

EnigmaToM: Improve LLMs' Theory-of-Mind Reasoning Capabilities with Neural Knowledge Base of Entity States

Hainiu Xu[♣] Siya Qi[♣] Jiazheng Li[♣] Yuxiang Zhou^{♣,♠}
Jinhua Du[♡] Caroline Catmur[♣] Yulan He^{♣,◇}

♣King's College London

♡Huawei London Research Centre

◇The Alan Turing Institute

♠Queen Mary University of London

{hainiu.xu, yulan.he}@kcl.ac.uk

Abstract

Theory-of-Mind (ToM), the ability to infer others' perceptions and mental states, is fundamental to human interaction but remains challenging for Large Language Models (LLMs). While existing ToM reasoning methods show promise with reasoning via perceptual perspective-taking, they often rely excessively on off-the-shelf LLMs, reducing their efficiency and limiting their applicability to high-order ToM reasoning. To address these issues, we present EnigmaToM, a novel neuro-symbolic framework that enhances ToM reasoning by integrating a Neural Knowledge Base of entity states (Enigma) for (1) a psychology-inspired *iterative masking* mechanism that facilitates accurate perspective-taking and (2) *knowledge injection* that elicits key entity information. Enigma generates structured knowledge of entity states to build spatial scene graphs for belief tracking across various ToM orders and enrich events with fine-grained entity state details. Experimental results on ToMi, HiToM, and FAN-ToM benchmarks show that EnigmaToM significantly improves ToM reasoning across LLMs of varying sizes, particularly excelling in high-order reasoning scenarios¹.

1 Introduction

Theory-of-Mind (ToM), the ability to understand that others have perceptions and mental states different from one's own, is fundamental to effective communication and social interaction (Premack and Woodruff, 1978; Apperly, 2010). ToM reasoning can be first-order, involving the understanding of another's mental state, or higher-order, requiring recursive thinking about others' beliefs. Higher-order ToM reasoning is particularly vital in real-world contexts such as negotiation (De Weerd et al., 2017). As Large Language Models (LLMs) become increasingly sophisticated in imitating human

¹The neural knowledge base Enigma can be downloaded via <https://huggingface.co/SeacowX/Enigma>. Code and data are available at <https://github.com/seacowx/EnigmaToM>.

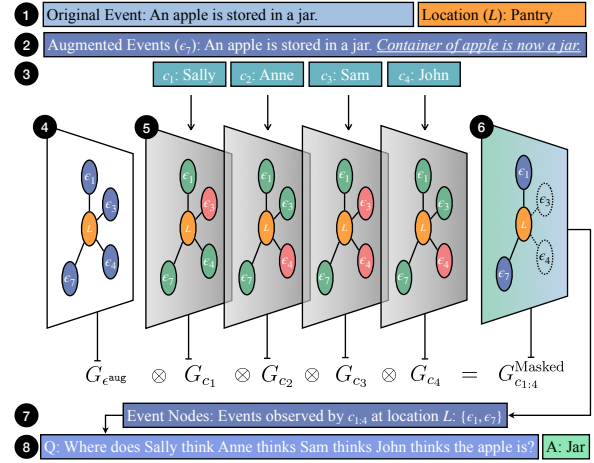


Figure 1: Example use-case of EnigmaToM framework in fourth-order ToM reasoning. An *event* (1) is enriched by adding information about entity-of-interests (*italic text* in (2)) derived from Enigma. Characters (3) are extracted using an off-the-shelf NER model. Spatial scene graphs (4) and (5) are constructed for perspective-taking through a masking mechanism (5 → 6). Event nodes are retrieved to construct character-centric event sequence (7), which is used for the final QA (8).

interactions, a plethora of studies have investigated LLMs' abilities to conduct ToM reasoning. While early studies show that LLMs exhibit traces of ToM capabilities (Bubeck et al., 2023; Kosinski, 2023), follow-up works impugn the robustness of such capabilities by showing that LLMs' ToM reasoning is often superficial (Sap et al., 2022; Ullman, 2023; Shapira et al., 2024).

A vital prerequisite for human ToM reasoning is *perceptual perspective-taking* (referred to as "perspective-taking" thereafter), which is the process of inferring the perception of other characters (Davis, 1983; Harwood and Farrar, 2006). In the case of ToM reasoning with LLMs, perspective-taking alleviates the reasoning burden of LLMs by identifying events that are observable by a given character and removing unobservable ones.

Centered around perspective-taking, numerous methods have been proposed. SimulatedToM (Wilf

et al., 2024) and Discrete World Models (DWM) (Huang et al., 2024) perform perspective-taking by directly prompting LLMs. While one may appreciate these methods’ simplicity, the quality of perspective-taking is largely dependent on the capability of LLMs. SymbolicToM, TimeToM, and PerceptToM took a neuro-symbolic approach. TimeToM (Hou et al., 2024) and PerceptToM (Jung et al., 2024) utilize temporal and perceptual information of events to derive characters’ perception by extracting common timestamps or perceived characters. However, accurately extracting perceived timestamps or perceivers becomes difficult as the length or complexity of the event trajectory increases. The most relevant work to ours is SymbolicToM, where perspective-taking is conducted by maintaining multiple belief graphs (Sclar et al., 2023). However, SymbolicToM constructs belief graphs using less powerful models including WANLI (Liu et al., 2022) and OpenIE (Stanovsky et al., 2018), limiting its generalizability to ToM tasks that involve complicated events. Further, as noted by Sclar et al. (2023), SymbolicToM lacks efficiency as the depth of ToM reasoning increases (see §3.4 for analysis).

Given the need for accurate and efficient perspective-taking in ToM reasoning, we introduce **Entity-Guided Masking** (EnigmaToM), a neuro-symbolic framework enhancing LLMs’ ToM reasoning (Figure 1). Perspective-taking relies on reasoning about event implications, where information about the states of key entities is crucial (Zhang et al., 2023). EnigmaToM employs a Neural Knowledge Base (Enigma) to generate structured entity-state information (§3.1). This entity-state information supports spatial scene graph construction for perspective-taking (§3.3) and event elicitation through knowledge injection (§3.2). Experiment results show that EnigmaToM improves the ToM reasoning capabilities of a range of LLMs. Furthermore, the iterative masking mechanism, grounded by theories from psychology (Arslan et al., 2017), guarantees the efficacy of EnigmaToM across ToM reasoning of varying orders.

We summarize our contributions as follows:

1. We introduce EnigmaToM, a neuro-symbolic framework for ToM reasoning that leverages a Neural Knowledge Base of Entity States to improve LLMs’ ToM reasoning capabilities.
2. Through the iterative masking mechanism, EnigmaToM conducts effective perspective-taking while greatly reducing the number of

character belief graphs that need to be tracked, thereby improving the efficiency in high-order ToM reasoning.

3. EnigmaToM improves LLMs’ ToM reasoning, especially for higher-order cases. Analysis show that EnigmaToM improves LLMs’ ToM reasoning ability up to the fourth order.

2 Related Work

Knowledge Base of Commonsense Knowledge in Natural Language Efforts to construct commonsense knowledge bases have a long history. Early work includes CyC, ConceptNet, and DBPedia (Lenat, 1995; Liu and Singh, 2004; Lehmann et al., 2015). Rashkin et al. (2018) introduced Event2Mind, an event-based knowledge graph that captures characters’ intentions and reactions. Subsequently, Sap et al. (2019) introduced ATOMIC, a commonsense knowledge graph that models if-then relationships for simple events. To explore more complex events, Tandon et al. (2020) introduced OpenPI, a dataset for entity state tracking in procedures. OpenPI was extended to OpenPI2.0 by introducing entity saliency scores and entity canonicalization (Zhang et al., 2024). Parallel efforts have developed neural models, including a GRU-based encoder-decoder model for Event2Mind (Rashkin et al., 2018), a decoder-only Transformer called COMET for ConceptNet and ATOMIC (Bosselut et al., 2019), and fine-tuned GPT-2 for OpenPI (Tandon et al., 2020).

Benchmarking LLMs’ ToM Reasoning Capabilities Many ToM benchmarks are inspired by the False Beliefs test (Wimmer and Perner, 1983), including event-based benchmarks such as ToMi (Le et al., 2019), HiToM (Wu et al., 2023), BigToM (Gandhi et al., 2024), and OpenToM (Xu et al., 2024), and dialogue-based datasets such as FAN-ToM (Kim et al., 2023). Based on the Smarties Test (Gopnik and Astington, 1988), Adv-CSFB (Shapira et al., 2024) and ToMChallenges (Ma et al., 2023) assess LLMs’ ability to reason about unexpected contents and unexpected transfers. ToMBench (Chen et al., 2024) and EPITOME (Jones et al., 2023) contain a suite of ToM tasks that go beyond False Beliefs and Smarties Test. MMTOM-QA extends ToM evaluation to multimodality (Jin et al., 2024) and InformativeBench evaluates ToM in multi-agent settings (Liu et al., 2024).

Improving LLMs’ ToM Reasoning Capabilities Methods for improving LLMs’ ToM reasoning ca-

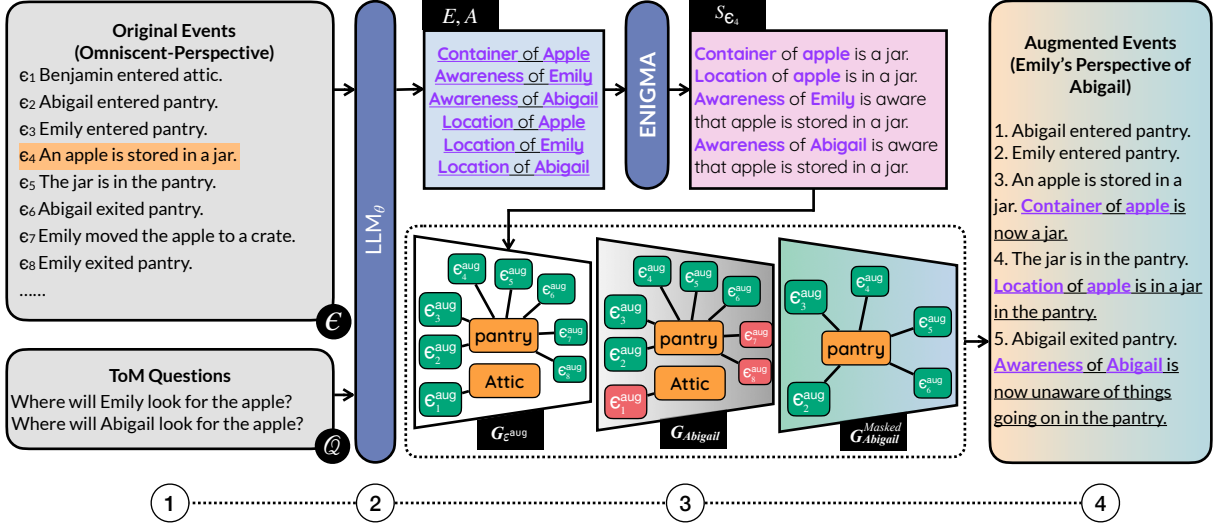


Figure 2: An overview of the EnigmaToM framework. In the graphs shown at bottom of ③, ● nodes denotes observed events while ● nodes denotes unobserved events. See detailed explanations in §3.

pabilities have focused on *perspective-taking*. SymbolicToM conducts perspective-taking via belief graphs (Sclar et al., 2023). SimulatedToM (Wilf et al., 2024) and DWM (Huang et al., 2024) conduct perspective-taking by prompting. DWM additionally prompts LLMs to infer the world state after a group of events. TimeToM utilizes the temporal order of events to conduct perspective-taking (Hou et al., 2024). PerceptToM does perspective-taking by prompting LLMs to infer perceivers of each event (Jung et al., 2024). For multimodal ToM, methods like NIPE and BIP-ALM leverage Bayesian Inverse Planning (Ying et al., 2023; Jin et al., 2024), with environments (e.g. 2D grids or videos) providing strong perspective-taking signals through observable trajectories.

3 The EnigmaToM Framework

Before presenting the EnigmaToM framework, we define the general setup of ToM reasoning tasks.

ToM Task Setup We focus on the widely studied ToM task of *reasoning about false beliefs* (Wimmer and Perner, 1983), which is typically formulated as QA tasks. Formally, given a context consisting of a sequence of events, $\mathcal{E} = \{\epsilon_i\}_{i=1}^n$, which involves multiple characters, $\mathcal{C} = \{c_j\}_{j=1}^m$, and a query regarding the belief of a particular character, $q_c, c \in \mathcal{C}$, the goal is to derive the most likely belief, b_c , from all potential beliefs, \mathcal{B}_c :

$$b_c^* = \arg \max_{b \in \mathcal{B}_c} \mathbb{P}(b | \mathcal{E}, q_c, c \in \mathcal{C}) \quad (1)$$

Further, \mathcal{E} can be concise events as seen in the ToMi dataset (Le et al., 2019) or utterances as seen

in the FANToM dataset (Kim et al., 2023). Beyond directly querying a character’s beliefs about the environment, one can also probe their beliefs regarding other characters’ perceptions, thereby enabling the assessment of higher-order ToM reasoning.

The EnigmaToM Framework Figure 2 provides an overview of our framework, we use circled number (③) to refer to components in the figure. At the core of EnigmaToM is a Neural Knowledge Base (NKB) of Entity States (Enigma). Given a sequence of events (①. \mathcal{E}) and the corresponding questions (①. \mathcal{Q}), EnigmaToM first leverages a chosen LLM (②) to identify key entities (e.g., characters and important objects) and their attributes relevant to ToM reasoning (Top left of ③). Enigma then produces state information for these entities after each event (Top right of ③, §3.1). With the entity state knowledge, EnigmaToM first conducts *Knowledge Injection* (referred to as "KI" thereafter) to enrich the original events by adding relevant fine-grained entity state details (§3.2). Among the entity state knowledge, spatial information of characters is used to conduct perspective-taking through an *Iterative Masking* mechanism (referred to as "IM" thereafter. Bottom of ③, §3.3). The modified events are provided to the LLM for final answers via zero-shot prompting (④). By offloading much of the ToM reasoning process to the symbolic IM component via perspective-taking, EnigmaToM reduces LLMs’ reasoning burden.

3.1 The Enigma Neural Knowledge Base

NKBs such as COMET are trained on a large corpus of structured knowledge in a sequence-

to-sequence manner (Bosselut et al., 2019). Following this approach, we fine-tuned a Llama3.1-8B (Dubey et al., 2024) model to function as our NKB (Zheng et al., 2024).² For training, we used OpenPI2.0 (Zhang et al., 2024), which consists of 25,600 human-annotated entity state changes derived from WikiHow articles. OpenPI2.0 was selected over ATOMIC and Event2Mind as it contains more complex events and entity states. Alternatively, as LLMs become increasingly adept at commonsense reasoning, they can serve as an NKB of entity states via prompting (Hwang et al., 2021). We denote the trained (T) and prompt-based (P) NKB as Enigma^T and Enigma^P, respectively.

To query the NKB, we adopt an *entity-attribute-guided* approach which contains two steps. In Step 1, given a sequence of events, \mathcal{E} , a set of ToM questions, \mathcal{Q} , and a chosen LLM parameterized by θ , we obtain a set of entities of interest, $E = \{e\}_{i=1}^n$, and their corresponding attributes, $A = \{a\}_{j=1}^m$, by zero-shot prompting:

$$E, A = \text{LLM}_\theta(\rho(\mathcal{E}, \mathcal{Q})) \quad (2)$$

where ρ denotes the prompt template (see Appendix D for details of the prompt). Then in Step 2, given an event, $\epsilon \in \mathcal{E}$, a set of entities of interest E and their corresponding attributes A , we query Enigma to retrieve the state of the entities after event ϵ :

$$s_\epsilon = \bigoplus_{i=1}^n \bigoplus_{j=1}^m \text{Enigma}(e_i, a_j, \epsilon) \quad \forall \epsilon \in \mathcal{E}, e_i \in E, a_j \in A \quad (3)$$

where \oplus denotes concatenation.

3.2 Knowledge Injection (KI) with Enigma

In prior studies, perspective-taking was regarded as filtering out events unobserved by a given character, yielding a subset $\mathcal{E}'_c \subseteq \mathcal{E}$. We argue that beyond event filtering, perspective-taking should enhance LLMs' comprehension of events. Fine-grained entity state knowledge is crucial for event reasoning (Zhang et al., 2023) but often omitted due to reporting bias (Shwartz and Choi, 2020). To address this, we propose a knowledge injection mechanism, KI, that utilizes Enigma to enrich observable events with fine-grained entity state information. In the first step of KI, a chosen LLM is used to infer key entities, E , and attributes, A , based on a given sequence of events, \mathcal{E} , and a set of ToM questions, \mathcal{Q} ,

²See Appendix B for details of the fine-tuning process.

(Equation 2). We then query Enigma with the recognized entities and their attributes to obtain their state information at each event (Equation 3). We exclude spatial information of characters, \mathcal{S}_c^p , as this will be handled in the subsequent masking process (§3.3). Given a sequence of events, $\mathcal{E} = \{\epsilon\}_{i=1}^n$, we augment it by injecting entity state knowledge, resulting in the sequence \mathcal{E}^{aug} :

$$\mathcal{E}^{\text{aug}} = \bigoplus_{i=1}^n \epsilon_i \oplus \hat{s}_{\epsilon_i}, \text{ where } \hat{s}_{\epsilon_i} = s_{\epsilon_i} \setminus s_c^p \quad (4)$$

where \oplus denotes concatenation. As fine-grained entity state knowledge is often omitted in events due to reporting bias (Shwartz and Choi, 2020), this mechanism compensates for the lost information. More importantly, by providing state information of key entities, KI reinforces LLMs' understanding of the observed events.

3.3 Perspective-Taking (IM) with Enigma

Studies in psychology have shown that people's beliefs about others' mental states rely only on information available to themselves³ (Arslan et al., 2017). Building on this insight, we assume that characters interpret others' beliefs through the lens of their own mental states, which allows us to employ *Iterative Masking* (IM) to facilitate efficient and accurate ToM reasoning across various order.

Perspective-taking with Enigma is accomplished by constructing spatial scene graphs and performing *Iterative Masking* (IM) using constructed graphs. Specifically, we obtain spatial information, \mathcal{S}_c^p , by querying Enigma about the location (*attr*) of a specific character (*ent*), c , using Equation (3). Spatial scene graphs are constructed based on spatial information to represent the detailed locations where each event takes place as perceived by a given character. The nodes of the scene graph represent events and locations, while the edges denote the "isin" relationship, specifying the location where each event takes place.

During IM, we first construct a character-oblivious spatial scene graph, $G_{\mathcal{E}^{\text{aug}}}$, which documents the *location* of each augmented event from an omniscient perspective. We then construct character-centric spatial scene graphs, G_c , that capture event locations from the perspective of each character. We introduce a null node, \emptyset , which indicates that the location of the current event is

³For instance, "Anne's belief about Sally's mental state" depends only on information available to Anne, i.e. events witnessed by Anne herself.

unknown to the character. During IM, the null node serves as a "mask" to exclude the event nodes, which are unobserved by the character, from $G_{\mathcal{E}^{\text{aug}}}$ (see Figure 1 and Figure 2). For high-order ToM reasoning, $G_{\mathcal{E}^{\text{aug}}}$ is masked sequentially by the order of characters in the belief chain⁴:

$$G_{c_{1:k}}^{\text{masked}} = G_{\mathcal{E}^{\text{aug}}} \bigotimes_{j=1}^k G_{c_j} \quad (5)$$

where \bigotimes represents the masking operation, and k corresponds to the ToM-order. The observable events of character $c_{1:k}$ with injected entity state knowledge can be constructed as:

$$\mathcal{E}_{c_{1:k}}^{\text{aug}} = V_{G_{c_{1:k}}^{\text{masked}}}^{\epsilon} \quad (6)$$

where $V_{G_{c_{1:k}}^{\text{masked}}}^{\epsilon}$ represents event nodes in $G_{c_{1:k}}^{\text{masked}}$. In the case of high-order ToM reasoning, $\mathcal{E}_{c_{1:k}}^{\text{aug}}$ is obtained by iteratively applying the belief of characters. As such, $\mathcal{E}_{c_{1:k}}^{\text{aug}}$ effectively encapsulates the beliefs of all characters in the belief chain. This allows us to transform the high-order ToM question to that of first-order. For instance, reasoning about "Sally's belief about Anne's belief" without EnigmaToM requires first inferring Sally's perceived world state, which then serves as the basis for modeling Anne's belief. With EnigmaToM, such nested dependencies and recursive reasoning are handled by the IM mechanism. Consequently, under $\mathcal{E}_{\text{Sally}, \text{Anne}}^{\text{aug}}$, deriving Sally's belief is sufficient to answer the original second-order Theory of Mind (ToM) question. Illustrative examples and further details on ToM order reduction are provided in Appendix C. We present illustrative examples and details of ToM order reduction in Appendix C.

3.4 Efficiency of EnigmaToM

The IM mechanism of EnigmaToM addresses the intractability of high-order ToM reasoning faced by SymbolicToM (Sclar et al., 2023). Due to the asymmetry of ToM modeling⁵, enumerating all possible mental states for characters across all ToM orders is a permutation problem. Suppose a ToM reasoning question involves m characters and the ToM order goes up to k^{th} -order, the worst-case complexity of constructing belief graphs in SymbolicToM is $\mathcal{O}\left(\sum_{i=1}^k \frac{m!}{(m-i)!}\right)$. In

⁴For instance, the masked spatial scene graph for "Sally's belief of Anne's mental state" is $G_{\mathcal{E}^{\text{aug}}} \otimes G_{\text{Sally}} \otimes G_{\text{Anne}}$.

⁵For example, in second-order ToM, Anne's belief of Sally's mental state is not equivalent to Sally's belief of Anne's mental state.

Dataset	O	Unit	#Units	#Qs
ToMi ⁶ (Le et al., 2019)	2	E	9.85	114
HiToM (Wu et al., 2023)	4	E	26.49	614
FANToM (Kim et al., 2023)	2	U	23.14	577

Table 1: Summary of datasets. **O**: highest ToM order tested. **Unit**: type of event sequence. "E" for event and "U" for utterance. **#Units**: avg. units per sequence. **#Qs**: avg. number of questions per sampled subset. Examples from each dataset can be found in Appendix A.

contrast, EnigmaToM constructs one spatial scene graph, $G_{\mathcal{E}^{\text{aug}}}$, which encapsulates omniscient spatial information, and m character-centric spatial scene graphs. Hence, the worst-case complexity for constructing spatial scene graphs in EnigmaToM is $\tilde{T}(m, k) = \mathcal{O}(m)$, which is linear with respect to the number of characters and independent of the ToM order k . We illustrate the difference in complexity in Appendix E.

4 Experiments

EnigmaToM is evaluated on three widely used ToM benchmarks (Table 1) and compared against the following generic and ToM-specific methods:

CoT (Wei et al., 2022) boosts LLMs' reasoning capabilities by prompting LLMs to explicitly list out their reasoning process.

SimToM (Wilf et al., 2024) conducts perspective-taking by directly querying the LLMs about the mental states of characters.

TimeToM[†] (Hou et al., 2024) leverage the temporal information of events to conduct perspective-taking. The final answer is obtained using a multi-perspective belief-solving prompt.

DWM (Huang et al., 2024) conducts perspective-taking by partitioning the events into chunks and querying the LLMs about characters' mental states after each chunk.

PerceptToM[†] (Jung et al., 2024) conducts perspective-taking by querying the LLMs about the characters' awareness of the events.

To ensure a fair comparison with established methods, we conduct controlled experiments by controlling the format and answer space of all ToM questions. In addition, we follow a realistic setting of ToM reasoning by using only the sequence of events and ToM questions from each dataset.

⁶We use the disambiguated ToMi (Sclar et al., 2023) from <https://github.com/msclar/symbolictom>.

[†]Official implementation is not available at the time of experiments (Sept-Dec, 2024). We implemented this method using prompts from the corresponding paper.

		Qwen2.5-7B	Llama3.1-8B	Gemma2-9B	Gemma2-27B	Llama3.3-70B ^{4bit}	Qwen2.5-72B ^{4bit}	GPT-4o
ToMi	Vanilla	0.722±0.045	0.647±0.011	0.741±0.037	0.715±0.048	0.767±0.015	0.717±0.034	0.767±0.041
	CoT	0.724±0.026	0.739±0.025	0.676±0.035	0.537±0.056	0.741±0.032	0.767±0.033	0.769±0.029
	SimToM	0.642±0.022	0.600±0.020	0.710±0.034	0.684±0.015	0.712±0.018	0.749±0.020	0.749±0.018
	TimeToM	0.567±0.024	0.630±0.019	0.681±0.028	0.587±0.036	0.739±0.021	0.865±0.018	0.723±0.016
	DWM	0.686±0.023	0.644±0.033	0.718±0.028	0.707±0.045	0.735±0.016	0.762±0.051	0.739±0.049
	PerceptToM	0.720±0.038	0.695±0.025	0.676±0.029	0.749±0.017	0.738±0.032	0.809±0.033	0.790±0.023
	Enigma ^p	0.706±0.044	0.738±0.056	0.865±0.031	0.833±0.018	0.828±0.012	0.839±0.014	0.847±0.030
	Enigma ^t	0.825±0.030	0.796±0.023	0.814±0.020	0.804±0.050	0.787±0.024	0.837±0.024	0.795±0.036
HiToM	Vanilla	0.378±0.013	0.333±0.015	0.471±0.009	0.527±0.018	0.534±0.008	0.456±0.012	0.521±0.006
	CoT	0.441±0.007	0.304±0.021	0.474±0.008	0.535±0.018	0.537±0.011	0.481±0.011	0.527±0.005
	SimToM	0.402±0.009	0.368±0.024	0.473±0.012	0.549±0.018	0.569±0.005	0.536±0.018	0.571±0.003
	TimeToM	0.316±0.010	0.462±0.013	0.302±0.012	0.302±0.013	0.623±0.006	0.415±0.013	0.633±0.008
	DWM	0.444±0.020	0.367±0.019	0.485±0.012	0.488±0.018	0.564±0.010	0.560±0.009	0.580±0.018
	PerceptToM	0.393±0.019	0.342±0.011	0.440±0.009	0.562±0.007	0.588±0.010	0.548±0.016	0.580±0.018
	Enigma ^p	0.508±0.012	0.477±0.005	0.555±0.010	0.576±0.004	0.696±0.007	0.605±0.007	0.733±0.017
	Enigma ^t	0.457±0.005	0.431±0.010	0.446±0.008	0.478±0.004	0.518±0.011	0.473±0.010	0.626±0.020
FANToM	Vanilla	0.400±0.015	0.429±0.022	0.485±0.016	0.553±0.011	0.486±0.022	0.532±0.025	0.476±0.020
	CoT	0.398±0.014	0.438±0.014	0.470±0.019	0.556±0.007	0.494±0.028	0.521±0.024	0.453±0.014
	SimToM	0.413±0.012	0.440±0.015	0.427±0.009	0.574±0.010	0.620±0.025	0.516±0.014	0.502±0.016
	TimeToM	0.252±0.020	0.260±0.012	0.299±0.011	0.300±0.021	0.580±0.017	0.409±0.026	0.404±0.016
	DWM	0.429±0.013	0.470±0.027	0.433±0.023	0.562±0.017	0.473±0.021	0.543±0.014	0.465±0.028
	PerceptToM	0.408±0.023	0.407±0.026	0.504±0.006	0.611±0.009	0.527±0.011	0.573±0.016	0.521±0.006
	Enigma ^p	0.445±0.026	0.442±0.018	0.439±0.023	0.462±0.014	0.515±0.020	0.450±0.013	0.531±0.015
	Enigma ^t	0.487±0.018	0.545±0.036	0.530±0.012	0.582±0.028	0.610±0.021	0.574±0.031	0.553±0.011

Table 2: Main results of EnigmaToM in comparison with existing methods on ToMi, HiToM, and FANToM datasets. Accuracy means and variances are calculated based on 5 runs, which used 5 different subsets of the corresponding dataset. The best and second best results are highlighted in **bold** and underline respectively.

Auxiliary information such as character names is obtained using an off-the-shelf NER model⁸.

Question Formatting We formulate ToMi as a free-form generation task where the model is instructed to choose between two possible answers. We formulate HiToM as a multiple-choice task as in the original paper (Wu et al., 2023). FANToM contains both free-form generation and multiple-choice questions. We follow the question formatting instructions in the original paper (Kim et al., 2023). For efficient and accurate parsing of LLM responses, we follow the convention of (Huang et al., 2024), instructing LLMs to wrap answers within the special <answer> and </answer> tokens. As introduced in §3.3, the recursive modeling of mental states in high-order ToM questions has been addressed by the IM mechanism, which allows us to transform high-order ToM questions into first-order questions. Similarly, TimeToM leverages temporal information to conduct symbolic modeling of high-order ToM (Hou et al., 2024). We apply such transformation when evaluating with TimeToM and EnigmaToM (Appendix C).

⁸<https://huggingface.co/dslim/bert-large-NER>

Towards Robust Evaluation To ensure robust evaluation, we construct 5 subsets for each dataset by sampling data points using commonly used random seeds[†]. Each subset of ToMi and HiToM contains 100 event sequences, whereas each subset of FANToM contains 50 multi-round dialogues. The number of QA pairs in each subset is shown in Table 1. We report both the mean accuracy and its variance based on the 5 runs.

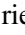
We evaluate each method using various instruction-tuned LLMs, including Llama3.1-8B, Llama3.3-70B[‡], Qwen2.5-7B, Qwen2.5-72B[‡], Gemma2-9B, Gemma2-27B, and GPT-4o (Dubey et al., 2024; Yang et al., 2024; Gemma, 2024; OpenAI, 2024). To ensure reproducibility, all experiments are done using zero-shot prompting with greedy decoding and a temperature of 0. LLM inference is carried out using  vLLM on 2 NVIDIA A100^{80GB} GPUs (Kwon et al., 2023).

Table 2 shows the main results of EnigmaToM in comparison with existing methods on ToMi,

[†]We use 12, 42, 96, 2012, and 2024 as random seeds.

[‡]Loaded in 4bit using BitsandBytes (Dettmers et al.) with weights from <https://huggingface.co/unsloth>.

HiToM, and FANToM datasets. In general, we see that EnigmaToM brings improvements in accuracy across all datasets and most LLMs. Specifically, Enigma^P outperforms other methods on ToMi and HiToM, while Enigma^T achieves superior performance on FANToM. EnigmaToM is particularly effective with smaller LLMs. For instance, Enigma^T boosts Qwen2.5-7B to exceed the zero-shot performance of Qwen2.5-72B^{4bit}. Further, results from the HiToM dataset demonstrate that EnigmaToM is particularly effective in high-order ToM reasoning. We analyze the effectiveness of EnigmaToM in tackling high-order ToM reasoning in §5.1. Moreover, results from Table 2 show that Enigma^P performs better on event-based datasets (ToMi and HiToM) while Enigma^T is more effective on a dialogue-based dataset (FANToM). We investigate such a discrepancy in §5.2 and §5.3.

5 Analysis

5.1 High-Order ToM Reasoning

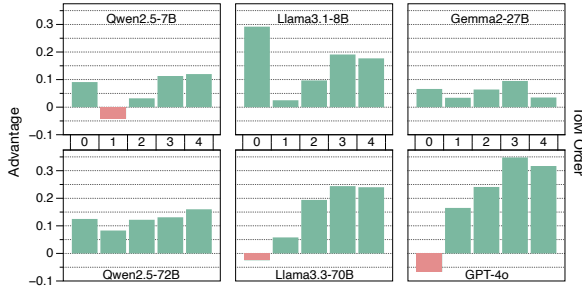


Figure 3: Relative advantage of EnigmaToM on HiToM dataset with respect to ToM order.

To assess the effectiveness of EnigmaToM in high order ToM reasoning, we analyze its performance on the HiToM dataset, which consists of ToM questions requiring reasoning up to the fourth order. We compute the relative advantage of EnigmaToM with Enigma^P over the zero-shot vanilla prompting baseline. From Figure 3, we observe that EnigmaToM improves mean accuracy across all orders of ToM reasoning, with notable effectiveness in higher-order ToM reasoning. Specifically, results from the Qwen2.5 and Llama3 families demonstrate that EnigmaToM has an increasing advantage as the order of ToM reasoning increases. For the third- and fourth-order ToM reasoning, EnigmaToM achieves an average improvement of 0.160 ± 0.003 and 0.148 ± 0.004 respectively, across all models compared to the baseline. We observe similar trends on ToMi and FANToM albeit they only contain ToM questions up to the second

		Llama3.3-70B ^{4bit}	Qwen2.5-72B ^{4bit}	GPT-4o
ToMi	Enigma ^P	0.828 \pm 0.012	0.839 \pm 0.014	0.847 \pm 0.030
	Enigma ^T	0.787 \pm 0.024	0.837 \pm 0.024	0.795 \pm 0.036
	w/o KI	0.834 \pm 0.067	0.845 \pm 0.026	0.811 \pm 0.028
	w/o IM	0.693 \pm 0.014	0.655 \pm 0.039	0.674 \pm 0.002
	w/o KI, IM	0.767 \pm 0.015	0.717 \pm 0.034	0.767 \pm 0.041
HiToM	Enigma ^P	0.696 \pm 0.007	0.605 \pm 0.007	0.733 \pm 0.017
	Enigma ^T	0.518 \pm 0.011	0.473 \pm 0.010	0.626 \pm 0.020
	w/o KI	0.726 \pm 0.004	0.632 \pm 0.003	0.751 \pm 0.004
	w/o IM	0.460 \pm 0.013	0.423 \pm 0.008	0.442 \pm 0.006
	w/o KI, IM	0.534 \pm 0.008	0.456 \pm 0.012	0.521 \pm 0.006
FANToM	Enigma ^P	0.515 \pm 0.020	0.450 \pm 0.013	0.531 \pm 0.015
	Enigma ^T	0.610 \pm 0.021	0.574 \pm 0.031	0.553 \pm 0.011
	w/o KI	0.607 \pm 0.018	0.542 \pm 0.036	0.539 \pm 0.012
	w/o IM	0.500 \pm 0.021	0.477 \pm 0.017	0.470 \pm 0.013
	w/o KI, IM	0.486 \pm 0.002	0.532 \pm 0.025	0.476 \pm 0.020

Table 3: Ablation study of EnigmaToM on ToMi, HiToM, and FANToM datasets. "w/o KI" indicates without *entity state knowledge injection*. "w/o IM" denotes without *perspective-taking via iterative masking*. ● Improved and ● decreased results are highlighted.

order. See Appendix F for complete results and analysis on all three datasets.

5.2 Ablation Study

To understand the effectiveness of each component of EnigmaToM, we conduct an ablation study by (1) keeping the injected knowledge but removing the masking-based perspective-taking mechanism (directly using \mathcal{E}^{aug} as context); and (2) conducting perspective-taking without knowledge injection (applying Equation 5 with $G_{\mathcal{E}}$ instead of $G_{\mathcal{E}^{\text{aug}}}$).

Enigma for Perspective Taking As shown in Table 3, removing the IM mechanism results in an average accuracy drop of -0.165 on ToMi, -0.172 on HiToM, and -0.103 on FANToM. This suggests that the iterative masking mechanism is effective in perspective-taking and crucial for EnigmaToM to achieve boosted performance in ToM reasoning (See Table A2 for complete results).

Enigma for Knowledge Injection Compared to IM for perspective-taking, entity state knowledge injection is less critical. On ToMi and HiToM, its removal slightly reduces Gemma2-27B’s performance on ToMi but improves performance for all other LLMs on both benchmarks, further highlighting IM’s effectiveness in perspective-taking. However, for FANToM, entity state knowledge is indispensable, as excluding it results in performance drops across all LLMs. For ToMi and HiToM, we hypothesize that larger LLMs are better at handling reporting bias. This aligns with the results shown in Table 2, where Enigma^P surpasses Enigma^T as LLM

Dataset	Model	Precision	Recall	F1-Score
TMi	Llama3.3-70B ^{4bit}	0.859	0.968	0.910
FTM	Llama3.3-70B ^{4bit}	0.880	0.970	0.923

Table 4: Performance analysis of key entity recognition in ToMi (TMi) and FANToM (FTM) datasets using Llama3.3-70B^{4bit}. See Appendix J for detailed description of the evaluation process.

	Model	Relevance	Accuracy	Avg. #Token
TMi	Enigma ^T _{8B}	0.847	0.807	7.665
	Enigma ^T _{70B}	0.870	0.860	9.740
FTM	Enigma ^T _{8B}	0.880	0.773	8.973
	Enigma ^T _{70B}	0.880	0.700	30.517

Table 5: Performance analysis of Enigma^T on ToMi (TMi) and FANToM (FTM) datasets. See Appendix J for detailed description of the evaluation process.

size increases, meaning that the fine-grained information about the state of the entity and its causal relationships with events are encapsulated more effectively in the larger LLMs. In such cases, potential inaccuracies in injected entity-state knowledge outweigh its benefits in addressing reporting bias, leading to decreased performance. In the case of FANToM, the dialogue-based nature of the dataset makes useful information sparser than in event-based datasets. Here, knowledge injection serves a different role: rather than primarily addressing reporting bias, it compresses important information from utterances into entity-state representation, effectively reducing LLMs’ workload in identifying crucial information. See Appendix H for examples.

5.3 Effectiveness of LLMs and Enigma

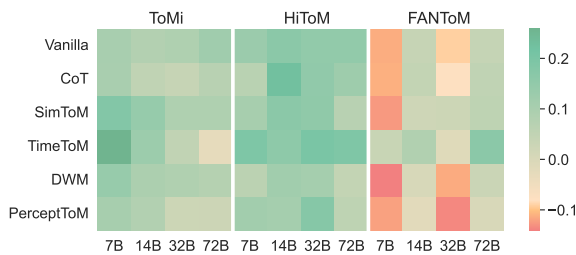


Figure 4: Relative advantage of EnigmaToM on ToMi, HiToM, and FANToM datasets. We use Enigma^P as the pivot method for ToMi and HiToM and Enigma^T for FANToM. Exact mean accuracies are shown in Table A1. Model sizes shown in x-axis are Qwen2.5 models.

The effectiveness of EnigmaToM is contingent upon the capabilities of both the Enigma neural knowledge base and the LLM deployed in the framework. In this section, we conduct analysis

of EnigmaToM with an aim to explore the following two questions: (1) Does EnigmaToM benefit LLMs of larger size? and (2) How effective is our Enigmeneural knowledge base and does scaling Enigmalead to increased performance?

We raise the first question by hypothesizing that perspective-taking, albeit the challenges posed by its multi-hop nature, could become solvable by more capable LLMs. Empirically speaking, the capability of LLMs positively correlates to their number of parameters. To eliminate potential confounding factors, we analyze the effectiveness of LLMs from the Qwen2.5 family with sizes ranging from 7B to 72B (Yang et al., 2024).

In the second question, we aim to examine the effectiveness of Enigma^T. Trained using data from OpenPI2.0, we wish to investigate how well can the knowledge encapsulated in Enigma^T be transferred to aid ToM reasoning. Aside from the performance of Enigma^T, we also explore the effectiveness of scaling of Enigma^T. In addition to the Enigma^T used in previous experiments, which is trained using a Llama3.1-8B model, we trained another Enigma^T based on Llama3.3-70B, which we denote as Enigma^T_{70B}. Experiments with Enigma^T_{70B} are carried out following the same procedure described in §4.

Scaling of Base LLMs We compute the relative advantage of EnigmaToM by calculating the difference in mean accuracy between EnigmaToM and the most performant baseline methods (see Table 2). We use Enigma^P as the pivot method for ToMi and HiToM and Enigma^T for FANToM. Figure 4 shows two trends: (1) a slight diminishment in advantage on ToMi and (2) a gradual increase in advantage on FANToM. We attribute this to the differing difficulty levels of these two datasets. ToMi, which consists of short sequences of concise events, becomes easier to solve with large-scale LLMs. Conversely, FANToM, featuring long sequences of lengthy dialogues, remains challenging even for larger LLMs. HiToM, positioned between these two extremes with long sequences of concise events, shows that EnigmaToM has a consistent advantage regardless of the LLM sizes. This discrepancy in performance and model scaling effect between ToMi and FANToM aligns with the analysis in §4 and §5.3. These findings suggest that while prompting large-scale LLMs can potentially tackle ToM reasoning involving short event sequences (as in ToMi), ToM reasoning about lengthy event or dialogue sequences (as in HiToM and FANToM) can benefit from the

		Qwen2.5-7B	Llama3.1-8B	Gemma2-9B	Gemma2-27B	Llama3.3-70B ^{4bit}	Qwen2.5-72B ^{4bit}	GPT-4o
ToMi	Enigma ^P	0.706 \pm 0.044	0.738 \pm 0.056	0.865 \pm 0.031	0.833 \pm 0.018	0.828 \pm 0.012	0.839 \pm 0.014	0.847 \pm 0.030
	Enigma ^T _{8B}	0.825 \pm 0.030	0.796 \pm 0.023	0.814 \pm 0.020	0.804 \pm 0.050	0.787 \pm 0.024	0.837 \pm 0.024	0.795 \pm 0.036
	Enigma ^T _{70B}	0.837 \pm 0.017	0.835 \pm 0.026	0.847 \pm 0.008	0.854 \pm 0.023	0.844 \pm 0.018	0.860 \pm 0.015	0.872 \pm 0.026
HiToM	Enigma ^P	0.508 \pm 0.012	0.477 \pm 0.005	0.555 \pm 0.010	0.576 \pm 0.004	0.696 \pm 0.007	0.605 \pm 0.007	0.733 \pm 0.017
	Enigma ^T _{8B}	0.457 \pm 0.005	0.431 \pm 0.010	0.446 \pm 0.008	0.478 \pm 0.004	0.518 \pm 0.011	0.473 \pm 0.010	0.626 \pm 0.020
	Enigma ^T _{70B}	0.456 \pm 0.007	0.444 \pm 0.012	0.481 \pm 0.006	0.489 \pm 0.012	0.517 \pm 0.010	0.467 \pm 0.009	0.733 \pm 0.010
FANToM	Enigma ^P	0.445 \pm 0.026	0.442 \pm 0.018	0.439 \pm 0.023	0.462 \pm 0.014	0.515 \pm 0.020	0.450 \pm 0.013	0.531 \pm 0.015
	Enigma ^T _{8B}	0.487 \pm 0.018	0.545 \pm 0.036	0.530 \pm 0.012	0.582 \pm 0.028	0.610 \pm 0.021	0.574 \pm 0.031	0.553 \pm 0.011
	Enigma ^T _{70B}	0.479 \pm 0.032	0.517 \pm 0.041	0.491 \pm 0.039	0.529 \pm 0.028	0.537 \pm 0.011	0.494 \pm 0.023	0.539 \pm 0.012

Table 6: Results of scaling Enigma^T from Llama3.1-8B to Llama3.3-70B. ● Improved, ● Unchanged, and ● decreased results are highlighted in the corresponding color.

fine-grained entity state knowledge as well as the symbolic masking mechanism of EnigmaToM.

Effectiveness of LLMs in Recognizing Key Entity-Attributes Recognizing entities and their attributes (Equation 2) that are indispensable for answering the ToM questions posed is a critical pre-requisite for the effectiveness of EnigmaToM. On the one hand, failing to recognize a key entity will disable EnigmaToM to properly augment the events with critical information. On the other hand, erroneously identify an extraneous entity will lead to inclusion of redundant information, which will prolong the context and increase the reasoning burden of the LLM. To evaluate the quality of key entities and attributes extracted by LLMs, we manually labeled 300 entities and attributes identified using Llama3.3-70B^{4bit}. Evaluation results from Table 4 suggests that LLMs are more than competent in identifying key entities and attributes. With a F1-score of 0.910 on the ToMi dataset and 0.923 on the FANToM dataset, it is safe to conclude that the vast majority of entities identified by LLMs are indeed vital to answering the ToM questions posed.

Effectiveness of Enigma in Generating Entity State Information To understand the effectiveness of scaling Enigma, we trained two Enigma^T models using the same OpenPI2.0 dataset. With Llama3.1-8B as the base model, we trained Enigma^T_{8B}. Further, with Llama3.3-70B, we trained Enigma^T_{70B}¹¹. Table 6 shows that there is an obvious discrepancy between the scaling effect of Enigma^T: Enigma^T_{70B} consistently outperforms Enigma^T_{8B} in ToMi and HiToM while underperforming Enigma^T_{8B} in FANToM.

- **Relevance:** whether the entity and attribute contribute to answering the ToM question. This evaluates the same aspects of EnigmaToM as the precision scores shown in Table 4

- **Accuracy:** whether the entity state can be inferred from the given context.

From Table 5, we see that the relevance scores of both ToMi and FANToM exceed 80%, indicating that Llama3.3-70B is capable of identifying entities and attributes useful for ToM reasoning. Entity state length and accuracy are closely correlated. Enigma^T_{70B} produces a more articulated response compared to its counterpart. Specifically, scaling Enigma^T brings 5.3% improvement in accuracy on ToMi while increasing response length by only 2.075 tokens. In contrast, FANToM experiences a significant 21.544 token increase in response length, which reduces Enigma^T’s efficiency as an information compressor and leads to greater hallucination, resulting in a 7.3% drop in accuracy. We provide demonstration examples in Appendix I.

6 Conclusion

In this work, we introduced EnigmaToM, a neuro-symbolic framework designed to enhance the ToM reasoning capabilities of LLMs. By leveraging a Neural Knowledge Base of Entity States through an iterative masking mechanism and knowledge injection, EnigmaToM accomplishes the bulk of ToM reasoning via perspective-taking through symbolic reasoning, which alleviates LLMs’ reasoning burden. Experimental results across multiple benchmarks demonstrate that EnigmaToM outperforms existing methods, particularly excelling in high-order ToM reasoning scenarios. Our analysis highlights the effectiveness of the iterative masking mechanism in maintaining strong performance across varying depths of ToM reasoning, as well as the critical role of fine-grained entity state knowledge in compressing key information in complex event sequences (as in FANToM). Furthermore, the framework’s efficiency and scalability make it a promising solution for addressing the computational challenges associated with high-order ToM reasoning tasks.

¹¹Training details are provided in Appendix B.

Limitations

ToM Reasoning Beyond Character Perception

EnigmaToM tackles ToM reasoning of characters' beliefs based on their perceptions. While we believe that reasoning about characters' perceptions serve as a cornerstone for all types of ToM reasoning, future work may explore methods to facilitate real-world ToM reasoning about characters' emotions, intentions, desires, and their inherent subjectivity (Zhou et al., 2025).

Neural Knowledge Base EnigmaToM relies on access to a Neural Knowledge Base (NKB) to retrieve entity-state information for answering ToM questions. While Table 5 shows that Enigma is capable of producing accurate entity-state information, it can be further improved (e.g. full-parameter fine-tuning instead of LoRA). Further, expanding the NKB to incorporate richer entity-state details, including emotional, temporal, and causal relationships, would be beneficial for ToM reasoning about high-level information.

Error Propagation While experiments demonstrate the effectiveness of the IM mechanism, it is prone to error propagation. In the case of high-order ToM reasoning, applying a wrong mask in the iterative masking process will lead to the event being erroneously excluded and vice versa. Additionally, in cases requiring complex reasoning about non-linear or intertwined event dependencies, the symbolic Iterative Masking (IM) mechanism may need to be enhanced.

Ethics Statement

This study aims to enhance LLMs' ToM reasoning by improving the accuracy and efficiency of perceptual perspective-taking, ultimately optimizing their effectiveness in communication. ToM reasoning is essential for enhancing LLMs' ability to interact with humans (e.g., in chatbots) or other LLMs (e.g., in multi-agent systems). The evaluation datasets used in this study have been peer-reviewed and widely adopted in previous research. However, these datasets may introduce issues such as cultural bias and often lack demographic information. Future research could incorporate auxiliary data, such as demographic and personality traits, to improve representativeness across diverse ethnic and cultural backgrounds.

Acknowledgments

This work was supported in part by the UK Engineering and Physical Sciences Research Council (EPSRC) through an iCASE award with Huawei London Research Centre and a Turing AI Fellowship (grant no. EP/V020579/1, EP/V020579/2).

References

- Ian Apperly. 2010. *Mindreaders: the cognitive basis of "theory of mind"*. Psychology Press.
- Burcu Arslan, Niels A Taatgen, and Rineke Verbrugge. 2017. Five-year-olds' systematic errors in second-order false belief tasks are due to first-order theory of mind strategy selection: a computational modeling study. *Frontiers in psychology*, 8:275.
- Antoine Bosselut, Hannah Rashkin, Maarten Sap, Chaitanya Malaviya, Asli Celikyilmaz, and Yejin Choi. 2019. **COMET: Commonsense transformers for automatic knowledge graph construction**. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4762–4779, Florence, Italy. Association for Computational Linguistics.
- Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrkke, Eric Horvitz, Ece Kamar, Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, et al. 2023. Sparks of artificial general intelligence: Early experiments with gpt-4. *arXiv preprint arXiv:2303.12712*.
- Zhuang Chen, Jincenzi Wu, Jinfeng Zhou, Bosi Wen, Guanqun Bi, Gongyao Jiang, Yaru Cao, Mengting Hu, Yunghwei Lai, Zexuan Xiong, and Minlie Huang. 2024. **ToMBench: Benchmarking theory of mind in large language models**. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 15959–15983, Bangkok, Thailand. Association for Computational Linguistics.
- Mark H Davis. 1983. Measuring individual differences in empathy: Evidence for a multidimensional approach. *Journal of personality and social psychology*, 44(1):113.
- Harmen De Weerd, Rineke Verbrugge, and Bart Verheij. 2017. Negotiating with other minds: the role of recursive theory of mind in negotiation with incomplete information. *Autonomous Agents and Multi-Agent Systems*, 31:250–287.
- Tim Dettmers, Mike Lewis, Sam Shleifer, and Luke Zettlemoyer. 8-bit optimizers via block-wise quantization. In *International Conference on Learning Representations*.
- Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. 2024. Qlora: Efficient finetuning of quantized llms. *Advances in Neural Information Processing Systems*, 36.

- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.
- Kanishk Gandhi, Jan-Philipp Fränken, Tobias Gerstenberg, and Noah Goodman. 2024. Understanding social reasoning in language models with language models. *Advances in Neural Information Processing Systems*, 36.
- Gemma. 2024. Gemma 2: Improving open language models at a practical size. *arXiv preprint arXiv:2408.00118*.
- Alison Gopnik and Janet W Astington. 1988. Children’s understanding of representational change and its relation to the understanding of false belief and the appearance-reality distinction. *Child development*, pages 26–37.
- Michelle D Harwood and M Jeffrey Farrar. 2006. Conflicting emotions: The connection between affective perspective taking and theory of mind. *British Journal of Developmental Psychology*, 24(2):401–418.
- Guiyang Hou, Wenqi Zhang, Yongliang Shen, Linjuan Wu, and Weiming Lu. 2024. [TimeToM: Temporal space is the key to unlocking the door of large language models’ theory-of-mind](#). In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 11532–11547, Bangkok, Thailand. Association for Computational Linguistics.
- Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2022. [LoRA: Low-rank adaptation of large language models](#). In *International Conference on Learning Representations*.
- X. Angelo Huang, Emanuele La Malfa, Samuele Marro, Andrea Asperti, Anthony G. Cohn, and Michael J. Wooldridge. 2024. [A notion of complexity for theory of mind via discrete world models](#). In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 2964–2983, Miami, Florida, USA. Association for Computational Linguistics.
- Jena D Hwang, Chandra Bhagavatula, Ronan Le Bras, Jeff Da, Keisuke Sakaguchi, Antoine Bosselut, and Yejin Choi. 2021. (comet-) atomic 2020: On symbolic and neural commonsense knowledge graphs. In *Proceedings of the AAAI conference on artificial intelligence*, volume 35, pages 6384–6392.
- Chuangyang Jin, Yutong Wu, Jing Cao, Jiannan Xiang, Yen-Ling Kuo, Zhiting Hu, Tomer Ullman, Antonio Torralba, Joshua Tenenbaum, and Tianmin Shu. 2024. [MMToM-QA: Multimodal theory of mind question answering](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 16077–16102, Bangkok, Thailand. Association for Computational Linguistics.
- Cameron Robert Jones, Sean Trott, and Ben Bergen. 2023. Epitome: Experimental protocol inventory for theory of mind evaluation. In *First Workshop on Theory of Mind in Communicating Agents*.
- Chani Jung, Dongkwan Kim, Jiho Jin, Jiseon Kim, Yeon Seonwoo, Yejin Choi, Alice Oh, and Hyunwoo Kim. 2024. [Perceptions to beliefs: Exploring precursory inferences for theory of mind in large language models](#). In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 19794–19809, Miami, Florida, USA. Association for Computational Linguistics.
- Hyunwoo Kim, Melanie Sclar, Xuhui Zhou, Ronan Bras, Gunhee Kim, Yejin Choi, and Maarten Sap. 2023. [FANToM: A benchmark for stress-testing machine theory of mind in interactions](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 14397–14413, Singapore. Association for Computational Linguistics.
- Michal Kosinski. 2023. Theory of mind may have spontaneously emerged in large language models.
- Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. 2023. Efficient memory management for large language model serving with pagedattention. In *Proceedings of the ACM SIGOPS 29th Symposium on Operating Systems Principles*.
- Matthew Le, Y-Lan Boureau, and Maximilian Nickel. 2019. [Revisiting the evaluation of theory of mind through question answering](#). 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)*, pages 5872–5877, Hong Kong, China. Association for Computational Linguistics.
- Jens Lehmann, Robert Isele, Max Jakob, Anja Jentzsch, Dimitris Kontokostas, Pablo N Mendes, Sebastian Hellmann, Mohamed Morsey, Patrick Van Kleef, Sören Auer, et al. 2015. Dbpedia—a large-scale, multilingual knowledge base extracted from wikipedia. *Semantic web*, 6(2):167–195.
- Douglas B Lenat. 1995. Cyc: A large-scale investment in knowledge infrastructure. *Communications of the ACM*, 38(11):33–38.
- Alisa Liu, Swabha Swayamdipta, Noah A. Smith, and Yejin Choi. 2022. [WANLI: Worker and AI collaboration for natural language inference dataset creation](#). In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 6826–6847, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Hugo Liu and Push Singh. 2004. [Conceptnet — a practical commonsense reasoning tool-kit](#). *BT Technology Journal*, 22:211–226.

- Wei Liu, Chenxi Wang, YiFei Wang, Zihao Xie, Rennai Qiu, Yufan Dang, Zhuoyun Du, Weize Chen, Cheng Yang, and Chen Qian. 2024. Autonomous agents for collaborative task under information asymmetry. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*.
- Xiaomeng Ma, Lingyu Gao, and Qihui Xu. 2023. [ToM-Challenges: A principle-guided dataset and diverse evaluation tasks for exploring theory of mind](#). In *Proceedings of the 27th Conference on Computational Natural Language Learning (CoNLL)*, pages 15–26, Singapore. Association for Computational Linguistics.
- OpenAI. 2024. Gpt-4o system card. *arXiv preprint arXiv:2410.21276*.
- David Premack and Guy Woodruff. 1978. Does the chimpanzee have a theory of mind? *Behavioral and brain sciences*, 1(4):515–526.
- Hannah Rashkin, Maarten Sap, Emily Allaway, Noah A. Smith, and Yejin Choi. 2018. [Event2Mind: Commonsense inference on events, intents, and reactions](#). In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 463–473, Melbourne, Australia. Association for Computational Linguistics.
- Maarten Sap, Ronan Le Bras, Emily Allaway, Chandra Bhagavatula, Nicholas Lourie, Hannah Rashkin, Brendan Roof, Noah A Smith, and Yejin Choi. 2019. Atomic: An atlas of machine commonsense for if-then reasoning. In *Proceedings of the AAAI conference on artificial intelligence*, volume 33, pages 3027–3035.
- Maarten Sap, Ronan Le Bras, Daniel Fried, and Yejin Choi. 2022. [Neural theory-of-mind? on the limits of social intelligence in large LMs](#). In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 3762–3780, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Melanie Sclar, Sachin Kumar, Peter West, Alane Suhr, Yejin Choi, and Yulia Tsvetkov. 2023. [Minding language models’ \(lack of\) theory of mind: A plug-and-play multi-character belief tracker](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 13960–13980, Toronto, Canada. Association for Computational Linguistics.
- Natalie Shapira, Mosh Levy, Seyed Hossein Alavi, Xuhui Zhou, Yejin Choi, Yoav Goldberg, Maarten Sap, and Vered Shwartz. 2024. [Clever hans or neural theory of mind? stress testing social reasoning in large language models](#). In *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2257–2273, St. Julian’s, Malta. Association for Computational Linguistics.
- Vered Shwartz and Yejin Choi. 2020. [Do neural language models overcome reporting bias?](#) In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 6863–6870, Barcelona, Spain (Online). International Committee on Computational Linguistics.
- Gabriel Stanovsky, Julian Michael, Luke Zettlemoyer, and Ido Dagan. 2018. [Supervised open information extraction](#). In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 885–895, New Orleans, Louisiana. Association for Computational Linguistics.
- Niket Tandon, Keisuke Sakaguchi, Bhavana Dalvi, Dheeraj Rajagopal, Peter Clark, Michal Guerquin, Kyle Richardson, and Eduard Hovy. 2020. [A dataset for tracking entities in open domain procedural text](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 6408–6417, Online. Association for Computational Linguistics.
- Tomer Ullman. 2023. Large language models fail on trivial alterations to theory-of-mind tasks. *arXiv preprint arXiv:2302.08399*.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837.
- Alex Wilf, Sihyun Lee, Paul Pu Liang, and Louis-Philippe Morency. 2024. [Think twice: Perspective-taking improves large language models’ theory-of-mind capabilities](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 8292–8308, Bangkok, Thailand. Association for Computational Linguistics.
- Heinz Wimmer and Josef Perner. 1983. Beliefs about beliefs: Representation and constraining function of wrong beliefs in young children’s understanding of deception. *Cognition*, 13(1):103–128.
- Yufan Wu, Yinghui He, Yilin Jia, Rada Mihalcea, Yulong Chen, and Naihao Deng. 2023. [Hi-ToM: A benchmark for evaluating higher-order theory of mind reasoning in large language models](#). In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 10691–10706, Singapore. Association for Computational Linguistics.
- Hainiu Xu, Runcong Zhao, Lixing Zhu, Jinhua Du, and Yulan He. 2024. [OpenToM: A comprehensive benchmark for evaluating theory-of-mind reasoning capabilities of large language models](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 8593–8623, Bangkok, Thailand. Association for Computational Linguistics.

An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxi Yang, Jingren Zhou, Junyang Lin, Kai Dang, Keming Lu, Keqin Bao, Kexin Yang, Le Yu, Mei Li, Mingfeng Xue, Pei Zhang, Qin Zhu, Rui Men, Runji Lin, Tianhao Li, Tingyu Xia, Xingzhang Ren, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yu Wan, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zihan Qiu. 2024. Qwen2.5 technical report. *arXiv preprint arXiv:2412.15115*.

Lance Ying, Katherine M Collins, Megan Wei, Cedegao E Zhang, Tan Zhi-Xuan, Adrian Weller, Joshua B Tenenbaum, and Lionel Wong. 2023. The neuro-symbolic inverse planning engine (nipe): Modeling probabilistic social inferences from linguistic inputs. In *First Workshop on Theory of Mind in Communicating Agents*.

Li Zhang, Hainiu Xu, Abhinav Kommula, Chris Callison-Burch, and Niket Tandon. 2024. [OpenPI2.0: An improved dataset for entity tracking in texts](#). In *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 166–178, St. Julian’s, Malta. Association for Computational Linguistics.

Li Zhang, Hainiu Xu, Yue Yang, Shuyan Zhou, Weiqiu You, Manni Arora, and Chris Callison-Burch. 2023. [Causal reasoning of entities and events in procedural texts](#). In *Findings of the Association for Computational Linguistics: EACL 2023*, pages 415–431, Dubrovnik, Croatia. Association for Computational Linguistics.

Yaowei Zheng, Richong Zhang, Junhao Zhang, Yanhan Ye, and Zheyang Luo. 2024. [LlamaFactory: Unified efficient fine-tuning of 100+ language models](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 3: System Demonstrations)*, pages 400–410, Bangkok, Thailand. Association for Computational Linguistics.

Yuxiang Zhou, Hainiu Xu, Desmond C Ong, Petr Slovak, and Yulan He. 2025. Modeling subjectivity in cognitive appraisal with language models. *arXiv preprint arXiv:2503.11381*.

A Data Examples

See below for examples of event sequences/dialogues and ToM questions from ToMi (Box A), HiToM (Box A), and FANToM (Box A) datasets.

ToMi Example

Benjamin entered the crawlspace.
Abigail entered the crawlspace.
Emily entered the crawlspace.
The t-shirt is in the cupboard.
The cupboard is in the crawlspace.
Abigail exited the crawlspace.
Emily moved the t-shirt to the basket.
The basket is in the crawlspace.
Benjamin hates the coat.
Emily exited the crawlspace.
Abigail entered the crawlspace.

<0th-Order (Factual) Question>

Where is the t-shirt in the begining?

<1st-Order ToM Question>

Where will Abigail search for the t-shirt?

<2nd-Order ToM Question>

Where does Benjamin think that Abigail search for the t-shirt?

HiToM Example

William, Lily, Aiden, Emma and Isla entered the porch.
The melon is in the green bathtub.
William made no movements and stayed in the porch for 1 minute.
William exited the porch.
Lily made no movements and stayed in the porch for 1 minute.
Lily exited the porch.
Aiden moved the melon to the blue pantry.
Lily likes the green bucket.
Aiden exited the porch.
Emma moved the melon to the green bucket.
Emma exited the porch.
Isla moved the melon to the green bathtub.
Isla exited the porch.
William, Lily, Aiden, Emma and Isla entered the waiting room. Emma, Isla and Aiden entered the basement.
The watermelon is in the blue suitcase.
Emma made no movements and stayed in the basement for 1 minute. Emma exited

the basement.

Isla moved the watermelon to the blue bottle.

Isla exited the basement.

Aiden made no movements and stayed in the basement for 1 minute.

Aiden exited the basement.

Emma, Isla and Aiden entered the waiting room.

Aiden likes the red pantry.

William, Isla, Aiden and Lily entered the front yard.

The beans is in the red bottle.

William made no movements and stayed in the front yard for 1 minute.

Isla likes the green cupboard.

William exited the front yard.

Isla moved the beans to the green bottle.

Isla exited the front yard.

Aiden made no movements and stayed in the front yard for 1 minute.

Aiden exited the front yard.

Lily made no movements and stayed in the front yard for 1 minute.

Lily exited the front yard.

William, Isla, Aiden and Lily entered the waiting room.

<0th-Order (Factual) Question>

Where is the melon really?

<1st-Order ToM Question>

Where does William really think the melon is?

<2nd-Order ToM Question>

Where does Lily think William thinks the melon is?

<3rd-Order ToM Question>

Where does Emma think Lily thinks William thinks the melon is?

<4th-Order ToM Question>

Where does Isla think Emma thinks Lily thinks William thinks the melon is?

FANToM Example

Armani: Hi Troy, it's nice to meet you.
What's been your experience in maintaining

good mental health while in a relationship?

Troy: Hey Armani! I've found that the most important thing for me is understanding that I need to take care of my own mental health first. I look at it like the whole oxygen mask in an airplane situation- you have to secure your own before helping someone else.

Armani: That's an interesting perspective, it's all about maintaining individual wellness before being able to fully contribute to a relationship. I've always subscribed to the idea of communication being an integral part of it too. Being open about my mental health issues with my partner has always helped to build understanding.

Troy: I definitely see the merit in that too. It can be hard to open up about these things sometimes, especially if the other person doesn't fully understand.

Armani: Absolutely, I think it's important for both partners to constantly educate themselves on each other's mental health issues. It not only encourages empathy but also helps in mitigating unnecessary tensions.

Troy: You're right. From my experience, I've found that maintaining a healthy work-life balance is also essential. Stress from work can really take a toll on my mental health, and it's hard to keep that stress from affecting my relationships.

Armani: I completely agree, Troy. Ignoring the effects of work-related stress on our mental wellbeing can have dire consequences on our relationships. Just as we wouldn't like to bring home a flu virus, we shouldn't be infecting our home with stress either. It's all about creating boundaries.

Troy: Absolutely. It's great to find someone else who understands the importance of maintaining mental health while in a relationship. It's definitely a balance, but it's worth it in the long run.

Armani: Couldn't agree more Troy. I'm glad we could have this open conversation about such an important topic. The more we talk, the more we can break the stigma surrounding mental health issues.

Cynthia: Hello Troy, Armani. The two of you have been engaged in quite a meaningful conversation, it seems.

Armani: Hi Cynthia, yes indeed. We've been discussing the importance of good mental health maintenance in a relationship.

Cynthia: Such a crucial topic! In my experience, clear actionable boundaries have played a very big role in mental wellness. Making sure me and my partner are on the same page about our needs and wants can genuinely de-stress the environment.

Troy: Couldn't agree more, Cynthia. The clear establishment of boundaries helps a lot in maintaining a harmonious balance.

Armani: Definitely, Cynthia. It's such a simple concept yet so often overlooked. People sometimes shy away from setting boundaries, afraid it might upset the other person. But it's needed for mutual respect and understanding.

Cynthia: Yes, Armani. And it's these boundaries that create stronger communication channels. So, when I'm experiencing a difficult time with my mental health, I find it easier to express to my partner.

Troy: That really hits home, Cynthia. It's incredibly liberating to be able to express ourselves without fear of judgement.

Armani: Absolutely, Troy! And the thing is, this whole conversation really highlights the importance of communication. Everything, from understanding personal mental health, setting boundaries, to dealing with stress, involves communicating effectively.

Cynthia: Here's to hoping that more people can learn and implement these practices in their relationships. Good mental health is so important, and talking about it openly like this is a great step in the right direction.

<0th-Order (Factual) Question>

What did Armani and Troy discuss as a preventative measure against work-related stress affecting their relationships?

<1st-Order ToM Question>

What does Cynthia believe are the necessary steps suggested by Armani and Troy for dealing with mental health issues in a relationship?

<2nd-Order ToM Question>

What does Armani believe about Cynthia's belief regarding the necessary steps suggested by Armani and Troy for dealing with mental health issues in a relationship? Where does Benjamin think that Abigail search for the t-shirt?

B Training of Enigma^T

We train Enigma^T based on Llama3.1-8B and Llama3.3-70B using OpenPI2.0 dataset. The training of Llama3.1-8B is done using Low-Rank Adaptation (Hu et al., 2022) and the training of Llama3.3-70B is done using Quantized Low-Rank Adaptation (Dettmers et al., 2024). The rank of the decomposed matrix is set to LoRA-Rank = 32. We use a batch size of 64 and a gradient accumulation of 2 steps, leading to an effective batch size of 128. The learning rate is initially set to 5e-5 and adjusted using a Cosine Annealing with a warm-up ratio of 0.01. Each model is trained for 3 epoches. The training is done on 2 NVIDIA A100^{80GB} GPUs using LlamaFactory (Zheng et al., 2024).

We modified the formatting convention used in OpenPI (Tandon et al., 2020). In the original formatting, the output is formatted as

[Attribute] of [Entity] is [Previous State]
before and [Current State] afterwards.

In this formulation, the model is tasked to predict both the previous state as well as the current state after the event has taken place. We modify this formulation by removing the previous state

and only asking the model to predict the current state:

[Attribute] of [Entity] becomes [State]

this modification is made to alleviate model's reasoning burden and to focus on deriving the effect the the event.

The events are provided to the model in a cumulative fashion. For instance, when generating entity state knowledge for the i -th event, the model is provided with the events from the first event to the i -th event.

C Transforming High-order ToM Question to First-order

As discussed in Hou et al. (2024), perspective-taking reduces the reasoning depth of high-order ToM questions. Therefore, we can transform high-order ToM questions into first-order questions after perspective-taking. We first justify how EnigmaToM enables reduction in ToM reasoning order and then provide a description of the question transformation process.

C.1 ToM Order Reduction

ToM order reduction is possible thanks to the IM mechanism, which carries out the high-order, multi-hop ToM reasoning process symbolically using spatial scene graphs. We use an example to elicit the necessity of IM mechanism.

Consider the following event sequence:

- William, Lily, Aiden, Emma and Isla entered the porch.
- The melon is in the green bathtub.
- Aiden moved the melon to the blue pantry.
- Lily likes the green bucket.
- Aiden exited the porch.
- Lily made no movements and stayed in the porch for 1 minute.
- Lily exited the porch.
- Emma moved the melon to the green bucket.
- Emma exited the porch.
- Isla moved the melon to the green bathtub.
- Isla exited the porch.
- William moved the melon to the red bucket.
- William exited the porch.
- William, Lily, Aiden, Emma, and Isla entered the waiting room.

Consider the following third-order ToM question:

Where does Emma think Lily thinks William thinks the melon is?

Answer: blue pantry

The above question can be reduced to the following first-order ToM question:

Where does William think the melon is?

Answer: red bucket

Notice that the answers to the third-order and first-order ToM question are different. This is because answering the original question requires three reasoning hops:

Step 1 Infer Emma’s belief about the environment.

Step 2 Given Emma’s belief, infer Lily’s belief about the environment.

Step 3 Given Emma’s belief of Lily’s belief, infer William’s belief about the environment.

In contrast, a first-order ToM question only requires a single reasoning hop:

Step 1 Infer William’s belief about the environment.

To alleviate the multi-hop reasoning burden from LLMs, we utilize the spatial scene graphs with the IM mechanism. Each character’s spatial scene graph can be treated as a reasoning hop. The aggregates of the scene graphs through IM effectively performs multi-hop reasoning symbolically. This allows us to transform a high-order ToM question into a simpler one while preserving accuracy.

For example, after IM, irrelevant events are masked (striked through), leaving only those representing Emma’s belief of Lily’s belief of William’s belief:

- William, Lily, Aiden, Emma and Isla entered the porch.
- The melon is in the green bathtub.
- Aiden moved the melon to the blue pantry.
- Lily likes the green bucket.
- Aiden exited the porch.
- Lily made no movements and stayed in the porch for 1 minute.
- Lily exited the porch.
- ~~Emma moved the melon to the green bucket.~~
- ~~Emma exited the porch.~~
- ~~Isla moved the melon to the green bathtub.~~
- ~~Isla exited the porch.~~
- ~~William moved the melon to the red bucket.~~

- ~~William exited the porch.~~
- ~~William, Lily, Aiden, Emma, and Isla entered the waiting room.~~

By reducing the original third-order ToM question into that of first-order, we are able to obtain the correct answer:

Original Third-order ToM Question

Where does Emma think Lily thinks William thinks the melon is?

Reduced First-order ToM Question

Where does William think the melon is?

Answer: blue pantry

C.2 Question Transformation

We transform the original high-order ToM questions to first-order by prompting LLMs with 5-shot demonstrations for each dataset. The transformed questions are used for QA with EnigmaToM as well as TimeToM. See below for an illustrative example of transforming a third-order ToM question into a first-order question:

3rd-Order \rightarrow 1st-Order ToM Question:

Where does Emma think Lily thinks William thinks the melon is?

\Rightarrow

Where does William think the melon is?

D Prompts

The following prompt is used to infer key entities and attributes from an event sequence and ToM question as described in §3.1 and Equation 2.

Infer Key Entities and Attributes

<Events>
{ {indexed narrative} }

<Questions>
{ {question list} }

Based on the list of <questions>, extract at most five entities and their attributes that are needed for answering the questions. Note that one entity could corresponds to multiple attributes. List the entities and their attributes. For instance, if the "location of tie", "placement of tie", and "color of crate" are important for answering the questions, the response should be formatted as follows:

<entities>
- location of tie

- placement of tie
 - color of crate
 </entities>

You must include at least one entity that is not a person. Only extract entities directly mentioned in the questions, do not make any further inference. Do not include any entities that indicate a time or point in time. First briefly reason about the content of the events and the questions and then provide a comprehensive list of at most five entities and their attributes with the following format:

<entities>
 - attribute of entity
 - attribute of entity
 ...
 </entities>

We use the following prompt to identify locations appeared in the event sequence. The location information is used as anchor when constructing the spatial scene graphs to map the locations produced by Enigma to a common location space.

Extract Locations from Events

<Events>
 {{indexed narrative}}

What are the rooms mentioned in these events? List all the rooms in the following format:

- Room1
 - Room2
 ...

Please exclude entities in which people cannot enter. Each narrative must contain at least one room and your answer must include at least one room. Provide your answer as bullet points without any explanation.

The following prompt is used to acquire the entity state knowledge by prompting LLMs (Enigma^P).

Generate Entity State Knowledge

<Events>
 {{indexed narrative}}

<Entity-of-Interest>
 {{eoi list}}

Given the list of events and entities-of-interest, track the state of the attribute of entities throughout the events. Generate state of each attribute of entities as a list in the following format:

- [Event Index]: [Entity Attribute] becomes [State]
 - [Event Index]: [Entity Attribute] becomes [State]
 ...

Determine spatial information according to the following instructions:

- All of the location changes, if exist, are explicitly stated in the events.
 - If the event does not state that a character left a room, assume that the character remains in the previous location.

Generate the answers exactly as instructed without any explanation or note. Only generate the event indices where there is a entity state change, omit other events.

E Visualization of Complexity

In §3.4, we provide a general analysis of the complexity of SymbolicToM and EnigmaToM with respect to the number of belief graphs need to be constructed. Here, we provide a visualization to better demonstrate the efficiency of the IM-based perspective-taking method of EnigmaToM. Figure A1 shows the number of belief graphs need to be constructed for SymbolicToM and EnigmaToM with respect to the number of characters (Part a) or the order of ToM reasoning (Part b). When plotting against the number of character, we fix the ToM reasoning order to be $k = 2$, which requires a minimum of 2 characters. When plotting against the Tom reasoning order, we fix the number of characters to be $m = 5$, which supports up to 5^{th} -order acyclic ToM reasoning.

		Qwen2.5-7B	Qwen2.5-14B	Qwen2.5-32B	Qwen2.5-72B ^{4bit}
ToMi	VNL	0.722 \pm 0.045	0.747 \pm 0.053	0.720 \pm 0.036	0.717 \pm 0.034
	CoT	<u>0.724\pm0.026</u>	<u>0.773\pm0.030</u>	0.769 \pm 0.029	0.767 \pm 0.033
	SToM	0.642 \pm 0.022	0.686 \pm 0.042	0.721 \pm 0.029	0.749 \pm 0.020
	TTToM	0.567 \pm 0.024	0.699 \pm 0.027	0.758 \pm 0.019	0.865\pm0.018
	DWM	0.686 \pm 0.023	0.732 \pm 0.033	0.723 \pm 0.053	0.762 \pm 0.051
	PToM	0.720 \pm 0.038	0.743 \pm 0.050	0.783 \pm 0.021	0.809 \pm 0.033
	Enigma^P	0.706 \pm 0.044	0.736 \pm 0.019	0.828\pm0.009	<u>0.839\pm0.014</u>
	Enigma^T	0.825\pm0.030	0.826\pm0.027	<u>0.809\pm0.027</u>	0.837 \pm 0.024
HiToM	VNL	0.378 \pm 0.013	0.389 \pm 0.014	0.524 \pm 0.007	0.456 \pm 0.012
	CoT	0.441 \pm 0.007	0.330 \pm 0.020	0.523 \pm 0.009	0.481 \pm 0.011
	SToM	0.402 \pm 0.009	0.389 \pm 0.025	0.520 \pm 0.007	0.536 \pm 0.018
	TTToM	0.316 \pm 0.010	0.397 \pm 0.017	0.473 \pm 0.011	0.415 \pm 0.013
	DWM	0.444 \pm 0.020	0.436 \pm 0.019	<u>0.566\pm0.011</u>	<u>0.560\pm0.009</u>
	PToM	0.393 \pm 0.018	<u>0.446\pm0.017</u>	0.500 \pm 0.011	0.548 \pm 0.016
	Enigma^P	0.508\pm0.012	0.554\pm0.013	0.674\pm0.011	0.605\pm0.007
	Enigma^T	<u>0.457\pm0.005</u>	0.414 \pm 0.013	0.504 \pm 0.013	0.473 \pm 0.010
FANToM	VNL	0.400 \pm 0.015	0.522 \pm 0.004	0.496 \pm 0.022	0.532 \pm 0.025
	CoT	0.398 \pm 0.014	0.516 \pm 0.011	0.482 \pm 0.030	0.521 \pm 0.024
	SToM	0.413 \pm 0.012	0.537 \pm 0.027	0.379 \pm 0.024	0.516 \pm 0.014
	TTToM	0.252 \pm 0.020	0.476 \pm 0.019	0.419 \pm 0.018	0.409 \pm 0.026
	DWM	0.429 \pm 0.013	0.558 \pm 0.012	<u>0.519\pm0.022</u>	0.543 \pm 0.014
	PToM	0.408 \pm 0.023	0.579\pm0.014	0.542\pm0.019	<u>0.573\pm0.016</u>
	Enigma^P	<u>0.445\pm0.026</u>	0.474 \pm 0.020	0.407 \pm 0.025	0.450 \pm 0.013
	Enigma^T	0.487\pm0.018	<u>0.560\pm0.030</u>	0.405 \pm 0.026	0.574\pm0.031

Table A1: Results of the scaling experiments. This table shows the exact mean accuracy of each method on ToMi, HiToM, and FANToM with base LLMs of different sizes. The results shown in this table are used to generate Figure 4.

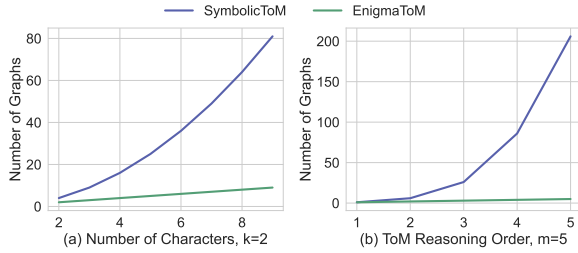


Figure A1: Visualization of the complexity of SymbolicToM and EnigmaToM with respect to the number of belief graphs need to be constructed.

F Detailed Results of High-Order ToM Reasoning with EnigmaToM

Complete Results As discussed in §5.1, we observe that EnigmaToM improves LLMs’ capability of conducting high-order ToM reasoning using the HiToM dataset. Although the ToM reasoning order of ToMi and FANToM are limited to second-order, we observe a similar trend where EnigmaToM brings considerable improvements in the case of second-order ToM reasoning. See Fig-

ure A2 for detailed results.

Ablation Study In §5.2, we analyzed the efficacy of the IM and KI components of EnigmaToM. Here, we conduct analysis on the impact of IM and KI with respect to the ToM reasoning order. Specifically, we compute the advantage of removing the IM or KI component from EnigmaToM as

$$\text{Adv} = \text{Acc}_{\text{w/o component}} - \text{Acc}_{\text{w/ component}} \quad (7)$$

where a positive advantage score means that the model performed better without the component and vice versa. Table A3, Table A4, and Table A5 shows the result of the ablation study. Aligning with our findings in §5.2, IM is indispensable to all ToM reasoning tasks, reiterating the importance of perspective-taking in ToM reasoning. On the other hand, KI is critical for the dialogue-based FANToM dataset while less effective for the event-based ToMi and HiToM datasets. This finding illustrates that Enigma’s functionality as an information compressor by compressing key information as entity state knowledge is substantial when tackling

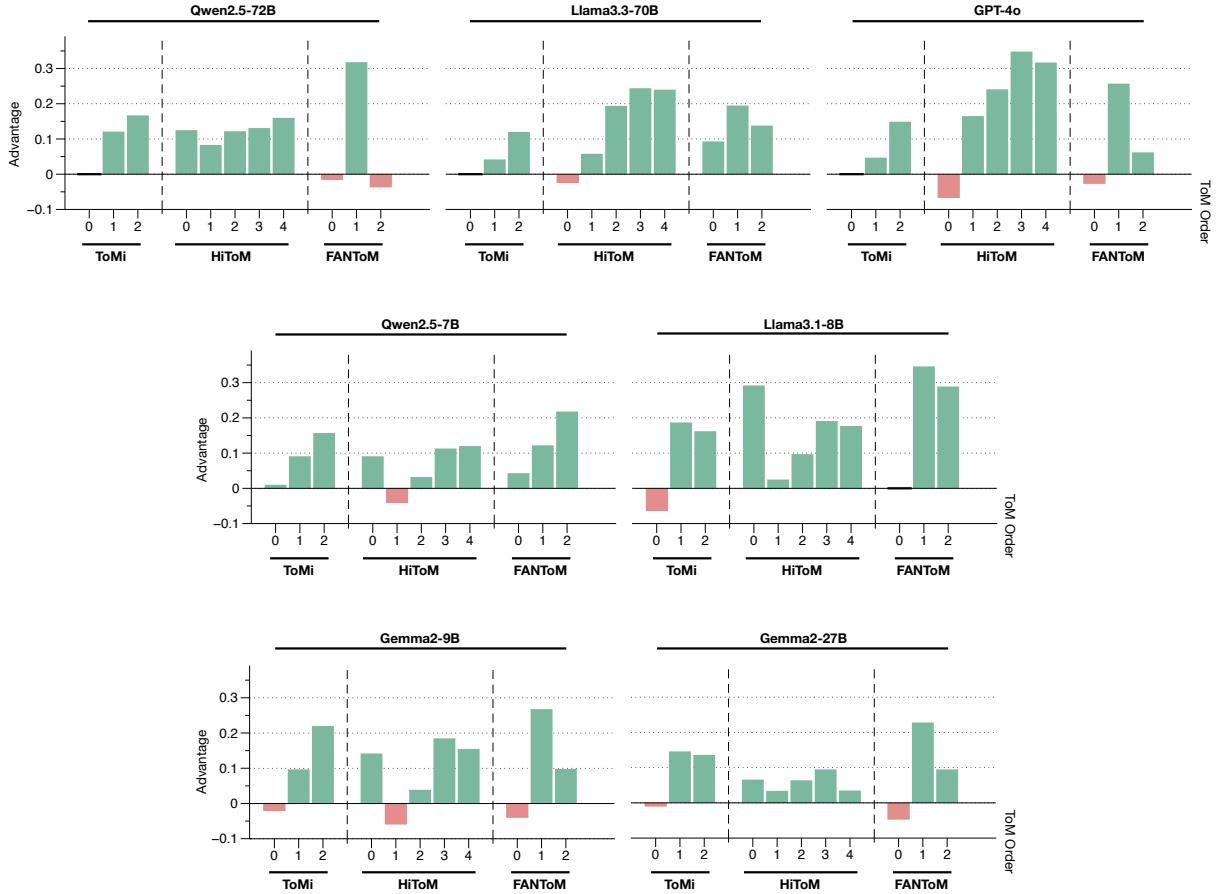


Figure A2: The complete results of the advantage of EnigmaToM with respect the ToM order on ToMi, HiToM, and FANToM datasets.

		Qwen2.5-7B	Llama3.1-8B	Gemma2-9B	Gemma2-27B	Llama3.3-70B ^{4bit}	Qwen2.5-72B ^{4bit}	GPT-4o
ToMi	Enigma ^p	0.706 \pm 0.044	0.738 \pm 0.056	0.865 \pm 0.031	0.833 \pm 0.018	0.828 \pm 0.012	0.839 \pm 0.014	0.847 \pm 0.030
	Enigma ^r	0.825 \pm 0.030	0.796 \pm 0.023	0.814 \pm 0.020	0.804 \pm 0.050	0.787 \pm 0.024	0.837 \pm 0.024	0.795 \pm 0.036
	w/o KI	0.805 \pm 0.018	0.799 \pm 0.022	0.842 \pm 0.023	0.825 \pm 0.031	0.834 \pm 0.067	0.845 \pm 0.026	0.811 \pm 0.028
	w/o IM	0.623 \pm 0.021	0.597 \pm 0.040	0.686 \pm 0.019	0.658 \pm 0.027	0.693 \pm 0.014	0.655 \pm 0.039	0.674 \pm 0.002
HiToM	Enigma ^p	0.508 \pm 0.012	0.477 \pm 0.005	0.555 \pm 0.010	0.576 \pm 0.004	0.696 \pm 0.007	0.605 \pm 0.007	0.733 \pm 0.017
	Enigma ^r	0.457 \pm 0.005	0.431 \pm 0.010	0.446 \pm 0.008	0.478 \pm 0.004	0.518 \pm 0.011	0.473 \pm 0.010	0.626 \pm 0.020
	w/o KI	0.511 \pm 0.007	0.463 \pm 0.006	0.571 \pm 0.012	0.614 \pm 0.007	0.726 \pm 0.004	0.632 \pm 0.003	0.751 \pm 0.004
	w/o IM	0.380 \pm 0.004	0.341 \pm 0.015	0.406 \pm 0.006	0.408 \pm 0.013	0.460 \pm 0.013	0.423 \pm 0.008	0.442 \pm 0.006
FANToM	Enigma ^p	0.445 \pm 0.026	0.442 \pm 0.018	0.439 \pm 0.023	0.462 \pm 0.014	0.515 \pm 0.020	0.450 \pm 0.013	0.531 \pm 0.015
	Enigma ^r	0.487 \pm 0.018	0.545 \pm 0.036	0.530 \pm 0.012	0.582 \pm 0.028	0.610 \pm 0.021	0.574 \pm 0.031	0.553 \pm 0.011
	w/o KI	0.487 \pm 0.022	0.544 \pm 0.033	0.530 \pm 0.017	0.579 \pm 0.021	0.607 \pm 0.018	0.542 \pm 0.036	0.539 \pm 0.012
	w/o IM	0.436 \pm 0.028	0.448 \pm 0.024	0.426 \pm 0.011	0.478 \pm 0.008	0.500 \pm 0.021	0.477 \pm 0.017	0.470 \pm 0.013

Table A2: Results of the ablation study of EnigmaToM on ToMi, HiToM, and FANToM datasets. "w/o KI" indicates without *entity state knowledge injection*. "w/o IM" denotes without *perspective-taking via iterative masking*.

ToM reasoning with sparse information (as in daily dialogue). Further, removing KI leads to larger performance degradation in LLMs of smaller size, further highlighting that small LLMs are less capable of dealing with reporting bias and prefer to have explicit state information regarding key entities in events.

G Detailed Results of Scaling Experiment

The exact mean accuracies of EnigmaToM and baseline methods on ToMi, HiToM, and FANToM datasets using LLMs of varying sizes are shown in Table A1.

	ToM Order	Qwen2.5-7B	Llama3.1-8B	Gemma2-9B	Gemma2-27B	Llama3.3-70B ^{4bit}	Qwen2.5-72B ^{4bit}	GPT-4o
w/o KI	<i>Zeroth</i>	+0.008	+0.016	+0.008	−0.000	+0.008	−0.000	−0.000
	<i>First</i>	−0.019	−0.019	−0.000	+0.005	+0.014	−0.000	−0.060
	<i>Second</i>	−0.004	−0.004	−0.039	−0.022	+0.017	−0.025	−0.030
w/o IM	<i>Zeroth</i>	+0.008	−0.024	−0.000	−0.000	+0.008	−0.000	+0.000
	<i>First</i>	−0.127	−0.226	−0.161	−0.283	−0.273	−0.025	−0.164
	<i>Second</i>	−0.025	−0.124	−0.253	−0.224	−0.206	−0.218	−0.269

Table A3: Results of the ablation study of EnigmaToM on ToMi with respect to ToM order. We display the advantage scores computed using Equation 7. "w/o KI" indicates without *entity state knowledge injection*. "w/o IM" denotes without *perspective-taking via iterative masking*.

	ToM Order	Qwen2.5-7B	Llama3.1-8B	Gemma2-9B	Gemma2-27B	Llama3.3-70B ^{4bit}	Qwen2.5-72B ^{4bit}	GPT-4o
w/o KI	<i>Zeroth</i>	−0.017	+0.075	+0.058	+0.005	+0.058	+0.042	+0.075
	<i>First</i>	−0.008	−0.017	+0.008	−0.017	+0.016	−0.017	+0.023
	<i>Second</i>	+0.045	+0.039	+0.026	+0.013	+0.019	+0.006	+0.002
	<i>Third</i>	−0.018	+0.018	+0.018	+0.065	+0.054	+0.030	0.009
	<i>Fourth</i>	+0.040	+0.040	+0.017	+0.046	+0.040	+0.029	+0.008
w/o IM	<i>Zeroth</i>	−0.000	−0.000	−0.000	−0.000	−0.000	−0.008	+0.016
	<i>First</i>	−0.019	−0.100	−0.092	−0.150	−0.259	−0.175	−0.229
	<i>Second</i>	−0.142	−0.084	−0.207	−0.193	−0.323	−0.271	−0.403
	<i>Third</i>	−0.214	−0.143	−0.238	−0.232	−0.268	−0.261	−0.365
	<i>Fourth</i>	−0.125	−0.103	−0.211	−0.205	−0.257	−0.177	−0.316

Table A4: Results of the ablation study of EnigmaToM on HiToM with respect to ToM order. We display the advantage scores computed using Equation 7. "w/o KI" indicates without *entity state knowledge injection*. "w/o IM" denotes without *perspective-taking via iterative masking*.

	ToM Order	Qwen2.5-7B	Llama3.1-8B	Gemma2-9B	Gemma2-27B	Llama3.3-70B ^{4bit}	Qwen2.5-72B ^{4bit}	GPT-4o
w/o KI	<i>Zeroth</i>	−0.032	−0.054	−0.031	−0.044	−0.053	−0.058	−0.010
	<i>First</i>	−0.308	−0.255	−0.265	−0.392	−0.237	−0.205	+0.005
	<i>Second</i>	−0.162	−0.212	−0.113	−0.230	−0.128	−0.119	−0.005
w/o IM	<i>Zeroth</i>	−0.012	−0.003	−0.025	−0.027	−0.041	−0.035	−0.025
	<i>First</i>	−0.335	−0.367	−0.343	−0.370	−0.327	−0.418	−0.348
	<i>Second</i>	−0.178	−0.215	−0.126	−0.218	−0.144	−0.073	+0.000

Table A5: Results of the ablation study of EnigmaToM on FANToM with respect to ToM order. We display the advantage scores computed using Equation 7. "w/o KI" indicates without *entity state knowledge injection*. "w/o IM" denotes without *perspective-taking via iterative masking*.

H Examples of Knowledge Injection

We demonstrate how KI is accomplished with an event sequence from ToMi. The important entities and attributes are generated using Llama3.3-70B^{4bit} and the entity state knowledge is generated by Enigma^T_{8B}.

Events from ToMi

- 1: Benjamin entered the crawlspace.
- 2: Abigail entered the crawlspace.
- 3: Emily entered the crawlspace.
- 4: The t-shirt is in the cupboard.
- 5: The cupboard is in the crawlspace.
- 6: Abigail exited the crawlspace.

- 7: Emily moved the t-shirt to the basket.
- 8: The basket is in the crawlspace.
- 9: Benjamin hates the coat.
- 10: Emily exited the crawlspace.
- 11: Abigail entered the crawlspace.

Entity State Generated by Enigma^T_{8B}

Event1:

- Location of Benjamin is in the crawlspace.

Event2:

- Location of Abigail is in the crawlspace.

Event3:

- Location of Emily is in the crawlspace.

Event4:

- Location of T-shirt is in cupboard.

Event5:

- Location of cupboard is in crawlspace.

Event6:

- Location of Abigail is outside the crawlspace.

Event7:

- Location of T-shirt is in basket.
- Content of basket is T-shirt.
- Content of cupboard is empty.

Event8:

- Location of basket is in the crawlspace.

Event9:

None

Event10:

- Location of Emily is outside the crawlspace.

Event11:

- Location of Abigail is in the crawlspace.

Augmented ToMi Events

- 1: Benjamin entered the crawlspace.
- 2: Abigail entered the crawlspace.
- 3: Emily entered the crawlspace.
- 4: The t-shirt is in the cupboard.
- Location of T-shirt is in the cupboard.
- 5: The cupboard is in the crawlspace.
- Location of cupboard is in the crawlspace.
- 6: Abigail exited the crawlspace.
- 7: Emily moved the t-shirt to the basket.
- Location of T-shirt is in basket. - Content of basket is T-shirt.
- Content of cupboard is empty.

- 8: The basket is in the crawlspace.

- Location of basket is in the crawlspace.

- 9: Benjamin hates the coat.

- 10: Emily exited the crawlspace.

- 11: Abigail entered the crawlspace.

I Examples of Entity State Knowledge Generated by Enigma^T_{8B} and Enigma^T_{70B}

Entity states generated by Enigma^T_{70B} is more informative compared to Enigma^T_{8B} . From the annotated entity states generated by Enigma^T_{8B} and Enigma^T_{70B} in both ToMi and FANToM, we observe that Enigma^T_{70B} is much more eloquent compared to its counterparts. Such eloquence makes the entity state knowledge more informative in ToMi:

Entity State for ToMi**Example Event 1**

Amelia moved the belt to the pantry.

 Enigma^T_{8B}

Amelia's knowledge of the belt's location is known

 Enigma^T_{70B}

Amelia's knowledge of the belt's location is that the belt is in the pantry

Example Event 2

Lucas exited the living room.

 Enigma^T_{8B}

Location of Lucas as known to Charlotte is absent

 Enigma^T_{70B}

Location of Lucas as known to Charlotte is outside the living room

However, in FANToM, the increased eloquence of Enigma^T_{70B} leads to more hallucination, resulting in more erroneous information and reducing its functionality as an information compressor (contents that can be entailed from the utterance is marked in ● green whereas the irrelevant contents are marked in ● red):

Entity State for FANToM

Example Event 1

Breanna: Long-distance, yes, but not for as long as you both. It was a few months, and things didn't work out for us. But talking about relationships, it brings back memories of my first crush.

Enigma^T_{8B}

Breanna's knowledge of longdistance relationship methods is known.

Enigma^T_{70B}

Breanna's knowledge of longdistance relationship methods is limited to her past personal experience of a few months, with an unsuccessful outcome, and she has been exposed to Nickolas and Raelynn's methods, including the importance of trust, understanding, commitment, and regular communication, such as scheduling regular talks and quick calls.

Example Event 2

Shayla: It was. But I made it through. It was a journey that brought me face-to-face with my own strength. It made me realize the importance of mental health and inspired me to help others who are going through something similar.

Enigma^T_{8B}

Significance of nurturing mental wellbeing is important.

Enigma^T_{70B}

Significance of nurturing mental wellbeing is heightened, as it is recognized as a crucial aspect of overall health, particularly in the aftermath of traumatic experiences like abusive relationships, and is seen as a key factor in building resilience and facilitating recovery.

line was provided to the annotators containing a description of the annotation task and examples that cover all possible label combinations.

The data provided to the annotators are identical to the ones used in EnigmaToM. The set of key entities (E) and their attributes (A) are generated using Llama3.3-70B^{4bit}. Hence, the relevance score reflects Llama3.3-70B^{4bit}'s capability of generating entities and attributes that are important to answering the corresponding ToM questions. Although it is not shown in the example, we provide the complete set of questions and cumulative events for each data point. In other words, given an event sequence $\{\epsilon\}_{i=1,n}$ and a set of ToM questions Q , we provided $\epsilon_{1:k}$ and Q when annotating the k^{th} event and its entity states.

The detailed annotation guideline is shown in the following pages.

J Evaluation of Enigma

As discussed in §5.3 and Table 5, we manually labeled 300 instances of entity state information from ToMi and FANToM datasets. The annotation was carried out by three graduate students (100 annotations each) majoring in computer science at a prestigious university. Since the annotation was conducted for entity state information, which only contains commonsense knowledge, no training was provided to the annotators. An annotation guide-

Annotation Guideline for Evaluation of Neural Knowledge Base

Task Formulation

In this task, you will be responsible for evaluating the quality of entity state information generated by two Neural Knowledge Base. You will be given a set of **Questions** in the question column, an **event** or **utterance** in the event column and a **entity state information** generated from a neural knowledge bank in the XXX_states column.

You will be annotating the following two aspects of the generated entity state information:

1. **Relevance**: whether the generated entity state knowledge is helpful with answering the questions.
2. **Accuracy**: whether the entity state can be entailed from the given event/utterance.

For both tasks, choose **Yes** if you believe that the entity state is relevant/accurate and choose **No** otherwise.

Note that all of the entity state follows a semi-structured formulation:

- **[Attribute]** of **[Entity]** is **[Current State]**

In the **Relevance** annotation task, you only need to pay attention to the **[Attribute]** of **[Entity]**.

In the **Accuracy** annotation task, you need to consider the complete entity state description.

Examples

Question: Where does Abigail think that Benjamin searches for the pumpkin?

Event: Abigail moved the pumpkin to the bathtub.

State: *location of the pumpkin is in bathtub*

In this example, we see that the Location of the Pumpkin is important in answering the question. **Therefore, we mark the relevance as Yes.** We see that there exists an entailment relationship between Event and State. **Therefore, we mark the accuracy as Yes.**

Question: Where does Abigail think that Benjamin searches for the pumpkin?

Event: Chloe entered the workshop.

State: *location of Chloe is in workshop*

In this example, we see that the Location of the Chloe is not related to the question. **Therefore, we mark the relevance as No.**

We see that there exists an entailment relationship between Event and State. **Therefore, we mark the accuracy as Yes.**

Question: Where does Abigail think that Benjamin searches for the pumpkin?

Event: Abigail moved the pumpkin to the bathtub.

State: *location of Abigail is in bathtub*

In this example, we see that the Location of the Abigail is important in answering the question. **Therefore, we mark the relevance as Yes.**

However, it is unlikely that Abigail needs to enter the bathtub to place the pumpkin. **Therefore, there does not exist an entailment relationship between Event and State. Therefore, we mark the accuracy as No.**

Question: Where does Abigail think that Benjamin searches for the pumpkin?

Event: Chloe entered the workshop.

State: *location of Chloe is in kitchen*

In this example, we see that the Location of the Chloe is not related to the question. **Therefore, we mark the relevance as No.**

We see that there does not exist an entailment relationship between Event and State. **Therefore, we mark the accuracy as No.**
