# Mining Logical Event Schemas From Pre-Trained Language Models

# Lane Lawley and Lenhart Schubert

University of Rochester Department of Computer Science

{llawley, schubert}@cs.rochester.edu

#### **Abstract**

We present NESL (the Neuro-Episodic Schema Learner), an event schema learning system that combines large language models, FrameNet parsing, a powerful logical representation of language, and a set of simple behavioral schemas meant to bootstrap the learning process. In lieu of a pre-made corpus of stories, our dataset is a continuous feed of "situation samples" from a pre-trained language model, which are then parsed into FrameNet frames, mapped into simple behavioral schemas, and combined and generalized into complex, hierarchical schemas for a variety of everyday scenarios. We show that careful sampling from the language model can help emphasize stereotypical properties of situations and de-emphasize irrelevant details, and that the resulting schemas specify situations more comprehensively than those learned by other systems.

## 1 Introduction

Work on the task of event schema acquisition is largely separable into two main camps: the *symbolic* camp, with its structurally rich but infamously brittle representations (Lebowitz, 1980; Norvig, 1987; Mooney, 1990); and the *statistical* camp, which utilizes complex models and vast amounts of data to produce large numbers of conceptually varied event schemas at the cost of representational richness and control over what schemas are learned (Chambers and Jurafsky, 2008; Pichotta and Mooney, 2016; Wanzare et al., 2017).

In an attempt to bridge these camps together, we introduce the Neuro-Episodic Schema Learner (NESL), a composite pipeline model bringing together (1) a large, pre-trained language model; (2) word vector embedding techniques; (3) a neural FrameNet parsing and information extraction system; (4) a formal logical semantic representation of English; and (5) a hierarchical event schema framework with extraordinary expressive power.

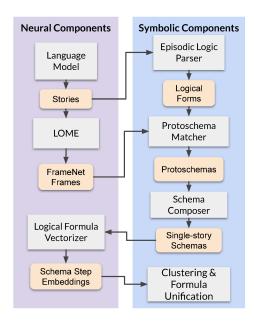


Figure 1: A diagram of the NESL schema learning pipeline. Peach-colored rounded rectangles represent data, while gray rectangles represent NESL's neural and symbolic components.

Besides the enriched schema framework, a key contribution of NESL is the idea of *latent schema sampling* (LSS), in which a pre-trained language model is induced, via prompt-engineering, into acting as a distribution over stories, parameterized by an implicit *latent schema* coarsely established by the prompt. By finding commonalities between multiple samples from this distribution and discarding infrequent details, NESL attempts to generate more accurate, less noisy schemas. In addition to eliminating NESL's need for its own training corpus, LSS allows NESL to generate schemas for user-provided situation descriptions on demand, greatly increasing the control over what sorts of event schemas are generated.

The remainder of this paper is organized into a description of our chosen semantic representation and event schema framework (Section 2); a description of each of NESL's components (Section 3); and a description of future evaluation work we intend to complete in the near future (Section 4).

## 2 Schema Model

## 2.1 Episodic Logic

We use Episodic Logic (EL) (Hwang and Schubert, 1993) as our semantic representation of stories and of schema formulas. EL is an intensional semantic representation that is both well-suited to inference tasks and structurally similar in its surface form to English. A unique feature of EL is its treatment of episodes (events, situations) as first-class individuals, which may be characterized by arbitrary formulas. For a formula  $\phi$  and an episode denoted by E,  $(\phi ** E)$  expresses that E is characterized by  $\phi$ ; i.e., E is an "episode of"  $\phi$  occurring. In the EL schema in Figure 2, for example, each formula in the STEPS section implicitly characterizes an episode (not displayed). The "header" formula at the top of that figure also characterizes an episode variable, ?E. EL can temporally relate episodes using predicates implementing the Allen Interval Algebra (Allen, 1983), allowing for complex and hierarchical situation descriptions.

#### 2.2 EL Schemas

Past approaches to statistical schema learning have largely represented schemas as sequences of lexical event tuples (Chambers and Jurafsky, 2008; Pichotta and Mooney, 2016). Seeking a richer representation, we adopt the rich, EL-based schema framework presented by Lawley et al. (2021), henceforth referred to in this paper as *EL schemas*. EL schemas are section-based: the main two sections, STEPS and ROLES, enumerate the temporal events ("steps") that make up the schema, and the type and relational constraints on the schema's participants, respectively (see Figure 2).

Designed as a suitable representation of humancentric events, EL schemas can also specify preconditions, postconditions, arbitrary temporal relationships between steps, and the *goals* of individual participants in the schema. All schema participants are represented as typed variables, all sharing a scope within the same schema, and formulas may include any number of variables as arguments. EL schemas also allow for recursive *nesting*: a schema may be embedded as a step in another schema, and implicitly expanded to check constraints or generate inferences.

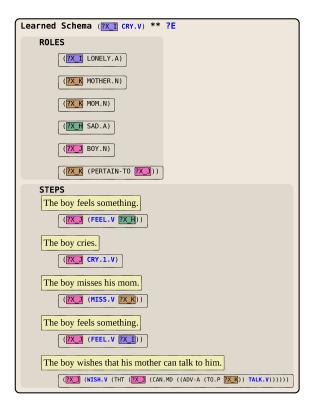


Figure 2: A schema generated by NESL from a single GPT-J 6B story sample. Note that the second step's verb predicate, CRY . 1 . V, contains a number: this identifies it as the header of a unique *protoschema* instance that was matched to the story.

For more information on the EL schema framework, see (Lawley et al., 2019) and (Lawley et al., 2021).

## 2.2.1 Protoschemas

The schema system we use here, due to Lawley et al. (2021), is designed to acquire schemas in the manner a very young child might. In his theory of the origin of intelligence in children, Jean Piaget hypothesized that event generalization, in babies, "always proceeds from the undifferentiated schema to the individual and to the general, combined and complementary" (Piaget and Cook, 1952). While we don't assert any particular theory of child development in this work, we follow the spirit of Piaget's claim and propose a system wherein complex schemas are learned from simple, universal ones.

As a way of modeling general actions a very young child would be likely to understand, we start with a handwritten corpus of several dozen "protoschemas", such as "X helps Y with action A", "X eats food F to alleviate hunger", etc.; the aim of this schema learning framework is to first match

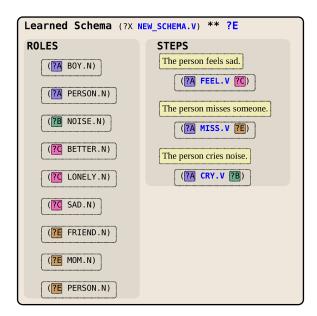


Figure 3: A schema generated by NESL combining multiple single-story schemas, such as the one shown in Figure 2. Note that some details from the single-story schema have been removed, and that variables have multiple possible types, sometimes conflicting; displaying these as certainty-weighted disjunctions is intended for a future version of NESL.

these protoschemas to stories, and then build progressively more complicated schemas from them.

The schemas built using protoschema matches may add complexity to them in two dimensions: the *compositional* dimension, in which new schemas temporally compose sequences of matched protoschemas to describe a new situation; and the *taxonomic* dimension, in which a protoschema is conceptually narrowed by additional type and relational constraints on its participants, such as "eating something" being narrowed to "eating an apple from a tree".

Protoschemas provide a rich semantic basis for understanding complex situations in terms of basic human behaviors, and this benefit is also reflected in an immediate practical convenience: protoschemas can be employed as a way of *canonicalizing* multiple distinct linguistic phrasings of the same basic action type. For more information on how we use protoschemas for this in NESL, see Section 3.2.

## 3 Learning Pipeline

NESL learns schemas using a multi-step pipeline, briefly described in order here, and further expounded upon in following subsections:

- 1. Using a procedure we call *latent schema sampling*, N short stories are sampled from the same topic-parameterized distribution defined by a language model and a task-specific prompt. (Section 3.1)
- 2. The N stories are then parsed into Episodic Logical Form (ELF), the formal semantic representation underlying the event schemas.
- 3. Simple protoschemas are then matched to each of the N EL-parsed stories with the help of LOME, a state-of-the-art neural FrameNet parser, whose identified FrameNet frames are mapped to corresponding EL protoschemas. (Section 3.2)
- 4. For each of the N stories, all identified protoschemas, and all unmatched ELF episodes and type formulas from the story, are composed into a *single-story schema*, in which constants are abstracted to variables and type-constrained, and events are related with a time graph. (Section 3.3)
- 5. The N single-story schemas are generalized into a single schema, incorporating common details and excising specious, incidental information from the entire set. (Section 3.4)

# 3.1 Latent Schema Sampling

Any story may have a generic schema formed from it. However, stories often contain incidental details: the sequence of going to school, taking an exam, and waving hello to a friend should likely have its third event discarded in a suitably general "school" schema. We adopt the hypothesis that the language model can generate stories according to a distribution implicitly parameterized by one or more latent schemas, from which it may deviate, but according to which it abstractly "selects" subsequent events. By inducing the language model into this distribution via prompt-engineering, and then sampling from it, we hypothesize that high-probability events will occur frequently across samples, and that incidental details will occur less frequently. By generating schemas for each sampled individual story, and generalizing them based on shared events and properties, we may approximate the language model's latent schemas and encode them into interpretable Episodic Logic schemas.

We implement latent schema sampling (LSS) with a sequence of two passes through the GPT-J

6B language model (Wang and Komatsuzaki, 2021) <sup>1</sup>, each pass performing a different task induced by a "prompt-engineering" approach, in which the language model is not fine-tuned, but instructed in natural language to perform a task on some input via its context window. The two language model-based tasks of LSS are described below, and illustrated in the top half of Figure 4.

## 3.1.1 Topic Extraction

To ensure that the stories generated by LSS share common and salient topics, and that the story topics conform to a degree of conceptual simplicity that will work well with the child-like protoschemas and the early-stage EL parsing pipeline, we estimate the set of topics from a known collection of simple stories. We assemble a collection of short, five-sentence stories by filtering the ROCStories dataset (Mostafazadeh et al., 2016), taking the 500 stories with the highest proportion of words shared with a classic collection of child-level stories (McGuffey, 1901).

We created a few-shot task prompt for GPT-J 6B to extract one to three simple topics, such as "going to a store" or "playing baseball", for a given story. We inserted each filtered ROCStory into this prompt and saved the generated topics for use in the next step.

## 3.1.2 Story Sampling

Adopting the hypothesis that story topics encode latent schemas instantiated by the story, we created a second few-shot prompt for the task of story generation given an extracted topic. For each topic generated in the previous step, we sample N stories from the language model with a temperature setting of 0.4 and a repetition penalty of 0.11. Sampled stories are filtered through a blacklist of 375 inappropriate words (Inflianskas, 2019), and stories caught in the content filter are "re-rolled" until N content-appropriate stories have been generated.

# 3.2 Protoschema Matching as FrameNet Parsing

To begin forming schemas, we bootstrap by matching general behavioral *protoschemas* to the stories. Protoschema matching is complicated by several issues, including the large number of actions that may evoke a general protoschema (e.g. walking,

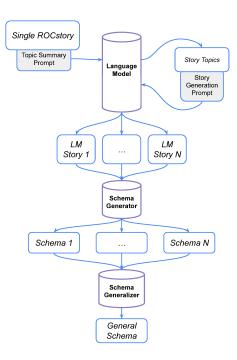


Figure 4: A diagram of the Latent Schema Sampling (LSS) procedure used to generate stories from a fixed-topic distribution.

running, driving, riding, and flying are all kinds of self-motion), and the large number of phrasings that express the same action. To help address these concerns, we first observe that the FrameNet project (Baker et al., 1998) contains many conceptual frames with protoschema analogs, such as self-motion, ingestion, and possession frames. While FrameNet frames lack hierarchical, compositional, and formal internal semantics, and are thus not suitable for the mechanistic inferences these schemas are meant to enable, these frames may be *mapped* to analogous protoschemas as a way to leverage existing work on FrameNet frame parsing.

Using LOME (Xia et al., 2021), a state-of-the-art neural FrameNet parsing system, we identify FrameNet frames in sampled stories. The invoking actions and roles of the frames are given as spans of text, which we reduce to single tokens using a dependency parser, selecting the first token in each span with a NSUBJ, DOBJ, or POBJ tag. The token indices are then aligned with those in the Episodic Logic parser to identify individuals in the logical domain whose type predicates were derived from those tokens. Finally, the frame, now with EL domain individuals as its roles, is mapped to a corresponding protoschema with hand-written mapping rules.

This FrameNet-based protoschema matching

<sup>&</sup>lt;sup>1</sup>We chose GPT-J 6B due to its number of parameters; at the time of this work, we believe it to be the largest publicly-downloadable auto-regressive language model.

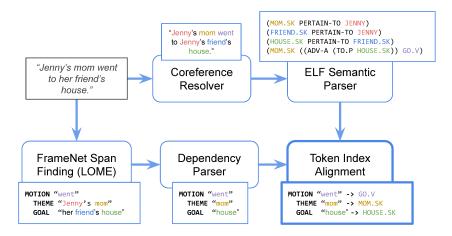


Figure 5: A presentation of the process by which FrameNet parsing with LOME is used for protoschema matching. The coreference resolver is provided by AllenNLP (Gardner et al., 2017), and the dependency parser is provided by spaCy (Honnibal and Montani, 2017).

process is illustrated, with an example, in Figure 5.

## 3.3 Single-Story Schema Formation

After parsing a story into Episodic Logic and using LOME to find protoschema matches in the story, we create a *single-story schema* to encapsulate the story's events and participants. This schema will be general in that the story's specific individuals will be abstracted to variables, and also in that the story's exact verbiage will be "normalized" into the semantic representations of EL and protoschemas. However, the schema will be non-general in that all of the story's details, without regard to what "usually" happens in the latent schema that generated it, will be kept, and only removed during multi-schema generalization (see Section 3.4).

The process of single-story schema generation is fairly simple:

- 1. The event formulae of the story are ordered as steps in a schema. If any of those steps matched to protoschemas, we instead substitute the header of that protoschema as the step; the full specification of the nested protoschema will be output separately, and may be freely expanded.
- All individual constants that are arguments to any of the verbs of those steps are replaced with variables.
- 3. All type and relational constraints on those variables are extracted from the EL parse of the story, as well as from any recursively nested protoschema steps, and enumerated in the ROLES section of the schema.

4. A verbalization of the schema, generated with the GPT-2 model described in Section 3.5, is fed to a GPT-J 6B model with a prompt designed for single-sentence story summaries. The summary sentence is then parsed back into EL, and its arguments are aligned, based on type, with participants in the schema. This formula is then used as the "header" episode of the schema, as seen in the generated header of the schema in Figure 2.

## 3.4 Multi-Story Schema Generalization

Once N stories have been sampled and N corresponding schemas obtained from them, those schemas must be generalized into a single schema. The remainder of this section describes the four main steps of the multi-schema generalization process:

- 1. Schema step clustering: similar steps of each of the N un-merged schemas are generalized together using vector embeddings. (Sections 3.4.1 and 3.4.2)
- 2. Step argument co-reference resolution: the verb arguments of the events of each general step are linked to one another based on co-reference information in the N un-merged schemas. (Section 3.4.3)
- 3. Temporal step ordering: the general schema's steps are put into a partial "happens-before" order. (Section 3.4.4)
- 4. Occurrence frequency filtering: infrequent details across the N un-merged schemas are assumed to be unimportant to the latent schema

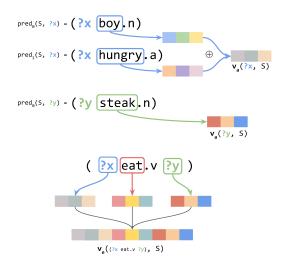


Figure 6: The procedure for embedding EL formulas as vectors. word2vec embeddings for each argument's lexical predicates are averaged together to form vector representations of each argument. These are concatenated, in subject-verb-object order, with verb predicates, to form vector representations of event formulas. Here,  $\oplus$  denotes an element-wise vector mean operation.

and discarded from the general schema approximating it. (Section 3.4.5)

Figure 3 shows a schema generalized from multiple single-story schemas, such as the one shown in Figure 2.

#### 3.4.1 Logical Formula Vector Embeddings

When generalizing several schemas sampled from the same "latent schema", we would like to combine and generalize non-identical, but functionally similar, steps, e.g. "the boy eats cake" and "the girl eats pie".

To do this, we define a function  $v_a(a, S)$ , which takes an argument symbol a and a schema S and returns the element-wise mean of the word vectors of all lexical predicates in S applied to argument a:

$$v_a(a,S) = \bigoplus_i [v_w(\operatorname{pred}_i(S,a))]$$

We also define a function  $v_s(\phi, S)$ , which takes a step formula  $\phi$  and a schema S and returns vector concatenation of  $\phi$ 's verb predicate and the vector representations of each of  $\phi$ 's arguments:

$$v_s(\phi, S) = v_w(\phi_{\text{verb}}) + + + [v_a(\phi_{\text{arg}_i}, S)]$$

Figure 6 illustrates this vectorization process using three "argument type" formulas as the values of the *pred* function. The lexical word vector function,

 $v_w$ , is implemented by word2vec.  $\oplus$  and # refer to element-wise mean and vector concatenation operations, respectively.

## 3.4.2 Formula Vector Clustering

After vector embeddings have been created for the EL formulas for each step of each of the N sampled schemas, we form clusters of similar formula embeddings. Because the vector elements for the formula's verb and for each argument are independent, clusters can correlate not only similar actions, like "boy draws tree" and "boy sketches tree", but also similar argument types, like "boy eats cake" and "girl eats pie". When suitably similar steps have been clustered from across each of the sampled schemas, the clusters form the basis for steps in a new, generalized schema.

Given N schemas, each generated from one of N sampled stories, we denote the sets of step vectors for schema 0 < i < N as  $ST_i = \{\vec{s}_{i,0},...,\vec{s}_{i,M}\}$ . For each step vector  $\vec{s}_{i,a}$ , we construct a list L of each step vector in each other schema,  $\vec{s}_{j\neq i,b}$ , and sort the elements of L in descending order according to their cosine similarity with  $\vec{s}_{i,a}$ . Steps suitably similar to  $\vec{s}_{i,a}$  are defined as the set of all elements  $L_{i\leq k}$  where  $k = \underset{k}{\operatorname{argmax}}[sim(L_k) - sim(L_{k+1})]$ , that is, all steps prior to the largest drop in cosine similarity values in the ordered list.

Once each step s has an associated cluster  $L_s$  of suitably similar steps, symmetricity is enforced by merging all similarity clusters with shared elements: each  $L_{s_1}$  and  $L_{s_2}$  are set to  $L_{s_1} \cup L_{s_2}$  if and only if  $L_{s_1} \cap L_{s_2} \neq \emptyset$ . When this is done, each final cluster  $L_s$  represents one step in the generalized schema.

#### 3.4.3 Schema Slot Co-Reference Resolution

After forming clusters of similar steps across N schemas, we would like to use these clusters as abstract steps in the merged, general schema. However, some clustered steps may not have specific instances in some schemas, even if they should ultimately share arguments with other steps in those schemas. Furthermore, specious co-reference may occur in LM-generated stories or in imperfect ELF parses, e.g., a person incorrectly being equivocated with their dog in the parsed logical domain; we would like to exclude such co-references from our merged, general schema on the basis of their (hoped-for) infrequency among the N samples.

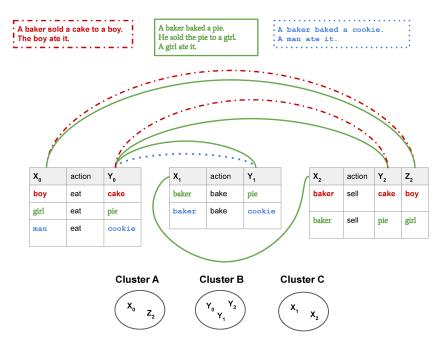


Figure 7: An illustration of the multigraph argument co-reference technique used to merge multiple sampled schemas into one general schema. Each action-argument table represents a cluster of steps obtained by the vector clustering approach in Section 3.4.2. Lines between pairs of argument columns correspond to an argument co-reference in one specific story instance. Each story instance is assigned a unique line pattern, color, and font, for easy identification. The transitive closures of the multigraph's edge relation form the clusters seen at the bottom of the figure; these clusters will form the variables of the eventual merged, generalized schema, whose abstracted steps will be, in no particular order, (A eat B), (C bake B), and (C sell B A). Note that this final sell action is only present in two of the three story instances.

To address these concerns, we perform argument co-reference resolution across step clusters by constructing a multigraph G=(V,E) where V is the set of all unique argument positions across all clusters, and a single edge between cluster arguments  $\arg_i \in V$  and  $\arg_j \in V$  signifies that, in at least one of the N sampled schemas, formulas from each cluster existed and shared an argument value in those two positions. One such edge between any two vertices may exist for each schema being merged.

After construction, the multigraph is reduced to a standard graph, G' = (V', E') with weighted edges: the set of all edges  $E_{i,j}$  between each pair of vertices  $\arg_i$  and  $\arg_j$  is converted to a single edge  $E'_{i,j} \in E'$  whose weight is given by:

$$W_{E_{i,j}'} = \frac{|E_{i,j}|}{|S(\arg_i) \cap S(\arg_j)|}$$

where  $S(v_x)$  is the set of all schemas with a step found in the step cluster for which  $\arg_x$  is an argument. Informally, this weight is the ratio of the number of schemas in which the argument coreference was detected to the number of schemas in which it  $could\ have\ been$  detected. To remove

specious, infrequent co-references,  $E^\prime$  is filtered to remove edges with  $W_{E_{i,j}^\prime} < 0.25.$ 

Argument clusters are formed by the transitive closure of the edge relation given by E', and each argument cluster then forms a final variable for the merged, general schema. The types for each of the final variables are taken from the union of all possible types across all schema instances, and each type predication is assigned a certainty score proportional to its frequency across all N schemas.

An example of the multigraph and final argument clusters generated by this process is illustrated in Figure 7.

## 3.4.4 Temporal Sorting of Schema Steps

Once steps from specific schema instances have been clustered into general steps, and shared arguments have been created across these general steps, we must finally derive their temporal order in the general schema. Episodes in Episodic Logic may be continuous intervals, and EL schemas support complex temporal relations between the episodes characterized by their steps. The schema step intervals, however, are slightly simplified from the full semantics of EL, and represented as a "time graph"

specifying before- and after-relationships between the start and end times of each episode.

Much like the argument co-reference resolution done in Section 3.4.3, we solve this with a *temporal* multi-graph. To build the graph, we further simplify the temporal model by assuming that step episodes never overlap and are defined solely by their start times. This assumption is desirable for its simplicity and computational efficiency, and suitable because the current version of the EL parser makes the same assumption. However, this algorithm could, in theory, be extended to operate on both start and end times.

We define the temporal multi-graph  $G^T=(V^T,E^T)$  such that the vertices  $V^T$  are the start times of each step in the merged, general schema, and, for all steps i and j in the general schema, one edge exists between each  $V_i^T$  and  $V_j^T$  for each occurrence of an instance of  $V_i^T$  happening before an instance of  $V_j^T$  in the same un-merged schema.

Similarly to how edges were assigned weights in the argument co-reference graph based on the ratio of the number of edges to the number of possible edges, we say that, in the general schema, step i happens before step j if and only if:

$$|E_{i,j}^T| > \frac{|S(i) \cap S(j)|}{2}$$

Informally, step i happens before step j if and only if the majority of un-merged schemas that contain both also have a happens-before edge between them.

## 3.4.5 Occurrence Frequency Filtering

Recall that the idea behind latent schema *sampling* is to filter out events that are unimportant to a core schema by exploiting their low frequency of generation by a large language model. Therefore, to finish our general schema, we must finally remove general steps that do not occur in enough of the N sampled schemas to distinguish themselves from "noise". We currently define this threshold for "enough" as  $\frac{N}{3}$ : at least one third of sampled schemas must contain an instance of a general step for the step to remain in the general schema.

# 3.5 Verbalization and Rendering

Once finalized, we post-process learned general schemas for human readability. To verbalize the formal ELF representations of the schema steps, we first apply rule-based transductions to serialize the formula's lexical EL predicates into a pseudo-English representation. Then, using a pre-trained, fine-tuned, 774M-parameter GPT-2 model (Radford et al., 2019), we convert these pseudo-English symbol sequences into proper English. Using the Huggingface Transformers library (Wolf et al., 2020), we fine-tuned this GPT-2 model on 1,200 pairs, manually annotated by a research assistant, for this task <sup>2</sup>.

After verbalizing the steps, we render the general schemas into an HTML representation for human review, automatically color-coding the variables with maximum mutual contrast for enhanced readability. An example of a verbalized and rendered schema is shown in Figure 2.

#### 4 Future Evaluation

This project is a work in progress; although NESL can generate one general schema every 10 minutes, and has generated several hundred to date, qualitative evaluation of the generated schemas has not yet been performed. Imminently, we intend to carry out two human-judged studies: one evaluating the quality of *inferences* generated by the learned schemas when given unseen stories, and another evaluating the quality of the schemas themselves.

Logical schemas enable consistent, structured, and interpretable inferences about novel text, by matching pieces of the text to pieces of the schema, replacing schema variables with entities from the story, and treating other formulas in the schema that use those newly-filled variables as inferences. It is crucial that we demonstrate the inferential capacity of these learned schemas, and as future work, we will be sourcing suitable inference datasets, designing a means of presenting inferences to human judges, and collecting quality evaluations.

In addition to inferences, we would like to evaluate whether the schemas we obtain are both *topically cohesive*, i.e., focused descriptions of one kind of situation; and *interesting*, i.e., capable of generating useful and novel inferences about situations, rather than obvious or redundant ones.

Using our GPT-2 verbalization model and schema rendering software, described in Section 3.5, to make our schemas and inferences readable to untrained human judges, we intend to immediately move forward with the design and execution

<sup>&</sup>lt;sup>2</sup>Our work on GPT-2 for formula verbalization is preliminary; as we continue this work, a more robust annotation protocol may be required, employing multiple annotators.

of these quality evaluation studies to complete the work.

## 5 Conclusion

We have described NESL, a hybrid neural and formal-logical schema learning system with which we aim to combine a richly structured schema representation, a human learning-inspired approach to schema acquisition, the linguistically flexible sentence understanding characteristic of neural NLP systems, and the large amount of knowledge contained in large language models.

By bringing a large and varied number of components from across the literature to bear, we have shown that general, semantically rich, and seemingly sensical schemas may be extracted from large language models and represented interpretably. In the immediate future, we also plan to demonstrate that these schemas are useful for inference tasks.

## References

- James F. Allen. 1983. Maintaining knowledge about temporal intervals. *Commun. ACM*, 26(11):832–843.
- Collin F. Baker, Charles J. Fillmore, and John B. Lowe. 1998. The berkeley framenet project. In *Proceedings* of the 36th Annual Meeting of the Association for Computational Linguistics and 17th International Conference on Computational Linguistics Volume 1, ACL '98/COLING '98, pages 86–90. Association for Computational Linguistics.
- Nathanael Chambers and Dan Jurafsky. 2008. Unsupervised learning of narrative event chains. *Proceedings of ACL-08: HLT*, pages 789–797.
- Matt Gardner, Joel Grus, Mark Neumann, Oyvind Tafjord, Pradeep Dasigi, Nelson F. Liu, Matthew Peters, Michael Schmitz, and Luke S. Zettlemoyer. 2017. Allennlp: A deep semantic natural language processing platform.
- Matthew Honnibal and Ines Montani. 2017. spaCy 2: Natural language understanding with Bloom embeddings, convolutional neural networks and incremental parsing. To appear.
- Chung Hee Hwang and Lenhart K Schubert. 1993. Episodic logic: A situational logic for natural language processing. *Situation Theory and its Applications*, 3:303–338.
- Roman Inflianskas. 2019. profanity-filter: A python library for detecting and filtering profanity. https://github.com/rominf/profanity-filter/blob/master/profanity\_filter/data/en\_profane\_words.txt.

- Lane Lawley, Gene L. Kim, and Lenhart Schubert. 2019. Towards natural language story understanding with rich logical schemas. In *Proceedings of the IWCS workshop on Natural Language and Computer Science*.
- Lane Lawley, Benjamin Kuehnert, and Lenhart K. Schubert. 2021. Learning general event schemas with episodic logic. In *NALOMA*.
- Michael Lebowitz. 1980. Generalization and Memory in an Integrated Understanding System. Ph.D. thesis, Yale University, New Haven, CT, USA. AAI8109800.
- William Holmes McGuffey. 1901. *The New McGuffey First Reader*. American Book Company.
- Raymond J Mooney. 1990. A general explanationbased learning mechanism and its application to narrative understanding. Morgan Kaufmann.
- Nasrin Mostafazadeh, Nathanael Chambers, Xiaodong He, Devi Parikh, Dhruv Batra, Lucy Vanderwende, Pushmeet Kohli, and James Allen. 2016. A corpus and cloze evaluation for deeper understanding of commonsense stories. In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 839–849, San Diego, California. Association for Computational Linguistics.
- Peter Norvig. 1987. Inference in text understanding. In *AAAI*, pages 561–565.
- Jean Piaget and Margaret Cook. 1952. *The origins of intelligence in children*, volume 8. International Universities Press New York.
- Karl Pichotta and Raymond J Mooney. 2016. Learning statistical scripts with 1stm recurrent neural networks. In *Thirtieth AAAI Conference on Artificial Intelligence*.
- Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners.
- Ben Wang and Aran Komatsuzaki. 2021. GPT-J-6B: A 6 Billion Parameter Autoregressive Language Model. https://github.com/ kingoflolz/mesh-transformer-jax.
- Lilian Wanzare, Alessandra Zarcone, Stefan Thater, and Manfred Pinkal. 2017. Inducing script structure from crowdsourced event descriptions via semi-supervised clustering. In *Proceedings of the 2nd Workshop on Linking Models of Lexical, Sentential and Discourse-level Semantics*, pages 1–11.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le

Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. 2020. Transformers: State-of-the-art natural language processing. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 38–45, Online. Association for Computational Linguistics.

Patrick Xia, Guanghui Qin, Siddharth Vashishtha, Yunmo Chen, Tongfei Chen, Chandler May, Craig Harman, Kyle Rawlins, Aaron Steven White, and Benjamin Van Durme. 2021. LOME: Large ontology multilingual extraction. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: System Demonstrations*, pages 149–159, Online. Association for Computational Linguistics.

# A Learned Schema Examples

# A.2 From the prompt "going to the library"

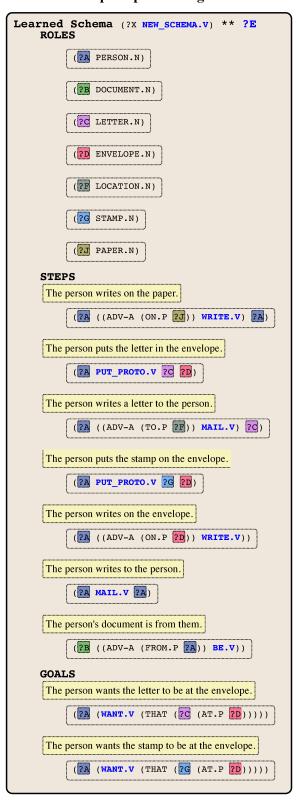
The following are additional examples of schemas learned by NESL.

## A.1 From the prompt "going fishing"

```
Learned Schema (?X NEW_SCHEMA.V) ** ?E ROLES
           (PB SEAFOOD.N)
           (?G PERSON.N)
           (?K FOOD.N)
           (?L SUPPER.N)
           (?M LOT.N)
     The person cooks the sea food for the food.
           (PG ((ADV-A (FOR.P PK)) COOK.V) PB)
     The person cleans the sea food.
           (PG CLEAN_PROTO.V PB)
     The person cooks the sea food.
           (2G COOK.V 2B)
     The person eats sea food.
          (?G EAT_PROTO.V ?B)
     The person goes somewhere.
           (?G GO.V)
     The person catches the sea food.
           ( CATCH_PROTO.V PB)
     The person catches the lot.
           (?G CATCH.V ?M)
     The person wants the sea food to be not dirty.
           ( ( ( WANT. V (THAT (NOT ( PB DIRTY. A)))))
     The person wants to not be hungry.
           ( C (WANT.V (THAT (NOT ( C HUNGRY.A)))))
```

```
Learned Schema (?x NEW_SCHEMA.V) ** ?E
ROLES
          (PA BOOK.N)
          (PE PERSON.N)
          (PI LIBRARY.N)
          (PJ ENTITY.N)
          (PL BOY.N)
           (PLO LOCATION.N)
          (PLX LOCATION.N)
          (?T EVENT.N)
     The person looks for a book.
          (PE ((ADV-A (FOR.P 2A)) LOOK_PROTO.V))
     The person reads a book about an event.
          (2E ((ADV-A (ABOUT.P 2T)) READ_PROTO.V) 2A)
     The person likes the book.
          (PE LIKE_PROTO.V PA)
     The person finds a book.
          (?E FIND_PROTO.V ?A)
     The person goes to the library.
          (PE ((ADV-A (FROM.P 2J)) GO_PROTO.V) PI)
     The person reads a book.
          (PE READ.V PA)
     GOALS
     The person wants to find the book.
          (ZE (WANT.V (KA (FIND.V ZA))))
     The person wants to possess a book.
          (PE (WANT.V (KA (POSSESS.V PA))))
     The person wants to be at the library.
          (PE (WANT.V (KA ((ADV-A (AT.P PI)) BE.V))))
```

## A.3 From the prompt "mailing a letter"



## **B** Language Model Prompts

# **B.1** Story Topic Summarization Prompt

This is the story generation prompt provided to GPT-J to obtain latent story schemas, or topics, from given stories. The format string at the end, %s,

is replaced with the input story, and a completion from the language model is truncated at a newline to obtain a topic. The prompt stories and latent schemas were written by research assistants and chosen based on subjective downstream schema acquisition quality on a development set of schema topics estimated from ROCstories.

```
Story:
Tom loved playing baseball.
He had a big game.
He was up to hit.
He hit a long drive.
He made a run and won the game.
Core story schemas:
playing a game
playing baseball
winning a game
Story:
The man took a shower.
The hot water went cold.
He still had soap in his hair.
He washed his hair quickly.
He was shivering when he got out of the
   shower.
Core story schemas:
bathing
taking a shower
Story:
Oran bought binoculars.
He took them outside.
He saw birds.
He watched them.
They became his friends.
Core story schemas:
watching animals
using binoculars
buying something
Story:
Emma went to school.
She studied math.
She ate lunch.
Her teacher gave her a lot of homework.
Later, she went home.
Core story schemas:
going to school
studying something
Story:
A farmer got up in the morning.
He put his boots on.
He went outside.
```

```
He milked the cow.
He went back to bed.
Core story schemas:
farming
milking a cow
Story:
My son is a little child.
He ran outside to play.
His friend was out there with him.
They played together with sticks.
My son came in from outside.
Core story schemas:
children playing
playing with a friend
playing outside
Story:
I went to my door yesterday.
I saw there was a new book.
It came right to me.
I was pretty happy about that.
I couldn't wait to read it.
Core story schemas:
getting mail
reading a book
Story:
The hedge started to grow.
Spring came around.
The hedge started to bud flowers.
The flowers grew.
The roses were very beautiful.
Core story schemas:
springtime
plants growing
Story:
%S
Core story schemas:
```

## **B.2** Story Generation Prompt

This is the story generation prompt provided to GPT-J to obtain a story given a topic as input. The format string at the end, %s, is replaced with the input topic, and a completion from the language model is truncated at a double-newline to obtain a story. The prompt stories were written by research assistants and chosen based on subjective downstream schema acquisition quality on a development set of schema topics estimated from ROCstories.

```
Stories about baseball:
Tom loved playing baseball.
He had a big game.
He was up to hit.
He hit a long drive.
He made a run and won the game.
Bob went to see a baseball game.
The players had nice bats.
The players swung at the ball.
One player hit the ball.
He hit a home run.
Jenny was playing baseball.
She took a bat and got ready.
She swung her bat at the ball.
She hit the ball.
She won the game.
Stories about showers:
The man took a shower.
The hot water went cold.
He still had soap in his hair.
He washed his hair quickly.
He was shivering when he got out of the
   shower.
Jenny took a shower.
She used soap to wash her body.
She washed her hair.
The water was warm.
She dried off with a towel.
Jack was dirty.
He needed to get clean.
He took a long shower.
The shower water was very hot.
He used plenty of soap.
Stories about plants:
Jessie loved plants.
She had plants in her apartment.
She watered the plants every day.
Her favorite plant was her fern.
Jessie wanted to buy more plants.
Alan bought a plant at the store.
The plant died.
He bought another plant.
He watered it.
It didn't die.
Plants are usually green.
Some plants are different colors.
Sometimes people keep plants in their
   houses.
They water those plants.
People like plants.
Stories about school:
Emma went to school.
She studied math.
She ate lunch.
```

Her teacher gave her a lot of homework. Later, she went home.

Jason was at school.

He ate lunch in the cafeteria.

After lunch, he went to class.

His teacher taught him about math.

He went home and ate dinner.

Abhishek loved to go to school.
His teacher gave him fun homework.
He finished his homework and gave it back to the teacher.
The teacher said Abhishek did a good job

The teacher gave him a good grade.

Stories about %s: