

Automated title and abstract screening for scoping reviews using the GPT-4 Large Language Model

Introduction

A scoping review is a relatively novel type of literature review that aims to map the key concepts and existing activity within an area of research (Arksey and O'Malley 2005). Like systematic reviews, scoping reviews typically use rigorous, transparent (Pham et al. 2014), and sometimes pre-registered methods for gathering and synthesising evidence, and increasingly use formal frameworks for both performing and reporting reviews (Micah D. J. Peters et al. 2021). Scoping reviews can inform future systematic reviews or primary research in the same area (Sutton et al. 2019). However, they differ from systematic reviews in aiming to describe the breadth of coverage of the available literature rather than research findings in depth (Arksey and O'Malley 2005).

Frameworks for performing a scoping review typically involve defining a research area or question, searching bibliographic databases for potentially relevant published material ('sources'), screening these sources to identify those relevant to the area or question, and systematically extracting and reporting data from the relevant sources (Institute 2015; Arksey and O'Malley 2005). The screening stage will usually involve initial screening of source titles and abstracts against pre-determined inclusion and exclusion criteria, followed by screening of source full texts, with both steps performed in replicate by at least two human reviewers (Micah D. J. Peters et al. 2020; Pham et al. 2014). Because database searches can return many hundreds or thousands of potentially relevant sources, these screening steps can require intensive human time and effort. Many software methods have been proposed or used to support or partially automate source screening for scoping reviews, including text mining to prioritise relevant sources for human screening (Shemilt et al. 2014; Howard et al. 2016; Chai et al. 2021), automated clustering and labelling of sources to support human decision-making (Stansfield, Thomas, and Kavanagh 2013), and 'crowdsourcing' screening to untrained workers via online platforms (Mortensen et al. 2017). A similar but more extensive set of methods have been developed and employed for systematic reviews (Khalil, Ameen, and Zarnegar 2022; Gates et al. 2019) for which the process of source screening is broadly comparable.

Since the release of the first Generative Pre-trained Transformer (GPT) Large Language Model (LLM) by OpenAI (San Francisco, California, United States of America) in 2018 (Radford et al. 2018), transformer-based LLMs and the GPT lineage in particular have seen rapid and widespread adoption for a range of automation tasks. Broadly, these models generate a probabilistically weighted list of ‘tokens’ (parts of text such as letter combinations and punctuation) to continue or complete some input text (a ‘prompt’), having been trained to do so by practising such predictions on large human-written corpora. When this generative process is iterated, it allows for a range of applications involving analysis and production of text, such as summarising documents, generating fiction in a particular genre or style, or conversing with humans (OpenAI 2023).

While LLMs are not yet widely used to screen sources for literature reviews, early work suggests they may perform well in this role. Guo et al. (Guo et al. 2023) reported the use of a GPT-lineage model (they do not specify which, though their published code suggests OpenAI’s ‘gpt-3.5-turbo’ model) to screen 24,307 titles and abstracts from five systematic reviews and one scoping review, achieving weighted average sensitivity of 76% and specificity of 91% when compared to human reviewers. Their approach involves giving the model a brief prompt instructing it to take on the persona of a researcher screening titles and abstracts, followed by a source’s title and abstract as well as the inclusion and exclusion criteria. The model is instructed to respond with a decision to include or exclude the source, and the process is iterated across the full set of sources to be screened. Syriani et al. (Syriani, David, and Kumar 2023) similarly reported the use of ‘gpt-3.5-turbo’ to screen titles and abstracts for a systematic review and achieved sensitivities of above 70%. They also systematically evaluated prompts given to the LLM to identify a prompt that performed best at the screening task; their chosen prompt, like that of Guo et al., placed the LLM in the role of an academic reviewer.

Both of these approaches made use of a single, fixed text prompt template, which the LLM then completes with additional text representing its response (the decision to include or exclude a source), a method sometimes called ‘zero-shot prompting’. Recent work has identified a number of methods which can be superior to zero-shot prompting when using LLMs for tasks that require complex or multi-step reasoning. These methods include ‘chain-of-thought prompting’ (Wei et al. 2022), in which a complex task is broken down into a series of intermediate steps, and the ‘tree of thoughts’ strategy (Yao et al. 2023), in which multiple parallel chains of thought are generated, compared, and integrated. The LLM is induced to follow these complex reasoning strategies either by being given examples of multi-step reasoning on similar tasks, or by being lead through the process with a series of intermediate prompts.

In this paper, I introduce a package for the R programming language (R Core Team 2023) called GPTscreenR that implements a chain-of-thought based approach to using GPT-4 for scoping review title and abstract screening, and evaluate its performance by comparison to human reviewers. The purpose of

this package is to assist and augment rather than replace human reviewers in performing scoping reviews. This package represents the first LLM-based screening tool designed specifically for scoping reviews. Further, this study provides the first report on the accuracy of LLM-based screening using the most recent iteration of the GPT model lineage, GPT-4, and using the recently developed chain-of-thought approach.

Methods

The GPTscreenR package

GPTscreenR is an R (R Core Team 2023) package released under the MIT open source licence. The source code is available for download from GitHub at <https://github.com/wilkox/GPTscreenR>. At the time of writing the most recent package version was 0.0.2, and the results presented in this paper were obtained using this version.

GPTscreenR consists of two main components. The first is a set of internal functions for interfacing with the OpenAI Application Programming Interface (API), which allows for ‘conversations’ with the GPT-4 LLM, as well as functions for representing and manipulating those conversations. These internal functions are designed to be model-agnostic, so that future versions of the package or users with particular needs can use different GPT models. The OpenAI API requires an OpenAI account, and OpenAI charges fees for use of the API. In order to access the API, GPTscreenR requires a secret key to be registered prior to source screening, and instructions for doing so are provided in the package documentation and on loading of the package in R if the key has not been correctly registered.

The second component is a set of user-facing functions to perform source screening with GPT-4. The `review_description()` function assists in generating a text description of the review’s objectives and inclusion and exclusion criteria, using the Population, Concept, and Context (PCC) framework (Micah D. J. Peters et al. 2020) for defining the review’s inclusion and exclusion criteria. The use of this function is optional, and users may instead choose to provide a description of the review and criteria for source selection using any framework or format they see fit.

The `screen_source()` function performs the main task of the package. This function mediates an conversation with GPT-4 in which chain-of-thought prompting (Wei et al. 2022) is used to guide GPT-4 through screening a source title and abstract against the study inclusion criteria. The template for this conversation is given in Fig. 1. The OpenAI API defines a conversation as a series of messages, each of which originates from one of three roles: *system*, representing an authoritative voice that can instruct GPT-4 on its task and behaviour; *user*, representing a human user that can interact with GPT-4; and *assistant*, representing the responses generated by GPT-4. In `screen_source()`,

the *system* role gives GPT-4 general instructions, while the *user* role provides the user-written review description and the source title and abstract.

The specific phrase ‘work step by step’ is derived from the ‘let’s think step by step’ prompt phrase which significantly improves LLM performance on multi-step reasoning tasks with a zero-shot prompt (Kojima et al. 2022), adapted to this chain-of-thought approach.

GPT-4 is instructed to summarise the inclusion criteria for the scoping review, compare the title and abstract against these summarised criteria, and make a final recommendation on whether to include or exclude the source (Fig. 1). Following the chain-of-thought approach, examples of summarised inclusion criteria and of comparing the source to these criteria are given for GPT-4 to use as templates or exemplars. This approach was chosen after noting that a major source of error when attempting to screen sources with zero-shot, one-shot, or few-shot prompts (i.e. with a single prompt and no, one, or a few examples) was that GPT-4 would fail to consider important inclusion criteria. As an example, Fig 2a presents a conversation with GPT-4 using a zero-shot prompt. The screening task in this example is intentionally adversarial and designed to lead the model towards making an error. In order to correctly recommend exclusion of the source, GPT-4 must notice that the review is looking for research on therapy alpaca interventions, but that the source reports on a therapy camel intervention. The presence of multiple other inclusion criteria which are met by the source, as well as the mention of alpacas in the source abstract, serve as distractors. In this example, GPT-4 incorrectly recommends inclusion. If the conversation is then continued to draw GPT-4’s attention to the error, it is able to identify and correct it (Fig 2b), suggesting that the error arises from a failure of GPT-4 to properly consider the relevant inclusion criterion rather than an inability to do so. Using the chain-of-thought approach overcomes this problem without the need for human intervention and correction (Fig 2c). GPT-4 identifies ‘The source reports on a study conducted in Residential [*sic*] aged care facilities that have used therapy alpacas as a part of their programme of care for residents’ as an inclusion criterion, and correctly assesses that this is the only inclusion criterion not met. GPT-4 then correctly recommends exclusion of the source.

The `screen_source()` function returns a list comprising the complete transcript of the conversation with GPT-4 and GPT-4’s final recommendation as a logical (Boolean) value. The conversation transcript can be used to interrogate cases where GPT-4 may have returned an incorrect or unexpected result. The package also provides a function `screen_sources()`, which applies `screen_source()` iteratively across a data frame of sources. `screen_sources()` caches screening outcomes to a file as they accumulate, so that screening can be split across multiple sessions and recover from interruptions.

System: You are being used to help researchers perform a scoping review. You are not interacting directly with a user.

A scoping review is a type of systematic review used to map the published scholarship on a topic. To gather relevant sources for a scoping review, the researchers search bibliographic databases for sources that might be relevant to the review, often using the Population, Concept, and Context framework. The titles and abstracts of sources that are found in this search are then screened against the review inclusion and exclusion criteria.

Your task is to screen a single source against the study criteria. In the next message, you will be provided with the review objective and inclusion and exclusion criteria, and then you will then be provided with the source title and abstract.

To screen the source, you must work step by step. First, generate a numbered list of statements that summarise the inclusion and exclusion criteria for the scoping review, including the Population, Concept, and Context if they are provided. The statements should be clear, comprehensive and complete. Any source for which all the statements are true is a source that meets the inclusion criteria. As a template, here are some example statements (these are a generic set of examples that are not related to the current scoping review):

1. The source reports the results of a randomised control trial
2. The source reports the results of a study in which:
 - 2a. The participants were all male; AND
 - 2b. The participants were all aged between 18 and 74 inclusive
3. The source reports the results of a study conducted in the European Union.

Aspects of the inclusion criteria with multiple elements should be broken down into separate points where possible. For example, instead of:

1. The source reports on a study of men who live in the European Union.

You should instead say:

1. The source reports on a study of people who are:
 - 1a. Male; and
 - 1b. Living in the European Union.

Then, continue to work step by step. Refer back to the set of statements you developed summarising the inclusion criteria. For each statement, decide whether or not the statement is true for the source described by the title and abstract. You must select from the following permitted responses: TRUE, FALSE, LIKELY TRUE, LIKELY FALSE, or NOT APPLICABLE. No other response is permitted. It is normal for the title and abstract to not have enough information to make a clear decision for every statement. There is a natural and normal amount of ambiguity in this process. For these situations, give your best guess, making use of your general knowledge, and deciding LIKELY TRUE or LIKELY FALSE. Responses like UNCLEAR or NOT ENOUGH INFORMATION are not permitted. After giving your response, give a one sentence explanation for your response. For example:

1. TRUE. The abstract describes the study design as a randomised control trial.
 - 1a. TRUE. The abstract mentions that all the participants were male.
 - 1b. TRUE. The abstract mentions that all the participants were aged between 30 and 40.
2. LIKELY TRUE. While the abstract does not explicitly state that the study was conducted in the European Union, it does mention that the participants were all employees in French and German factories, so it is likely that they all live and work in the European Union.

Finally, consider your decisions on whether the title and abstract meet the conclusion criteria. Overall, is it likely true that the source meets the inclusion criteria? End your response with a single word on a new line, either INCLUDE or EXCLUDE, representing your recommendation on whether the source is likely to meet the inclusion criteria. The response must end with a line containing only one of these two words; any other reply will cause the automatic parsing of your response to fail, which will be troublesome for the user. The response must not end with a blank line."

User: *The user-generated review description.*

GPT-4: *GPT-4 replies with a summary of the inclusion criteria, assessment of the source against each criterion, and a final recommendation.*

Figure 1: Template for the conversation with GPT-4 mediated by the `screen_source()` function. Messages⁵ with variable content, including user-provided data as well as GPT-4's responses, are given in italics.

User: You are being used to help researchers perform a scoping review. A scoping review is a type of systematic review used to map the published scholarship on a topic. To gather relevant sources for a scoping review, the researchers search bibliographic databases for sources that match a selected Population, Concept, and Context (the inclusion criteria). The titles and abstracts of sources that are found in this search are then screened against the inclusion criteria.

This review has the following objective and inclusion criteria:

OBJECTIVE: This scoping review will examine the existing primary research on the role of therapy alpacas in enhancing the mental health and emotional well-being of elderly residential aged care facility residents

POPULATION: Elderly people living in residential aged care facilities. 'Elderly' is defined as 65 years of age or older.

CONTEXT: Residential aged care facilities that have used therapy alpacas as part of their programme of care for residents. 'Residential aged care facility' is defined as a residential setting intended as a long-term place of residence for elderly people which includes provision of support for activities of daily living (e.g. meal preparation, bathing, housekeeping) and nursing support (e.g. medication management). Such facilities will also typically offer other structured programmes and facilities to provide entertainment, diversion, and wellbeing. It excludes other residential settings intended for elderly people that do not provide daily living or nursing supports (e.g. independent living villages) or that are not long-term (e.g. hospitals or hospices).

CONCEPT: The impact of therapy alpaca programmes on stress reduction, emotional well-being, mental health, overall life satisfaction, or similar outcomes for residents.

Here is the title and abstract of the source for you to screen:

TITLE: Therapy Ungulate Visits Reduce Depression, Anxiety, and Stress in Australian Nursing Home Residents

ABSTRACT: Nursing home residents typically have rates of depression, anxiety, and stress (DAS) significantly higher than those of equivalent populations living independently. As well as reducing quality of life and increasing morbidity and mortality, DAS significantly increases the burden on both nursing homes and health systems in caring for the elderly. Recent pilot trials of the use of therapy ungulates (alpacas, llamas, and camels) have shown promising initial results in reducing DAS among nursing home residents. We conducted a randomised control trial involving 8,124 residents across 17 nursing homes in urban and rural Australia in which residents were randomised to either a weekly 2-hour therapy camel visit or equivalent time in a television viewing lounge. Residents receiving the intervention had significantly ($p < 0.01$) reduced levels of depression (mean DASS-21 score for depression reduction 8.2 points), anxiety (mean reduction 11 points) and stress (mean reduction 13 points) compared to the control arm. This study represents the first large-scale randomised trial of ungulate therapy for nursing home residents, and provides strong support for its effectiveness in reducing DAS.

Please respond with a single word, either INCLUDE or EXCLUDE, representing your recommendation

GPT-4: INCLUDE

Figure 2a: A conversation with GPT-4, using a zero-shot prompt to instruct the model to screen a single source for a hypothetical scoping review. This adversarial task has been intentionally designed to influence GPT-4 towards making an error.

User: This doesn't seem right. The inclusion criteria includes the concept of therapy alpaca programmes, while the abstract describes a study of a therapy camel programme. Can you explain your reasoning?

GPT-4: I apologize for the oversight. You're right. While the abstract mentions "therapy ungulates" (which includes alpacas, llamas, and camels), the specific intervention studied was therapy camel visits, not alpacas. Given the specific criteria you've set for the scoping review, this study does not directly address the role of therapy alpacas.

Based on this clarification, the correct recommendation would be:

EXCLUDE.

Figure 2b: A continuation of the conversation in Fig. 2a, in which GPT-4 is lead to recognise and correct the error.

Validation

To validate GPTscreenR's approach, six scoping reviews were identified from the Open Science Framework (OSF; <https://osf.io>) where the review inclusion criteria and the results of title and abstract screening were publicly available. A summary of the review characteristics is provided in Table 1. Small random subsets of screened sources from four of the reviews (**COVID**, **smartphones**, **solastalgia**, and **teachers**) were used during initial testing and refinement of the `screen_source()` function, while the full set or random subsets from all six reviews were used for final validation. Random subsets were used where the large number of sources made it prohibitive in time and cost to screen the full set. The total number of sources available for screening and the number used for validation from each review are given in Table 1.

Some of the reviews did not include the full abstract text of sources in the publicly available files, and where these abstracts could not be obtained from external databases these sources were excluded from validation. There were also many cases where missing, malformed, or duplicate data required either manual correction or exclusion of sources. The scoping review data, code used to prepare this data for validation, and code for calculating summary statistics are available in a reproducible form in the package repository on GitHub (<https://github.com/wilcox/GPTscreenR/tree/master/validation>).

The consensus human reviewer decision at the title and abstract screening level was used as the gold standard outcome. Sensitivity and specificity were calculated by comparing GPT-4's recommendation against this gold standard. Three of the scoping reviews (**COVID**, **solastalgia**, and **teachers**) included individual human reviewer decisions in addition to the final decision in their publicly available datasets, and these were used to calculate human intraobserver agreement (Cohen's kappa) using the R function `cohen.kappa()` from the **psych** package (William Revelle 2023). This was compared to human/GPT-4 agreement across all screened sources, calculated with the same method.

System: You are being used to help researchers perform a scoping review. You are not interacting directly with a user.

A scoping review is a type of systematic review used to map the published scholarship on a topic. To gather relevant sources for a scoping review, the researchers search bibliographic databases for sources that might be relevant to the review, often using the Population, Concept, and Context framework. The titles and abstracts of sources that are found in this search are then screened against the review inclusion and exclusion criteria.

Your task is to screen a single source against the study criteria. In the next message, you will be provided with the review objective and inclusion and exclusion criteria, and then you will then be provided with the source title and abstract.

To screen the source, you must work step by step. First, generate a numbered list of statements that summarise the inclusion and exclusion criteria for the scoping review, including the Population, Concept, and Context if they are provided. The statements should be clear, comprehensive and complete. Any source for which all the statements are true is a source that meets the inclusion criteria. As a template, here are some example statements (these are a generic set of examples that are not related to the current scoping review):

1. The source reports the results of a randomised control trial
2. The source reports the results of a study in which:
 - 2a. The participants were all male; AND
 - 2b. The participants were all aged between 18 and 74 inclusive
3. The source reports the results of a study conducted in the European Union.

Aspects of the inclusion criteria with multiple elements should be broken down into separate points where possible. For example, instead of:

1. The source reports on a study of men who live in the European Union.

You should instead say:

1. The source reports on a study of people who are:
 - 1a. Male; and
 - 1b. Living in the European Union.

Then, continue to work step by step. Refer back to the set of statements you developed summarising the inclusion criteria. For each statement, decide whether or not the statement is true for the source described by the title and abstract. You must select from the following permitted responses: TRUE, FALSE, LIKELY TRUE, LIKELY FALSE, or NOT APPLICABLE. No other response is permitted. It is normal for the title and abstract to not have enough information to make a clear decision for every statement. There is a natural and normal amount of ambiguity in this process. For these situations, give your best guess, making use of your general knowledge, and deriving LIKELY TRUE, or LIKELY FALSE. Responses like UNCLEAR or NOT ENOUGH INFORMATION are not permitted. After giving your response, give a one sentence explanation for your response. For example:

1. TRUE. The abstract describes the study design as a randomised control trial.
 - 1a. TRUE. The abstract mentions that all the participants were male.
 - 1b. TRUE. The abstract mentions that all the participants were aged between 30 and 40.
2. LIKELY TRUE. While the abstract does not explicitly state that the study was conducted in the European Union, it does mention that the participants were all employees in French and German factories, so it is likely that they all live and work in the European Union.

Finally, consider your decisions on whether the title and abstract meet the inclusion criteria. Overall, is it likely true that the source meets the inclusion criteria? End your response with a single word on a new line, either INCLUDE or EXCLUDE, representing your recommendation on whether the source is likely to meet the inclusion criteria. The response must end with a line containing only one of these two words; any other reply will cause the automatic parsing of your response to fail, which will be troublesome for the user. The response must not end with a blank line.'

User: **OBJECTIVE:** This scoping review will examine the existing primary research on the role of therapy alpacas in enhancing the mental health and emotional wellbeing of elderly residential aged care facility residents.

POPULATION: Elderly people living in residential aged care facilities. 'Elderly' is defined as 65 years of age or older.

CONTEXT: Residential aged care facilities that have used therapy alpacas as part of their programme of care for residents. 'Residential aged care facility' is defined as a residential setting intended as a long-term place of residence for elderly people which includes provision of support for activities of daily living (e.g. meal preparation, bathing, housekeeping) and nursing support (e.g. medication management). Such facilities will also typically offer other structured programmes and facilities to provide entertainment, diversion, and wellbeing. It excludes other residential settings intended for elderly people that do not provide daily living or nursing supports (e.g. independent living villages) or that are not long-term (e.g. hospitals or hospices).

CONCEPT: The impact of therapy alpaca programmes on stress reduction, emotional wellbeing, mental health, overall life satisfaction, or similar outcomes for residents.

TITLE: Therapy Ungulate Visits Reduce Depression, Anxiety, and Stress in Australian Nursing Home Residents

ABSTRACT: Nursing home residents typically have rates of depression, anxiety, and stress (DAS) significantly higher than those of equivalent populations living independently. As well as reducing quality of life and increasing morbidity and mortality, DAS significantly increases the burden on both nursing homes and health systems in caring for the elderly. Recent pilot trials of the use of therapy ungulates (alpacas, llamas, and camels) have shown promising initial results in reducing DAS among nursing home residents. We conducted a randomised control trial involving 8,124 residents across 17 nursing homes in urban and rural Australia in which residents were randomised to either a weekly 2-hour therapy camel visit or equivalent time in a television viewing lounge. Residents receiving the intervention had significantly ($p < 0.01$) reduced levels of depression (mean DASS-21 score for depression reduction 8.2 points), anxiety (mean reduction 11 points) and stress (mean reduction 13 points) compared to the control arm. This study represents the first large-scale randomised trial of ungulate therapy for nursing home residents, and provides strong support for its effectiveness in reducing DAS.

GPT-4: Inclusion and exclusion criteria:

1. The source reports on a study of people who are:
 - 1a. 65 years old or older; and
 - 1b. Living in residential aged care facilities; and
 - 1c. Not living in settings intended for elderly people that do not provide daily living or nursing supports or that are not long-term.
2. The source reports on a study conducted in Residential aged care facilities that have used therapy alpacas as a part of their programme of care for residents.
3. The source reports on the impact of therapy alpaca programmes on stress reduction, emotional wellbeing, mental health, overall life satisfaction, or similar outcomes for residents.

Screening of the source:

1. 1a. LIKELY TRUE. The ages of the participants are not mentioned specifically, but since they are referred to as 'nursing home residents', it can be assumed that most are likely 65 or older.
- 1b. TRUE. The participants were all living in nursing homes.
- 1c. LIKELY TRUE. The study was conducted in nursing homes, which generally include provision of support for activities of daily living and nursing support, making them likely to fit the definition given.
2. FALSE. While the utilising therapy animals, the animals used in this study were camels, not alpacas.
3. TRUE. The study reported on the impact of the therapy involving camels on depression, anxiety, and stress levels in nursing home residents.

RECOMMENDATION: EXCLUDE

Figure 2c: A conversation with GPT-4, instructing it to perform the same adversarial task presented in Fig. 2a, but using the chain-of-thought approach. This is not a continuation of the conversation presented in Fig. 2a and Fig. 2b but a new conversation. In this case, GPT-4 correctly recommends excluding the source. The model also correctly identifies that the source meets all of the relevant inclusion criteria except for the requirement that the source report on an alpaca intervention.

Comparison to zero-shot method

In order to directly compare the chain-of-thoughts approach to a zero-shot approach (i.e. an conversation consisting of a single prompt with no examples, followed by GPT-4’s response), the validation screening task was repeated using the prompt designed by Guo et al. (Guo et al. 2023), substituting the permitted responses `INCLUDE` and `EXCLUDE` for `included` and `excluded` respectively in order to maintain compatibility with GPTscreenR’s parsing of the response (Fig. 3). The number of sources used for zero-shot validation were the same as those used for the chain-of-thoughts method validation given in Table 2, although a different subset was selected by the randomisation process in cases where the full set was not used. The code used to prepare this data and calculate summary statistics was otherwise identical to that used for the chain-of-thoughts method validation and is available in the ‘zeroshot’ branch of the package repository on GitHub (<https://github.com/wilkox/GPTscreenR/tree/zeroshot>).

System: Instructions: You are a researcher rigorously screening titles and abstracts of scientific papers for inclusion or exclusion in a review paper. Use the criteria below to inform your decision. If any exclusion criteria are met or not all inclusion criteria are met, exclude the article. If all inclusion criteria are met, include the article. Only type “INCLUDE” or “EXCLUDE” to indicate your decision. Do not type anything else.

Title: *The source title*

Abstract: *The source abstract*

Inclusion criteria: *The inclusion criteria*

Exclusion criteria: *The exclusion criteria*

GPT-4: *GPT-4 replies with a decision.*

Figure 3: Zero-shot prompt used as a comparator for the chain-of-thoughts approach, derived from the approach of Guo et al. (Guo et al. 2023). Variable content, including user-provided data as well as GPT-4’s response, is given in italics.

Results

1,140 sources were screened from the six scoping reviews. GPTscreenR achieved a weighted average sensitivity of 0.72 and weighted average specificity of 0.87. For the three reviews that provided individual reviewer decisions, the weighted average Cohen’s kappa was 0.67, while the weighted average Cohen’s kappa between final human and GPT-4 decisions was 0.52.

The zero-shot method achieved a weighted average sensitivity of 0.78 and weighted average specificity of 0.84. The weighted average Cohen’s kappa between human and GPT-4 decisions was 0.52.

Discussion

GPTscreenR and its chain-of-thought approach were able to achieve a sensitivity and specificity similar to though not as high as those reported by Guo et al. (Guo et al. 2023). Also similarly to their report, Cohen’s kappa was higher for inter-human than human-GPT-4 agreement, suggesting that despite this generally good performance GPT-based methods still do not match the performance of a human reviewer when the consensus human decision is used as the gold standard. For reasons discussed further below, this is likely because of poorer performance relative to humans on more difficult cases rather than overall poorer performance.

Comparison to the zero-shot method

The replication of Guo et al.’s zero-shot prompt on the reviews used for validation in this study resulted in similar weighted average sensitivity (0.78) compared to their report (0.76), though somewhat poorer specificity (0.84 compared to 0.91). There are a number of possible factors contributing to this difference. Firstly, this study used the GPT-4 model while Guo et al. likely used the ‘gpt-3.5-turbo’ model, although it might be expected that the more advanced GPT-4 would generally perform better on the same task. Secondly, this study included only scoping reviews, while Guo et al. examined five systematic reviews and only one scoping review, although the reported sensitivity (1.00) and specificity (0.94) for that scoping review were higher than the weighted average for this study. Thirdly, validation with both the chain-of-thoughts and zero-shot methods was dependent on the availability of source and screening data from the included reviews, and it was noted that most of the reviews had issues with missing, malformed, missing, and duplicate data, including lacking the text of some or all source abstracts and sufficient information to retrieve them from public databases. While effort was made to rectify these issues, many sources from the included reviews could not be used for validation, and the differing quality of the datasets used in the two studies may have affected the performance of the zero-shot method. Finally, the different reviews used for validation in the two studies likely posed different levels of difficulty the zero-shot approach. If LLM-based approaches to source screening become more widely adopted, the performance of different methods can be more accurately determined across a growing sample of reviews.

Compared against the zero-shot method, the chain-of-thought method achieved poorer sensitivity but higher specificity. This result is consistent with rationale for selecting the chain-of-thought method, which was to reduce type I errors (false positives) of the type demonstrated in Fig. 2a. However, this came at the cost of higher type II error. It is not possible to retrospectively examine the reasoning process that lead GPT-4 to make an incorrect recommendation such as a false negative, as asking it to explain its reasoning post-hoc will result in a confabulated statistically likely response; as they have no direct access to their internal processes, LLMs are unable to meaningfully introspect.

Limitations of this approach

Compared to the prompts used by Guo et al. (Guo et al. 2023) and Syriani et al. (Syriani, David, and Kumar 2023), both the prompt used for this approach and the response generated by GPT-4 have substantially higher token counts. This results in both a longer time to screen a single source (typically 30-40 seconds) and a higher dollar cost, as the OpenAI API currently charges on a per-token basis. This may make GPTscreenR less appealing to some reviewers, particularly when screening very large source sets.

Validation of GPTscreenR was limited by a small number of scoping reviews with publicly available, high-quality data on human title and abstract screening decisions. Because a subset of these reviews were used for both testing and validation, there is a risk of over-fitting of the prompt to these particular reviews. The public release of GPTscreenR encourages users to contribute data from their own scoping reviews to support more accurately measuring GPTscreenR’s real-world performance, and to guide further refinement of the approach.

Conclusions

TODO

References

- Arksey, Hilary, and Lisa O’Malley. 2005. “Scoping studies: towards a methodological framework.” *International Journal of Social Research Methodology* 8 (1): 19–32. <https://doi.org/10.1080/1364557032000119616>.
- Chai, Kevin E. K., Robin L. J. Lines, Daniel F. Gucciardi, and Leo Ng. 2021. “Research Screener: a machine learning tool to semi-automate abstract screening for systematic reviews.” *Systematic Reviews* 10 (1): 93. <https://doi.org/10.1186/s13643-021-01635-3>.
- Gates, Allison, Samantha Guitard, Jennifer Pillay, Sarah A. Elliott, Michele P. Dyson, Amanda S. Newton, and Lisa Hartling. 2019. “Performance and usability of machine learning for screening in systematic reviews: a comparative evaluation of three tools.” *Systematic Reviews* 8 (1): 278. <https://doi.org/10.1186/s13643-019-1222-2>.
- Guo, Eddie, Mehul Gupta, Jiawen Deng, Ye-Jean Park, Mike Paget, and Christopher Naugler. 2023. “Automated Paper Screening for Clinical Reviews Using Large Language Models.” *arXiv*. <https://doi.org/10.48550/arxiv.2305.00844>.
- Howard, Brian E., Jason Phillips, Kyle Miller, Arpit Tandon, Deepak Mav, Mihir R. Shah, Stephanie Holmgren, et al. 2016. “SWIFT-Review: a text-mining workbench for systematic review.” *Systematic Reviews* 5 (1): 87. <https://doi.org/10.1186/s13643-016-0263-z>.

- Institute, The Joanna Briggs. 2015. *The Joanna Briggs Institute Reviewers' Manual 2015: Methodology for JBI Scoping Reviews*. The Joanna Briggs Institute.
- Khalil, Hanan, Daniel Ameen, and Armita Zarnegar. 2022. "Tools to support the automation of systematic reviews: a scoping review." *Journal of Clinical Epidemiology* 144: 22–42. <https://doi.org/10.1016/j.jclinepi.2021.12.005>.
- Kojima, Takeshi, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2022. "Large Language Models Are Zero-Shot Reasoners." *arXiv*. <https://doi.org/10.48550/arxiv.2205.11916>.
- Mortensen, Michael L., Gaelen P. Adam, Thomas A. Trikalinos, Tim Kraska, and Byron C. Wallace. 2017. "An exploration of crowdsourcing citation screening for systematic reviews." *Research Synthesis Methods* 8 (3): 366–86. <https://doi.org/10.1002/jrsm.1252>.
- OpenAI. 2023. "Introducing ChatGPT." "<https://openai.com/blog/chatgpt>".
- Peters, Micah D J, Casey Marnie, Andrea C Tricco, Danielle Pollock, Zachary Munn, Lyndsay Alexander, Patricia McInerney, Christina M Godfrey, and Hanan Khalil. 2020. "Updated Methodological Guidance for the Conduct of Scoping Reviews." *JBI Evidence Synthesis* 18 (10): 2119–26. <https://doi.org/10.11124/jbies-20-00167>.
- Peters, Micah D. J., Casey Marnie, Heather Colquhoun, Chantelle M. Garritty, Susanne Hempel, Tanya Horsley, Etienne V. Langlois, et al. 2021. "Scoping reviews: reinforcing and advancing the methodology and application." *Systematic Reviews* 10 (1): 263. <https://doi.org/10.1186/s13643-021-01821-3>.
- Pham, Mai T., Andrijana Rajić, Judy D. Greig, Jan M. Sargeant, Andrew Papadopoulos, and Scott A. McEwen. 2014. "A scoping review of scoping reviews: advancing the approach and enhancing the consistency." *Research Synthesis Methods* 5 (4): 371–85. <https://doi.org/10.1002/jrsm.1123>.
- R Core Team. 2023. *R: A Language and Environment for Statistical Computing*. Vienna, Austria: R Foundation for Statistical Computing. <https://www.R-project.org/>.
- Radford, Alec, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al. 2018. "Improving Language Understanding by Generative Pre-Training." "<https://openai.com/research/language-unsupervised>".
- Shemilt, Ian, Antonia Simon, Gareth J. Hollands, Theresa M. Marteau, David Ogilvie, Alison O'Mara-Eves, Michael P. Kelly, and James Thomas. 2014. "Pinpointing needles in giant haystacks: use of text mining to reduce impractical screening workload in extremely large scoping reviews." *Research Synthesis Methods* 5 (1): 31–49. <https://doi.org/10.1002/jrsm.1093>.
- Stansfield, Claire, James Thomas, and Josephine Kavanagh. 2013. "'Clustering' documents automatically to support scoping reviews of research: a case

- study.” *Research Synthesis Methods* 4 (3): 230–41. <https://doi.org/10.1002/jrsm.1082>.
- Sutton, Anthea, Mark Clowes, Louise Preston, and Andrew Booth. 2019. “Meeting the review family: exploring review types and associated information retrieval requirements.” *Health Information & Libraries Journal* 36 (3): 202–22. <https://doi.org/10.1111/hir.12276>.
- Syriani, Eugene, Istvan David, and Gauransh Kumar. 2023. “Assessing the Ability of ChatGPT to Screen Articles for Systematic Reviews.” *arXiv*.
- Wei, Jason, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou. 2022. “Chain of Thought Prompting Elicits Reasoning in Large Language Models.” *arXiv*. <https://doi.org/10.48550/arxiv.2201.11903>.
- William Revelle. 2023. *Psych: Procedures for Psychological, Psychometric, and Personality Research*. Evanston, Illinois: Northwestern University. <https://CRAN.R-project.org/package=psych>.
- Yao, Shunyu, Dian Yu, Jeffrey Zhao, Izhak Shafran, Thomas L. Griffiths, Yuan Cao, and Karthik Narasimhan. 2023. “Tree of Thoughts: Deliberate Problem Solving with Large Language Models.” *arXiv*. <https://arxiv.org/abs/2305.10601>.