

Full Automation of Goal-driven LLM Dialog Threads with And-Or Recursors and Refiner Oracles

Paul Tarau

University of North Texas

February 29, 2024

code at <https://github.com/ptarau/recursors/>

demo at <https://deepllm.streamlit.app/>

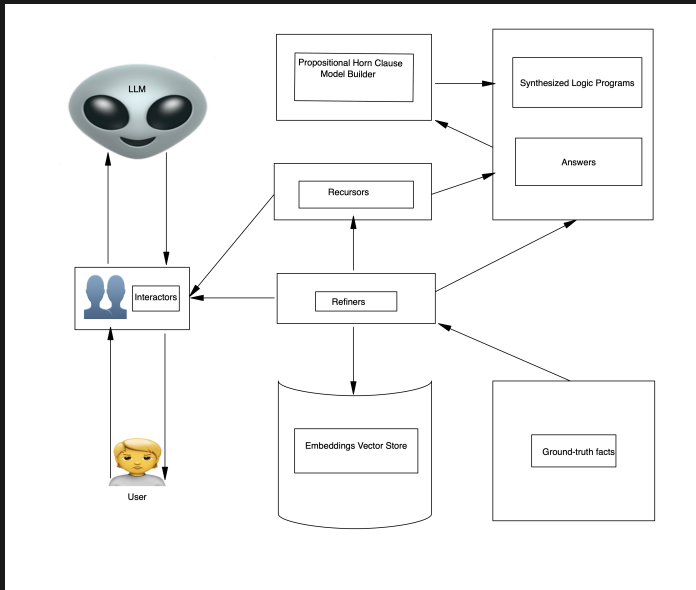
demo at <https://deep-auto-quests.streamlit.app/>

demo at <https://docdiver.streamlit.app/>

Overview

- we automate deep step-by step reasoning in an LLM dialog thread by recursively exploring alternatives (OR-nodes) and expanding details (AND-nodes) up to a given depth
- starting from a single succinct task-specific initiator we steer the automated dialog thread to stay focussed on the task by synthesizing a prompt that summarizes the depth-first steps taken so far
- semantic similarity to ground-truth facts or oracle advice from another LLM instance is used to restrict the search space and validate the traces of justification steps returned as focussed and trustable answers
- applications:
 - consequence predictions
 - causal explanations
 - recommendation systems
 - topic-focussed exploration of scientific literature

DeepLLM: System Architecture



DeepLLM: key components

- Active components (“Agents”): interactors, recursors, refiners.
 - Interactors manage input prompts and task breakdown
 - Recursors handle iterative exploration of subtasks
 - Refiners enhance clarity and relevance of LLM responses
- Resources:
 - Ground truth facts: sentences collected from online sources or local
 - Vector store: enables “semantic search” via embeddings of sentences

A first example: AGI as seen by GPT4

Prompter:

- ☒ scientific_concept_explorer
- ☐ consequence_prediction
- ☐ causal_inference
- ☐ step_by_step_guidance_to_achieve_a_goal
- ☐ recommender_system
- ☐ supporting_arguments_for_a_thesis
- ☐ cons_and_pros_for_a_thesis

Activate LLM!

Visualize relation graph!

Browse relation graph in new tab!

```
'artificial general intelligence!:-  
  'Machine learning algorithms';  
  'Cognitive computing capabilities';  
  'Autonomous decision-making processes'.  
'Machine learning algorithms':-  
  'Supervised Learning Models',  
  'Neural Network Architecture',  
  'Reinforcement Learning Techniques',  
  'Deep Learning Algorithms',  
  'Feature Extraction Methods'.  
'Cognitive computing capabilities':-  
  'Neural Network Models',  
  'Cognitive Architecture Integration',  
  'Machine Learning Algorithms',  
  'Semantic Knowledge Representation',  
  'Distributed Cognitive Systems'.  
'Autonomous decision-making processes':-  
  'Reinforcement Learning Algorithms',  
  'Probabilistic Graphical Models',  
  'Multi-Agent Systems',  
  'Decision-Theoretic Planning',  
  'Stochastic Control Policies'.  
'Supervised Learning Models'.  
'Neural Network Architecture'.  
'Reinforcement Learning Techniques'.
```

Figure: A Logic Program: AGI as seen by GPT4

Graph of Relations defining AGI

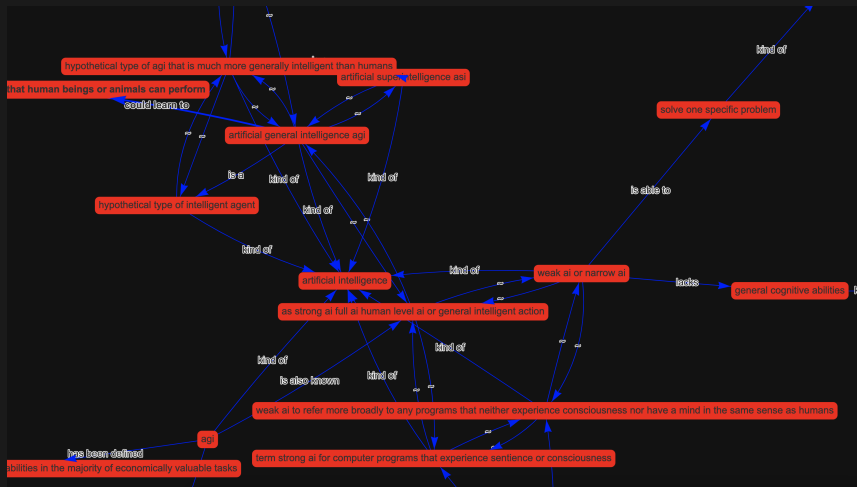


Figure: Relations defining AGI

Logic Programming : Propositional Horn Clauses in Prolog

```
pet :- dog.  
pet :- cat.  
pet :- snake.  
  
dog :- barks, walks, bites.  
  
cat :- purrs, walks, hisses.  
  
snake:- hisses, slither, bites.  
  
walks :- true.  
purrs  :- true.  
bites  :- true.  
hisses :- true.  
  
slither :- false.  
bites   :- false.  
barks   :- false.
```

```
?- pet.  
true
```

DeepLLM: Logic Programming, but done differently!

- SLD-resolution's clause selection via unification is replaced by LLM-driven dynamic clause head creation with an option of focusing by proximity of embeddings to ground truth facts
- as dialog units are sentences, the underlying logic is propositional
- client-side management (via the API) of the LLM's memory is based on the equivalent of a *goal stack* and a *goal trace* recording our steps on the current search path
- instead of variable bindings, answers are traces of justification steps clearly explaining where they are derived from
- when their depth-limit is reached, the items on the goal stack are interpreted as "abducibles", statements that can be hypothetically assumed and then checked against "integrity constraints"

Keeping the ground truth in mind

- our depth-bounded refinement steps support compilation of the dialog threads to a Horn Clause program to be explored with logic programming solvers
- modular, task specific, customizable prompt engineering primitives are aggregated together for “AND-step” and “OR-step” prompts
- normalized semantic similarity measures of embeddings can be made available when generating probabilistic logic programs
- sentences in authoritative documents or collections seen as "ground-truth facts" can be used to select abducibles via semantic similarity or advice of an LLM-based oracle

Overall, our approach exploits synergies between structured prompt engineering, logic-guided recursion over LLM queries and semantic search in an embeddings vector store.

Interactors

- the LLM API: same client-server interface for local and outside LLM
- prompt templates example:
 - AND-OR prompt patterns for causal reasoning

```
causal_prompter = dict(  
    name='causal',  
    and_p="""We need causal explanations in this context: "$context"  
    Generate 3-5 explanations of 2-4 words each for the causes of "$g".  
    Itemize your answer, one reason for "$g" per line.  
    No explanations needed, just the 2-4 words noun phrase,  
    nothing else.""",  
    or_p="""We need causal explanations in this context: "$context"  
    Generate 2-3 alternative explanations citing facts that might  
    cause "$g".  
    Itemize your answer, one noun phrase per line.  
    No explanations needed, just the noun phrase, nothing else.  
    """)
```

- OR-prompter expanding the head h in a series of alternatives a_1, \dots, a_n :
 - $h :- a_1.$
 - $h :- a_2.$
 - \dots
 - $h :- a_m.$
- more concisely expressed with a disjunctive body as:
 - $h :- a_1 ; a_2 ; \dots ; a_m.$
- the result of an AND-prompter:
 - $h :- b_1, b_2, \dots, b_n.$
- when a depth limit is reached, the remaining unexplored goals are considered as *abducibles*

- The Embeddings **Vector Store**
 - Filtering abducibles with semantic distance to **ground-truth facts**
- Refining decisions with **LLM-based oracles**
 - The LLM-based True/False Decider
 - The LLM-based Rater

Toward Trustable Generative AI

- **trustable AI**: requirement for wide adoption of today's LLMs in medical, educational, defense and several other business applications
- likely subject of upcoming government regulations
- \Rightarrow DeepLLM provides a principled approach toward trustable generative AI:
 - oracles enforcing consistency and factuality
 - focussed reasoning steps enforced by casting AND/OR steps into a propositional Horn clause program

What if $NP == coNP$?: exploring the consequences

Topic to explore:

Proof that NP equals coNP

Prompter:

- ☐ scientific_concept_explorer
- ☒ consequence_prediction
- ☐ causal_inference
- ☐ step_by_step_guidance_to_achieve_a_goal
- ☐ recommender_system
- ☐ supporting_arguments_for_a_thesis
- ☐ cons_and_pros_for_a_thesis

Activate LLM!

```
'Proof that NP equals coNP':-  
  'Revolution in complexity theory';  
  'Transformation of cryptographic systems';  
  'Paradigm shift in algorithm design'.  
'Revolution in complexity theory':-  
  'Fundamental shift in computational understanding',  
  'New algorithms development surge',  
  'Transformation of encryption systems',  
  'Paradigm shift in mathematical logic',  
  'Reshaping of computer science education'.  
'Transformation of cryptographic systems':-  
  'Enhanced security protocols',  
  'Advanced encryption techniques',  
  'Shift in cyber defense strategies',  
  'New computational complexities',  
  'Evolution of digital communication'.  
'Paradigm shift in algorithm design':-  
  'Advanced computational capabilities',  
  'Revolution in problem-solving methods',  
  'New era of software development',  
  'Unprecedented efficiency in data processing',  
  'Transformation in artificial intelligence algorithms'.
```

Figure: Exploring consequences of a counterfactual hypothesis

An application: DeepQA, Iterating on multiple follow-up questions

DeepQA: a DeepLLM App exploring self-generated follow-up questions

☐ Local LLM?

Enter your OPENAI_API_KEY:

OpenAI LLM

☒ GPT-4☐ GPT-3.5

Maximum depth

2

14

☒ Show generated Definite Clause Grammar?

Question to start with:

How can humans align superintelligent LLMs to human goals and values?

Activate LLM!

Q: How can humans align superintelligent LLMs to human goals and values?

A: One way to align superintelligent Language Models (LLMs) to human goals and values is through a process known as value alignment. This involves training these models on a diverse range of human texts, allowing them to learn and understand our values, ethics, and goals. The more diverse the training data, the more likely the model will align with a broad spectrum of human values.

Q: How can we ensure the diversity of the training data and avoid potential biases in the LLMs?

A: One way to ensure the diversity of the training data is to consciously include texts from a wide range of sources, cultures, and perspectives. This can help to represent a broad spectrum of human values and experiences. To avoid potential biases, it's important to critically evaluate the sources of the training data and to actively seek out and include underrepresented voices.

Q: How can humans align superintelligent LLMs to human goals and values?

A: One way to align superintelligent Language Models (LLMs) to human goals and values is through a process known as value alignment. This involves training these models on a diverse range of human texts, allowing them to learn and understand our values, ethics, and goals. The more diverse the training data, the more likely the model will align with a broad spectrum of human values.

Q: How can we ensure the diversity of the training data and avoid potential biases in the LLMs?

A: Another approach is to use algorithms that can detect and mitigate biases in the training data. These algorithms can help to identify patterns of bias and to adjust the training data accordingly, ensuring a more balanced representation of human values and goals.

Q: How can humans align superintelligent LLMs to human goals and values?

A: One way to align superintelligent Language Models (LLMs) to human goals and values is through a process known as value alignment. This involves training these models on a diverse range of human texts, allowing them to learn and understand our values, ethics, and goals. The more diverse the training data, the more likely the model will align with a broad spectrum of human values.

Q: How can we ensure the diversity of the training data and avoid potential biases in the LLMs?

A: A third method is to involve a diverse group of people in the process of curating and reviewing the training data. This can help to ensure a wide range of perspectives and to identify potential biases that might be overlooked by a more homogenous group.

Q: How can humans align superintelligent LLMs to human goals and values?

Share ☆ ↻

Manage

How DeepLLM enhances a user's LLM Interactions?

- Reduces complexity of prompt engineering
- Automates decomposition of tasks into simpler sub-tasks
- Refines prompts for more accurate user intent capture
- Generates dynamic responses for a wide range of queries
- Improves user experience with concise and relevant information

Future Work

- Granularity Refiners
 - higher granularity: one can work with sentences/statements instead of noun phrases
 - lower granularity: SVO triplets
- Question generators
 - asking LLMs to generate follow-up questions, recursively (done, see DeepQA)
- Diversifiers and Harmonizers
 - to restrict unwanted “hallucinatory” generation twists *diversifiers* for OR-nodes and *harmonizers* for AND-nodes can be expressed as additional integrity constraints
 - diversifier: no two OR-nodes are too close semantically
 - harmonizer: no two AND-nodes are too far semantically
 - both could be implemented with help of semantic distances in the embeddings vector store
- Extended Implementation Techniques
 - fixpoint computation with torch tensors (actually already working!)
 - extend the power of the underlying logic language: ASP, probabilistic LP

Conclusion

- automation of deep step-by step reasoning in LLM dialog threads up to a given depth, by recursively descending from a single succinct initiator phrase, while staying focussed on the task at hand
- we have made LLMs function as de facto logic programming engines that mimic Horn Clause resolution
- instead of trying to parse sentences into logic formulas, we have accommodated our logic engine to fit the natural language reasoning patterns LLMs have been trained on
- semantic similarity to ground-truth facts and oracle advice from another LLM instance has been used to restrict the search space and validate the traces of justification steps returned as answers
- focussed, controllable output, enabling:
 - deep investigations into details of specific scientific domains
 - expert-level causal reasoning or consequence predictions
- a method to extract hallucination-free focussed knowledge expressed as a logic program that encapsulates trustable AI in a clearly expressed and easily verifiable framework

Questions?

Links, ready to try out!

code at <https://github.com/ptarau/recursors/>

demo at <https://deepllm.streamlit.app/>

demo at <https://deep-auto-quests.streamlit.app/>

demo at <https://docdiver.streamlit.app/>