

# Exploring Alternative Approaches to Language Modeling for Learning from Data and Knowledge

Yuxin Zi, Kaushik Roy, Vignesh Narayanan, Amit Sheth

Artificial Intelligence Institute, University of South Carolina  
yzi@email.sc.edu, kaushikr@email.sc.edu, vignar@sc.edu, amit@sc.edu

## Abstract

Despite their extensive application in language understanding tasks, large language models (LLMs) still encounter challenges including hallucinations - occasional fabrication of information - and alignment issues - lack of associations with human-curated world models (e.g., intuitive physics or common-sense knowledge). Moreover, the black-box nature of LLMs presents significant obstacles in training them effectively to achieve desired behaviors. In particular, modifying the concept embedding spaces of LLMs can be highly intractable. This process involves analyzing the implicit impact of such adjustments on the myriad parameters within LLMs and the resulting inductive biases. We propose a novel architecture that wraps powerful function approximation architectures within an outer, interpretable read-out layer. This read-out layer can be scrutinized to explicitly observe the effects of concept modeling during the training of the LLM. Our method stands in contrast with gradient-based implicit mechanisms, which depend solely on adjustments to the LLM parameters and thus evade scrutiny. By conducting extensive experiments across both generative and discriminative language modeling tasks, we evaluate the capabilities of our proposed architecture relative to state-of-the-art LLMs of similar sizes. Additionally, we offer a qualitative examination of the interpretable read-out layer and visualize the concepts it captures. The results demonstrate the potential of our approach for effectively controlling LLM hallucinations and enhancing the alignment with human expectations.

## 1 Introduction

Language modeling involves extracting extensive patterns from vast amounts of data and representing these patterns in high-dimensional vector spaces (Zhou et al. 2023). These vector spaces enable us to assess the similarity or dissimilarity between different concepts by calculating the distances between their embedded vectors. For instance, the embeddings of various fruits such as apples, grapes, and watermelons will be located within a small distance (e.g.,  $\epsilon$ ) from each other. As a result, the common attributes (e.g., being fruits) or unique characteristics (e.g., nutritional value) of these concepts are implicitly encoded within the parameters of the language models that embed them. The implicit representation of concept features by language model parameters,

like fruits providing natural sugars, presents challenges in aligning the model’s parameter space to achieve desired outcomes (Huang et al. 2023).

Consider a scenario where we train an LLM to capture concepts of *apples*, *grapes*, and *watermelons* for a specific task, such as designing a diet plan focused on high levels of antioxidants. Given that *watermelons* do not possess significant levels of antioxidants, the LLM should distinguish the concept of *watermelon* from fruits that are high in antioxidants. A knowledge graph (KG) provides explicit information about the nutritional properties of fruits. Therefore, instead of relying on implicit representations by the LLM’s parameters, existing work enforces categorizations within the LLM’s concept embedding space by integrating vector representations of KG concepts during training (Sarzynska-Wawer et al. 2021). Nevertheless, these methods either incorporate KG concept embeddings as part of the input to the LLMs, or combine them into the parametric space of these black-box models. Consequently, they result in changes to the concept space of LLMs that are impossible to interpret.

To address this challenge, we introduce an outer interpretable read-out layer to enable the explicit observation of modifications to a language model’s concept embedding space. This layer is added as the last layer after the LLM’s parametric architectural layers. In the subsequent sections, we develop and formalize the proposed approach. We then evaluate our approach on benchmark language modeling and understanding tasks. The findings demonstrate that our method performs comparably with state-of-the-art LLMs of similar size while offering a cleaner and more interpretable way to understand the LLM’s concept representations.

## 2 The Interpretable Read-Out Layer

Let  $g_\theta(x) : x \rightarrow \mathbb{R}^d$  be a function parameterized by  $\theta$ , which embeds a concept  $x$  as a  $d$  dimensional embedding vector. Let  $f_\beta : g_\theta(x) \rightarrow y, y \in \mathbb{R}^C$  be a linear function  $\beta^T g_\theta(x)$  which maps the embedding for  $x$  to a set of  $C$  target outcomes, for example,  $C$  target classes for sentiment classification or  $C$  possible next concepts (words or tokens) for language generation. The similarity of concepts in an embedding space depends on the target task. That is, after  $f_\beta$  is trained for a target task, we define two concepts  $x_i$  and  $x_j$  are similar if  $\|g_\theta(x_i) - g_\theta(x_j)\| \leq \epsilon$ , where  $\epsilon$  is some small number, and  $\|\cdot\|$  is an appropriate distance metric.

We introduce a non-parametric *interpretable read-out layer* over  $g_\theta(x)$ , denoted by  $\Phi(x)$ , which characterizes the concept  $x$  using  $g_\theta(x)$  and other “support” concepts before computing  $f_\beta$ . In the next section, we define such “support” concepts in the interpretable read-out layer, and explain their roles for explicit concept modeling. Specifically, we delve into the core underpinnings of language modeling and subsequently examine how fundamental operations of state-of-the-art language models, such as transformers, formalize these underpinnings (Vaswani et al. 2017). We then generalize these operations to demonstrate that adding the interpretable read-out layer is the logical next step for endowing the model with interpretability.

### A Closer Look at Language Modeling, the Transformer Architecture and Defining the Interpretable Read-Out Layer $\Phi(x)$

While our discussion has focused on concept embeddings of single words, such as *apple*, it is essential to note that a LLM also embeds *concept-phrases* of arbitrary length, like *a red apple with vitamin C*. The key idea is that when computing for the LLM’s concept embedding space, a set of co-occurring concepts that constitute a concept phrase serve as each other’s “support”. Intuitively, the co-occurring concepts *apple* and *vitamin C*, serve to support the interpretation of the individual tokens (*apple*, *vitamin* and *C*) as belonging together, indicating their presumed proximity in the LLM’s embedding space. If the concept *apple* were co-occurring with *Steve Jobs* in the concept phrase, the concept *Steve Jobs* would support a different interpretation of *apple* (i.e., embeddings for *apple*, *vitamin*, and *C* would not be close in this case). The core operation in a transformer, the self-attention operator, leverages this idea of “support” during computing for concept and concept phrase embeddings.

Let  $T$  be the set of all possible concepts in a language. Let  $X = [x_1, x_2, \dots, x_N]$  denote a concept phrase, which is an ordered list of concepts such that each  $x_n \in X$  is an element in the set  $T$ . Thus the self-attention for the concept  $x_i$ , denoted by  $\mathbf{SA}_{x_i}$  is computed as  $\sum_j \sigma(\frac{q_i^T k_j}{\sqrt{d}}) v_j = \sum_j \left( \frac{\exp(\frac{q_i^T k_j}{\sqrt{d}})}{\sum_j \exp(\frac{q_i^T k_j}{\sqrt{d}})} \right) v_j$ . The  $q_i$ ,  $k_j$ , and  $v_j$  are vectors computed using  $\mathbf{W}_{q\theta}(x_i)$ ,  $\mathbf{W}_{k\theta}(x_j)$ , and  $\mathbf{W}_{v\theta}(x_j)$ , respectively ( $g_\theta$  is as defined in Section 2). The matrices  $\mathbf{W}_q$ ,  $\mathbf{W}_k$ , and  $\mathbf{W}_v$  are square matrices of dimension  $d \times d$ . Notice that the self-attention computation for  $x_i$  depends on  $j$  other co-occurring “support” concepts. We now generalize this idea to obtain the interpretable read-out layer.

$$\begin{aligned} \mathbf{SA}_{x_i} &= \sum_j \sigma(\frac{q_i^T k_j}{\sqrt{d}}) v_j = \sum_j \left( \frac{\exp(\frac{q_i^T k_j}{\sqrt{d}})}{\sum_j \exp(\frac{q_i^T k_j}{\sqrt{d}})} \right) v_j \\ &= \sum_j \left( \frac{\exp(\frac{\|q_i\|^2 - \|q_i - k_j\|^2 + \|k_j\|^2}{2\sqrt{d}})}{\sum_j \exp(\frac{\|q_i\|^2 - \|q_i - k_j\|^2 + \|k_j\|^2}{2\sqrt{d}})} \right) v_j \end{aligned} \quad (1)$$

Then, we substitute  $q_i = -q'_i = -\mathbf{W}_{q\theta}(x_i)$ ,  $k_j = -k'_j = -\mathbf{W}_{k\theta}(x_j)$ , and  $v_j = \mathbf{W}_{v\theta}(x_j)$ . We also observe that

$\exp(\frac{-\|l-m\|^2}{\sqrt{d}})$  is the Gaussian kernel with bandwidth  $\sqrt{d}$ , denoted by the inner product  $\phi(l)^T \phi(m)$ , where  $\phi(\cdot)$  is the infinite-dimensional map (Gaussian kernel is the infinite-dimensional inner product). Thus, we rewrite Eq (1) as:

$$\begin{aligned} \mathbf{SA}_{x_i} &= \sum_j \left( \frac{C \|\phi(\mathbf{W}_{q\theta}(x_i))^T \phi(\mathbf{W}_{k\theta}(x_j))\|^2}{C \sum_j \|\phi(\mathbf{W}_{q\theta}(x_i))^T \phi(\mathbf{W}_{k\theta}(x_j))\|^2} \right) v_j, \\ v_j &= \mathbf{W}_{v\theta}(x_j), \quad C = \phi(0)^T \phi(0) \end{aligned} \quad (2)$$

$$\begin{aligned} &= \sum_j \left( \frac{\Phi(\mathbf{W}_{q\theta}(x_i))^T \Phi(\mathbf{W}_{k\theta}(x_j))}{\sum_j \Phi(\mathbf{W}_{q\theta}(x_i))^T \Phi(\mathbf{W}_{k\theta}(x_j))} \right) v_j, \\ \Phi(\cdot)^T \Phi(\cdot) &= \|\phi(\cdot)^T \phi(\cdot)\|^2 \end{aligned} \quad (3)$$

We replace  $\|\phi(\cdot)^T \phi(\cdot)\|^2$  with a new inner product (kernel)  $\Phi(\cdot)^T \Phi(\cdot)$  as the product of two kernels is still a kernel. In this way, we have introduced a read-out layer over  $g_\theta$  and the notion of “support” concepts  $x_j$  for characterizing the concept  $x_i$ . In the next section, we will explain how we enable the explicit observation of language modeling outcomes by incorporating the interpretable read-out layer  $\Phi(\cdot)$  and the “support” concepts. It is important to note that although earlier methods have demonstrated the feasibility of constructing a kernel to instantiate the self-attention operation, they have yet to explicitly derive the specific form containing the Gaussian kernel as in this work (Tsai et al. 2019; Chowdhury et al. 2021). Furthermore, our objective is not to offer an alternative kernel re-formulation, but to achieve a formulation through which we can explicitly interpret the LLM’s concept embeddings using “support” concepts, as we will detail in the following sections.

### Language Modeling by Leveraging the Interpretable Read-Out Layer

**Support Concepts** Recall the target task of designing a diet plan rich in antioxidants from Section 1, using just the fruits - *apples*, *grapes*, and *watermelons*. We can obtain “support” concepts as paths from external KGs. For example, the KGs can have the triples *Apple is\_a Fruit*, and *Apple has Antioxidants*, which can be reformulated as the path *Antioxidants has<sup>-1</sup> Apple is\_a Fruit* using inverse relationships. Similarly, for *grape*, we have the path *Antioxidants has<sup>-1</sup> Grape is\_a Fruit*. Thus, we can characterize the concepts *Apple* and *Grape* as being *fruits rich in antioxidants*, an appropriate description of a fruit category for the target task. Thus in this example,  $x_i$  in Eq (3) is the concept *apple* and the  $x_j$  are the KG paths, *Antioxidants has<sup>-1</sup> Apple is\_a Fruit*, and *Antioxidants has<sup>-1</sup> Grape is\_a Fruit*. Figure 1 depicts this process.

**Layering the Interpretable Read-Out Layer over Different Parametric Function Approximation Architectures for Language Modeling** In traditional parametric language modeling, the self-attention operations such as described in Eq (1) are layered on top of one another (e.g., 12 layers in BERT (Devlin et al. 2018)). In our approach, we

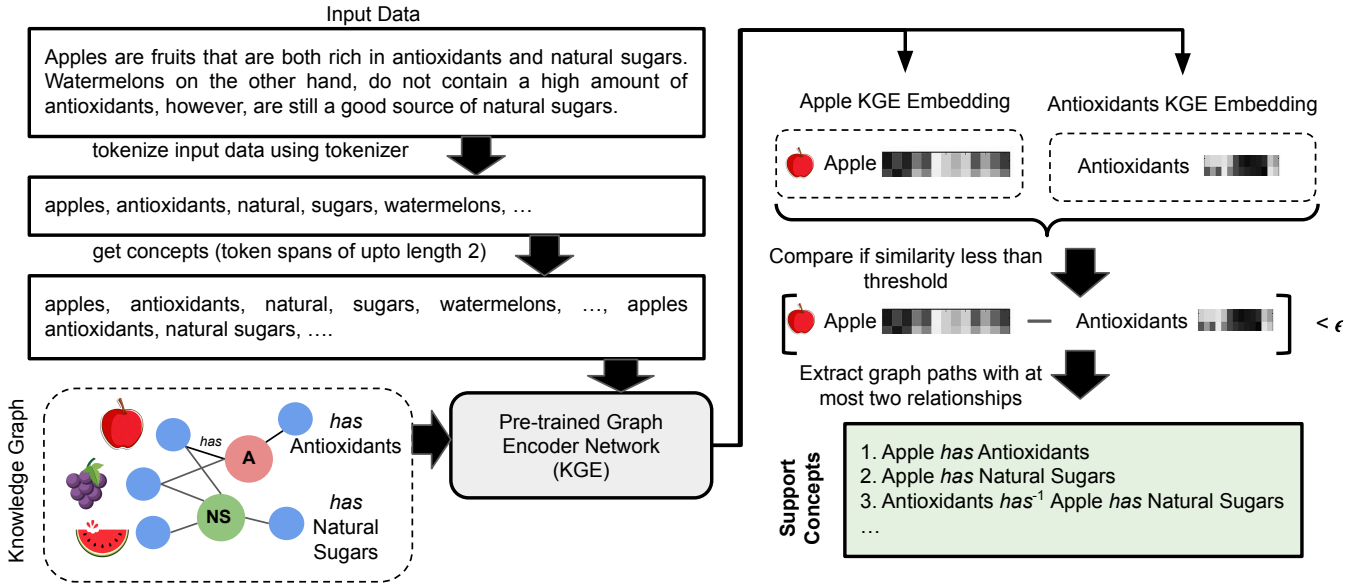


Figure 1: Process for deriving supportive concepts from data and an external KG. Initially, the data undergoes tokenization. Subsequently, token spans of length  $K$  are identified as potential concept candidates for retrieving graph paths from the KG (where  $K$  is set to 2 in the illustration). Graph paths are extracted by evaluating vector distances between concepts and graph nodes. These graph paths then serve as the resulting supportive concepts.

instead lay the read-out layer  $\Phi(\cdot)$  over the parametric architecture  $g_\theta$  (e.g., could be a 12-layer transformer architecture). This allows us to leverage the high-capacity function approximation capabilities of complex parametric architectures in  $g_\theta$ , while still retaining the explicit interpretation of concept descriptions using “support” concepts, as explained in the previous section (Section 2). Furthermore, we can now experiment with different parametric architectures for  $g_\theta$  to achieve optimal language modeling performance. Figure 2 illustrates the idea of layering the interpretable read-out layer on top of parametric architectures and how the layer’s output can be interpreted.

### 3 Definitions

In following sections, we introduce definitions required for explaining our method and experiments, then elaborate on the choices for the parametric function approximation  $g_\theta$  and the interpretable read-out layer  $\Phi(\cdot)$ .

#### Data Structure Definitions

**Knowledge Graphs** We formally define a KG and its paths (paths in the next section), as we use KG paths as “support” concepts in our experiments. A KG is denoted by  $KG(\mathbf{V}, \mathbf{E}, \mathbf{L})$ , where the sets  $\mathbf{V}$  and  $\mathbf{E}$  are the vertices and edges of the graph, respectively. The set  $\mathbf{L}$  consists of relationships represented by the edges  $e(v_1, v_2) \in \mathbf{E}$ ,  $v_1, v_2 \in \mathbf{V}$ . The relationships represented by the edges  $e(v_1, v_2) \in \mathbf{E}$  are given by a set of  $N$  boolean-valued functions:

$$\mathbf{L} = \{r_{e(v_1, v_2)} : y_e \in \{0, 1\}^N \mid e \in \mathbf{E}\} \quad (4)$$

Here 1 and 0 denotes whether the relationship  $r_{e(v_1, v_2)}$  between vertices  $v_1, v_2 \in V$  holds or not. We use  $L(e)$

to denote the value  $y_e$  associated with the relationship  $r_{e(v_1, v_2)}$ . We use  $subject(e)$  and  $object(e)$  to denote the incident vertices  $v_1, v_2$  for the edge  $e(v_1, v_2)$ . Figure 1 (bottom-left) shows an example of a KG where the vertices  $\mathbf{V}$  is the set  $\{apple, watermelon, grape, antioxidants, natural\ sugar\}$ , the edges  $\mathbf{E}$  is the set  $\{has(apple, antioxidants), has(apple, natural\ sugar), \dots\}$ , and  $\mathbf{L}$  is the set  $\{apple\ has\ antioxidants : 1, apple\ has\ natural\ sugar : 1, \dots\}$ . Note that the set  $\mathbf{L}$  is not completely specified, i.e., only the relationships that hold ( $L(e) = 1$ ) are stored, and the ones that are not in  $\mathbf{L}$  are assumed not to hold ( $L(e) = 0$ ).

**Knowledge Graph Paths** Given a KG denoted by  $KG(\mathbf{V}, \mathbf{E}, \mathbf{L})$ , a  $K$  length path  $p_K(v_l, v_m)$  between two vertices  $v_l, v_m \in \mathbf{V}$ , is a sequence of edges  $e_1, e_2, \dots, e_K \in \mathbf{E}$ , such that  $L(e_k) = 1 \wedge subject(e_1) = v_l \wedge object(e_K) = v_m \wedge object(e_k) = subject(e_{k+1})$ ,  $k \in \{1, 2, \dots, K\}$ . We have shown examples of KG paths in Section 2.

#### Task Definitions

We experiment with generative (text generation) and discriminative (classification) modeling tasks. We describe these tasks and how a KG is used to solve the tasks formally in the following subsections.

**Generative Modeling** In generative modeling, given a context sequence of  $H$  previous concepts  $[x_1, \dots, x_H]$ , the task is next-concept prediction, i.e., predicting the concept  $x_{H+1}$  from among a predefined vocabulary of  $C$  concepts. The forward pass for this computation is written using a

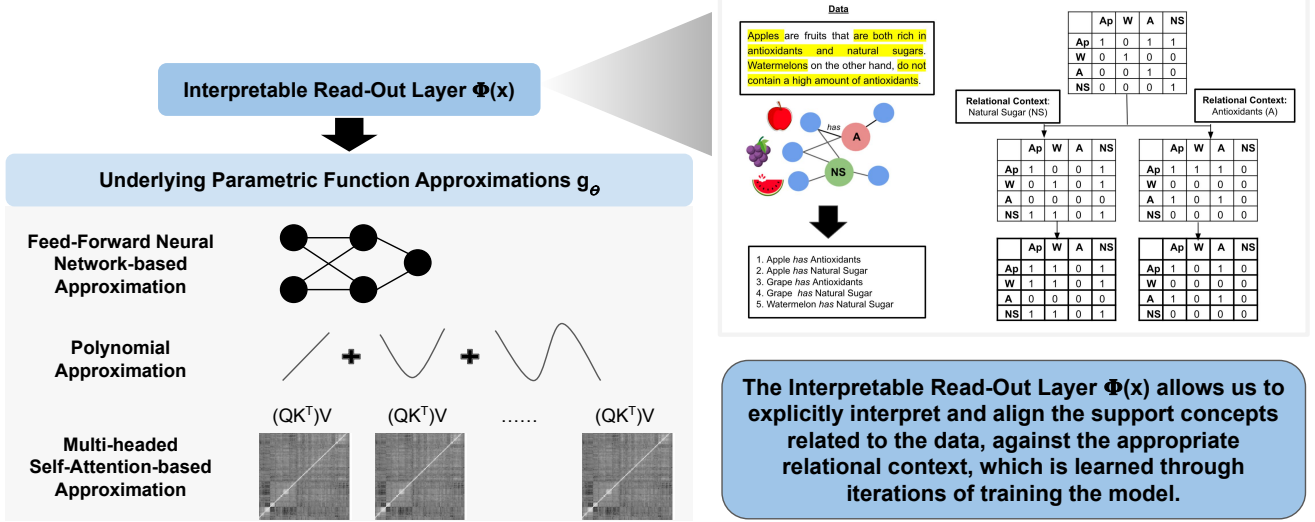


Figure 2: Shows different possible parametric function choices for  $g_\theta$  under the interpretable read-out layer  $\Phi(x)$ , all of which have universal approximation properties. The figure also illustrates how the read-out layer can be interpreted as defining the appropriate relational context for the target task of designing a diet plan rich in antioxidants (Section 1), using *apples*, *grapes*, and *watermelons*. The relevant data-specific context is captured in  $\Phi(\cdot)$  throughout the training iterations of the model. This is depicted as an adjacency matrix (formed by utilizing thresholded dot product values calculated by Equation (2)).

modification to  $f_\beta$  (introduced in Section 2) and Eq 3 as:

$$\begin{aligned}
 & f_\beta(x_{H+1} \mid [x_1, \dots, x_H]) \\
 &= \beta \sum_j \left( \frac{(\mathbf{W}_q g_\theta(x_i))^T (\mathbf{W}_k g_\theta(x_j))}{\sum_j (\mathbf{W}_q g_\theta(x_i))^T (\mathbf{W}_k g_\theta(x_j))} \right) \mathbf{W}_v g_\theta(x_j), \\
 & x_j \in \{x_1, \dots, x_H\}
 \end{aligned} \tag{5}$$

Here the “support” concepts are the set of concepts in context sequence  $[x_1, \dots, x_H]$ . Note here that we set  $\Phi(\cdot)$  to be the identity function, thus reducing the inner product in Eq (3) to the standard scalar product in Eq (5), as this is what we experiment with. Crucially, this change still retains the explicit interpretation described in Section 2 as the inner-product is still the outermost layer over  $g_\theta$ . We leave the exploration of different  $\Phi(\cdot)$  for future work.

**Generative Modeling using Knowledge Graph paths as “support” concepts  $x_j$  in Eq (5)** For a concept  $x$  and KG denoted by  $KG(\mathbf{V}, \mathbf{E}, \mathbf{L})$ , we define all KG paths of up to length  $K$  corresponding to concept  $x$  as

$$\begin{aligned}
 & P_K(x) \\
 &= \{p_K(v_l, v_m) \mid k \leq K \ v_l, v_m \in \mathbf{V}, \\
 & \|g^{KGE}(x) - g^{KGE}(v_m)\| \leq \varepsilon\}
 \end{aligned} \tag{6}$$

where  $\|g^{KGE}(x) - g^{KGE}(v_m)\| \leq \varepsilon$  denotes if embeddings for  $v_m$  and concept  $x$  are “close enough” using a KG embedding (KGE) model  $g^{KGE}$ . Let  $\mathbb{I}_{KG}(x, x_j)$  be an indicator function that returns 1 or 0 indicating whether the “support” concept  $x_j \in \Phi(x)$  or not, i.e., if the “support” concept is a path in the KG. Thus, to incorporate KG paths,

we make a slight modification to Eq (5) as follows:

$$\begin{aligned}
 & f_\beta(x_{H+1} \mid [x_1, \dots, x_H]) \\
 &= \beta \sum_j \left( \frac{(\mathbf{W}_q g'_\theta(x_i))^T (\mathbf{W}_k g'_\theta(x_j))}{\sum_j (\mathbf{W}_q g'_\theta(x_i))^T (\mathbf{W}_k g'_\theta(x_j))} \right) \mathbf{W}_v g'_\theta(x_j) \\
 & g'_\theta = \begin{cases} g_\theta & \text{if } \mathbb{I}_{KG}(x_{H+1}, x_j) = 0 \\ g^{KGE} & \text{if } \mathbb{I}_{KG}(x_{H+1}, x_j) = 1, \end{cases} \\
 & x_j \in \{x_1, \dots, x_H\} \cup P_K(x_{H+1})
 \end{aligned} \tag{7}$$

Here the “support” concepts include both the set of concepts  $[x_1, \dots, x_H]$  and the KG paths in  $P_K(x_i)$ .

**Discriminative Modeling** In discriminative modeling, given a dataset of  $D$  data points, i.e., concept set and label pairs  $(l_d = \{x_1, x_2, \dots, x_H\}, m_d = c)$ , where the labels  $c$  are from among a predefined set of labels  $\mathcal{C}$ , the task is to predict the label  $c$  using  $f_\beta$ . Let  $y$  denote the predicted variable for the class label  $c$ . The forward pass computation for the logit corresponding to label  $c \in \mathcal{C}$  is as follows:

$$\begin{aligned}
 & f_\beta(y = c \mid \{x_1, \dots, x_H\}) \\
 &= \beta \sum_j \left( \frac{(\mathbf{W}_q g_\theta(x_i))^T (\mathbf{W}_k g_\theta(x_j))}{\sum_j (\mathbf{W}_q g_\theta(x_i))^T (\mathbf{W}_k g_\theta(x_j))} \right) \mathbf{W}_v g_\theta(x_j), \\
 & x_j \in \{l_1 \cup l_2 \cup \dots \cup l_D\} \setminus \{x_1, \dots, x_H\}
 \end{aligned} \tag{8}$$

Thus the “support” concepts for a datapoint during discriminative learning are concepts from all the other  $D - 1$  datapoints in the dataset. Extending the “support” concepts to the discriminative modeling case involves including the union of

all the KG paths  $P_K(x_j)$  corresponding to the support  $x_j$ . To incorporate KG paths, we modify  $g_\theta$  as in Eq 7.

### Defining Approximation Architectures for $g_\theta$

As illustrated in Figure 2, we experiment with different approximation architectures, all of which have universal approximation properties. They are feed-forward neural networks, polynomial approximations, and the multi-headed self-attention architecture.

**Feed-Forward Neural Network** First, we experiment with a  $Z$  layer feed-forward neural network described by:

$$\begin{aligned} e_x &= \mathbf{E}x, \mathbf{E} \in \mathbb{R}^d \\ p_x &= e_x + \mathbf{P}_e x, \mathbf{P}_e \in \mathbb{R}^d \\ z_x &= \max(\mathbf{W}_z^T p_x, 0), \mathbf{W}_1 \in \mathbb{R}^{d \times d_1}, \mathbf{W}_{z \setminus \{1, \mathbf{Z}\}} \in \mathbb{R}^{d_z \times d_{z+1}} \\ \mathbf{W}_z &\in \mathbb{R}^{d_{z-1} \times d}, z \in \{1, 2, \dots, Z\} \end{aligned} \quad (9)$$

The matrices  $\mathbf{E}$ ,  $\mathbf{P}_e$ , and the  $\mathbf{W}_z$  are the trainable weights in the network (the embedding matrix, the position encoding matrix, and the network weights and biases). We layer multiple feed-forward structures as described in Eq (9) (12 layers in our experiments, each layer with its own trainable weights) to obtain deeper approximation architectures. From each lower layer to the upper layer, we extract the last  $d$  dimensional column of  $z_x$  output at that lower layer. Finally, to obtain the  $d$  dimensional vector corresponding to  $g_\theta(x)$ , we extract the last column of  $z_x$  from the final layer.

**Polynomial Approximation** Next, we experiment with a polynomial approximation where we compute powers of  $x$  up to order  $J$ . This is described by the equations

$$\begin{aligned} e_x &= \mathbf{E}x, \mathbf{E} \in \mathbb{R}^d \\ p_x &= e_x + \mathbf{P}_e x, \mathbf{P}_e \in \mathbb{R}^d \\ z_x &= \{p_x^1, p_x^2, \dots, p_x^J\} \end{aligned} \quad (10)$$

Each of the  $\{p_x^1, p_x^2, \dots, p_x^J\}$  is computed dimension-wise (i.e., each of the  $d$  dimensions of  $x$  is raised to the power  $j$ ). The matrices  $\mathbf{E}$  and  $\mathbf{P}_e$  are the trainable weights in the network (the embedding and position encoding matrices). Once again, we layer multiple such structures to obtain deeper approximation architectures. From each lower layer, we take the average of the polynomial powers at that layer to obtain a  $d$  dimensional vector to pass to the layer above it. Finally, to obtain the  $d$  dimensional vector corresponding to  $g_\theta(x)$ , we take the average of the polynomial powers at the last layer.

### Multiheaded Self-Attention

$$\begin{aligned} e_x &= \mathbf{E}x, \mathbf{E} \in \mathbb{R}^d \\ p_x &= e_x + \mathbf{P}_e x, \mathbf{P}_e \in \mathbb{R}^d \\ q_x^a, k_x^a, v_x^a &= \{\mathbf{W}_q^a x, \mathbf{W}_k^a x, \mathbf{W}_v^a x \mid a \in \{1, 2, \dots, A\}\} \\ z_x &= \left\{ \sigma \left( \frac{(q_x^a)^T k_x^a}{\sqrt{d}} \right) v_x^a \mid a \in \{1, 2, \dots, A\} \right\} \end{aligned} \quad (11)$$

The matrices  $\mathbf{E}$ ,  $\mathbf{P}_e$ , and the  $\mathbf{W}_q^a$ ,  $\mathbf{W}_k^a$ , and  $\mathbf{W}_v^a$  are the trainable weights of the network.  $A$  denotes the number of attention heads. Again, multiple self-attention blocks are layered, as in the other two cases. From each lower layer, we take the average of the elements of  $z_x$  at that layer to obtain a  $d$  dimensional vector to pass to the layer above it. The final  $d$  dimensional vector corresponding to  $g_\theta(x)$  is obtained from the average of the elements in  $z_x$  from the last layer.

## 4 Experiments

In this section, we describe our hyperparameter configurations and experiments for generative and discriminative modeling in Sections 4, 4 and 4, respectively. Due to space concerns, the table and figure captions contain the discussion about all the experiments. *We provide the GLUE leaderboard result for context for the numbers in the results tables. However, it should be noted that the leaderboard models are up to 10 times larger than the models implemented in this paper. We will compare larger models in future work when the models have finished training.*

### Hyperparameter Configurations

For all our experiments, we use a single A100 GPU. For the feed-forward neural network described in Section 3, we set  $d = 384$  (chosen by tuning from the set  $\{200, 384, 768\}$ ), and  $d_z = 4000$  (chosen by tuning from the set  $\{500, 1000, 2000, 4000\}$ ). For the polynomial approximation described in Section 3, we set  $d = 384$  (chosen by tuning from the set  $\{200, 384, 768\}$ ), and  $J = 5$  (chosen by tuning from the set  $\{2, 3, 4, 5\}$ ). Finally, for the multi-headed-self-attention-based network described in Section 3, we set  $d = 384$  (chosen by tuning from the set  $\{200, 384, 768\}$ ), and  $A = 12$  (chosen by tuning from the set  $\{4, 8, 12\}$ ). These form the basic units, and to get deep architectures, we layer them 4, 4, and 6 times for the feed-forward approximation, polynomial approximation, and self-attention-based approximation, respectively (chosen by tuning from the set  $\{4, 6, 12\}$ ). We also include layer normalization between the layers. We use a train-validation split of 80-20 for all our experiments, and all reported results are evaluation loss scores.

### Generative Modeling

**Context Size, Tokenization, and Batch Size:** For all models, the context size is set to 1024 (chosen by tuning from the set  $\{256, 512, 1024\}$ ). Batch size is set to 32 (chosen by tuning from the set  $\{8, 16, 32\}$ ). For tokenization, we use the GPT-2 tokenizer<sup>1</sup>, which consists of 50,257 tokens ( $C$  in Section 3). **Embeddings for the  $g^{KGE}$  from Eq (7):** We use ConceptNet Numberbatch Embeddings<sup>2</sup> and EWISE WordNet Embeddings<sup>3</sup> for ConceptNet and WordNet, respectively.

**Parameter Initializations:** All the parameter matrices for all three methods are randomly initialized. Tables 1 and 2 show the results.

<sup>1</sup>[https://huggingface.co/docs/transformers/model\\_doc/gpt2](https://huggingface.co/docs/transformers/model_doc/gpt2)

<sup>2</sup><https://github.com/commonsense/conceptnet-numberbatch>

<sup>3</sup><https://github.com/malllabiisc/EWISE>

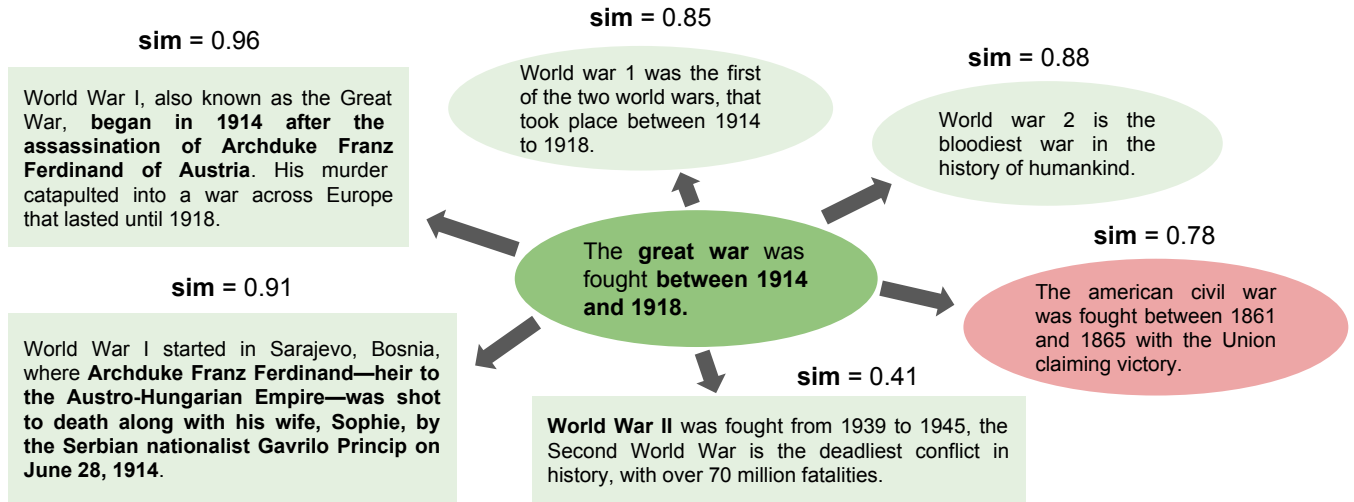


Figure 3: Interpretability results on a sentence similarity example. Since the read-out layer includes explicit inner products, we visualize and highlight in bold, the inner products from different models for a test example comparing sentence similarities against an anchor text (center). The squares green boxes represent the  $OLM_{pf}$  model, the green oval boxes represent the  $OLM_{mha}$  models, and the pink oval boxes represent the  $OLM_{nn}$  model. A significant advantage to using language modeling using our method is that we can directly use it to visualize the quality of the outputs. We find here also that a lot of the highlights correspond to relationships from ConceptNet (e.g., austria, 1914, austro-hungarian, great war, etc.), showing explicitly that KGs benefit performance in the discriminative modeling case.

**Datasets and Knowledge Graphs** We experiment with two text generation tasks. One, we train models for text generation in the style of Shakespeare by using the tiny-Shakespeare dataset<sup>4</sup>, which consists of 338, 025 tokens. We use a train-evaluation split of 80% and 20%, respectively. Second, we train models for autocomplete (next-word prediction) using the OpenWebText dataset<sup>5</sup>, which consists of  $\sim 9$  Billion (9, 040, 017, 095) tokens. For the first task, we used the KGs WordNet<sup>6</sup> and ConceptNet<sup>7</sup>, respectively. The relationships across both include: *Antonym, Distinct-From, EtymologicallyRelatedTo, LocatedNear, RelatedTo, SimilarTo, Synonym, AtLocation, CapableOf, Causes, CausesDesire, CreatedBy, DefinedAs, DerivedFrom, Desires, Entails, ExternalURL, FormOf, HasA, HasContext, HasFirst-Subevent, HasLastSubevent, HasPrerequisite, HasProperty, InstanceOf, IsA, MadeOf, MannerOf, MotivatedByGoal, ObstructedBy, PartOf, ReceivesAction, SenseOf, SymbolOf, and UsedFor* (Speer, Chin, and Havasi 2017).

**Results** Table 1 shows the results on the tiny-Shakespeare dataset. We see that multi-headed-self-attention architecture takes only 5 minutes and converges to the least evaluation loss. Across the board, including KGs results in worse performance. The polynomial fit converges the fastest in terms of the number of epochs but the slowest in terms of the number of minutes. Finally, we see that the difference in the evaluation losses with and without KG is significantly smaller

using the polynomial approximation method.

Model	w/o KG	/w KG	#Ep	Mins
Feed-forward Network	3.2	7.1	300	10
Polynomial Approximation	2.5	<b>4.2</b>	<b>100</b>	22
Multi-head Attention	<b>1.55</b>	6.8	200	5
GPT-2-Large-Fine-Tuned	2.2	6.7	250	<b>4</b>

Table 1: Results on the tiny-Shakespeare dataset

Table 2 shows the results on the OpenWebText dataset. We see similar trends as in the tiny-Shakespeare dataset. The multi-headed-self-attention architecture takes the least amount of days and achieves the least evaluation loss. Once again, including KGs consistently results in worse performance. The polynomial fit converges the fastest in terms of the number of epochs but the slowest in terms of the number of days. The polynomial approximation also again shows the least difference in the evaluation losses with and without KGs compared to the other two methods. Note that GPT-2 is fine-tuned, and our models are trained from scratch.

### Discriminative Modeling

We experiment with General Language Understanding Evaluation (GLUE) benchmark Tasks - STS (Semantic Textual Similarity Benchmark), MNLI (Multi-genre Natural Language Inference), QNLI (Question Answering Natural Language Inference), WNLI (Winograd Natural Language Inference), RTE (Recognizing Textual Entailment), and QQP (Quora Question Pairs) (Wang et al. 2018).

<sup>4</sup><https://raw.githubusercontent.com/karpathy/char-rnn/master/data/tinyshakespeare/input.txt>

<sup>5</sup><https://skylion007.github.io/OpenWebTextCorpus/>

<sup>6</sup><https://wordnet.princeton.edu/>

<sup>7</sup><https://conceptnet.io/>



Model	w/o KG	/w KG	#Ep	Days
Feed-forward Network	5.2	10.3	300	45
Polynomial Approximation	3.5	<b>6.2</b>	<b>200</b>	60
Multi-head Attention	<b>2.8</b>	9.5	300	32
GPT-2-Large-Fine-Tuned	3.2	8.7	250	<b>2</b>

Table 2: Results on the OpenWebText dataset

Table 3 shows the results. Here, OLM denotes our language model, with the neural network (nn), polynomial-fit (pf), and multiheaded (mha) self-attention architecture. We find that across the board, adding “support” concepts from the KGs improves scores significantly. Interestingly, the polynomial fit exhibits the best performance among the options for  $g_\theta$ .

Model	STS	QQP	QNLI	WNLI	MNLI	RTE
(GLUE <sub>LEADER</sub> )	93.5	<b>90.9</b>	<b>96.7</b>	97.9	<b>92.5</b>	<b>93.6</b>
OLM <sub>nn</sub> w/o KG	85.71	86.11	89.9	89.5	75.3	85.1
OLM <sub>nn</sub> /w KG	89.2	90.2	90.5	90.2	82.1	90.3
OLM <sub>pf</sub> w/o KG	90.89	86.41	92.3	90.11	88.53	90.4
OLM <sub>pf</sub> /w KG	<b>93.55</b>	90.51	95.56	<b>98.7</b>	92.08	92.3
OLM <sub>mha</sub> w/o KG	88.7	86.2	90.3	91.2	86.3	87.3
OLM <sub>mha</sub> /w KG	90.5	88.8	93.6	97.9	90.8	90.56

Table 3: Results on the GLUE Benchmark tasks.

## Interpretability

We perform interpretability analysis for a few examples from the test set, specifically for the task of sentence similarity. Figure 3 shows an example comparing different sentences talking about World War 1, World War 2, and the American Civil War.

## 5 Future Work

### Additional Forms of Support Concepts

In Section 2, we explored the application of KGs in relation to support concepts. Moving forward, we will incorporate support concepts from Instructing Tuning datasets as described in the subsequent two paragraphs.

**Support Concepts from Instruction Tuning Datasets**  
Consider the (Prompt, Instruction) pair, (Prompt: *Given information <data>, give me a food item rich in antioxidants containing apples*, Instruction: *Here is a food item rich in antioxidants containing apples - A fruit salad with the fruits apples and grapes*), where an example of <data> is shown in Figure 2.  $x_i$  is the concept *apple*,  $x_j$  can be the “support” concepts in Instruction that affects the distribution of LLM’s generated output given  $x_i$ .

**A Note on Instruction-based Fine-Tuning of LLMs**  
Instruction-based fine-tuning of LLMs using Reinforcement Learning with Human Feedback (RLHF) updates a policy function that uses proximal policy gradient-based methods to modify the LLM’s output distribution (Ouyang et al. 2022). In the case of RLHF, the policy function for modifying the LLM can be seen as an interpretable read-out layer over the LLM. More formally, let  $\pi : (x_i, x_j) \rightarrow [0, 1]$  denote a distribution over different “support” concepts  $x_j$ ’s in the Instruction given concept  $x_i$  in the Prompt.

The quantity  $\left( \frac{\Phi(\mathbf{W}_{\mathbf{q}g_\theta}(x_i))^T \Phi(\mathbf{W}_{\mathbf{k}g_\theta}(x_j))}{\sum_j \Phi(\mathbf{W}_{\mathbf{q}g_\theta}(x_i))^T \Phi(\mathbf{W}_{\mathbf{k}g_\theta}(x_j))} \right)$  in Eq (2) also defines such a distribution. Thus RLHF achieves the objective of reducing the divergence between distributions  $\pi$  and  $\left( \frac{\Phi(\mathbf{W}_{\mathbf{q}g_\theta}(x_i))^T \Phi(\mathbf{W}_{\mathbf{k}g_\theta}(x_j))}{\sum_j \Phi(\mathbf{W}_{\mathbf{q}g_\theta}(x_i))^T \Phi(\mathbf{W}_{\mathbf{k}g_\theta}(x_j))} \right)$ , for example, reducing the KL-divergence in RLHF using proximal policy gradient-based methods.

## Efficient Inner Product Implementations

Since we derive an explicit inner product or kernel formulation in the interpretable read-out layer, we can exploit high-quality sub-quadratic approximations to kernels to significantly speed up learning and inference (Poli et al. 2023; Hallgren 2021).

## 6 Conclusion

In this paper, we explore alternative perspectives on language modeling and demonstrate its comparable effectiveness to parametric approaches. Additionally, we illustrate its added advantage of facilitating easy visualization and interpretation. Consequently, the methodologies outlined in this paper hold promise for implementing mitigation strategies concerning observed adverse effects in language models, such as hallucinations and alignment issues.

## Acknowledgements

This research is partially supported by NSF Award 2335967 “EAGER: Knowledge-guided neurosymbolic AI with guardrails for safe virtual health assistants” (Sheth et al. 2021, 2022; Sheth, Roy, and Gaur 2023).

## References

- Chowdhury, S. P.; Solomou, A.; Dubey, A.; and Sachan, M. 2021. On learning the transformer kernel. *arXiv preprint arXiv:2110.08323*.
- Devlin, J.; Chang, M.-W.; Lee, K.; and Toutanova, K. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.
- Hallgren, F. 2021. Kernel PCA with the Nyström method. *arXiv preprint arXiv:2109.05578*.
- Huang, S.; Dong, L.; Wang, W.; Hao, Y.; Singhal, S.; Ma, S.; Lv, T.; Cui, L.; Mohammed, O. K.; Liu, Q.; et al. 2023. Language is not all you need: Aligning perception with language models. *arXiv preprint arXiv:2302.14045*.

Ouyang, L.; Wu, J.; Jiang, X.; Almeida, D.; Wainwright, C.; Mishkin, P.; Zhang, C.; Agarwal, S.; Slama, K.; Ray, A.; et al. 2022. Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems*, 35: 27730–27744.

Poli, M.; Massaroli, S.; Nguyen, E.; Fu, D. Y.; Dao, T.; Bacus, S.; Bengio, Y.; Ermon, S.; and Ré, C. 2023. Hyena Hierarchy: Towards Larger Convolutional Language Models. *arXiv preprint arXiv:2302.10866*.

Sarzynska-Wawer, J.; Wawer, A.; Pawlak, A.; Szymanowska, J.; Stefaniak, I.; Jarkiewicz, M.; and Okruszek, L. 2021. Detecting formal thought disorder by deep contextualized word representations. *Psychiatry Research*, 304: 114135.

Sheth, A.; Gaur, M.; Roy, K.; and Faldu, K. 2021. Knowledge-intensive language understanding for explainable ai. *IEEE Internet Computing*, 25(5): 19–24.

Sheth, A.; Gaur, M.; Roy, K.; Venkataraman, R.; and Khandelwal, V. 2022. Process knowledge-infused ai: Toward user-level explainability, interpretability, and safety. *IEEE Internet Computing*, 26(5): 76–84.

Sheth, A.; Roy, K.; and Gaur, M. 2023. Neurosymbolic Artificial Intelligence (Why, What, and How). *IEEE Intelligent Systems*, 38(3): 56–62.

Speer, R.; Chin, J.; and Havasi, C. 2017. Conceptnet 5.5: An open multilingual graph of general knowledge. In *Proceedings of the AAAI conference on artificial intelligence*, volume 31.

Tsai, Y.-H. H.; Bai, S.; Yamada, M.; Morency, L.-P.; and Salakhutdinov, R. 2019. Transformer dissection: a unified understanding of transformer’s attention via the lens of kernel. *arXiv preprint arXiv:1908.11775*.

Vaswani, A.; Shazeer, N.; Parmar, N.; Uszkoreit, J.; Jones, L.; Gomez, A. N.; Kaiser, Ł.; and Polosukhin, I. 2017. Attention is all you need. *Advances in neural information processing systems*, 30.

Wang, A.; Singh, A.; Michael, J.; Hill, F.; Levy, O.; and Bowman, S. R. 2018. GLUE: A multi-task benchmark and analysis platform for natural language understanding. *arXiv preprint arXiv:1804.07461*.

Zhou, C.; Li, Q.; Li, C.; Yu, J.; Liu, Y.; Wang, G.; Zhang, K.; Ji, C.; Yan, Q.; He, L.; et al. 2023. A comprehensive survey on pretrained foundation models: A history from bert to chatgpt. *arXiv preprint arXiv:2302.09419*.