

AI Report

We are highly confident this text is

AI Generated

AI Probability

100%

This number is the probability that the document is AI generated, not a percentage of AI text in the document.

Plagiarism



The plagiarism scan was not run for this document. Go to gptzero.me to check for plagiarism.

Untitled Document - 11/29/2025

Anonymous User

This paper proposes ContextNav, a novel Chain-of-Thought (CoT)-driven navigation framework designed for multimodal embodied question answering (Embodied QA). The core idea is to augment a large multimodal model (LMM) with an explicit reasoning chain that can be dynamically adapted to changing visual and textual contexts.

The authors introduce a new benchmark called CEN (Contextual Embodied Navigation), which requires agentic navigation in 3D environments (based on AI2-THOR) conditioned on both linguistic queries and evolving visual observations. ContextNav is shown to outperform baselines such as VLN-BERT, RecCoT, and classical map-based agents, with analysis showing improvements in long-horizon reasoning and policy transferability.

The paper focuses on something quite underexplored--how to make vision-language agents actually reason and act over time, rather than just respond to single inputs. The focus on embodied QA with CoT-style reasoning feels fresh and highly relevant.

Most CoT papers just use reasoning chains to get better answers. Here, it's used to drive navigation--deciding where to go next, based on what's already seen and reasoned about. That's a neat twist on the usual setup.

I like that the authors didn't just propose a model--they also introduced a new benchmark (CEN) that requires both visual and linguistic context to solve. This feels realistic and useful to the broader community.

The method outperforms several strong baselines (VLN-BERT, RecCoT, etc.), especially on long-horizon tasks. The gains are consistent, and not just on cherry-picked examples.

The paper includes helpful qualitative visualizations of how CoT steps evolve, how memory is used, and what decisions are made. This makes the agent's reasoning process transparent and auditable.

The paper combines CoT + memory, but there's no solid ablation to tell which part really matters. It would help to show performance with and without the memory module.

Since the system heavily relies on the underlying vision-language model (e.g., for perception, grounding, reasoning), it'd be helpful to test or at least discuss what happens if you switch to a smaller or weaker LMM.

Some examples suggest that the model "imagines" object positions it hasn't actually seen. Are there ways to ground the reasoning better or prevent hallucinations?

CoT is generated at every step, which might be slow in real-world settings. But there's no mention of how fast inference is, or whether it can run in real-time. It's a bit hard to reproduce the method without more info. There's no pseudocode, no clear prompt templates, and no info on model size, inference cost, or environment config. Some of this could go in the appendix.

Could you provide an ablation where you disable the memory module? How much does it actually contribute compared to CoT alone? Which specific LMM is used, and did you try smaller or open-source models? How much does performance depend on having a "strong" LMM?

Are there known cases where the CoT causes the agent to make wrong or speculative moves due to hallucinated beliefs? How long does one inference step take? Can this be deployed in a near real-time scenario? Could the benchmark (CEN) be extended to tasks beyond navigation--like multi-room exploration, tool use, or question answering with object manipulation?

Would you consider releasing pseudocode or a Colab to help others reproduce the results? If these questions are addressed in detail, especially around memory vs. CoT contribution and scalability, I'd definitely consider raising my score.

 Sentences that are likely AI-generated.

FAQs

What is GPTZero?

GPTZero is the leading AI detector for checking whether a document was written by a large language model such as ChatGPT. GPTZero detects AI on sentence, paragraph, and document level. Our model was trained on a large, diverse corpus of human-written and AI-generated text with support for English, Spanish, French, German, and other languages. To date, GPTZero has served over 10 million users around the world, and works with over 100 organizations in education, hiring, publishing, legal, and more.

When should I use GPTZero?

Our users have seen the use of AI-generated text proliferate into education, certification, hiring and recruitment, social writing platforms, disinformation, and beyond. We've created GPTZero as a tool to highlight the possible use of AI in writing text. In particular, we focus on classifying AI use in prose. Overall, our classifier is intended to be used to flag situations in which a conversation can be started (for example, between educators and students) to drive further inquiry and spread awareness of the risks of using AI in written work.

Does GPTZero only detect ChatGPT outputs?

No, GPTZero works robustly across a range of AI language models, including but not limited to ChatGPT, GPT-5, GPT-4, GPT-3, Gemini, Claude, and AI services based on those models.

What are the limitations of the classifier?

The nature of AI-generated content is changing constantly. As such, these results should not be used to punish students. We recommend educators to use our behind-the-scene [Writing Reports](#) as part of a holistic assessment of student work. There always exist edge cases with both instances where AI is classified as human, and human is classified as AI. Instead, we recommend educators take approaches that give students the opportunity to demonstrate their understanding in a controlled environment and craft assignments that cannot be solved with AI. Our classifier is not trained to identify AI-generated text after it has been heavily modified after generation (although we estimate this is a minority of the uses for AI-generation at the moment). Currently, our classifier can sometimes flag other machine-generated or highly procedural text as AI-generated, and as such, should be used on more descriptive portions of text.

I'm an educator who has found AI-generated text by my students. What do I do?

Firstly, at GPTZero, we don't believe that any AI detector is perfect. There always exist edge cases with both instances where AI is classified as human, and human is classified as AI. Nonetheless, we recommend that educators can do the following when they get a positive detection: Ask students to demonstrate their understanding in a controlled environment, whether that is through an in-person assessment, or through an editor that can track their edit history (for instance, using our [Writing Reports](#) through Google Docs). Check out our list of [several recommendations](#) on types of assignments that are difficult to solve with AI.

Ask the student if they can produce artifacts of their writing process, whether it is drafts, revision histories, or brainstorming notes. For example, if the editor they used to write the text has an edit history (such as Google Docs), and it was typed out with several edits over a reasonable period of time, it is likely the student work is authentic. You can use GPTZero's Writing Reports to replay the student's writing process, and view signals that indicate the authenticity of the work.

See if there is a history of AI-generated text in the student's work. We recommend looking for a long-term pattern of AI use, as opposed to a single instance, in order to determine whether the student is using AI.