

# KALM: KNOWLEDGE-AWARE INTEGRATION OF LOCAL, DOCUMENT, AND GLOBAL CONTEXTS FOR LONG DOCUMENT UNDERSTANDING

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

With the advent of pre-trained language models (LMs), increasing research efforts have been focusing on infusing commonsense and domain-specific knowledge to prepare LMs for downstream tasks. These works attempt to leverage knowledge graphs, the *de facto* standard of symbolic knowledge representation, along with pre-trained LMs. While existing approaches leverage external knowledge, it remains an open question how to jointly incorporate knowledge graphs representing varying contexts—from local (e.g., sentence), to document-level, to global knowledge—to enable knowledge-rich and interpretable exchange across these contexts. Such rich contextualization can be especially beneficial for long document understanding tasks since standard pre-trained LMs are typically bounded by the input sequence length. In light of these challenges, we propose **KALM**, a **Knowledge-Aware Language Model** that jointly leverages knowledge in local, document-level, and global contexts for long document understanding. KALM first encodes long documents and knowledge graphs into the three knowledge-aware context representations. It then processes each context with context-specific layers, followed by a “context fusion” layer that facilitates interpretable knowledge exchange to derive an overarching document representation. Extensive experiments demonstrate that KALM achieves state-of-the-art performance on three long document understanding tasks across 6 datasets/settings. Further analyses reveal that the three knowledge-aware contexts are complementary and they all contribute to model performance, while the importance and information exchange patterns of different contexts vary with respect to different tasks and datasets.

## 1 INTRODUCTION

Pre-trained language models (LMs) have become the dominant paradigm in NLP research, while knowledge graphs (KGs) are the *de facto* standard of symbolic knowledge representation. Recent advances in knowledge-aware NLP focus on combining the two paradigms (Wang et al., 2021b; Zhang et al., 2021; He et al., 2021), infusing encyclopedic (Vrandečić & Krötzsch, 2014; Pellissier Tanon et al., 2020), commonsense (Speer et al., 2017), and domain-specific (Feng et al., 2021a; Chang et al., 2020) knowledge with LMs. Knowledge-grounded models achieved state-of-the-art performance in tasks including question answering (Sun et al., 2022), commonsense reasoning (Kim et al., 2022; Liu et al., 2021), and social text analysis (Zhang et al., 2022; Hu et al., 2021).

Prior approaches to infusing LMs with knowledge typically focused on three hitherto orthogonal directions: incorporating knowledge related to local (e.g., sentence-level), document-level, or global context. **Local** context approaches argue that sentences mention entities, and the external knowledge of entities, such as textual descriptions (Balachandran et al., 2021; Wang et al., 2021b) and metadata (Ostapenko et al., 2022), help LMs realize they are more than just tokens. **Document-level** context approaches argue that core idea entities are repeatedly mentioned throughout the document, while related concepts might be discussed in different paragraphs. These methods attempt to leverage entities and knowledge across paragraphs with techniques such as document graphs (Feng et al., 2021a; Zhang et al., 2022; Hu et al., 2021). **Global** context approaches argue that unmentioned yet connecting entities help connect the dots for knowledge-based reasoning, thus knowledge graph subgraphs are encoded with graph neural networks alongside textual content (Zhang et al., 2021;

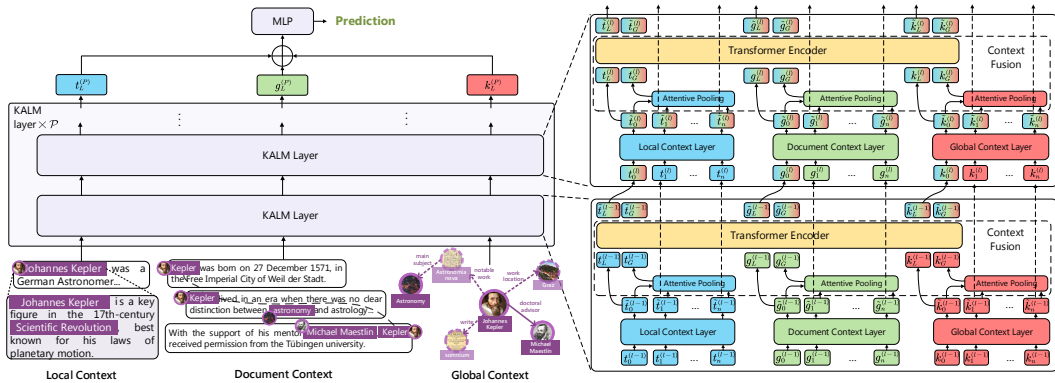


Figure 1: Overview of KALM, which encodes long documents and knowledge graphs into local, document, and global contexts while enabling interpretable information exchange across contexts.

Yasunaga et al., 2021). However, despite their individual pros and cons, how to integrate the three document contexts in a knowledge-aware and interpretable way remains an open problem.

Controlling for varying scopes of knowledge and context representations could benefit numerous language understanding tasks, especially those centered around long documents. Bounded by the inherent limitation of input sequence length, existing knowledge-aware LMs are mostly designed to handle short texts (Wang et al., 2021b; Zhang et al., 2021). However, processing long documents containing thousands of tokens (Beltagy et al., 2021) requires attending to varying document contexts, disambiguating long-distance co-referring entities and events, and more.

In light of these challenges, we propose **KALM**, a **Knowledge-Aware Language Model** for long document understanding. Specifically, KALM first derives three context- and knowledge-aware representations from the long input document and an external knowledge graph: the local context represented as raw text, the document-level context represented as a document graph, and the global context represented as a knowledge graph subgraph. KALM layers then encode each context with context-specific layers, followed by our proposed novel ContextFusion layers to enable knowledge-rich and interpretable information exchange across the three knowledge-aware contexts. A unified document representation is then derived from context-specific representations that also interact with other contexts. An illustration of the proposed KALM is presented in Figure 1.

While KALM is a general method for long document understanding, we evaluate the model on three tasks across six datasets/settings that are particularly sensitive to broader contexts and external knowledge: political perspective detection, misinformation detection, and roll call vote prediction. Extensive experiments demonstrate that KALM outperforms pre-trained LMs, task-agnostic knowledge-aware baselines, and strong task-specific baselines on all six datasets. In ablation experiments, we further establish KALM’s ability to enable information exchange, better handle long documents, and improve data efficiency. In addition, KALM and the proposed ContextFusion layers reveal and help interpret the roles and information exchange patterns of different contexts.

## 2 KALM METHODOLOGY

### 2.1 PROBLEM DEFINITION

Let  $\mathbf{d} = \{d_1, \dots, d_n\}$  denote a natural language document with  $n$  paragraphs, where each paragraph contains a sequence of  $n_i$  tokens  $\mathbf{d}_i = \{w_{i1}, \dots, w_{in_i}\}$ . Knowledge-aware long document understanding assumes the access to an external knowledge graph  $KG = (\mathcal{E}, \mathcal{R}, \mathbf{A}, \epsilon, \varphi)$ , where  $\mathcal{E} = \{e_1, \dots, e_N\}$  denotes the entity set,  $\mathcal{R} = \{r_1, \dots, r_M\}$  denotes the relation set,  $\mathbf{A}$  is the adjacency matrix where  $a_{ij} = k$  indicates  $(e_i, r_k, e_j) \in KG$ ,  $\epsilon(\cdot) : \mathcal{E} \rightarrow \text{str}$  and  $\varphi(\cdot) : \mathcal{R} \rightarrow \text{str}$  map the entities and relations to their textual descriptions.

Given pre-defined document labels, knowledge-aware natural language understanding aims to learn document representations and classify  $\mathbf{d}$  into its corresponding label with the help of  $KG$ .

## 2.2 KNOWLEDGE-AWARE CONTEXTS

We hypothesize that a holistic representation of long documents should incorporate contexts and relevant knowledge at three levels: the local context (e.g., a sentence with descriptions of mentioned entities), the broader document context (e.g., a long document with cross-paragraph entity reference structure), and the global/external context represented as external knowledge (e.g., relevant knowledge base subgraphs). Each of the three contexts uses different granularities of external knowledge, while existing works fall short of jointly integrating the three types of representations. To this end, KALM firstly employs different ways to introduce knowledge in different levels of contexts.

**Local context.** Represented as the raw text of sentences and paragraphs, the local context models the smallest unit in long document understanding. Prior works attempted to add sentence metadata (e.g., tense, sentiment, topic) (Zhang et al., 2022), adopt sentence-level pre-training tasks based on KG triples (Wang et al., 2021b), or leverage knowledge graph embeddings along with textual representations (Hu et al., 2021). While these methods were effective, in the face of LM-centered NLP research, they are ad-hoc add-ons and not fully compatible with existing pre-trained LMs. As a result, KALM proposes to directly concatenate the textual descriptions of entities  $\epsilon(e_i)$  to the paragraph if  $e_i$  is mentioned. In this way, the original text is directly augmented with the entity descriptions, informing the LM that entities such as "Kepler" are more than mere tokens and help to combat the spurious correlations of pre-trained LMs (McMilin, 2022). For each augmented paragraph  $d'_i$ , we adopt pre-trained  $\text{LM}(\cdot)$  and mean pooling to extract a paragraph representation. We use BART (Lewis et al., 2020) as  $\text{LM}(\cdot)$  without further notice. We also add a fusion token at the beginning of the paragraph sequence for information exchange across contexts. After processing all  $n$  paragraphs, we obtain the local context representation  $\mathbf{T}^{(0)}$  as follows:

$$\mathbf{T}^{(0)} = \{\mathbf{t}_0^{(0)}, \dots, \mathbf{t}_n^{(0)}\} = \{\theta_{rand}, \text{LM}(d'_1), \dots, \text{LM}(d'_n)\} \quad (1)$$

where  $\theta_{rand}$  denotes a randomly initialized vector of the fusion token in the local context and the superscript  $(0)$  indicates the 0-th layer.

**Document-level context.** Represented as the structure of the full document, the document-level context is responsible for modeling cross-paragraph entities and knowledge on a document level. While existing works attempted to incorporate external knowledge in documents via document graphs (Feng et al., 2021a; Hu et al., 2021), they fall short of leveraging the overlapping entities and concepts between paragraphs that underpin the reasoning of long documents. To this end, we propose **knowledge coreference**, a simple and effective mechanism for modeling text-knowledge interaction on the document level. Specifically, a document graph with  $n + 1$  nodes is constructed, consisting of one fusion node and  $n$  paragraph nodes. If paragraph  $i$  and  $j$  both mention entity  $e_k$  in the external knowledge base, node  $i$  and  $j$  in the document graph are then connected with relation type  $k$ . In addition, the fusion node is connected to every paragraph node with a super-relation. As a result, we obtain the adjacency matrix of the document graph  $\mathbf{A}^g$ . Paired with the knowledge-guided GNN to be introduced in Section 2.3, knowledge coreference enables the information flow across paragraphs guided by external knowledge. Node feature initialization of the document graph is as follows:

$$\mathbf{G}^{(0)} = \{\mathbf{g}_0^{(0)}, \dots, \mathbf{g}_n^{(0)}\} = \{\theta_{rand}, \text{LM}(d_1), \dots, \text{LM}(d_n)\} \quad (2)$$

**Global context.** Represented as external knowledge graphs, the global context is responsible for leveraging unseen entities and facilitating KG-based reasoning. Existing works mainly focused on extracting knowledge graph subgraphs (Yasunaga et al., 2021; Zhang et al., 2021) and encoding them alongside document content. Though many tricks are proposed to extract and prune knowledge graph subgraphs, in KALM, we employ a straightforward approach: for all mentioned entities in the long document, KALM merges their  $k$ -hop neighborhood to obtain a knowledge graph subgraph. Specifically, we use  $k = 2$  following previous works (Zhang et al., 2021; Vashishth et al., 2019), striking a balance between KB structure and computational efficiency while KALM could support any  $k$  settings. A fusion entity is then introduced and connected with every other entity, resulting in a connected graph. In this way, KALM cuts back on the preprocessing for modeling global knowledge

and better preserve the information in the KG. Knowledge graph embedding methods (Bordes et al., 2013) are then adopted to initialize node features of the KG subgraph:

$$\mathbf{K}^{(0)} = \{\mathbf{k}_0^{(0)}, \dots, \mathbf{k}_{|\rho(\mathbf{d})|}^{(0)}\} = \{\theta_{rand}, \text{KGE}(e_1), \dots, \text{KGE}(e_{|\rho(\mathbf{d})|})\} \quad (3)$$

where  $\text{KGE}(\cdot)$  denotes the knowledge graph embeddings trained on the original KG,  $|\rho(\mathbf{d})|$  indicates the number of mentioned entities identified in document  $\mathbf{d}$ . We use TransE (Bordes et al., 2013) to learn KB embeddings and use them for  $\text{KGE}(\cdot)$ , while the knowledge base embeddings are not adapted in the KALM learning process.

### 2.3 KALM LAYERS

After obtaining the local, document-level, and global context representations of long documents, we employ KALM layers to learn document representations. Specifically, each KALM layer consists of three context-specific layers to process each context. A ContextFusion layer is then adopted to enable the knowledge-rich and interpretable information exchange across the three contexts.

#### 2.3.1 CONTEXT-SPECIFIC LAYERS

**Local context layer.** The local context is represented as a sequence of vectors extracted from the knowledge-enriched text with the help of pre-trained LMs. We adopt transformer encoder layers (Vaswani et al., 2017) to encode the local context:

$$\tilde{\mathbf{T}}^{(\ell)} = \{\tilde{\mathbf{t}}_0^{(\ell)}, \dots, \tilde{\mathbf{t}}_n^{(\ell)}\} = \phi\left(\text{TrmEnc}(\{\mathbf{t}_0^{(\ell)}, \dots, \mathbf{t}_n^{(\ell)}\})\right) \quad (4)$$

where  $\phi(\cdot)$  denotes non-linearity,  $\text{TrmEnc}$  denotes the transformer encoder layer, and  $\tilde{\mathbf{t}}_0^{(\ell)}$  denotes the transformed representation of the fusion token.

**Document-level context layer.** The document-level context is represented as a document graph based on knowledge coreference. To better exploit the entity-based relations in the document graph, we propose a knowledge-aware GNN architecture to enable **knowledge-guided message passing** on the document graph:

$$\tilde{\mathbf{G}}^{(\ell)} = \{\tilde{\mathbf{g}}_0^{(\ell)}, \dots, \tilde{\mathbf{g}}_n^{(\ell)}\} = \text{GNN}\left(\{\mathbf{g}_0^{(\ell)}, \dots, \mathbf{g}_n^{(\ell)}\}\right) \quad (5)$$

where  $\text{GNN}(\cdot)$  denotes the proposed knowledge-guided graph neural networks as follows:

$$\tilde{\mathbf{g}}_i^{(\ell)} = \phi\left(\alpha_{i,i} \Theta \mathbf{g}_i^{(\ell)} + \sum_{j \in \mathcal{N}(i)} \Theta \mathbf{g}_j^{(\ell)}\right) \quad (6)$$

where  $\alpha_{i,j}$  denotes the knowledge-guided attention weight and is defined as follows:

$$\alpha_{i,j} = \frac{\exp\left(\text{ELU}(\mathbf{a}^T [\Theta \mathbf{g}_i^{(\ell)} \parallel \Theta \mathbf{g}_j^{(\ell)} \parallel \Theta f(\text{KGE}(a_{ij}^g))])\right)}{\sum_{k \in \mathcal{N}(i)} \exp\left(\text{ELU}(\mathbf{a}^T [\Theta \mathbf{g}_i^{(\ell)} \parallel \Theta \mathbf{g}_k^{(\ell)} \parallel \Theta f(\text{KGE}(a_{ik}^g))])\right)} \quad (7)$$

where  $\tilde{\mathbf{g}}_0^{(\ell)}$  denotes the transformed representation of the fusion node,  $\mathbf{a}$  and  $\Theta$  are learnable parameters,  $a_{ij}^g$  is the  $i$ -th row and  $j$ -th column value of adjacency matrix  $\mathbf{A}^g$  of the document graph,  $\text{ELU}$  denotes the exponential linear unit activation function (Clevert et al., 2015), and  $f(\cdot)$  is a learnable linear layer. The term of  $\Theta f(\text{KGE}(a_{ij}^g))$  is responsible for enabling the knowledge-guided message passing on the document graph, enabling KALM to incorporate the entity and concept patterns in different paragraphs and their document-level interactions.

**Global context layer.** The global context is represented as a relevant knowledge graph subgraph. We follow previous works and adopt GATs (Veličković et al., 2018) to encode the global context:

$$\tilde{\mathbf{K}}^{(\ell)} = \{\tilde{\mathbf{k}}_0^{(\ell)}, \dots, \tilde{\mathbf{k}}_{|\rho(d)|}^{(\ell)}\} = \text{GAT}\left(\{\mathbf{k}_0^{(\ell)}, \dots, \mathbf{k}_{|\rho(d)|}^{(\ell)}\}\right) \quad (8)$$

where  $\tilde{\mathbf{k}}_0^{(\ell)}$  denotes the transformed representation of the fusion entity.

### 2.3.2 CONTEXTFUSION LAYER

The local, document, and global contexts model external knowledge within sentences, across the document, and beyond the document. These contexts are closely connected and a robust long document understanding method should reflect their interactions. Existing approaches mostly leverage only one or two of the contexts (Wang et al., 2021b; Feng et al., 2021a; Zhang et al., 2022), falling short of jointly leveraging the three knowledge-aware contexts. In addition, they mostly adopted direct concatenation or MLP layers (Zhang et al., 2022; 2021; Hu et al., 2021), falling short of enabling context-specific information to flow across contexts in a knowledge-rich and interpretable manner. As a result, we propose the ContextFusion layer to tackle these challenges. We firstly take a local perspective and extract the representations of the fusion tokens, nodes, and entities in each context:

$$\left[\mathbf{t}_L^{(\ell)}, \mathbf{g}_L^{(\ell)}, \mathbf{k}_L^{(\ell)}\right] = \left[\tilde{\mathbf{t}}_0^{(\ell)}, \tilde{\mathbf{g}}_0^{(\ell)}, \tilde{\mathbf{k}}_0^{(\ell)}\right] \quad (9)$$

We then take a global perspective and use the fusion token/node/entity as the query to conduct attentive pooling  $\text{ap}(\cdot, \cdot)$  across all other tokens/nodes/entities in each context:

$$\left[\mathbf{t}_G^{(\ell)}, \mathbf{g}_G^{(\ell)}, \mathbf{k}_G^{(\ell)}\right] = \left[\text{ap}(\tilde{\mathbf{t}}_0^{(\ell)}, \{\tilde{\mathbf{t}}_i^{(\ell)}\}_{i=1}^n), \text{ap}(\tilde{\mathbf{g}}_0^{(\ell)}, \{\tilde{\mathbf{g}}_i^{(\ell)}\}_{i=1}^n), \text{ap}(\tilde{\mathbf{k}}_0^{(\ell)}, \{\tilde{\mathbf{k}}_i^{(\ell)}\}_{i=1}^n)\right] \quad (10)$$

where attentive pooling  $\text{ap}(\cdot, \cdot)$  is defined as follows:

$$\text{ap}(\mathbf{q}, \{\mathbf{k}_i\}_{i=1}^n) = \sum_{i=1}^n \frac{\exp(\mathbf{q} \cdot \mathbf{k}_i)}{\sum_{j=1}^n \exp(\mathbf{q} \cdot \mathbf{k}_j)} \mathbf{k}_i \quad (11)$$

In this way, the fusion token/node/entity in each context serve as the information exchange portals. We then use a transformer encoder layer to enable information exchange across the contexts:

$$\left[\tilde{\mathbf{t}}_L^{(\ell)}, \tilde{\mathbf{g}}_L^{(\ell)}, \tilde{\mathbf{k}}_L^{(\ell)}, \tilde{\mathbf{t}}_G^{(\ell)}, \tilde{\mathbf{g}}_G^{(\ell)}, \tilde{\mathbf{k}}_G^{(\ell)}\right] = \phi\left(\text{TrmEnc}\left(\left[\mathbf{t}_L^{(\ell)}, \mathbf{g}_L^{(\ell)}, \mathbf{k}_L^{(\ell)}, \mathbf{t}_G^{(\ell)}, \mathbf{g}_G^{(\ell)}, \mathbf{k}_G^{(\ell)}\right]\right)\right) \quad (12)$$

As a result,  $\tilde{\mathbf{t}}_L^{(\ell)}$ ,  $\tilde{\mathbf{g}}_L^{(\ell)}$ , and  $\tilde{\mathbf{k}}_L^{(\ell)}$  are the representations of the fusion token/node/entity that incorporates information from other contexts. We formulate the output of the  $l$ -th layer as follows:

$$\mathbf{T}^{(\ell+1)} = \{\tilde{\mathbf{t}}_L^{(\ell)}, \tilde{\mathbf{t}}_1^{(\ell)}, \dots, \tilde{\mathbf{t}}_n^{(\ell)}\}, \mathbf{G}^{(\ell+1)} = \{\tilde{\mathbf{g}}_L^{(\ell)}, \tilde{\mathbf{g}}_1^{(\ell)}, \dots, \tilde{\mathbf{g}}_n^{(\ell)}\}, \mathbf{K}^{(\ell+1)} = \{\tilde{\mathbf{k}}_L^{(\ell)}, \tilde{\mathbf{k}}_1^{(\ell)}, \dots, \tilde{\mathbf{k}}_n^{(\ell)}\} \quad (13)$$

Our proposed ContextFusion layer is interactive since it enables the information to flow across different document contexts, instead of direct concatenation or hierarchical processing. The ContextFusion layer is interpretable since the attention weights in  $\text{TrmEnc}(\cdot)$  could provide insights into the roles and importance of each document context, which will be further explored in Section 3.3. To the best of our knowledge, KALM is the first work to jointly consider the three levels of document context and enable interpretable information exchange across document contexts.

Table 1: Model performance on three tasks and six datasets. Acc, MaF, miF, and BAcc denote accuracy, macro-averaged F1-score, micro-averaged F1-score, and balanced accuracy. Certain task-specific models did not report standard deviation in the original paper.

Task	Dataset	Metric	Task SOTA	Best LM	Knowledge-Aware LMs					KALM	
					KELM	KnowBERT	Joshi et al.	KGAP	GreaseLM		GreaseLM+
PDD	SemEval	Acc	89.90 ( $\pm 0.6$ )	86.99 ( $\pm 1.9$ )	86.40 ( $\pm 2.3$ )	84.73 ( $\pm 3.4$ )	81.88 ( $\pm 2.1$ )	87.73 ( $\pm 1.8$ )	86.64 ( $\pm 1.5$ )	85.66 ( $\pm 1.8$ )	<b>91.45</b> ( $\pm 0.8$ )
		MaF	86.11 ( $\pm 1.1$ )	80.62 ( $\pm 3.8$ )	83.98 ( $\pm 1.0$ )	75.72 ( $\pm 5.3$ )	77.15 ( $\pm 3.8$ )	82.00 ( $\pm 3.1$ )	80.32 ( $\pm 3.0$ )	77.23 ( $\pm 4.1$ )	<b>87.65</b> ( $\pm 1.2$ )
	Allsides	Acc	87.17 ( $\pm 0.2$ )	68.71 ( $\pm 4.3$ )	80.71 ( $\pm 2.4$ )	66.56 ( $\pm 0.7$ )	80.88 ( $\pm 2.1$ )	83.65 ( $\pm 1.3$ )	80.23 ( $\pm 1.2$ )	82.16 ( $\pm 5.5$ )	<b>87.26</b> ( $\pm 0.2$ )
		MaF	86.72 ( $\pm 0.3$ )	65.39 ( $\pm 5.7$ )	79.74 ( $\pm 2.7$ )	58.81 ( $\pm 0.5$ )	79.73 ( $\pm 2.3$ )	82.92 ( $\pm 1.4$ )	79.17 ( $\pm 1.2$ )	80.81 ( $\pm 7.1$ )	<b>86.79</b> ( $\pm 0.2$ )
MD	SLN	miF	89.17	88.17 ( $\pm 0.6$ )	84.11 ( $\pm 0.6$ )	78.67 ( $\pm 3.2$ )	82.72 ( $\pm 5.1$ )	92.17 ( $\pm 1.2$ )	73.83 ( $\pm 0.9$ )	88.17 ( $\pm 0.8$ )	<b>94.22</b> ( $\pm 1.2$ )
		MaF	94.18 ( $\pm 1.1$ )	88.46 ( $\pm 4.9$ )	82.80 ( $\pm 1.3$ )	79.80 ( $\pm 2.0$ )	83.98 ( $\pm 3.7$ )	92.30 ( $\pm 0.9$ )	75.20 ( $\pm 0.8$ )	88.64 ( $\pm 0.6$ )	<b>94.18</b> ( $\pm 1.1$ )
	LUN	miF	69.05	60.10 ( $\pm 1.7$ )	59.28 ( $\pm 2.1$ )	59.66 ( $\pm 1.1$ )	58.57 ( $\pm 3.4$ )	65.52 ( $\pm 2.3$ )	56.54 ( $\pm 1.5$ )	64.29 ( $\pm 2.4$ )	<b>71.28</b> ( $\pm 1.7$ )
		MaF	68.26	58.57 ( $\pm 2.1$ )	57.30 ( $\pm 1.6$ )	59.19 ( $\pm 1.3$ )	56.73 ( $\pm 4.0$ )	63.94 ( $\pm 2.9$ )	55.75 ( $\pm 1.6$ )	62.65 ( $\pm 3.7$ )	<b>69.82</b> ( $\pm 1.2$ )
RCVP	Random	BAcc	90.33	89.94 ( $\pm 0.2$ )	89.13 ( $\pm 1.1$ )	86.72 ( $\pm 0.9$ )	92.43 ( $\pm 0.5$ )	77.98 ( $\pm 0.5$ )	89.99 ( $\pm 1.5$ )	91.01 ( $\pm 0.2$ )	<b>92.36</b> ( $\pm 0.4$ )
		MaF	84.92	86.10 ( $\pm 0.7$ )	84.76 ( $\pm 2.0$ )	79.33 ( $\pm 2.4$ )	89.64 ( $\pm 0.6$ )	68.11 ( $\pm 6.0$ )	84.72 ( $\pm 3.0$ )	87.29 ( $\pm 0.3$ )	<b>89.33</b> ( $\pm 0.4$ )
	Time-based	BAcc	89.92	90.40 ( $\pm 0.8$ )	90.80 ( $\pm 0.2$ )	87.07 ( $\pm 0.9$ )	92.63 ( $\pm 1.6$ )	77.90 ( $\pm 0.6$ )	88.21 ( $\pm 2.7$ )	91.69 ( $\pm 0.1$ )	<b>94.46</b> ( $\pm 0.4$ )
		MaF	84.35	85.21 ( $\pm 2.1$ )	86.62 ( $\pm 0.4$ )	78.90 ( $\pm 1.9$ )	89.31 ( $\pm 2.4$ )	70.81 ( $\pm 4.6$ )	79.73 ( $\pm 7.4$ )	87.95 ( $\pm 0.3$ )	<b>91.97</b> ( $\pm 0.5$ )

## 2.4 LEARNING AND INFERENCE

After a total of  $\mathcal{P}$  KALM layers, we obtain the final document representation as  $[\tilde{\mathbf{t}}_L^{(\mathcal{P})}, \tilde{\mathbf{g}}_L^{(\mathcal{P})}, \tilde{\mathbf{k}}_L^{(\mathcal{P})}]$ . Given the document label  $a \in \mathcal{A}$ , the label probability is formulated as  $p(a|\mathbf{d}) \propto \exp(\text{MLP}_a([\tilde{\mathbf{t}}_L^{(\mathcal{P})}, \tilde{\mathbf{g}}_L^{(\mathcal{P})}, \tilde{\mathbf{k}}_L^{(\mathcal{P})}]))$ . We then optimize KALM with the cross entropy loss function. At inference time, the predicted label is  $\text{argmax}_a p(a|\mathbf{d})$ .

## 3 EXPERIMENTS

### 3.1 EXPERIMENT SETTINGS

**Tasks and Datasets.** We propose KALM, a general method for knowledge-aware long document understanding. We evaluate KALM on three tasks that especially benefit from external knowledge and broader context: political perspective detection, misinformation detection, and roll call vote prediction. We follow previous works to adopt SemEval (Kiesel et al., 2019) and Allsides (Li & Goldwasser, 2019) for political perspective detection, LUN (Rashkin et al., 2017) and SLN (Rubin et al., 2016) for misinformation detection, and the 2 datasets proposed in Mou et al. (2021) for roll call vote prediction. For external KGs, we follow existing works to adopt the KGs in KGAP (Feng et al., 2021a), CompareNet (Hu et al., 2021), and ConceptNet (Speer et al., 2017) for the three tasks.

**Baseline methods.** We compare KALM with three types of baseline methods for holistic evaluation: pre-trained LMs, task-agnostic knowledge-aware methods, and task-specific models. For pre-trained LMs, we evaluate RoBERTa (Liu et al., 2019b), Electra (Clark et al., 2019), DeBERTa (He et al., 2020), BART (Lewis et al., 2020), and LongFormer (Beltagy et al., 2020) on the three tasks. For task-agnostic baselines, we evaluate KGAP (Feng et al., 2021a), GreaseLM (Zhang et al., 2021), and GreaseLM+ on the three tasks. Task-specific models are introduced in the following sections. For pre-trained LMs, task-agnostic methods, and KALM, we run each method five times and report the average performance and standard deviation. For task-specific models, we compare with the results originally reported since we follow the exact same experiment settings and data splits.

### 3.2 MODEL PERFORMANCE

We present the performance of task-specific methods, pre-trained language models, task-agnostic knowledge-aware baselines, and KALM in Table 1. We select the best-performing task-specific baseline (Task SOTA) and pre-trained language model (BestLM), while the full results are available in Tables 4, 5, and 6 in the appendix. Table 1 demonstrates that:

- KALM consistently outperforms all task-specific models, pre-trained language models, and knowledge-aware methods on all three tasks and six datasets/settings. Statistical significance tests in Section C.5 further demonstrates KALM’s superiority over existing models.
- Knowledge-aware language models generally outperform pre-trained LMs, which did not incorporate external knowledge bases in the pre-training process. This suggests that incorporating external knowledge bases could enrich document representations and boost downstream task performance.

Table 2: Ablation study of the three document contexts and the ContextFusion layer. The local, document, and global contexts all contribute to model performance, while the ContextFusion layer is better than existing strategies at enabling information exchange across contexts.

Task	Dataset	Metric	Ours	Remove Context			Substitute ContextFusion		
			KALM	w/o local	w/o document	w/o global	MInt	concat	sum
PDD	SemEval	Acc	<b>91.45</b> ( $\pm 0.8$ )	83.55 ( $\pm 0.8$ )	83.57 ( $\pm 1.1$ )	84.11 ( $\pm 0.9$ )	81.91 ( $\pm 0.9$ )	83.52 ( $\pm 1.8$ )	83.21 ( $\pm 1.0$ )
		MaF	<b>87.65</b> ( $\pm 1.2$ )	74.25 ( $\pm 1.3$ )	76.13 ( $\pm 2.0$ )	74.92 ( $\pm 1.8$ )	70.47 ( $\pm 3.6$ )	74.27 ( $\pm 4.0$ )	73.59 ( $\pm 2.1$ )
	Allsides	Acc	<b>87.26</b> ( $\pm 0.2$ )	83.72 ( $\pm 4.0$ )	82.88 ( $\pm 5.1$ )	80.59 ( $\pm 6.3$ )	83.08 ( $\pm 4.0$ )	83.27 ( $\pm 4.2$ )	83.50 ( $\pm 3.5$ )
		MaF	<b>86.79</b> ( $\pm 0.2$ )	83.10 ( $\pm 4.2$ )	81.86 ( $\pm 6.2$ )	78.98 ( $\pm 8.1$ )	82.39 ( $\pm 4.2$ )	82.28 ( $\pm 5.3$ )	82.64 ( $\pm 4.0$ )
MD	SLN	MiF	<b>94.22</b> ( $\pm 1.2$ )	80.94 ( $\pm 5.5$ )	83.50 ( $\pm 5.7$ )	83.94 ( $\pm 4.7$ )	86.33 ( $\pm 2.1$ )	82.67 ( $\pm 9.2$ )	79.89 ( $\pm 6.3$ )
		MaF	<b>94.18</b> ( $\pm 1.1$ )	82.95 ( $\pm 4.4$ )	85.55 ( $\pm 4.4$ )	85.65 ( $\pm 3.4$ )	86.79 ( $\pm 1.9$ )	85.26 ( $\pm 6.2$ )	82.71 ( $\pm 4.1$ )
	LUN	MiF	<b>71.28</b> ( $\pm 1.7$ )	41.13 ( $\pm 5.8$ )	50.18 ( $\pm 6.3$ )	57.94 ( $\pm 4.1$ )	48.78 ( $\pm 6.3$ )	53.52 ( $\pm 6.5$ )	63.27 ( $\pm 4.0$ )
		MaF	<b>69.82</b> ( $\pm 1.2$ )	35.95 ( $\pm 7.3$ )	47.27 ( $\pm 7.3$ )	55.58 ( $\pm 4.6$ )	44.11 ( $\pm 9.0$ )	48.98 ( $\pm 7.9$ )	61.86 ( $\pm 4.4$ )
RCVP	Random	BAcc	<b>92.36</b> ( $\pm 0.3$ )	91.29 ( $\pm 2.4$ )	91.35 ( $\pm 0.4$ )	91.34 ( $\pm 0.5$ )	92.14 ( $\pm 0.5$ )	91.82 ( $\pm 0.8$ )	91.18 ( $\pm 1.5$ )
		MaF	89.33 ( $\pm 0.4$ )	88.16 ( $\pm 2.5$ )	87.81 ( $\pm 0.8$ )	88.50 ( $\pm 0.4$ )	<b>89.35</b> ( $\pm 0.7$ )	89.01 ( $\pm 1.0$ )	88.19 ( $\pm 1.6$ )
	Time-based	BAcc	<b>94.46</b> ( $\pm 0.4$ )	93.58 ( $\pm 1.4$ )	93.47 ( $\pm 0.5$ )	93.91 ( $\pm 0.5$ )	93.06 ( $\pm 1.7$ )	92.37 ( $\pm 2.2$ )	93.06 ( $\pm 1.0$ )
		MaF	<b>91.97</b> ( $\pm 0.5$ )	90.60 ( $\pm 2.1$ )	90.73 ( $\pm 0.6$ )	91.29 ( $\pm 0.5$ )	90.06 ( $\pm 2.4$ )	88.56 ( $\pm 4.5$ )	90.21 ( $\pm 1.1$ )

- GreaseLM+ outperforms GreaseLM by adding the global context, which suggests the importance of jointly leveraging the three document contexts. KALM further introduces interpretable information exchange across contexts through the ContextFusion layer and achieves state-of-the-art performance. We further investigate the importance of three document contexts and the ContextFusion layer in Section 2.3.2.

### 3.3 CONTEXT EXCHANGE STUDY

By jointly modeling three document contexts and employing the ContextFusion layer, KALM facilitates interpretable exchange across the three document contexts. We conduct an ablation study to examine whether the contexts and the ContextFusion layer are essential in the KALM architecture. Specifically, we remove the three contexts one at a time and change the ContextFusion layer into MInt (Zhang et al., 2021), concatenation, and sum. We present the results in Table 2, which demonstrates that:

- All three levels of document contexts, local, document, and global, contribute to model performance. These results substantiate the necessity of jointly leveraging the three document contexts for long document understanding.
- When substituting our proposed ContextFusion layers with three existing combination strategies, MInt (Zhang et al., 2021), direct concatenation, and summation, performance drops are observed across multiple datasets. This suggests that the proposed ContextFusion layer successfully boost model performance by enabling information exchange across contexts.

In addition to boosting model performance, the ContextFusion layer enables the interpretation of how different document contexts contribute to document understanding. We calculate the average of attention weights’ absolute values of the multi-head attention in the TrmEnc( $\cdot$ ) layer of ContextFusion and illustrate in Figure 2. It is demonstrated that the three contexts’ contribution and information exchange patterns vary with respect to datasets and KALM layers. Specifically, local and global contexts are important for the LUN dataset, document and global contexts are important for the task of roll call vote prediction, and the SLN dataset equally leverages the three contexts. However, for the task of political perspective detection, the importance of the three aspects varies with the depth of KALM layers. This is especially salient on SemEval, where KALM firstly takes a view of the whole document, then draws from both local and document-level contexts, and closes by leveraging global knowledge to derive an overall document representation.

In summary, the ContextFusion layer in KALM successfully identifies the relative importance and information exchange patterns of the three contexts, providing insights into how KALM arrives at the conclusion and which context information should be the focus of future research. We further demonstrate that the role and importance of each context changes as training progresses in Section C.2 in the appendix.

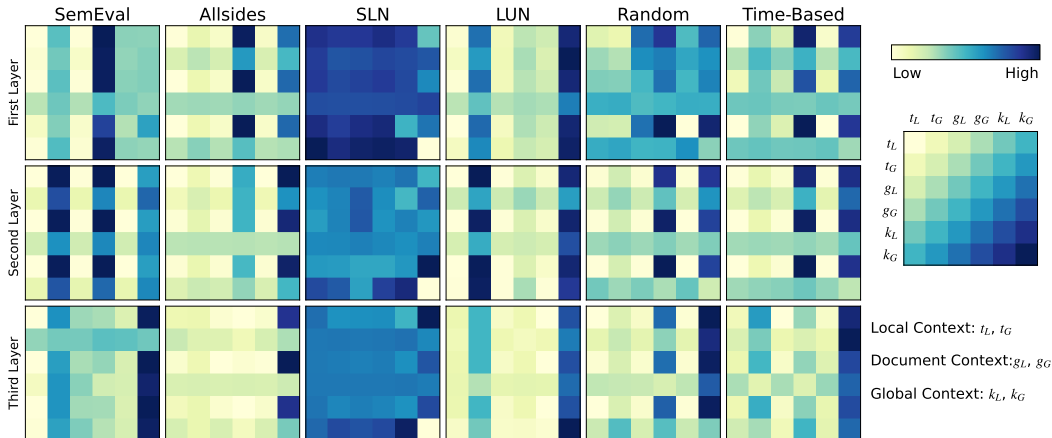


Figure 2: Interpreting the roles of the three contexts with attention maps in the ContextFusion layer.  $t_L, t_G, g_L, g_G, k_L, k_G$  denote the context representations in equations (9) and (10), so that the first two columns indicate how the local context attends to information in other contexts, the next two columns for the document-level context, and the last two columns for the global context.

### 3.4 LONG DOCUMENT STUDY

KALM complements the scarce literature in knowledge-aware long document understanding. In addition to more input tokens, the understanding of long documents often relies on more knowledge reference and knowledge reasoning. To examine whether KALM indeed improved in the face of longer documents and more external knowledge, we illustrate the performance of KALM and competitive baselines with respect to document length and knowledge intensity in Figure 3. Specifically, we use the number of mentioned entities to represent knowledge intensity and the number of sentences to represent document length, mapping each data point onto a two-dimensional space. It is illustrated that while baseline methods are prone to mistakes when the document is long and knowledge is rich, KALM alleviates this issue and performs better in the top-right corner. We further analyze KALM and more baseline methods’ performance on long documents with great knowledge intensity in Figure 6 in the appendix.

## 4 RELATED WORK

Knowledge graphs are playing an increasingly important role in language models and NLP research. Commonsense (Speer et al., 2017; Ilievski et al., 2021; Bosselut et al., 2019; West et al., 2022; Li et al., 2022a) and domain-specific KGs (Feng et al., 2021a; Li et al., 2022b; Gyori et al., 2017) serve as external knowledge to augment pre-trained LMs, which achieves state-of-the-art performance on question answering (Zhang et al., 2021; Yasunaga et al., 2021; Mitra et al., 2022; Bosselut et al., 2021; Oguz et al., 2022; Feng et al., 2022; Heo et al., 2022; Ma et al., 2022), social text analysis (Hu et al., 2021; Zhang et al., 2022), commonsense reasoning (Kim et al., 2022; Jung et al., 2022; Amayuelas et al., 2021; Liu et al., 2022), and text generation (Rony et al., 2022). These approaches (Lu et al., 2022; Zhang et al., 2019; Yu et al., 2022b; Sun et al., 2020; Yamada et al., 2020; Qiu et al., 2019b; Xie et al., 2022) could be mainly categorized by the three levels of the context where knowledge injection happens.

**Local** context approaches focus on entity mentions and external knowledge in individual sentences to enable fine-grained knowledge inclusion. A straightforward way is to encode KG entities with KG embeddings (Bordes et al., 2013; Lin et al., 2015; Cucala et al., 2021; Sun et al., 2018) and infuse the embeddings with language representations (Hu et al., 2021; Feng et al., 2021a; Kang et al., 2022). Later approaches focus on augmenting pre-trained LMs with KGs by introducing knowledge-aware training tasks and LM architectures (Wang et al., 2021b;a; Sridhar & Yang, 2022; Moiseev et al., 2022; Kaur et al., 2022; Hu et al., 2022; Arora et al., 2022; de Jong et al., 2021; Meng et al., 2021; He et al., 2021). Topic models were also introduced to enrich document representation learning



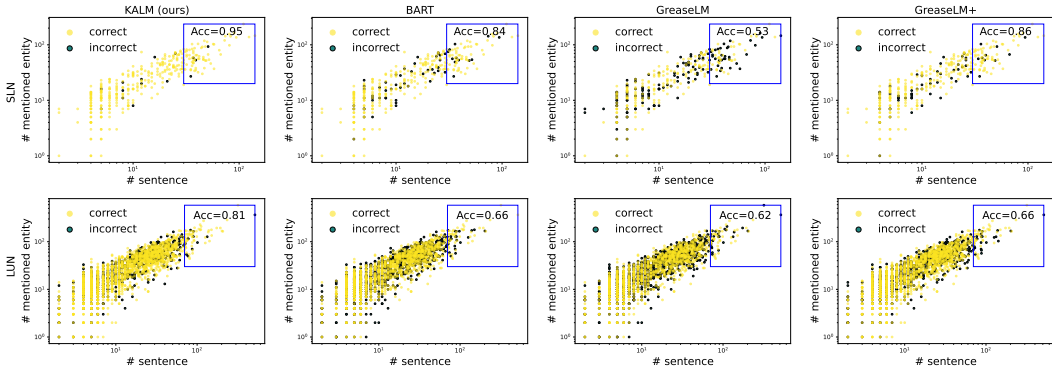


Figure 3: Error analysis of KALM and baseline methods. KALM successfully improves in the top-right corner, which represents documents with more sentences and more entailed knowledge.

(Gupta et al., 2018; Chaudhary et al., 2020; Wang et al., 2018). However, local context approaches fall short of leveraging inter-sentence and inter-entity knowledge, resulting in models that could not grasp the full picture of the text-knowledge interactions.

**Document-level** context approaches take a view of the whole document, jointly considering external knowledge across multiple sentences and paragraphs. The predominant way of achieving document-level knowledge understanding is through the construction of "document graphs" (Zhang et al., 2022), where the textual content, external knowledge bases, and other sources of information are encoded and represented as different components in graphs, often heterogeneous information networks (Hu et al., 2021; Feng et al., 2021a; Zhang et al., 2022; Yu et al., 2022a). Graph neural networks are then employed to learn graph representations, which fuses both textual information and external KGs. However, document-level context approaches fall short of preserving the original KG structure, resulting in models with reduced knowledge reasoning abilities.

**Global** context approaches focus on the external KG, extracting relevant KG subgraphs based on entity mentions of the long document. Pruned with certain mechanisms (Yasunaga et al., 2021) or not (Qiu et al., 2019a), these KG subgraphs are encoded with GNNs, and such representations are fused with LMs from simple concatenation (Hu et al., 2021) to deeper interactions (Zhang et al., 2021). However, global context approaches leverage external KGs in a stand-alone manner, falling short of enabling the dynamic integration of textual content and external KGs.

While existing approaches successfully introduced external KG to LMs, long document understanding poses new challenges to knowledge-aware NLP. Long documents possess greater knowledge intensity where more entities are mentioned, more relations are leveraged, and more reasoning is required to fully understand the nuances, while existing approaches are mostly designed for sparse knowledge scenarios. In addition, long documents also exhibit the phenomenon of knowledge co-reference, where central ideas and entities are reiterated throughout the document and co-exist in different levels of document contexts. In light of these challenges, we propose KALM to jointly leverage the local, document, and global contexts of long documents for knowledge incorporation.

## 5 CONCLUSION

In this paper, we propose KALM, a knowledge-aware long document understanding approach that introduces external knowledge to three levels of document contexts and enables interactive and interpretable exchange across them. Extensive experiments demonstrate that KALM achieves state-of-the-art performance on three tasks across six datasets. Our analysis shows that KALM provides insights into the roles and patterns of individual contexts, improves the handling of long documents with greater knowledge intensity, and has better data efficiency than existing works.

## REFERENCES

- Oshin Agarwal, Heming Ge, Siamak Shakeri, and Rami Al-Rfou. Knowledge graph based synthetic corpus generation for knowledge-enhanced language model pre-training. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pp. 3554–3565, 2021.
- Alfonso Amayuelas, Shuai Zhang, Xi Susie Rao, and Ce Zhang. Neural methods for logical reasoning over knowledge graphs. In *International Conference on Learning Representations*, 2021.
- Simran Arora, Sen Wu, Enci Liu, and Christopher Re. Metadata shaping: A simple approach for knowledge-enhanced language models. In *Findings of the Association for Computational Linguistics: ACL 2022*, pp. 1733–1745, Dublin, Ireland, May 2022.
- Vidhisha Balachandran, Bhuwan Dhingra, Haitian Sun, Michael Collins, and William W Cohen. Investigating the effect of background knowledge on natural questions. *NAACL-HLT 2021*, pp. 25, 2021.
- Iz Beltagy, Matthew E Peters, and Arman Cohan. Longformer: The long-document transformer. *arXiv preprint arXiv:2004.05150*, 2020.
- Iz Beltagy, Arman Cohan, Hannaneh Hajishirzi, Sewon Min, and Matthew E Peters. Beyond paragraphs: Nlp for long sequences. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies: Tutorials*, pp. 20–24, 2021.
- Steven Bird, Ewan Klein, and Edward Loper. *Natural language processing with Python: analyzing text with the natural language toolkit*. ” O’Reilly Media, Inc.”, 2009.
- Antoine Bordes, Nicolas Usunier, Alberto Garcia-Duran, Jason Weston, and Oksana Yakhnenko. Translating embeddings for modeling multi-relational data. *Advances in neural information processing systems*, 26, 2013.
- Antoine Bosselut, Hannah Rashkin, Maarten Sap, Chaitanya Malaviya, Asli Celikyilmaz, and Yejin Choi. Comet: Commonsense transformers for automatic knowledge graph construction. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 4762–4779, 2019.
- Antoine Bosselut, Ronan Le Bras, and Yejin Choi. Dynamic neuro-symbolic knowledge graph construction for zero-shot commonsense question answering. In *Proceedings of the 35th AAAI Conference on Artificial Intelligence (AAAI)*, 2021.
- David Chang, Ivana Balažević, Carl Allen, Daniel Chawla, Cynthia Brandt, and Andrew Taylor. Benchmark and best practices for biomedical knowledge graph embeddings. In *Proceedings of the 19th SIGBioMed Workshop on Biomedical Language Processing*, pp. 167–176, Online, July 2020.
- Yatin Chaudhary, Hinrich Schütze, and Pankaj Gupta. Explainable and discourse topic-aware neural language understanding. In *International Conference on Machine Learning*, pp. 1479–1488. PMLR, 2020.
- Kevin Clark, Minh-Thang Luong, Quoc V Le, and Christopher D Manning. Electra: Pre-training text encoders as discriminators rather than generators. In *International Conference on Learning Representations*, 2019.
- Djork-Arné Clevert, Thomas Unterthiner, and Sepp Hochreiter. Fast and accurate deep network learning by exponential linear units (elus). *arXiv preprint arXiv:1511.07289*, 2015.
- David Jaime Tena Cucala, Bernardo Cuenca Grau, Egor V Kostylev, and Boris Motik. Explainable gnn-based models over knowledge graphs. In *International Conference on Learning Representations*, 2021.
- Michiel de Jong, Yury Zemlyanskiy, Nicholas FitzGerald, Fei Sha, and William W Cohen. Mention memory: incorporating textual knowledge into transformers through entity mention attention. In *International Conference on Learning Representations*, 2021.

- Yupei Du, Qi Zheng, Yuanbin Wu, Man Lan, Yan Yang, and Meirong Ma. Understanding gender bias in knowledge base embeddings. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 1381–1395, Dublin, Ireland, May 2022.
- William Falcon and The PyTorch Lightning team. PyTorch Lightning, 3 2019.
- Shangbin Feng, Zilong Chen, Wenqian Zhang, Qingyao Li, Qinghua Zheng, Xiaojun Chang, and Minnan Luo. Kgap: Knowledge graph augmented political perspective detection in news media. *arXiv preprint arXiv:2108.03861*, 2021a.
- Shangbin Feng, Zhaoxuan Tan, Zilong Chen, Peisheng Yu, Qinghua Zheng, Xiao Chang, and Minnan Luo. Legislator representation learning with social context and expert knowledge. 2021b.
- Yue Feng, Zhen Han, Mingming Sun, and Ping Li. Multi-hop open-domain question answering over structured and unstructured knowledge. In *Findings of the Association for Computational Linguistics: NAACL 2022*, pp. 151–156, Seattle, United States, July 2022.
- Paolo Ferragina and Ugo Scaiella. Fast and accurate annotation of short texts with wikipedia pages. *IEEE software*, 29(1):70–75, 2011.
- Matthias Fey and Jan Eric Lenssen. Fast graph representation learning with pytorch geometric. *arXiv preprint arXiv:1903.02428*, 2019.
- Joseph Fisher, Dave Palfrey, Christos Christodoulopoulos, and Arpit Mittal. Measuring social bias in knowledge graph embeddings. *arXiv preprint arXiv:1912.02761*, 2019.
- Joseph Fisher, Arpit Mittal, Dave Palfrey, and Christos Christodoulopoulos. Debiasing knowledge graph embeddings. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 7332–7345, 2020.
- Sean M Gerrish and David M Blei. Predicting legislative roll calls from text. In *Proceedings of the 28th International Conference on Machine Learning, ICML 2011*, 2011.
- Pankaj Gupta, Yatin Chaudhary, Florian Buettner, and Hinrich Schuetze. texttvec: Deep contextualized neural autoregressive topic models of language with distributed compositional prior. In *International Conference on Learning Representations*, 2018.
- Benjamin M Gyori, John A Bachman, Kartik Subramanian, Jeremy L Muhlich, Lucian Galescu, and Peter K Sorger. From word models to executable models of signaling networks using automated assembly. *Molecular systems biology*, 13(11):954, 2017.
- Xu Han, Shulin Cao, Xin Lv, Yankai Lin, Zhiyuan Liu, Maosong Sun, and Juanzi Li. Openke: An open toolkit for knowledge embedding. In *Proceedings of the 2018 conference on empirical methods in natural language processing: system demonstrations*, pp. 139–144, 2018.
- Charles R Harris, K Jarrod Millman, Stéfan J Van Der Walt, Ralf Gommers, Pauli Virtanen, David Cournapeau, Eric Wieser, Julian Taylor, Sebastian Berg, Nathaniel J Smith, et al. Array programming with numpy. *Nature*, 585(7825):357–362, 2020.
- Lei He, Suncong Zheng, Tao Yang, and Feng Zhang. Klmo: Knowledge graph enhanced pretrained language model with fine-grained relationships. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pp. 4536–4542, 2021.
- Pengcheng He, Xiaodong Liu, Jianfeng Gao, and Weizhu Chen. Deberta: Decoding-enhanced bert with disentangled attention. In *International Conference on Learning Representations*, 2020.
- Yu-Jung Heo, Eun-Sol Kim, Woo Suk Choi, and Byoung-Tak Zhang. Hypergraph transformer: Weakly-supervised multi-hop reasoning for knowledge-based visual question answering. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 373–390, 2022.

- Linmei Hu, Tianchi Yang, Luhao Zhang, Wanjun Zhong, Duyu Tang, Chuan Shi, Nan Duan, and Ming Zhou. Compare to the knowledge: Graph neural fake news detection with external knowledge. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pp. 754–763, 2021.
- Shengding Hu, Ning Ding, Huadong Wang, Zhiyuan Liu, Jingang Wang, Juanzi Li, Wei Wu, and Maosong Sun. Knowledgeable prompt-tuning: Incorporating knowledge into prompt verbalizer for text classification. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 2225–2240, Dublin, Ireland, May 2022.
- Filip Ilievski, Pedro Szekely, and Bin Zhang. Cskg: The commonsense knowledge graph. In *European Semantic Web Conference*, pp. 680–696. Springer, 2021.
- Mandar Joshi, Kenton Lee, Yi Luan, and Kristina Toutanova. Contextualized representations using textual encyclopedic knowledge. *arXiv preprint arXiv:2004.12006*, 2020.
- Yong-Ho Jung, Jun-Hyung Park, Joon-Young Choi, Mingyu Lee, Junho Kim, Kang-Min Kim, and SangKeun Lee. Learning from missing relations: Contrastive learning with commonsense knowledge graphs for commonsense inference. In *Findings of the Association for Computational Linguistics: ACL 2022*, pp. 1514–1523, 2022.
- Minki Kang, Jinheon Baek, and Sung Ju Hwang. KALA: knowledge-augmented language model adaptation. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pp. 5144–5167, Seattle, United States, July 2022.
- Jivat Kaur, Sumit Bhatia, Milan Aggarwal, Rachit Bansal, and Balaji Krishnamurthy. LM-CORE: Language models with contextually relevant external knowledge. In *Findings of the Association for Computational Linguistics: NAACL 2022*, pp. 750–769, Seattle, United States, July 2022.
- Daphna Keidar, Mian Zhong, Ce Zhang, Yash Raj Shrestha, and Bibek Paudel. Towards automatic bias detection in knowledge graphs. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pp. 3804–3811, Punta Cana, Dominican Republic, November 2021.
- Johannes Kiesel, Maria Mestre, Rishabh Shukla, Emmanuel Vincent, Payam Adineh, David Corney, Benno Stein, and Martin Potthast. Semeval-2019 task 4: Hyperpartisan news detection. In *Proceedings of the 13th International Workshop on Semantic Evaluation*, pp. 829–839, 2019.
- Yu Jin Kim, Beong-woo Kwak, Youngwook Kim, Reinald Kim Amplayo, Seung-won Hwang, and Jinyoung Yeo. Modularized transfer learning with multiple knowledge graphs for zero-shot commonsense reasoning. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pp. 2244–2257, Seattle, United States, July 2022.
- Peter Kraft, Hirsh Jain, and Alexander M Rush. An embedding model for predicting roll-call votes. In *Proceedings of the 2016 conference on empirical methods in natural language processing*, pp. 2066–2070, 2016.
- Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pp. 7871–7880, 2020.
- Chang Li and Dan Goldwasser. Encoding social information with graph convolutional networks for political perspective detection in news media. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 2594–2604, 2019.
- Chang Li and Dan Goldwasser. Using social and linguistic information to adapt pretrained representations for political perspective identification. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pp. 4569–4579, 2021.

- Dawei Li, Yanran Li, Jiayi Zhang, Ke Li, Chen Wei, Jianwei Cui, and Bin Wang. C<sup>3</sup>KG: A Chinese commonsense conversation knowledge graph. In *Findings of the Association for Computational Linguistics: ACL 2022*, pp. 1369–1383, Dublin, Ireland, May 2022a.
- Zongren Li, Qin Zhong, Jing Yang, Yongjie Duan, Wenjun Wang, Chengkun Wu, and Kunlun He. Deepkg: an end-to-end deep learning-based workflow for biomedical knowledge graph extraction, optimization and applications. *Bioinformatics*, 38(5):1477–1479, 2022b.
- Paul Pu Liang, Chiyu Wu, Louis-Philippe Morency, and Ruslan Salakhutdinov. Towards understanding and mitigating social biases in language models. In *International Conference on Machine Learning*, pp. 6565–6576. PMLR, 2021.
- Yankai Lin, Zhiyuan Liu, Maosong Sun, Yang Liu, and Xuan Zhu. Learning entity and relation embeddings for knowledge graph completion. In *Twenty-ninth AAAI conference on artificial intelligence*, 2015.
- Jiacheng Liu, Alisa Liu, Ximing Lu, Sean Welleck, Peter West, Ronan Le Bras, Yejin Choi, and Hannaneh Hajishirzi. Generated knowledge prompting for commonsense reasoning. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 3154–3169, 2022.
- Liyuan Liu, Haoming Jiang, Pengcheng He, Weizhu Chen, Xiaodong Liu, Jianfeng Gao, and Jiawei Han. On the variance of the adaptive learning rate and beyond. *arXiv preprint arXiv:1908.03265*, 2019a.
- Ye Liu, Yao Wan, Lifang He, Hao Peng, and S Yu Philip. Kg-bart: Knowledge graph-augmented bart for generative commonsense reasoning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pp. 6418–6425, 2021.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*, 2019b.
- Yinquan Lu, Haonan Lu, Guirong Fu, and Qun Liu. Kelm: Knowledge enhanced pre-trained language representations with message passing on hierarchical relational graphs. In *ICLR 2022 Workshop on Deep Learning on Graphs for Natural Language Processing*, 2022.
- Kaixin Ma, Hao Cheng, Xiaodong Liu, Eric Nyberg, and Jianfeng Gao. Open domain question answering with a unified knowledge interface. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 1605–1620, Dublin, Ireland, May 2022.
- Emily McMilin. Selection bias induced spurious correlations in large language models. *arXiv preprint arXiv:2207.08982*, 2022.
- Ninareh Mehrabi, Pei Zhou, Fred Morstatter, Jay Pujara, Xiang Ren, and Aram Galstyan. Lawyers are dishonest? quantifying representational harms in commonsense knowledge resources. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pp. 5016–5033, 2021.
- Zaiqiao Meng, Fangyu Liu, Thomas Clark, Ehsan Shareghi, and Nigel Collier. Mixture-of-partitions: Infusing large biomedical knowledge graphs into bert. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pp. 4672–4681, 2021.
- Sayantana Mitra, Roshni Ramnani, and Shubhashis Sengupta. Constraint-based multi-hop question answering with knowledge graph. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies: Industry Track*, pp. 280–288, 2022.
- Fedor Moiseev, Zhe Dong, Enrique Alfonseca, and Martin Jaggi. SKILL: Structured knowledge infusion for large language models. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pp. 1581–1588, Seattle, United States, July 2022.

- Xinyi Mou, Zhongyu Wei, Lei Chen, Shangyi Ning, Yancheng He, Changjian Jiang, and Xuan-Jing Huang. Align voting behavior with public statements for legislator representation learning. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pp. 1236–1246, 2021.
- Moin Nadeem, Anna Bethke, and Siva Reddy. Stereoset: Measuring stereotypical bias in pretrained language models. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pp. 5356–5371, 2021.
- Barlas Oguz, Xilun Chen, Vladimir Karpukhin, Stan Peshterliev, Dmytro Okhonko, Michael Schlichtkrull, Sonal Gupta, Yashar Mehdad, and Scott Yih. UniK-QA: Unified representations of structured and unstructured knowledge for open-domain question answering. In *Findings of the Association for Computational Linguistics: NAACL 2022*, pp. 1535–1546, Seattle, United States, July 2022.
- Alissa Ostapenko, Shuly Wintner, Melinda Fricke, and Yulia Tsvetkov. Speaker information can guide models to better inductive biases: A case study on predicting code-switching. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 3853–3867, Dublin, Ireland, May 2022.
- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems*, 32, 2019.
- F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.
- Thomas Pellissier Tanon, Gerhard Weikum, and Fabian Suchanek. Yago 4: A reasonable knowledge base. In *European Semantic Web Conference*, pp. 583–596. Springer, 2020.
- Matthew E Peters, Mark Neumann, Robert L Logan IV, Roy Schwartz, Vidur Joshi, Sameer Singh, and Noah A Smith. Knowledge enhanced contextual word representations. *arXiv preprint arXiv:1909.04164*, 2019.
- Jay Pujara, Eriq Augustine, and Lise Getoor. Sparsity and noise: Where knowledge graph embeddings fall short. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pp. 1751–1756, Copenhagen, Denmark, September 2017.
- Delai Qiu, Yuanzhe Zhang, Xinwei Feng, Xiangwen Liao, Wenbin Jiang, Yajuan Lyu, Kang Liu, and Jun Zhao. Machine reading comprehension using structural knowledge graph-aware network. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pp. 5896–5901, 2019a.
- Delai Qiu, Yuanzhe Zhang, Xinwei Feng, Xiangwen Liao, Wenbin Jiang, Yajuan Lyu, Kang Liu, and Jun Zhao. Machine reading comprehension using structural knowledge graph-aware network. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pp. 5896–5901, 2019b.
- Hannah Rashkin, Eunsol Choi, Jin Yea Jang, Svitlana Volkova, and Yejin Choi. Truth of varying shades: Analyzing language in fake news and political fact-checking. In *Proceedings of the 2017 conference on empirical methods in natural language processing*, pp. 2931–2937, 2017.
- Md Rashad Al Hasan Rony, Ricardo Usbeck, and Jens Lehmann. DialoKG: Knowledge-structure aware task-oriented dialogue generation. In *Findings of the Association for Computational Linguistics: NAACL 2022*, pp. 2557–2571, Seattle, United States, July 2022.

- Victoria L Rubin, Niall Conroy, Yimin Chen, and Sarah Cornwell. Fake news or truth? using satirical cues to detect potentially misleading news. In *Proceedings of the second workshop on computational approaches to deception detection*, pp. 7–17, 2016.
- Robyn Speer, Joshua Chin, and Catherine Havasi. Conceptnet 5.5: An open multilingual graph of general knowledge. In *Thirty-first AAAI conference on artificial intelligence*, 2017.
- Rohit Sridhar and Diyi Yang. Explaining toxic text via knowledge enhanced text generation. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pp. 811–826, Seattle, United States, July 2022.
- Tianxiang Sun, Yunfan Shao, Xipeng Qiu, Qipeng Guo, Yaru Hu, Xuan-Jing Huang, and Zheng Zhang. Colake: Contextualized language and knowledge embedding. In *Proceedings of the 28th International Conference on Computational Linguistics*, pp. 3660–3670, 2020.
- Yueqing Sun, Qi Shi, Le Qi, and Yu Zhang. JointLK: Joint reasoning with language models and knowledge graphs for commonsense question answering. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pp. 5049–5060, Seattle, United States, July 2022.
- 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*, 2018.
- Zhaoxuan Tan, Zilong Chen, Shangbin Feng, Qingyue Zhang, Qinghua Zheng, Jundong Li, and Minnan Luo. Kracl: Contrastive learning with graph context modeling for sparse knowledge graph completion. *arXiv preprint arXiv:2208.07622*, 2022.
- Shikhar Vashishth, Soumya Sanyal, Vikram Nitin, and Partha Talukdar. Composition-based multi-relational graph convolutional networks. In *International Conference on Learning Representations*, 2019.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua Bengio. Graph attention networks. In *International Conference on Learning Representations*, 2018.
- Denny Vrandečić and Markus Krötzsch. Wikidata: a free collaborative knowledgebase. *Communications of the ACM*, 57(10):78–85, 2014.
- Ruize Wang, Duyu Tang, Nan Duan, Zhongyu Wei, Xuan-Jing Huang, Jianshu Ji, Guihong Cao, Daxin Jiang, and Ming Zhou. K-adapter: Infusing knowledge into pre-trained models with adapters. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pp. 1405–1418, 2021a.
- Wenlin Wang, Zhe Gan, Wenqi Wang, Dinghan Shen, Jiaji Huang, Wei Ping, Sanjeev Satheesh, and Lawrence Carin. Topic compositional neural language model. In *International Conference on Artificial Intelligence and Statistics*, pp. 356–365. PMLR, 2018.
- Xiaozhi Wang, Tianyu Gao, Zhaocheng Zhu, Zhengyan Zhang, Zhiyuan Liu, Juanzi Li, and Jian Tang. Kepler: A unified model for knowledge embedding and pre-trained language representation. *Transactions of the Association for Computational Linguistics*, 9:176–194, 2021b.
- Max Welling and Thomas N Kipf. Semi-supervised classification with graph convolutional networks. In *J. International Conference on Learning Representations (ICLR 2017)*, 2016.
- Peter West, Chandra Bhagavatula, Jack Hessel, Jena Hwang, Liwei Jiang, Ronan Le Bras, Ximing Lu, Sean Welleck, and Yejin Choi. Symbolic knowledge distillation: from general language models to commonsense models. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pp. 4602–4625, Seattle, United States, July 2022.

- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, et al. Transformers: State-of-the-art natural language processing. In *Proceedings of the 2020 conference on empirical methods in natural language processing: system demonstrations*, pp. 38–45, 2020.
- Tianbao Xie, Chen Henry Wu, Peng Shi, Ruiqi Zhong, Torsten Scholak, Michihiro Yasunaga, Chien-Sheng Wu, Ming Zhong, Pengcheng Yin, Sida I Wang, et al. Unifiedskg: Unifying and multi-tasking structured knowledge grounding with text-to-text language models. *arXiv preprint arXiv:2201.05966*, 2022.
- Ikuya Yamada, Akari Asai, Hiroyuki Shindo, Hideaki Takeda, and Yuji Matsumoto. Luke: Deep contextualized entity representations with entity-aware self-attention. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 6442–6454, 2020.
- Zichao Yang, Diyi Yang, Chris Dyer, Xiaodong He, Alex Smola, and Eduard Hovy. Hierarchical attention networks for document classification. In *Proceedings of the 2016 conference of the North American chapter of the association for computational linguistics: human language technologies*, pp. 1480–1489, 2016.
- Michihiro Yasunaga, Hongyu Ren, Antoine Bosselut, Percy Liang, and Jure Leskovec. Qa-gnn: Reasoning with language models and knowledge graphs for question answering. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pp. 535–546, 2021.
- Donghan Yu, Chenguang Zhu, Yuwei Fang, Wenhao Yu, Shuohang Wang, Yichong Xu, Xiang Ren, Yiming Yang, and Michael Zeng. KG-FiD: Infusing knowledge graph in fusion-in-decoder for open-domain question answering. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 4961–4974, Dublin, Ireland, May 2022a.
- Donghan Yu, Chenguang Zhu, Yiming Yang, and Michael Zeng. Jacket: Joint pre-training of knowledge graph and language understanding. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pp. 11630–11638, 2022b.
- Wenqian Zhang, Shangbin Feng, Zilong Chen, Zhenyu Lei, Jundong Li, and Minnan Luo. KCD: Knowledge walks and textual cues enhanced political perspective detection in news media. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pp. 4129–4140, Seattle, United States, July 2022.
- Xikun Zhang, Antoine Bosselut, Michihiro Yasunaga, Hongyu Ren, Percy Liang, Christopher D Manning, and Jure Leskovec. Greaselm: Graph reasoning enhanced language models. In *International Conference on Learning Representations*, 2021.
- Zhengyan Zhang, Xu Han, Zhiyuan Liu, Xin Jiang, Maosong Sun, and Qun Liu. Ernie: Enhanced language representation with informative entities. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 1441–1451, 2019.

## A LIMITATIONS

Our proposed KALM has two minor limitations:

- KALM relies on existing knowledge graphs to facilitate knowledge-aware long document understanding. While knowledge graphs are effective and prevalent tools for modeling real-world symbolic knowledge, they are often sparse and hardly exhaustive (Tan et al., 2022; Pujara et al., 2017). In addition, external knowledge is not only limited to knowledge graphs but also exists in textual, visual, and other symbolic forms. We leave it to future work on how to jointly leverage multiple forms and sources of external knowledge in document understanding.



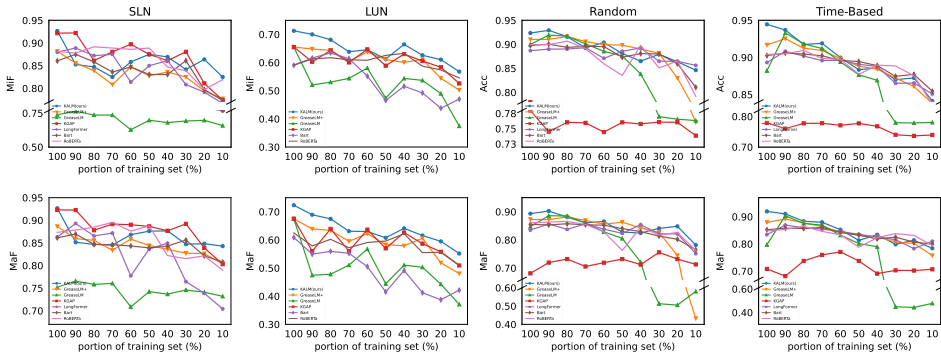


Figure 4: KALM and competitive baselines’ performance when training data decreases from 100% to 10%. KALM maintains steady performance with as little as 10% to 20% of training data, while baseline methods witness serious performance deterioration.

- KALM leverages TagMe (Ferragina & Scaiella, 2011) to identify entity mentions and build the three knowledge-aware contexts. While TagMe and other entity identification tools are effective, they are not 100% correct, resulting in potentially omitted entities and external knowledge. In addition, running TagMe on hundreds of thousands of long documents is time-consuming and resource-consuming even if processed in parallel. We leave it to future work on how to leverage knowledge graphs for long document understanding without explicitly using entity linking tools.

## B SOCIETAL IMPACT

KALM is a knowledge-aware long document understanding approach that jointly leverages pre-trained LMs and knowledge graphs on three levels of contexts. Consequently, KALM might exhibit many of the biases of the adopted language models (Liang et al., 2021; Nadeem et al., 2021) and knowledge graphs (Fisher et al., 2020; 2019; Mehrabi et al., 2021; Du et al., 2022; Keidar et al., 2021). As a result, KALM might leverage the biased and unethical correlations in LMs and KGs to arrive at conclusions. We encourage KALM users to audit its output before using it beyond the standard benchmarks. We leave it to future work on how to leverage knowledge graphs in pre-trained LMs with a focus on fairness and equity.

## C ADDITIONAL EXPERIMENTS

### C.1 DATA EFFICIENCY

Existing works argue that introducing knowledge graphs to NLP tasks could improve data efficiency and help alleviate the need for extensive training data (Zhang et al., 2022). By introducing knowledge to all three document contexts and enabling knowledge-rich context information exchange, KALM might be in a better position to tackle this issue. To examine whether KALM has indeed improved data efficiency, we compare the performance of KALM with competitive baselines when trained on partial training sets and illustrate the results in Figure 4. It is demonstrated that while performance did not change greatly with 30% to 100% training data, baseline methods witness significant performance drops when only 10% to 20% of data are available. In contrast, KALM maintains steady performance with as little as 10% of training data.

### C.2 CONTEXT EXCHANGE STUDY (CONT.)

In Section 3.3, we conducted an ablation study of the three knowledge-aware contexts and explored how the ContextFusion layer enables the interpretation of context contribution and information exchange patterns. It is demonstrated that the three contexts play different roles with respect to datasets and KALM layers. In addition, we explore whether the role and information exchange patterns of contexts change when the training progresses. Figure 5 illustrates the results with respect to training

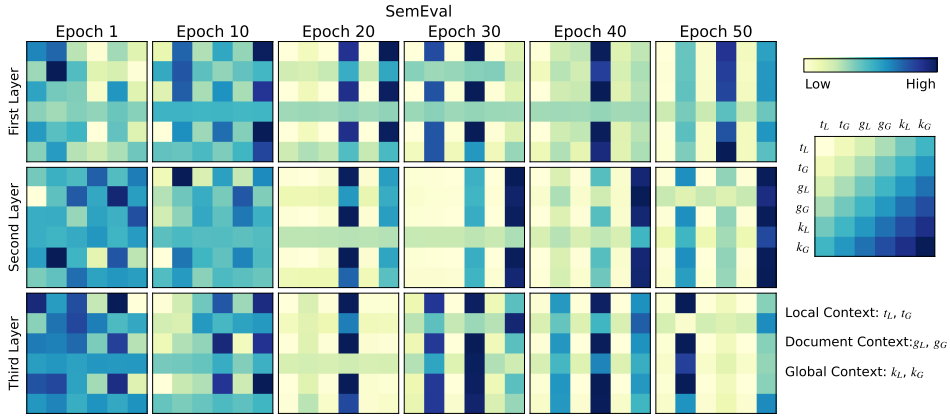


Figure 5: Interpreting the roles of the three contexts with respect to training progress on the SemEval dataset.  $t_L, t_G, g_L, g_G, k_L, k_G$  denote the context representations in equations (9) and (10), so that the first two columns indicate how the local context attends to information in other contexts, the next two columns for the document context, and the last two columns for the global context.

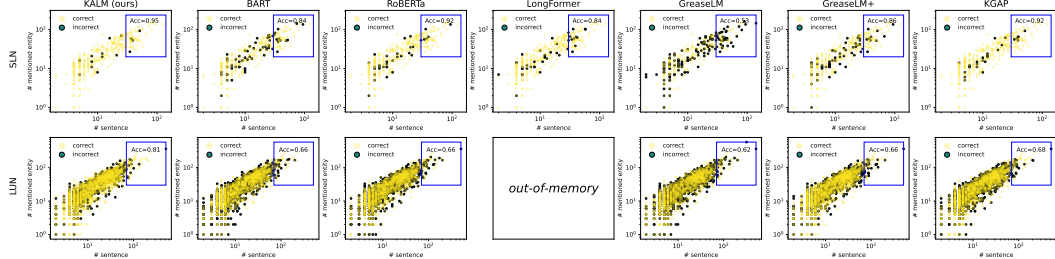


Figure 6: Error analysis of KALM and baselines. KALM successfully improves in the top-right corner, which represents documents with more sentences and more entailed knowledge.

epochs, which shows that the attention matrices started out dense and ended sparse, indicating that the role of different contexts is gradually developed through time.

### C.3 LONG DOCUMENT STUDY (CONT.)

We present error analysis with respect to document length and knowledge intensity on more baseline methods, including language models (RoBERTa, BART, LongFormer), knowledge-aware LMs (KGAP, GreaseLM, GreaseLM+), and our proposed KALM in Figure 6. Our conclusion still holds true: KALM successfully improves performance on documents that are longer and contain more external knowledge, which are positioned in the top-right corner of the figure.

### C.4 MANUAL ERROR ANALYSIS

We manually examined 20 news articles from the LUN misinformation detection dataset where KALM made a mistake. Several news articles focused on the same topic of marijuana legalization, and some others focused on international affairs such as the conflict in Iraq. These articles feature entities and knowledge that are much more recent such as "pot-infused products" and "ISIS jihadists", which are emerging concepts and generally not covered by existing knowledge graphs. We present the relevant sentences in Table 3. This indicates the need for more comprehensive, up-to-date, and temporal knowledge graphs that grow with the world.

Sample ID	Example Sentences
1853	... the legalization of recreational marijuana ... has created new markets for <b>pot-infused products</b> ... ... children who were taken to emergency departments due to accidental <b>THC</b> ingestion ...
1169	Mr. Kerry met with Iraqi foreign minister <b>Hoshiyar Zebari</b> about providing help in fighting the <b>ISIS jihadists</b> ... ... territory north and north-east of Baghdad where the predominantly <b>Sunni militants</b> have penetrated within ...

Table 3: Example sentences in the articles where KALM made a mistake. Emerging entities that are not covered by existing knowledge graphs are in **bold**.

Table 4: Model performance on the task of political perspective detection.

	Baseline	SemEval		Allsides	
		Acc	MaF	Acc	MaF
<b>task specific</b>	HLSTM	81.71	/	76.45	74.95
	MAN	86.21	84.33	85.00	84.25
	KCD	89.90 ( $\pm 0.6$ )	86.11 ( $\pm 1.1$ )	87.17 ( $\pm 0.2$ )	86.72 ( $\pm 0.3$ )
<b>language model</b>	RoBERTa	85.56 ( $\pm 1.6$ )	77.94 ( $\pm 3.5$ )	68.71 ( $\pm 4.3$ )	65.39 ( $\pm 5.7$ )
	Electra	78.87 ( $\pm 2.8$ )	62.85 ( $\pm 7.9$ )	63.14 ( $\pm 2.3$ )	58.24 ( $\pm 3.8$ )
	DeBERTa	86.99 ( $\pm 1.9$ )	80.62 ( $\pm 3.8$ )	67.86 ( $\pm 4.3$ )	63.50 ( $\pm 5.9$ )
	BART	86.62 ( $\pm 1.5$ )	79.87 ( $\pm 2.6$ )	60.56 ( $\pm 3.8$ )	54.64 ( $\pm 5.4$ )
	LongFormer	82.81 ( $\pm 2.3$ )	73.09 ( $\pm 4.5$ )	62.88 ( $\pm 3.0$ )	58.03 ( $\pm 4.6$ )
<b>task agnostic</b>	KELM	86.40 ( $\pm 2.3$ )	83.98 ( $\pm 1.0$ )	80.71 ( $\pm 2.4$ )	79.74 ( $\pm 2.7$ )
	KnowBERT-Wordnet	81.71 ( $\pm 5.5$ )	72.28 ( $\pm 6.7$ )	60.54 ( $\pm 0.4$ )	58.77 ( $\pm 0.6$ )
	KnowBERT-Wikidata	76.72 ( $\pm 3.0$ )	66.21 ( $\pm 5.0$ )	60.56 ( $\pm 0.7$ )	58.81 ( $\pm 0.5$ )
	KnowBERT-W+W	84.73 ( $\pm 3.4$ )	75.72 ( $\pm 5.3$ )	60.44 ( $\pm 0.3$ )	58.46 ( $\pm 0.5$ )
	Joshi et al.	81.88 ( $\pm 2.1$ )	77.15 ( $\pm 3.8$ )	80.88 ( $\pm 2.1$ )	79.73 ( $\pm 2.3$ )
	KGAP	87.73 ( $\pm 1.8$ )	82.00 ( $\pm 3.1$ )	83.65 ( $\pm 1.3$ )	82.92 ( $\pm 1.4$ )
	GreaseLM	86.64 ( $\pm 1.5$ )	80.32 ( $\pm 3.0$ )	80.23 ( $\pm 1.2$ )	79.17 ( $\pm 1.2$ )
	GreaseLM+	85.66 ( $\pm 1.8$ )	77.23 ( $\pm 4.1$ )	82.16 ( $\pm 5.5$ )	80.81 ( $\pm 7.1$ )
	KALM (Ours)	<b>91.45</b> ( $\pm 0.8$ )	<b>87.65</b> ( $\pm 1.2$ )	<b>87.26</b> ( $\pm 0.2$ )	<b>86.79</b> ( $\pm 0.2$ )

### C.5 SIGNIFICANCE TESTING

To examine whether KALM significantly outperforms baselines on the three tasks, we conduct one-way repeated measures ANOVA test for the results in Table 4, Table 5, and Table 6. It is demonstrated that the performance gain is significant on five of the six datasets, specifically SemEval (against the second-best KCD on Acc and MaF), SLN (against the second-best KGAP on MiF and MaRecall), LUN (against the second-best CompareNet on MiF, MaF and MaRecall), Random (against the second-best GreasesLM+ on BAcc and MaF), and Time-Based (against the second-best GreaseLM+ on BAcc and MaF).

### C.6 TASK-SPECIFIC MODEL PERFORMANCE

We present the full results for task-specific methods, pre-trained language models, knowledge-aware task-agnostic models, and KALM on the three tasks and six datasets/settings in Tables 4, 5, and 6.

### C.7 IS LOCAL CONTEXT ENOUGH?

Though long document understanding requires attending to a long sequence of tokens, it is possible that sometimes only one or two sentences would give away the label of the document. We examine this by removing the document-level and global contexts in KALM, leaving only the local context to simulate this scenario. Comparing the local-only variant with the full KALM, there are 14.78%, 10.53%, 8.21%, 4.85%, 1.4%, and 3.18% performance drops across the six datasets in terms of macro-averaged F1-score. As a result, it is necessary to go beyond local context windows in long document understanding.

Table 5: Model performance on the task of misinformation detection.

Baseline		SLN				LUN			
		MiF	MaPrecision	MaRecall	MaF	MiF	MaPrecision	MaRecall	MaF
task specific	Rubin et al.	/	88.00	82.00	/	/	/	/	/
	Rashkin et al.	/	/	/	/	/	/	/	65.00
	GCN + Attn	85.27	85.59	85.27	85.24	67.08	68.60	67.00	66.42
	GAT + Attn	84.72	85.65	84.72	84.62	66.95	68.05	66.86	66.37
	CompareNet	89.17	89.82	89.17	89.12	69.05	72.94	69.04	68.26
language model	RoBERTa	88.17 ( $\pm 0.6$ )	89.02 ( $\pm 1.8$ )	88.17 ( $\pm 0.6$ )	87.34 ( $\pm 1.2$ )	59.09 ( $\pm 1.7$ )	62.49 ( $\pm 2.6$ )	59.11 ( $\pm 1.6$ )	55.52 ( $\pm 1.5$ )
	Electra	75.44 ( $\pm 2.2$ )	83.22 ( $\pm 0.6$ )	75.44 ( $\pm 2.2$ )	67.53 ( $\pm 4.1$ )	60.10 ( $\pm 1.7$ )	63.26 ( $\pm 1.2$ )	60.11 ( $\pm 1.7$ )	58.57 ( $\pm 2.1$ )
	DeBERTa	86.89 ( $\pm 6.6$ )	89.43 ( $\pm 3.7$ )	86.89 ( $\pm 6.6$ )	88.46 ( $\pm 4.9$ )	57.62 ( $\pm 3.1$ )	64.03 ( $\pm 0.9$ )	57.63 ( $\pm 3.1$ )	52.24 ( $\pm 5.3$ )
	BART	86.06 ( $\pm 0.6$ )	86.13 ( $\pm 0.5$ )	86.06 ( $\pm 0.6$ )	86.12 ( $\pm 0.6$ )	59.05 ( $\pm 2.2$ )	60.89 ( $\pm 4.5$ )	59.07 ( $\pm 2.2$ )	54.18 ( $\pm 2.8$ )
	LongFormer	88.00 ( $\pm 2.5$ )	88.84 ( $\pm 1.5$ )	87.44 ( $\pm 2.5$ )	86.29 ( $\pm 3.4$ )	out-of-memory			
task agnostic	KELM	84.11 ( $\pm 0.6$ )	85.23 ( $\pm 0.7$ )	84.11 ( $\pm 0.6$ )	82.80 ( $\pm 1.3$ )	59.28 ( $\pm 2.1$ )	61.09 ( $\pm 2.8$ )	59.29 ( $\pm 2.1$ )	57.30 ( $\pm 1.6$ )
	KnowBERT-Wordnet	74.72 ( $\pm 3.3$ )	77.22 ( $\pm 1.8$ )	74.72 ( $\pm 3.3$ )	72.74 ( $\pm 8.5$ )	55.63 ( $\pm 1.8$ )	56.29 ( $\pm 2.0$ )	55.63 ( $\pm 1.8$ )	55.02 ( $\pm 1.7$ )
	KnowBERT-Wikidata	72.17 ( $\pm 2.5$ )	73.57 ( $\pm 0.6$ )	72.17 ( $\pm 2.5$ )	69.41 ( $\pm 6.9$ )	57.57 ( $\pm 0.5$ )	57.27 ( $\pm 0.6$ )	57.57 ( $\pm 0.5$ )	56.76 ( $\pm 0.6$ )
	KnowBERT-W+W	78.67 ( $\pm 3.2$ )	79.36 ( $\pm 3.1$ )	78.67 ( $\pm 3.2$ )	79.80 ( $\pm 0.9$ )	65.52 ( $\pm 2.3$ )	67.50 ( $\pm 1.6$ )	65.53 ( $\pm 2.3$ )	63.94 ( $\pm 2.0$ )
	Joshi et al.	92.72 ( $\pm 5.1$ )	84.95 ( $\pm 2.8$ )	83.37 ( $\pm 5.2$ )	83.98 ( $\pm 3.7$ )	58.57 ( $\pm 3.4$ )	62.56 ( $\pm 4.0$ )	58.59 ( $\pm 3.4$ )	56.73 ( $\pm 4.0$ )
	KGAP	92.17 ( $\pm 1.2$ )	92.67 ( $\pm 0.9$ )	92.17 ( $\pm 1.2$ )	92.30 ( $\pm 0.9$ )	65.52 ( $\pm 2.3$ )	67.50 ( $\pm 1.6$ )	65.53 ( $\pm 2.3$ )	63.94 ( $\pm 2.9$ )
	GreaseLM	73.83 ( $\pm 0.9$ )	74.33 ( $\pm 0.8$ )	73.83 ( $\pm 0.9$ )	75.20 ( $\pm 0.8$ )	56.54 ( $\pm 1.5$ )	58.12 ( $\pm 2.7$ )	56.55 ( $\pm 1.5$ )	55.75 ( $\pm 1.6$ )
	GreaseLM+	88.17 ( $\pm 0.8$ )	88.56 ( $\pm 0.6$ )	88.17 ( $\pm 0.8$ )	88.64 ( $\pm 0.6$ )	64.29 ( $\pm 2.4$ )	65.13 ( $\pm 2.7$ )	64.31 ( $\pm 2.4$ )	62.65 ( $\pm 3.7$ )
	KALM (Ours)	<b>94.22</b> ( $\pm 1.2$ )	<b>94.33</b> ( $\pm 1.1$ )	<b>94.22</b> ( $\pm 1.1$ )	<b>94.18</b> ( $\pm 1.1$ )	<b>71.28</b> ( $\pm 1.7$ )	<b>72.33</b> ( $\pm 2.7$ )	<b>71.29</b> ( $\pm 1.7$ )	<b>69.82</b> ( $\pm 1.2$ )

Table 6: Model performance on the task of roll call vote prediction.

Baseline		Random		Time-Based	
		BAcc	MaF	BAcc	MaF
task specific	ideal-point	86.46	80.02	/	/
	ideal-vector	87.35	80.15	81.95	75.49
	Vote	90.22	84.92	89.76	84.35
	PAR	90.33	/	89.92	/
language model	RoBERTa	89.94 ( $\pm 0.2$ )	86.10 ( $\pm 0.7$ )	90.40 ( $\pm 0.8$ )	84.78 ( $\pm 2.2$ )
	Electra	87.47 ( $\pm 0.3$ )	80.23 ( $\pm 0.7$ )	88.92 ( $\pm 0.4$ )	82.50 ( $\pm 1.7$ )
	DeBERTa	86.98 ( $\pm 0.4$ )	80.07 ( $\pm 1.2$ )	88.59 ( $\pm 0.1$ )	81.38 ( $\pm 1.0$ )
	BART	89.76 ( $\pm 0.5$ )	85.52 ( $\pm 0.6$ )	90.25 ( $\pm 0.6$ )	85.21 ( $\pm 2.1$ )
	LongFormer	88.69 ( $\pm 0.4$ )	83.52 ( $\pm 1.2$ )	89.32 ( $\pm 1.4$ )	83.42 ( $\pm 3.8$ )
task agnostic	KELM	89.13 ( $\pm 1.1$ )	84.76 ( $\pm 2.0$ )	90.80 ( $\pm 0.2$ )	86.62 ( $\pm 0.4$ )
	KnowBERT-Wordnet	86.72 ( $\pm 0.9$ )	79.33 ( $\pm 2.4$ )	86.92 ( $\pm 0.6$ )	78.90 ( $\pm 1.9$ )
	KnowBERT-Wikidata	85.98 ( $\pm 0.8$ )	78.48 ( $\pm 1.0$ )	86.45 ( $\pm 0.5$ )	78.21 ( $\pm 0.7$ )
	KnowBERT-W+W	85.75 ( $\pm 1.0$ )	78.70 ( $\pm 2.4$ )	87.07 ( $\pm 1.0$ )	78.42 ( $\pm 2.2$ )
	Joshi et al.	91.43 ( $\pm 0.5$ )	89.64 ( $\pm 0.6$ )	92.63 ( $\pm 1.6$ )	89.31 ( $\pm 2.4$ )
	KGAP	77.98 ( $\pm 0.5$ )	68.11 ( $\pm 6.0$ )	77.90 ( $\pm 0.6$ )	70.81 ( $\pm 4.6$ )
	GreaseLM	89.99 ( $\pm 1.5$ )	84.72 ( $\pm 3.0$ )	88.21 ( $\pm 2.7$ )	79.73 ( $\pm 7.4$ )
	GreaseLM+	91.01 ( $\pm 0.2$ )	87.29 ( $\pm 0.3$ )	91.69 ( $\pm 0.1$ )	87.95 ( $\pm 0.3$ )
KALM (Ours)	<b>92.36</b> ( $\pm 0.3$ )	<b>89.33</b> ( $\pm 0.4$ )	<b>94.46</b> ( $\pm 0.4$ )	<b>91.97</b> ( $\pm 0.5$ )	

## D EXPERIMENT DETAILS

### D.1 DATASET DETAILS

We present important dataset details in Table 7. We follow the exact same dataset settings and splits in previous works (Zhang et al., 2022; Hu et al., 2021; Feng et al., 2021b) for fair comparison.

### D.2 BASELINE DETAILS

We compare KALM with pre-trained language models, task-specific baselines, and task-agnostic knowledge-aware methods to ensure a holistic evaluation. In the following, we provide a brief description of each of the baseline methods. We also highlight whether one approach leverages knowledge graphs and the three document contexts in Table 8.

- **HLSTM** (Yang et al., 2016) is short for hierarchical long short-term memory networks. It was used in previous works (Li & Goldwasser, 2019; 2021) for political perspective detection.
- **MAN** (Li & Goldwasser, 2021) proposes to leverage social and linguistic information to design pre-training tasks and fine-tune on the task of political perspective detection.
- **KCD** (Zhang et al., 2022) proposes to leverage multi-hop knowledge reasoning with knowledge walks and textual cues with document graphs for political perspective detection.
- Rubin et al. (2016) proposes the SLN dataset and leverages satirical cues for misinformation detection.

Task	Dataset	# Document	# Class	Class Distribution	Document Length	Originally Proposed In
PPD	SemEval	645	2	407 / 238	793.00 $\pm$ 736.93	Kiesel et al. (2019)
	Allsides	10,385	3	4,164 / 3,931 / 2,290	1316.81 $\pm$ 2978.71	Li & Goldwasser (2019)
MD	SLN	360	2	180 / 180	551.32 $\pm$ 661.82	Rubin et al. (2016)
	LUN	51,854	4	10,745 / 14,797 / 7,692 / 18,620		Rashkin et al. (2017)
RCVP	random time-based	1,155	2	304,655 / 95,464	653.94 $\pm$ 424.32	Mou et al. (2021)

Table 7: Dataset statistics. The number of long documents and class distribution does not add up for RCVP since multiple legislators vote on the same legislation.

- Rashkin et al. (2017) proposes the LUN dataset and argues that misinformation detection should have more fine-grained labels than true or false.
- GCN (Welling & Kipf, 2016) and GAT (Veličković et al., 2018) are leveraged along with the attention mechanism by Hu et al. (2021) for misinformation detection on graphs.
- CompareNet (Hu et al., 2021) proposes to leverage knowledge graphs and compare the textual content to external knowledge for misinformation detection.
- Ideal-point (Gerrish & Blei, 2011) and ideal-vector (Kraft et al., 2016) propose to use 1d and 2d representations of political actors for roll call vote prediction.
- Vote (Mou et al., 2021) proposes to jointly leverage legislation text and the social network information for roll call vote prediction.
- PAR (Feng et al., 2021b) proposes to learn legislator representations with social context and expert knowledge for roll call vote prediction.
- RoBERTa (Liu et al., 2019b), Electra (Clark et al., 2019), DeBERTa (He et al., 2020), BART (Lewis et al., 2020), and LongFormer (Beltagy et al., 2020) are pre-trained language models. We use the pre-trained weights *roberta-base*, *electra-small-discriminator*, *deberta-v3-base*, *bart-base*, and *longformer-base-4096* in Huggingface Transformers (Wolf et al., 2020) to extract sentence representations, average across the whole document, and classify with softmax layers.
- KELM (Agarwal et al., 2021) proposes to generate synthetic pre-training corpora based on structured knowledge bases. In this paper, we further pre-trained the *roberta-base* checkpoint on the KELM synthetic corpus and report performance on downstream tasks.
- KnowBERT (Peters et al., 2019) is one of the first works to leverage external knowledge bases to enrich language representations. We used the three pre-trained models, KnowBERT-Wordnet, KnowBERT-Wikidata, and KnowBERT-W+W for document representation extraction and report performance on downstream tasks.
- Joshi et al. (2020) proposes to learn contextualized language representations by adding Wikipedia text to the input sequences and jointly learning text representations. This is similar to KALM’s setting with only the local context, where Wikipedia descriptions of entities are concatenated to input texts.
- KGAP (Feng et al., 2021a) proposes to construct document graphs to jointly encode textual content and external knowledge. Gated relational graph convolutional networks are then adopted for document representation learning.
- GreaseLM (Zhang et al., 2021) proposes to encode textual content with language model layers, encode knowledge graph subgraphs with graph neural networks and KG embeddings, and adopt MInt layers to fuse the two for question answering. In this paper, we implement GreaseLM by using MInt layers to fuse the local and global contexts.
- GreaseLM+ is our extended version of GreaseLM, which adds the document-level context while keeping the original MInt layer instead of our proposed ContextFusion layer.
- KALM is our proposed approach for knowledge-aware long document understanding. It jointly infuses the local, document-level, and global contexts with external knowledge graphs and adopts ContextFusion layers to derive an overarching document representation.

Table 8: Checklist of whether baselines leverage knowledge graphs and the three document contexts.

	Baseline	Knowledge	Local	Document	Global
<b>task specific</b>	HLSTM (Yang et al., 2016)	✗	✓	✓	✗
	MAN (Li & Goldwasser, 2021)	✗	✓	✓	✗
	KCD (Zhang et al., 2022)	✓	✓	✓	✗
	Rubin et al. (2016)	✗	✓	✓	✗
	Rashkin et al. (2017)	✗	✓	✓	✗
	GCN + Attn (Welling & Kipf, 2016)	✓	✓	✓	✗
	GAT + Attn (Veličković et al., 2018)	✓	✓	✓	✗
	CompareNet (Hu et al., 2021)	✓	✓	✓	✗
	ideal-point (Gerrish & Blei, 2011)	✗	✓	✗	✗
	ideal-vector (Kraft et al., 2016)	✗	✓	✗	✗
	Vote (Mou et al., 2021)	✗	✓	✓	✗
PAR (Feng et al., 2021b)	✓	✓	✓	✗	
<b>language model</b>	RoBERTa (Liu et al., 2019b)	✗	✓	✗	✗
	Electra (Clark et al., 2019)	✗	✓	✗	✗
	DeBERTa (He et al., 2020)	✗	✓	✗	✗
	BART (Lewis et al., 2020)	✗	✓	✗	✗
	LongFormer (Beltagy et al., 2020)	✗	✓	✓	✗
<b>task agnostic</b>	KELM (Agarwal et al., 2021)	✓	✓	✗	✗
	KnowBERT (Peters et al., 2019)	✓	✓	✗	✗
	Joshi et al. (2020)	✓	✓	✗	✗
	KGAP (Feng et al., 2021a)	✓	✗	✓	✗
	GreaseLM (Zhang et al., 2021)	✓	✓	✗	✓
	GreaseLM+ (ours)	✓	✓	✓	✓
<b>KALM (ours)</b>	✓	✓	✓	✓	

Hyperparameter	PPD		MD		RCVP	
	SemEval	Allsides	SLN	LUN	random	time-based
max epochs	50	25	3	5	100	
optimizer	RAdam (Liu et al., 2019a)					
seed LM	BART (Lewis et al., 2020)					
KB embedding	TransE (Bordes et al., 2013)					
dimension of hidden layers	512		512		128	
learning rate	1e-3		1e-3		1e-4	
weight decay	1e-5		1e-5		1e-5	
# KALM layers	2		2		2	
# attention heads	8		8		8	
dropout	0.5		0.5		0.5	
batch size	16		16		4	

Table 9: Hyperparameter settings of KALM.

### D.3 EVALUATION METRICS DETAILS

We adopted these evaluation metrics throughout the paper: Acc (accuracy), MaF (macro-averaged F1-score), MiF (micro-averaged F1-score), MaPrecision (macro-averaged precision), MaRecall (macro-averaged recall), and BAcc (balanced accuracy). These metrics are chosen based on which metrics are used in previous works regarding the three tasks.

### D.4 HYPERPARAMETER DETAILS

We present KALM’s hyperparameter settings in Table 9. We conduct hyperparameter searches for different datasets and report the best setups.

### D.5 WHERE DID THE NUMBERS COME FROM?

For task-specific baselines, we directly use the results reported in previous works (Zhang et al., 2022; Hu et al., 2021; Feng et al., 2021b) since we follow the same experiment settings and the comparison is thus fair. For pre-trained LMs and task-agnostic baselines, we run each method **five times** with different random seeds and report the average performance as well as standard deviation. Figure 4 is an exception, where we only run each method one time due to computing constraints.

### D.6 MORE EXPERIMENT DETAILS

We provide more details about the experiments that are worth further explaining.

- Table 6: We implement pre-trained LMs and task-agnostic baselines for roll call vote prediction by using them to learn representations of legislation texts, concatenate them with the legislator representations learned with PAR (Feng et al., 2021b), and adopt softmax layers for classification.
- Table 2: We remove each context by only applying ContextFusion layers to the other two context representations. We follow the implementation of MInt described in Zhang et al. (2021). We implement concat and sum by using the concatenation and summation of the three context representations as the overall document representation.
- Figure 2: The multi-head attention in the ContextFusion layer provides a  $6 \times 6$  attention weight matrix indicating how information flowed across different contexts. The six rows (columns) stand for the local view of the local context, the global view of the local context, the local view of the document-level context, the global view of the document-level context, the local view of the global context, and the global view of the global context, which are described in detail in Section 2.3.2. The values in each square are the average of the absolute values of the attention weights across all data samples in the validation set.

### D.7 COMPUTATIONAL RESOURCES DETAILS

We used a GPU cluster with 16 NVIDIA A40 GPUs, 1,988G memory, and 104 CPU cores for the experiments. Running KALM with the best parameters takes approximately 1.5, 16, 3, 4, 1, and 1 hour(s) for the six datasets (SemEval, Allsides, SLN, LUN, random, time-based).

### D.8 SCIENTIFIC ARTIFACT DETAILS

KALM is built with the help of many existing scientific artifacts, including TagMe (Ferragina & Scaiella, 2011), pytorch (Paszke et al., 2019), pytorch lightning (Falcon & The PyTorch Lightning team, 2019), transformers (Wolf et al., 2020), pytorch geometric (Fey & Lenssen, 2019), sklearn (Pedregosa et al., 2011), numpy (Harris et al., 2020), nltk (Bird et al., 2009), OpenKE (Han et al., 2018), and the three adopted knowledge graphs (Feng et al., 2021a; Hu et al., 2021; Speer et al., 2017). We commit to make our code and data publicly available upon acceptance to facilitate reproduction and further research.