the $\delta$-hyperbolicity of the dataset.

Since computing $\delta$-hyperbolicity can be computationally expensive we calculate it in batches. We sample 1500 points from the training dataset and compute $\delta_{rel}$. We repeat this process 5 times. For evaluation, we use the codebase of Ho et al. (2020), which is open-sourced[1]. Similarly to Misra et al. (2023), we evaluate the model performance with the Exact Match (EM) score.

### 5.1   Dataset Preparation

We use four datasets in a closed-book QA setting, where context was omitted to prioritize the reasoning capabilities of the model. The complete statistics for these datasets can be found in Tables 1 and 2. To ensure consistent evaluation across all datasets, we focus on the 2-hop questions.

**2WikiMultiHopQA**

(Ho et al., 2020), hereafter referred to as 2WikiHop for simplicity, consists of two-hop English questions constructed over a knowledge base containing 98,284 entities and 29 relations sourced

from WikiData (Vrandečić and Krötzsch, 2014).

Since the test splits of the 2WikiHop are private, the validation split was repurposed as the test set, with 10% of the training data reserved for validation. This adaptation mirrors the approach taken by Misra et al. (2023).

**MetaQA** consists of questions that can have multiple answers, different from 2WikiHop, where each question is associated with a single answer. For our study, we focused exclusively on a subset of the MetaQA dataset containing questions with a single possible answer to ensure consistency in the evaluation. In particular, MetaQA does not directly provide evidence for each question. To address this, we generated the necessary evidence for each question, as detailed in Appendix A.1. For this dataset, we used the official train/dev/test split[2].

**MLPQ** (Tan et al., 2023) consists of multilingual questions paired with corresponding language-specific knowledge graphs. For our study, we use the evidence (paths) of the dataset as a unified knowledge graph, resulting in a total of 51,401 entities and 72 relations. We specifically focus on files that contain English questions alongside potential French-language entities and relations. To maintain consistency during question parsing, French relations are translated into English. The dataset is divided into training, validation, and test sets with a ratio of 8:1:1, resulting in 57,283 questions for training, 7,160 for validation, and 7,161 for testing. Although the evidence parts are structured as triples, it is worth noting that due to the multilingual nature of the dataset, the tail of the first evidence may not always match the head of the subsequent evidence. To address this, we normalize by always selecting the English entity to construct the knowledge graph.

The Path Questions (**PQ**) dataset (Zhou et al., 2018a) is a QA dataset designed for multi-hop reasoning, leveraging entity relationships sourced from a knowledge base called Freebase (Bollacker et al., 2008). Our focus is on the 2-hop reasoning

---

[1]https://github.com/Alab-NII/2wikimultihop

[2]https://github.com/yuyuz/MetaQA

| Data | Model | 2WikiHop | MetaQA | MLPQ | PQ |
|------|-------|----------|--------|------|-----|
| Dev | Euclidean | 44.36 | 22.92 | 81.03 | 18.28 |
| Dev | Hyperbolic | **46.93** | **28.33** | **82.60** | **29.03** |
| Test | Euclidean | 14.88 | 19.76 | 72.10 | 11.90 |
| Test | Hyperbolic | **15.20** | **25.40** | **74.58** | **23.21** |

Table 3: Exact match scores for hopping prompt

| Data | Model | 2WikiHop | MetaQA | MLPQ | PQ |
|------|-------|----------|--------|------|-----|
| Dev | Euclidean | 88.60 | 95.51 | 97.08 | **100** |
| Dev | Hyperbolic | **89.34** | **95.65** | **97.14** | **100** |
| Test | Euclidean | 79.24 | **95.27** | 95.91 | **98.95** |
| Test | Hyperbolic | **80.11** | 95.07 | **96.79** | **98.95** |

Table 4: Exact match scores for parsing prompt

| Parsing | Hopping | 2WikiHop | MetaQA | MLPQ | PQ |
|---------|---------|----------|--------|------|-----|
| Euclidean | Euclidean | 13.39 | 19.20 | 72.59 | 12.04 |
| Hyperbolic | Euclidean | 13.56 | 19.08 | 72.74 | 12.04 |
| Euclidean | Hyperbolic | 13.40 | **24.74** | **73.48** | **23.04** |
| Hyperbolic | Hyperbolic | **13.65** | 24.72 | 73.40 | 22.51 |

Table 5: Exact match scores on test set for T5 with the additional Euclidean/hyperbolic layer for both prompts.

| Data | Model | 2WikiHop | MetaQA | MLPQ | PQ |
|------|-------|----------|--------|------|-----|
| Dev | Euclidean | 32.43 | 13.89 | 78.74 | 12.76 |
| Dev | Hyperbolic | **34.22** | **17.09** | **79.64** | **27.42** |
| Test | Euclidean | 10.24 | 12.08 | 70.36 | 11.90 |
| Test | Hyperbolic | **10.70** | **15.10** | **72.12** | **22.02** |

Table 6: Exact match scores for hopping stage when applying the additional layer without soft prompts.

subset, which comprises 1,908 questions, their corresponding answers, and the reasoning paths used to derive them. We adopt the same dataset split as Wang et al. (2024), which is an 8:1:1 ratio for training, dev, and test sets, respectively. However, contrary to Wang et al. (2024), we exclude the reasoning walks found in the dev and test splits during training, making the task more challenging since all supporting evidence present in the dev and test sets was also part of their training split.

We chose PQ over the similar PQL (Zhou et al., 2018b) dataset since PQ's smaller size allowed us to stay within our computational budget while capturing the same multi-hop reasoning patterns.

# 6 Results

In this section, we compare the two spaces in a variety of settings, provide an ablation study on the curvature, show that computationally both cases are similar, present the results of distance analysis in the two spaces, and give some insights into dataset difficulty.

## 6.1 Hyberbolic vs. Euclidean Layer

Table 3 presents the results across all datasets under prompt tuning for the hopping prompt. The results demonstrate that the hyperbolic layer outperforms the Euclidean counterpart across all datasets. For 2WikiHop, the hyperbolic layer increases the EM score from 44.36% (Euclidean) to 46.93%, reflecting a performance boost of 2.57%. Similarly, for MetaQA, the hyperbolic layer achieves an exact match (EM) score of 28.33%, outperforming the Euclidean layer, which achieves 22.92%, resulting in an improvement of 5.41%.

Interestingly, the smallest improvement occurs

in MLPQ, where the hyperbolic layer increases the EM score from 81.03% to 82.60%, marking a marginal gain of 1.57%. This limited improvement could be attributed to the already high baseline performance achieved by the Euclidean layer in this dataset. With less room for improvement, the hierarchical modeling advantages of the hyperbolic layer are less pronounced. Notably, MLPQ's predominantly linear knowledge-graph structure—with over 80% of nodes having an out-degree of one (see Figure 6)—constrains the benefits of hyperbolic space on this dataset. For the same reason, we can also see marginal improvements for the parsing prompt in Table 4. This contrasts with datasets like MetaQA and PQ (for the hopping prompt in Table 3), where a lower baseline provides more opportunities for substantial gains. More importantly, given that the parsing task is independent of the knowledge graph structure, it might not inherently benefit from the hierarchical properties of the hyperbolic space.

Table 5 compares all configurations of hyperbolic and Euclidean layers for the parsing and hopping stages. The results show that in most cases the use of hyperbolic space for the hopping stage gives the highest performance improvements. This finding is expected as the hopping stage is influenced by the knowledge graph hierarchies, which are better captured in hyperbolic space. Notably, while the Euclidean-hyperbolic configuration for parsing-hopping stages achieves the best performance in most cases, the hyperbolic-hyperbolic configuration follows closely, with only a marginal difference (ranging from 0.02 to 0.47 EM points). This suggests that while hyperbolic hopping sig-
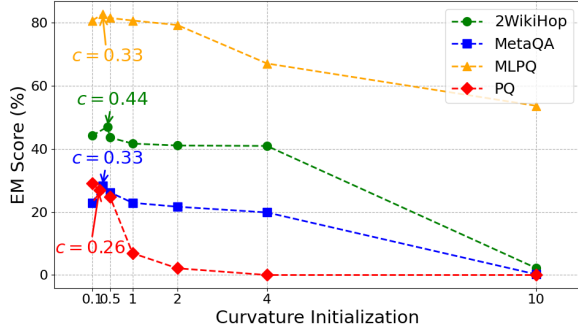
Figure 3: Curvature ablation for random walk training showing exact match score on dev sets. Initializing the curvature with or around $\delta$-hyperbolicity gives the highest EM score.

|  | **2WikiHop** | **MetaQA** | **MLPQ** | **PQ** |
|---|---|---|---|---|
| *Random Walk Dataset (Hopping Prompt)* | | | | |
| $\delta$ | $0.22_{\pm 0.012}$ | $0.25_{\pm 0.015}$ | $0.25_{\pm 0.013}$ | $0.28_{\pm 0.017}$ |
| $c$ | 0.44 | 0.33 | 0.33 | 0.26 |
| *Parsing Dataset (Parsing Prompt)* | | | | |
| $\delta$ | $0.29_{\pm 0.017}$ | $0.33_{\pm 0.018}$ | $0.29_{\pm 0.020}$ | $0.28_{\pm 0.019}$ |
| $c$ | 0.25 | 0.19 | 0.25 | 0.26 |

Table 7: Mean $\delta$-hyperbolicity and curvature values for random walks and parsing prompt data.

nificantly improves the performance, the choice of space for the parsing stage could also have a relatively small impact.

Table 6 presents the results of the additional Euclidean and hyperbolic layers on the random walk dev set without the use of soft prompts. The results demonstrate that the hyperbolic layer consistently outperforms the Euclidean layer, achieving higher scores on all datasets. This indicates that the performance observed is not dependent on soft prompting, as the hyperbolic layer exhibits superior results even in its absence.

## 6.2 Curvature Ablation

One hyperparameter of the hyperbolic layer is the curvature, which can be initialized arbitrarily. Figure 3 presents the results of the curvature ablation study with different initializations. A key takeaway from this study is that initialization plays a crucial role in model performance. Specifically, setting the curvature based on the relative $\delta$-hyperbolicity (see Section 3.3) of each dataset, as shown in Table 7, yields the best (or very close to it for PQ) across all datasets. In contrast, while smaller curvatures such as 0.1 and 1.0 still yield competitive results, increasing the curvature beyond 1.0 leads to a no-
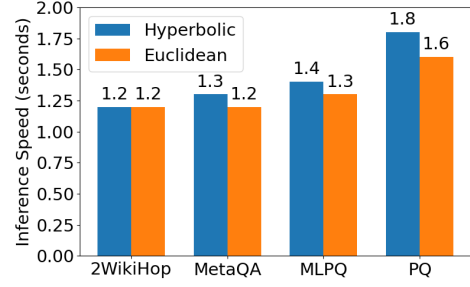


Figure 4: Average inference time per batch (=8) on the test data. The hyperbolic layer causes a negligible increase in inference time over the Euclidean layer.

| Hopping | 2WikiHop | MetaQA | MLPQ | PQ |
|---|---|---|---|---|
| First relation | 100 | 100 | 81.25 | 100 |
| Second relation | 100 | 100 | 69.01 | 99.46 |

Table 8: Percentage of cases where the geodesic distance in hyperbolic space between the source entity and its relations is larger than the Euclidean distance.

table degradation in EM scores. For instance, with a curvature of 10.0, the EM scores drop drastically to $2.21\%$ for 2WikiHop and $0.22\%$ for MetaQA, demonstrating that inappropriate curvature values can severely impact model effectiveness. These findings suggest that hyperbolic models benefit from curvature settings that reflect the structure of the data. Since hyperbolic space expands exponentially, setting the curvature to match a dataset's $\delta$-hyperbolicity allows the model to better reflect hierarchical relationships, thereby improving multihop reasoning accuracy. In contrast, Euclidean space lacks this adaptability, making it less effective when reasoning over complex data in knowledge graphs.

## 6.3 Computational Analysis

Another crucial aspect of this study is the computational complexity associated with hyperbolic layers. Despite improved performance, adding a hyperbolic layer introduces negligible time and memory overhead as shown in Figure 4. This observation is significant because it demonstrates that hyperbolic layers can achieve superior performance without significantly increasing the computational cost of the model. This makes hyperbolic layers a practical and efficient choice for tasks involving graph-structured data, such as multi-hop reasoning.

## 6.4 Embedding Distances

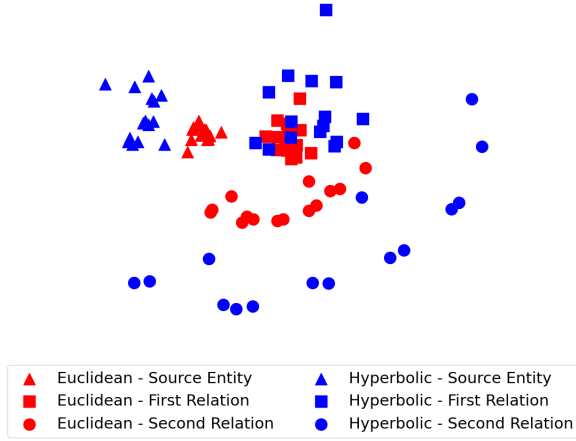To investigate how the distance between the source entities and their relations compare to each other in

Figure 5: Embeddings for 15 input samples from the MetaQA dataset, each structured as "source entity; first relation; second relation" in Euclidean versus Poincaré layer. The Euclidean embeddings use Euclidean distance while the hyperbolic embeddings use geodesic distance. Due to the exponential growth of the hyperbolic space, entities and relations can be more spread out as the paths become longer.
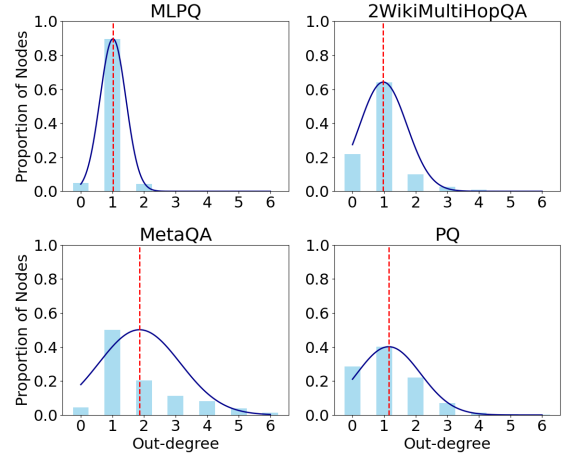


Figure 6: Distribution of out-going relations (out-degree) for each dataset. As the proportion of nodes with a degree of 2 or higher goes up, the complexity of the dataset also increases.

the Euclidean and hyperbolic layers, we looked into their embeddings in these layers. Table 8 presents a comparison between hyperbolic and Euclidean embeddings in terms of their geodesic and Euclidean distances. Specifically, it reports the percentage of cases where the geodesic distance in hyperbolic space is larger than the Euclidean distance for the Euclidean embeddings. The comparison is conducted for both the first and second relational hops with respect to the source entity.

For 2WikiHop, MetaQA, and PQ datasets, it is evident that the hyperbolic relation embeddings consistently exhibit a greater distance from the source entity in comparison to the Euclidean embeddings, as evidenced by almost 100% of the cases. However, MLPQ demonstrates a notable decrease in these percentages for both the first (81.25%) and the second (69.01%) relational hops. This behavior can be attributed to the structural characteristics inherent in the MLPQ knowledge graph. The fact that over 80% of its nodes have an out-degree of 1 (only one out-going relation as seen in Figure 6) indicates that MLPQ is predominantly a linear knowledge graph rather than a hierarchical one. The hyperbolic space is particularly beneficial for tree-like structures, where distances expand exponentially with branching. However, in a mostly linear graph, entities and relations are more evenly spaced out, meaning that Euclidean space can capture these relationships almost just as effectively.

Since MLPQ lacks significant tree-like expansion, its hyperbolic distances are not consistently larger than Euclidean distances, leading to significantly lower percentages compared to the other datasets.

The results confirm that for tree-like knowledge graphs, the hyperbolic geodesic distance typically exceeds the Euclidean distance. This outcome is expected due to the exponential expansion of hyperbolic space. Unlike Euclidean space, where distance scales linearly, hyperbolic space exhibits exponential growth, allowing entities and relations to be more sparsely distributed (Figure 5), which, in turn, makes it easier for the model to find relevant paths.

Such geometric properties yield a better theoretical justification for the performance gains seen in hyperbolic models. By allowing for increased spatial separation between entities along a path, hyperbolic space lessens interference between rival paths and improves the model's capability to learn effective reasoning multi-hop chains. This has a significant impact in scenarios where disambiguation between relation paths proves central—something Euclidean space has difficulty with given its linear growth and reduced ability to represent hierarchical branching. Hence, the benefit is not simply in numeric improvements, but rather in how the geometry changes the landscape of the embedding to better reflect the nature of reasoning problems.

## 6.5 Dataset Difficulty

Depending on the knowledge graph, datasets can have different levels of difficulty. Figure 6 presents

the proportion of nodes in each dataset with their out-going relations (out-degree) in their knowledge graphs. The MLPQ dataset is comparatively simpler, with over 80% of its nodes having an out-degree of 1. This characteristic significantly reduces the complexity of navigating the graph, as there is typically only one possible path from a given source node. In contrast, the MetaQA is more challenging, with only 50% of its nodes having an out-degree of 1 while more than 40% possess an out-degree of 2 or higher. The presence of multiple paths increases ambiguity, making traversal more complex and negatively impacting the performance, particularly in the random walk stage.

## 7   Conclusion and Future Work

We carried out a rigorous and careful investigation of using hyperbolic versus Euclidean representations in multi-hop reasoning and showed some of the advantages of the former compared to the latter.

Our experiments also confirm that initializing the curvature using the relative delta hyperbolicity of the dataset provides a robust and effective starting point for learning, ensuring that the model captures the hierarchical relationships within the data with greater accuracy. We also provided evidence for the hyperbolic geometry showing more effectiveness when the dataset has more hierarchical characteristics. These findings underscore the importance of understanding the structural properties of the data when selecting appropriate model architectures. In addition, our findings open several promising avenues for future research.

Given that the number of outgoing relations for each node in different knowledge graphs is not equal in many cases, future work could investigate a general manifold structure as well as other hyperbolic spaces for multi-hop reasoning.

While our current framework focuses on encoder-decoder models in a closed-book QA setting, future work should investigate the generalizability of hyperbolic representations across broader architectures and tasks. First, extending our approach to decoder-only language models would help assess the geometric advantages in other generative models. Second, applying hyperbolic reasoning layers to open-book QA tasks, where models retrieve and integrate external evidence, would clarify the interaction between geometric embedding space and retrieval-based reasoning. Finally, examining our method on a wider variety of datasets

and QA formats—including multi-lingual, noisy, or longer-hop reasoning datasets—will be critical to understanding the full scope of its effectiveness and limitations.

## 8   Limitations

Although our approach shows promising results, it has certain limitations: First, we focus exclusively on the closed-book QA setting, where no external context is provided to the model. This limitation inherently limits the amount of information available to answer questions, as the model relies solely on its trained knowledge. As a result, our approach may underperform compared to models that use additional context, such as open-book (Jiang et al., 2022; Feng et al., 2020; Xu et al., 2021) or retrieval-augmented (Shi et al., 2024a,b; Zhang et al., 2024b) methods, which can provide more relevant information during inference. Second, our experiments were conducted using a frozen model, where only a small number of parameters in the additional layer were fine-tuned. While this approach reduces computational cost and maintains efficiency, it may limit the ability of the models that require full fine-tuning for higher accuracy. In such cases, the impact of one hyperbolic layer with only one million trainable parameters might fade away compared to a billion parameters.

## 9   Ethics Statement

In this study, we exclusively employ pre-trained knowledge from the T5 model and the datasets utilized in our experimental setup. We have solely transformed the data as outlined in Section 5.1, without introducing or curating additional external knowledge sources. However, it is crucial to acknowledge that the datasets may contain biased, inaccurate, or incomplete information, which could influence the model's reasoning and outputs. Furthermore, these datasets may inadvertently include private or sensitive information that has not been explicitly identified, as addressing such concerns is beyond the scope of this study.

## Acknowledgments

## References

Ivana Balažević, Carl Allen, and Timothy Hospedales. 2019. Multi-relational poincaré graph embeddings. *Preprint*, arXiv:1905.09791.

Kurt Bollacker, Colin Evans, Praveen Paritosh, Tim Sturge, and Jamie Taylor. 2008. Freebase: a collaboratively created graph database for structuring human knowledge. In *Proceedings of the 2008 ACM SIGMOD international conference on Management of data*, pages 1247–1250.

Antoine Bordes, Nicolas Usunier, Alberto Garcia-Duran, Jason Weston, and Oksana Yakhnenko. 2013. Translating embeddings for modeling multi-relational data. In *Advances in Neural Information Processing Systems*, volume 26. Curran Associates, Inc.

Ines Chami, Adva Wolf, Da-Cheng Juan, Frederic Sala, Sujith Ravi, and Christopher Ré. 2020. Low-dimensional hyperbolic knowledge graph embeddings. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 6901–6914.

Ines Chami, Zhitao Ying, Christopher Ré, and Jure Leskovec. 2019. Hyperbolic graph convolutional neural networks. In *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc.

Nan Chen, Xiangdong Su, and Feilong Bao. 2024. Hyperbolic representations for prompt learning. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 8487–8492, Torino, Italia. ELRA and ICCL.

Bhuwan Dhingra, Christopher Shallue, Mohammad Norouzi, Andrew Dai, and George Dahl. 2018. Embedding text in hyperbolic spaces. In *Proceedings of the Twelfth Workshop on Graph-Based Methods for Natural Language Processing (TextGraphs-12)*, pages 59–69.

Aleksandr Ermolov, Leyla Mirvakhabova, Valentin Khrulkov, Nicu Sebe, and Ivan Oseledets. 2022. Hyperbolic vision transformers: Combining improvements in metric learning. In *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 7399–7409. IEEE.

Yanlin Feng, Xinyue Chen, Bill Yuchen Lin, Peifeng Wang, Jun Yan, and Xiang Ren. 2020. Scalable multi-hop relational reasoning for knowledge-aware question answering. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1295–1309.

Octavian-Eugen Ganea, Gary Bécigneul, and Thomas Hofmann. 2018. Hyperbolic neural networks. In *Proceedings of the 32nd International Conference on Neural Information Processing Systems*, pages 5350–5360.

Xanh Ho, Anh-Khoa Duong Nguyen, Saku Sugawara, and Akiko Aizawa. 2020. Constructing a multihop qa dataset for comprehensive evaluation of reasoning steps. In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 6609–6625, Barcelona, Spain (Online). International Committee on Computational Linguistics.

Zhengbao Jiang, Jun Araki, Haibo Ding, and Graham Neubig. 2022. Understanding and improving zero-shot multi-hop reasoning in generative question answering. In *Proceedings of the 29th International Conference on Computational Linguistics*, pages 1765–1775.

Valentin Khrulkov, Leyla Mirvakhabova, Evgeniya Ustinova, Ivan Oseledets, and Victor Lempitsky. 2020. Hyperbolic image embeddings. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 6418–6428.

Prodromos Kolyvakis, Alexandros Kalousis, and Dimitris Kiritsis. 2019. Hyperkg: Hyperbolic knowledge graph embeddings for knowledge base completion. *arXiv preprint arXiv:1908.04895*.

Ni Lao, Tom Mitchell, and William W. Cohen. 2011. Random walk inference and learning in a large scale knowledge base. In *Proceedings of the 2011 Conference on Empirical Methods in Natural Language Processing*, pages 529–539, Edinburgh, Scotland, UK. Association for Computational Linguistics.

Xi Victoria Lin, Richard Socher, and Caiming Xiong. 2018. Multi-hop knowledge graph reasoning with reward shaping. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 3243–3253, Brussels, Belgium. Association for Computational Linguistics.

Jiaxu Liu, Xinping Yi, and Xiaowei Huang. 2024. Deephgcn: Toward deeper hyperbolic graph convolutional networks. *IEEE Transactions on Artificial Intelligence*, 5(12):6172–6185.

Peter J Liu, Mohammad Saleh, Etienne Pot, Ben Goodrich, Ryan Sepassi, Lukasz Kaiser, and Noam Shazeer. 2018. Generating wikipedia by summarizing long sequences. In *International Conference on Learning Representations*.

Qi Liu, Maximilian Nickel, and Douwe Kiela. 2019. *Hyperbolic graph neural networks*. Curran Associates Inc., Red Hook, NY, USA.

Ang Ma, Yanhua Yu, Chuan Shi, Shuai Zhen, Liang Pang, and Tat-Seng Chua. 2024. Pmhr: Path-based multi-hop reasoning incorporating rule-enhanced reinforcement learning and kg embeddings. *Electronics*, 13(23):4847.

Kanishka Misra, Cicero Nogueira dos Santos, and Siamak Shakeri. 2023. Triggering multi-hop reasoning for question answering in language models using soft prompts and random walks. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 972–985, Toronto, Canada. Association for Computational Linguistics.

Sebastien Montella, Lina M Rojas Barahona, and Johannes Heinecke. 2021. Hyperbolic temporal knowledge graph embeddings with relational and time curvatures. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 3296–3308.

Chau Duc Minh Nguyen, Tim French, Wei Liu, and Michael Stewart. 2023. Cyle: Cylinder embeddings for multi-hop reasoning over knowledge graphs. In *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics*, pages 1736–1751.

Maximillian Nickel and Douwe Kiela. 2017. Poincaré embeddings for learning hierarchical representations. *Advances in neural information processing systems*, 30.

Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of machine learning research*, 21(140):1–67.

Shimizu Ryohei, Mukuta Yusuke, and Harada Tatsuya. 2021. Hyperbolic neural networks++. In *Proceedings of International Conference on Learning Representations (ICLR)*.

Ramit Sawhney, Shrey Pandit, Vishwa Shah, Megh Thakkar, and Shafiq Joty. 2024. AdaPT: A set of guidelines for hyperbolic multimodal multilingual NLP. In *Findings of the Association for Computational Linguistics: NAACL 2024*, pages 1757–1771, Mexico City, Mexico. Association for Computational Linguistics.

Michael Schlichtkrull, Thomas N Kipf, Peter Bloem, Rianne Van Den Berg, Ivan Titov, and Max Welling. 2018. Modeling relational data with graph convolutional networks. In *The semantic web: 15th international conference, ESWC 2018, Heraklion, Crete, Greece, June 3–7, 2018, proceedings 15*, pages 593–607. Springer.

Noam Shazeer and Mitchell Stern. 2018. Adafactor: Adaptive learning rates with sublinear memory cost. In *International Conference on Machine Learning*, pages 4596–4604. PMLR.

Yucheng Shi, Qiaoyu Tan, Xuansheng Wu, Shaochen Zhong, Kaixiong Zhou, and Ninghao Liu. 2024a. Retrieval-enhanced knowledge editing in language models for multi-hop question answering. In *Proceedings of the 33rd ACM International Conference on Information and Knowledge Management*, CIKM '24, page 2056–2066, New York, NY, USA. Association for Computing Machinery.

Zhengliang Shi, Shuo Zhang, Weiwei Sun, Shen Gao, Pengjie Ren, Zhumin Chen, and Zhaochun Ren. 2024b. Generate-then-ground in retrieval-augmented generation for multi-hop question answering. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 7339–7353, Bangkok, Thailand. Association for Computational Linguistics.

Zhiqing Sun, Zhi-Hong Deng, Jian-Yun Nie, and Jian Tang. Rotate: Knowledge graph embedding by relational rotation in complex space. In *International Conference on Learning Representations*.

Yiming Tan, Yongrui Chen, Guilin Qi, Weizhuo Li, and Meng Wang. 2023. Mlpq: A dataset for path question answering over multilingual knowledge graphs. *Big Data Research*, 32:100381.

Max van Spengler, Erwin Berkhout, and Pascal Mettes. 2023. Poincaré resnet. In *2023 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 5396–5405. IEEE Computer Society.

Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua Bengio. 2018. Graph attention networks. In *International Conference on Learning Representations*.

Denny Vrandečić and Markus Krötzsch. 2014. Wikidata: a free collaborative knowledgebase. *Communications of the ACM*, 57(10):78–85.

Guojia Wan, Shirui Pan, Chen Gong, Chuan Zhou, and Gholamreza Haffari. 2021. Reasoning like human: hierarchical reinforcement learning for knowledge graph reasoning. In *Proceedings of the Twenty-Ninth International Conference on International Joint Conferences on Artificial Intelligence*, pages 1926–1932.

Bin Wang, Fuyong Xu, Peiyu Liu, and Zhenfang Zhu. 2024. Hypermr: Hyperbolic hypergraph multi-hop reasoning for knowledge-based visual question answering. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 8505–8515.

Zhen Wang, Jianwen Zhang, Jianlin Feng, and Zheng Chen. 2014. Knowledge graph embedding by translating on hyperplanes. In *Proceedings of the AAAI conference on artificial intelligence*, volume 28.

Guanchen Xiao, Jinzhi Liao, Zhen Tan, Yiqi Yu, and Bin Ge. 2022. Hyperbolic directed hypergraph-based reasoning for multi-hop kbqa. *Mathematics*, 10(20):3905.

Weiwen Xu, Yang Deng, Huihui Zhang, Deng Cai, and Wai Lam. 2021. Exploiting reasoning chains for multi-hop science question answering. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 1143–1156.

Yi Xu, Dongsheng Wang, Bo Chen, Ruiying Lu, Zhibin Duan, Mingyuan Zhou, et al. 2022. Hyperminer: Topic taxonomy mining with hyperbolic embedding. *Advances in Neural Information Processing Systems*, 35:31557–31570.

Bishan Yang, Scott Wen-tau Yih, Xiaodong He, Jianfeng Gao, and Li Deng. 2015. Embedding entities and relations for learning and inference in knowledge bases. In *Proceedings of the International Conference on Learning Representations (ICLR) 2015*.

Jiahao Zhang, Haiyang Zhang, Dongmei Zhang, Liu Yong, and Shen Huang. 2024a. End-to-end beam retrieval for multi-hop question answering. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 1718–1731.

Xiaoming Zhang, Ming Wang, Xiaocui Yang, Daling Wang, Shi Feng, and Yifei Zhang. 2024b. Hierarchical retrieval-augmented generation model with rethink for multi-hop question answering. *arXiv e-prints*, pages arXiv–2408.

Mantong Zhou, Minlie Huang, and Xiaoyan Zhu. 2018a. An interpretable reasoning network for multi-relation question answering. In *Proceedings of the 27th International Conference on Computational Linguistics*, pages 2010–2022.

Mantong Zhou, Minlie Huang, and Xiaoyan Zhu. 2018b. An interpretable reasoning network for multi-relation question answering. In *Proceedings of the 27th International Conference on Computational Linguistics*, pages 2010–2022, Santa Fe, New Mexico, USA. Association for Computational Linguistics.

Xingchen Zhou, Peng Wang, Qiqing Luo, and Zhe Pan. 2021. Multi-hop knowledge graph reasoning based on hyperbolic knowledge graph embedding and reinforcement learning. In *Proceedings of the 10th International Joint Conference on Knowledge Graphs*, pages 1–9.

Anjie Zhu, Deqiang Ouyang, Shuang Liang, and Jie Shao. 2022. Step by step:: A hierarchical framework for multi-hop knowledge graph reasoning with reinforcement learning.

# A   Appendix

## A.1   Evidence Creation for MetaQA

The MetaQA dataset contains structured information about questions, including the source entity, the tail entity (answer), and the intermediate relations that connect them for each question. For instance, in the question:

*"What are the languages spoken in the films directed by [Joel Zwick]?"*

| Pair | Relation |
|------|----------|
| (movie, language) | in_language |
| (movie, year) | release_year |
| (movie, writer) | written_by |
| (movie, director) | directed_by |
| (movie, genre) | has_genre |
| (movie, actor) | starred_actors |
| (language, movie) | in_language_reversed |
| (year, movie) | release_year_reversed |
| (writer, movie) | written_by_reversed |
| (director, movie) | directed_by_reversed |
| (genre, movie) | has_genre_reversed |
| (actor, movie) | starred_actors_reversed |

Table 9: Pair to relation mapping of the MetaQA dataset

The source entity is *Joel Zwick*, and the answer would be *Greek*. Additionally, the dataset provides the intermediate relations forming the reasoning path from the source entity to the answer. For this example, the intermediate relations are represented as a path string:

*director_to_movie_to_language.*

To find the evidence, we first parse the path string into the pairs:

- (director, movie)

- (movie, language)

Each pair represents a segment of the reasoning path. These pairs are then mapped to their corresponding relations in the knowledge graph using the mapping defined in Table 9. For instance:

- (director, movie) → directed_by_reversed

- (movie, language) → in_language

Using the source entity, the intermediate relations, and the answer entity, we construct the complete entity-relation-entity-relation-entity chain for each question. This chain serves as the evidence for the reasoning process.

## A.2   Further Performance Comparison and Model Efficiency on the PQ Dataset

| | EM Score |
|------|----------|
| Ours with Euclidean | 89.01 |
| Ours with Hyperbolic | 94.24 |
| HyperMR (Wang et al., 2024) | **96.2** |

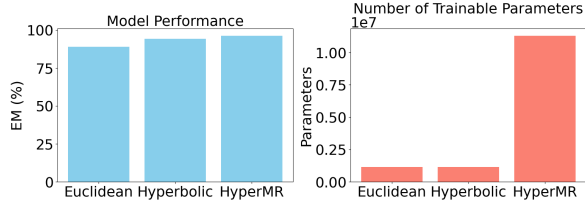Table 10: Comparison of our model with HyperMR on the PQ dataset.

Figure 7: On the left: EM score on PQ test set for our model with Euclidean/hyperbolic layer vs HyperMR Wang et al. (2024). On the right: Number of trainable parameters. Even though our approach has lower performance, our model only has approximately $\frac{1}{10}$ of the trainable parameters as their model.

In Table 10, we compare the test performance of our approach against the state-of-the-art hyperbolic model, HyperMR, on the PQ dataset. In this case, we incorporated reasoning paths into the training process to ensure a fair comparison. Our results indicate that the hyperbolic layer outperforms its Euclidean counterpart, improving accuracy from 89.01% to 94.24%. While our hyperbolic model performs slightly lower than HyperMR, it achieves this with only a fraction of the trainable parameters (one million compared to 11 million for HyperMR) as illustrated in Figure 7.