

Figure 1. Locator map of California to highlight the groundwater sustainability plan (GSP) areas reviewed in this study. The 5 sample GSPs reviewed in most detail are in dark green and the 49 other GSPs with rubrics that we analyzed are in light green. We attempted, but were unable, to analyze 9 GSPs due to technical issues with the report PDF or the human scored rubric (shown in yellow). The remaining GSPs shown in grey were reviewed by humans using different methods and did not have a standardized scoring rubric so they were not included in this study.

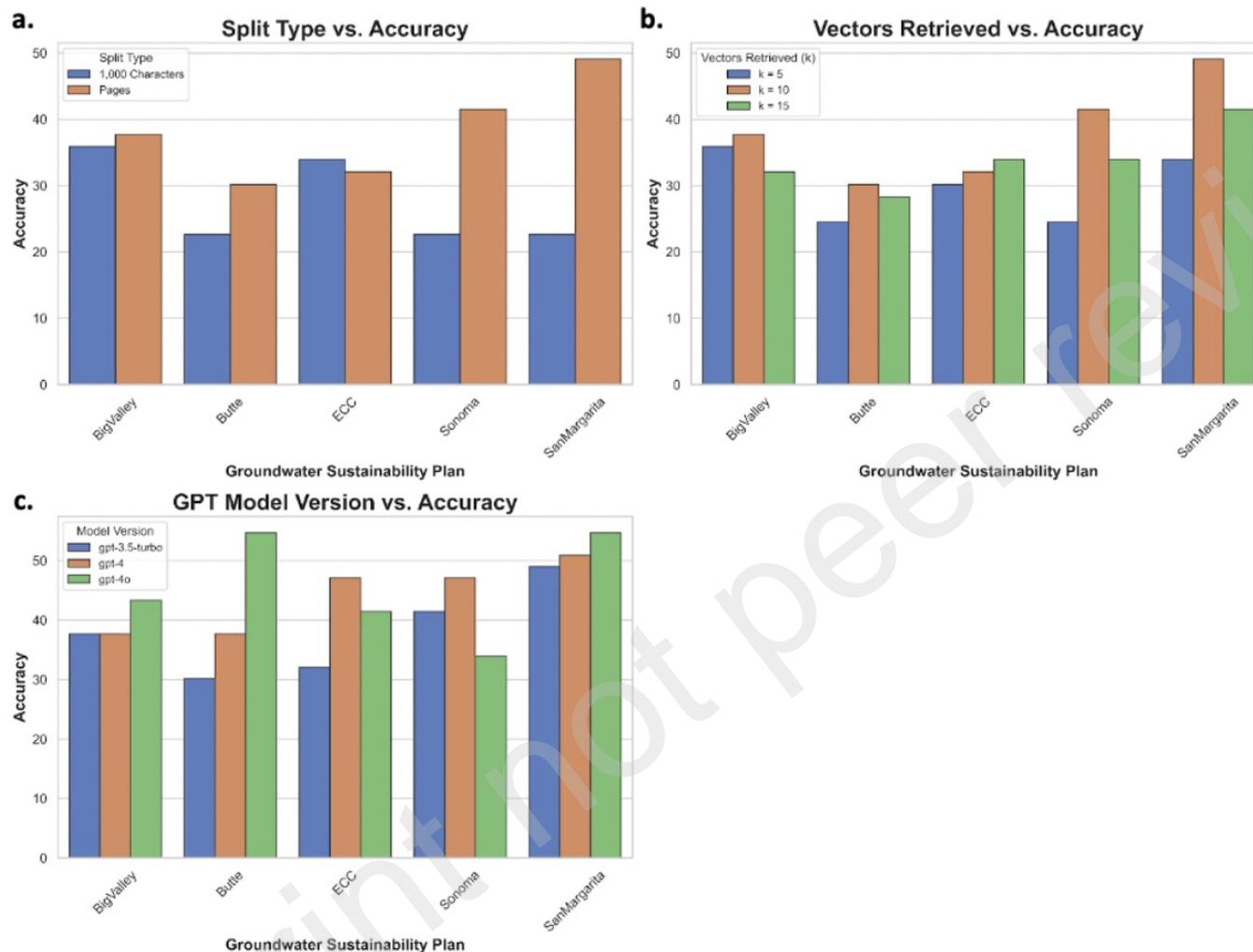


Figure 2: Accuracy (percent of questions answered correctly *100) comparison across model and hyperparameter choices for non-fine-tuned models. Panels a. and b. depict accuracy for experiments on split type to define vector chunk size and vector retrieval (chunk count) size using GPT-3.5 Turbo, while panel c. compares accuracy across model versions.

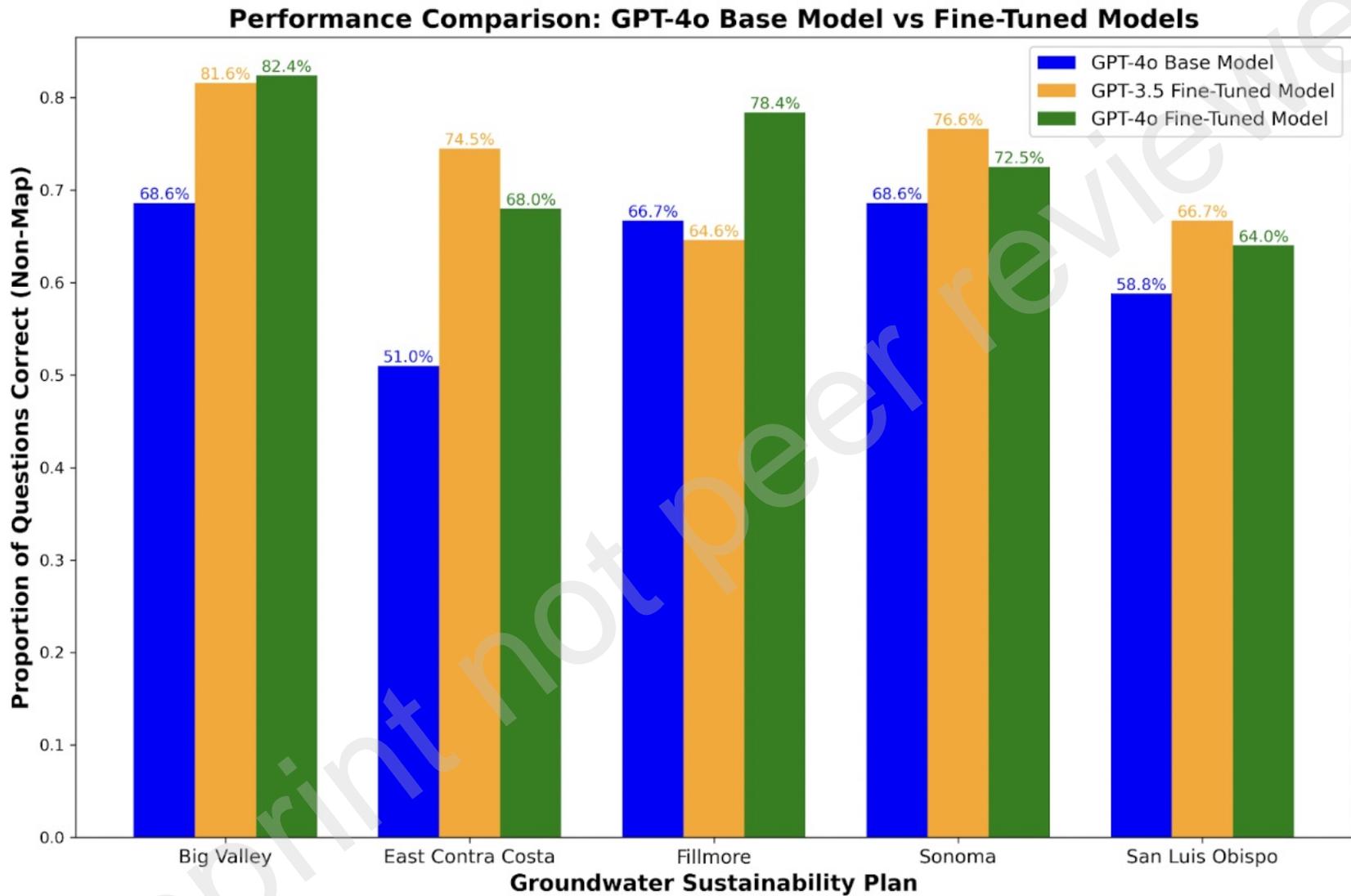


Figure 3: Performance Comparison: GPT-4o Base Model vs GPT-3.5 Fine-Tuned Model vs GPT-4o Fine-Tuned Model.

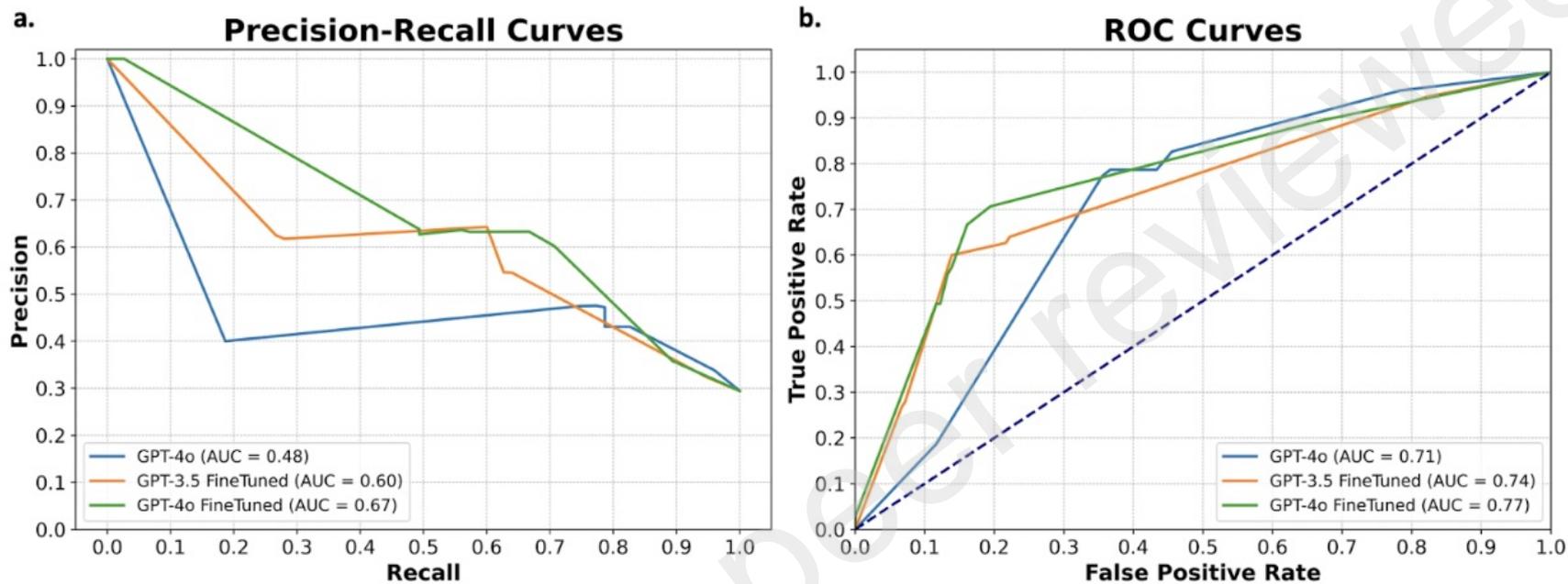


Figure 4: Precision-Recall Curves and ROC Curves for GPT-4o Base, GPT-3.5 Fine-Tuned and GPT-4o Fine-Tuned Models.

Improving Conservation Efficiency: Accelerating Groundwater Sustainability Plan Reviews

Using Large Language Models

Ryan Bernstein¹, Seneth Waterman¹,
Kirk R. Klausmeyer², Nicholas Murphy², Melissa M. Rohde^{3,4}, Cody Carroll¹

¹ *University of San Francisco*

² *The Nature Conservancy of California, San Francisco, CA*

³ *Rohde Environmental Consulting, LLC, Seattle, WA*

⁴ *SUNY College of Environmental Science and Forestry, Syracuse, NY*

Highlights

- LLMs accelerate conservation tasks like reviewing technically dense Groundwater Sustainability Plans.
- A fine-tuned GPT-4o model achieved 73% agreement with human reviewers in tests.
- Hybrid Human/LLM reviews can cut marginal GSP review time by a factor of 240.
- We developed an automated LLM pipeline to review plans for \$1.40 per plan.
- LLMs have the potential to *support but not replace* human environmental reviews.

Abstract

Background: The effective implementation of environmental policy relies on thorough review and public consultation, yet the analysis of lengthy, technical documents is a resource-intensive process that creates a significant barrier to entry for many civil society organizations and overburdened agencies. This analytical bottleneck can limit oversight and hinder the achievement of sustainability and environmental justice goals.

Objective & Methods: This study evaluates the potential for Large Language Models (LLMs) to augment human capacity for large-scale policy review. Using the mandated assessment of 108 Groundwater Sustainability Plans in California as a case study, we compared the results of a custom fine-tuned LLM against a comprehensive, multi-year review previously conducted by expert human analysts.

Results: Our method led to a significant acceleration of the review process. Our LLM performed an initial analysis of a plan in under two minutes, a 240-fold increase in speed over the eight-hour average for a human expert. The LLM's qualitative assessments achieved a 73% agreement rate with the human-led benchmark, demonstrating substantial utility in identifying key policy components and deficiencies.

Conclusion: LLMs represent a transformative tool for the science-policy interface, not as a replacement for human expertise but as a powerful accelerator in a human-in-the-loop system. By drastically reducing the initial effort of document review, this technology can enhance the capacity of organizations to participate in environmental decision-making, helping to overcome persistent bottlenecks in policy monitoring and evaluation. This approach has broad, international implications for improving the efficiency and equity of environmental governance.

Introduction

Advances in large language models (LLMs) offer promise for increasing efficiency in many industries. LLMs are capable of human-like reading, writing, and communication skills, making them valuable for tasks like document summarization, classification, and policy analysis (Radford et al., 2018; Vaswani et al., 2017; Hadi et al., 2023). In fields like education, medicine, and finance, LLMs have already shown their utility in reducing costs and enhancing productivity (Mello et al. 2023; Biswas, 2023; Wu et al., 2023; Thawkar et al., 2023; Goyal et al. 2023). LLMs have also passed the Bar Exam and can prove mathematical theorems at the olympiad level, though not always consistently (Bommarito et al. 2022; Frieder et al. 2024). The use of LLMs is increasingly being explored for sustainable development and environmental conservation, including applications in energy management, emission forecasting, waste management, pollution classification, environmental monitoring and biodiversity conservation (Grief et al. 2025; Ullah et al. 2024). These new tools offer exciting opportunities for progress on difficult social and environmental challenges.

One significant environmental challenge is the decline of freshwater species abundance across the globe (WWF, 2024). This issue is most apparent in regions where human needs for water conflict with ecosystem needs. Groundwater is a major source of water in semi-arid and arid regions, where climate change and growing human demand during droughts impact freshwater species. Groundwater-dependent ecosystems (GDEs), including rivers, wetlands, and springs, are crucial to maintain biodiversity, provide critical habitat, improve water quality, and store carbon (Eamus and Froend, 2006; Howard et al., 2023; Rohde et al., 2024; Yang et al., 2020). These ecosystems are highly sensitive to groundwater depletion, which has been exacerbated by unsustainable extraction practices (Rohde et al., 2021). Efficient groundwater management is pivotal for environmental sustainability and ecological well-being, particularly in heavily populated and agricultural regions like California.

California's Sustainable Groundwater Management Act (SGMA), enacted in 2014, requires Groundwater Sustainability Agencies (GSAs) to develop Groundwater Sustainability Plans (GSPs) for high- and medium-priority basins (Harrer et al., 2020). These GSAs are charged with both formulating and executing GSPs, and are required to meet groundwater sustainability goals by 2040. GSPs are vital for outlining groundwater management plans to avoid the six undesirable results of SGMA (chronic lowering of groundwater levels, reduction in groundwater storage, seawater intrusion, land subsidence, degradation of water quality, and depletion of interconnected surface water), and assessing the impacts of groundwater management on beneficial users (domestic, agricultural, municipal, environmental, tribes, and disadvantaged communities) (Harrer et al. 2020). GSP reviews by state regulators and interested parties are necessary to determine if GSAs have sufficiently met their regulatory requirements under SGMA; however they require significant time and financial investment by technical experts (Dumas Leslie, 2017). The evaluation of the 108 GSPs completed by The Nature Conservancy (TNC), partner organizations, and contractors during the public comment periods (issued by each GSA and the California Department of Water Resources) required significant financial and time investments to review and write comment letters (Perrone et al., 2023).

This study evaluates the ability of LLMs to review GSPs as accurately as human evaluators, through the lens of environmental review. Using a sample of 5 GSPs, we compared previous human-led evaluations to LLM-assisted reviews, considering accuracy and efficiency as metrics. By investigating different LLM configurations through experiments with prompt engineering, retrieval-augmented generation, and fine-tuning, we assessed the LLM's ability to replicate human decision-making with the goal of reducing contractor review hours and monetary costs. We then scaled the most promising method up to the full set of 56 GSPs with scoring rubrics that were technically feasible to review (Figure 1).

This work seeks to address the pressing need for scalable and cost-effective tools to facilitate and accelerate environmental review of resource management plans. By exploring the capabilities and limitations of LLMs in the context of GSPs, this study aims to illustrate their potential to support human-led environmental policy review and adaptive management to address the pressing environmental challenges facing the globe.

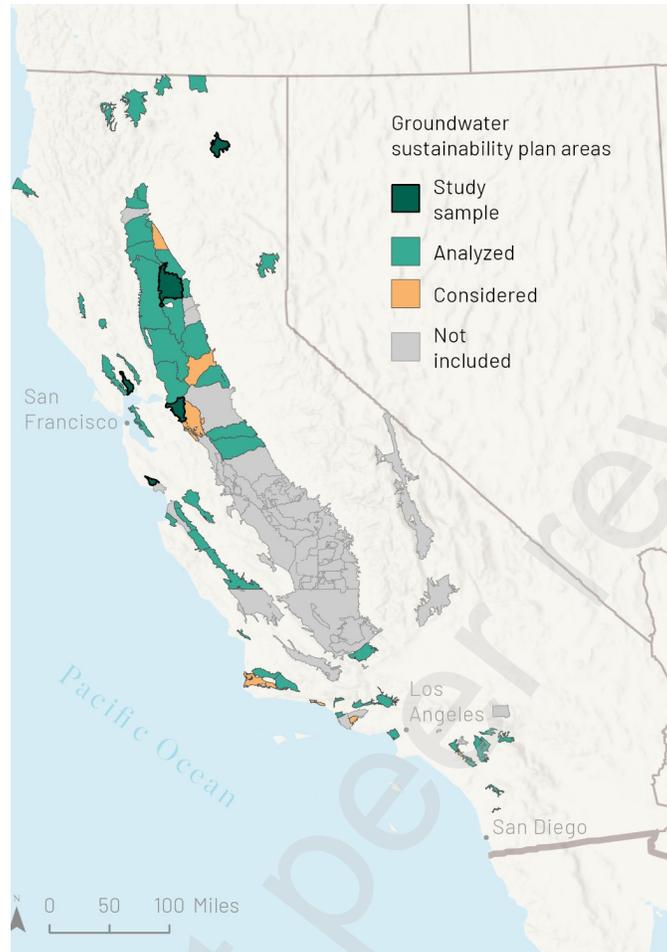


Figure 1. Locator map of California to highlight the groundwater sustainability plan (GSP) areas reviewed in this study. The 5 sample GSPs reviewed in most detail are in dark green and the 49 other GSPs with rubrics that we analyzed are in light green. We attempted, but were unable, to analyze 9 GSPs due to technical issues with the report PDF or the human scored rubric (shown in yellow). The remaining GSPs shown in grey were reviewed by humans using different methods and did not have a standardized scoring rubric so they were not included in this study.

Approach

We acquired a sample of 65 GSPs and their corresponding human-generated scoring rubrics from a review led by TNC (Perrone et al, 2023). For our experimentation, we looked at a randomly selected

subset of 5 representative GSPs (Big Valley, Butte, East Contra Costa, Sonoma, Santa Margarita) and their 5 corresponding GSP scoring rubrics (Figure 1). The sample GSPs come in the form of lengthy PDF documents, ranging from 212 to 459 pages long, with associated maps and scientific figures. For the scope of this paper, we only focused on leveraging the text of the PDFs and left the multimodal problem of incorporating image data for future research. GSP scoring rubrics each consist of 70 rows that have corresponding questions, one word answers as responses (“Yes/No/Somewhat”), relevant excerpts of text drawn from GSP entries, and references to relevant pages or sections of the plan. They are also divided into 5 Topics: Environment (17 entries), Disadvantaged Communities (12 entries), Climate Change (7 entries), Ecosystems (7 entries) and Tribes (6 entries). Seventeen of these rows were removed as they contained questions that necessitated the interpretation of map data for their answers, such as maps containing domestic well locations. Table 1 is an example of a row recreated from a rubric. It displays the evaluation criteria for Disadvantaged Community (DAC) stakeholder engagement in the GSP development process. The criteria include the level of engagement, spectrum of responses, and relevant search terms. The table also provides a detailed excerpt from the GSP document and the corresponding answer from the human reviewer, which in this case is classified as "Somewhat."

Table 1: Example response from the human-review of GSP rubric. This example is taken from the review of the Butte GSP.

| Topic | Criteria | Spectrum | Key Search Terms | Answer | Relevant Text from GSP |
|----------------------------------|---|--|--|----------|--|
| Disadvantaged Communities (DACs) | Does the GSP document how DAC stakeholders were given opportunities to engage in the GSP development process? If so, please describe the level of engagement (i.e., whether stakeholders were on an advisory committee, GSA board, working group) in the Notes. Answer according to level of engagement: Yes = Involve, | No = no reference to DACs, no reference to stakeholder engagement during development; Somewhat = no specific reference to DAC as a stakeholder, provides stakeholder engagement; Yes = references DACs and how they were engaged specifically during development beyond informing DACs | Disadvantaged Community, DAC, stakeholder, beneficial user | Somewhat | "Lassen and Modoc counties are fulfilling their unfunded, mandated roles as Groundwater Sustainability Agencies (GSAs) to develop this Groundwater Sustainability Plan (GSP) after exhausting its administrative challenges to the California Department of Water Resources' (DWR's) determination that Big Valley qualifies as a medium-priority basin. Both counties are disadvantaged, have declining populations, and have no ability to cover the costs of GSP development and implementation. These disadvantaged communities are on the losing end of the digital divide. While the GSAs made every attempt to conduct BVAC meetings with the ability for remote public participation, there were still major logistical and technical challenges with both conducting such meetings and members of the public participating. Those participants that had internet connectivity frequently could not hear |

| | | | | | |
|--|---|--|--|--|---|
| | Collaborate, Empower; Somewhat = Inform and Consult; No = No mention of engagement. | | | | or understand the dialogue in the Big Valley community venues and could not interact in the most effective way. However, the GSAs made the best of the circumstances and addressed all comments provided through the various means. The GSAs recognized the obstacles presented by the COVID pandemic early in the efforts to develop a GSP and were proactive in reaching out to both the Governor and Legislature to identify potential solutions. ...” |
|--|---|--|--|--|---|

Our main scientific inquiry is to test whether a LLM-based application programming interface (API) is capable of answering the questions in the scoring rubric in the same way as a human reviewer.

Our feature engineering consisted of altering the format of the provided questions in a way that an LLM could understand, while also ensuring that we maintained the integrity of the questions’ original content. The GSP criteria and spectrum were combined, the spectrum was converted to complete sentences, and all acronyms were defined. Table 2 illustrates an example of how we incorporated this feature engineering.

Table 2: Example of prompt formation using the Criteria and Spectrum columns found in the GSP evaluation documents to create a “Question” to feed to the LLM.

| Criteria | Spectrum | Question |
|---|--|--|
| Does the GSP document how DAC stakeholders were given opportunities to engage in the GSP development process? If so, please describe the level of engagement (e.g., whether stakeholders were on an advisory committee, GSA board, working group) in the Notes. Answer according to level of engagement: Yes = Involve, Collaborate, Empower; Somewhat = Inform and Consult; No = No mention of engagement. | No = no reference to DACs, no reference to stakeholder engagement during development; Somewhat = no specific reference to DAC as a stakeholder, provides stakeholder engagement; Yes = references DACs and how they were engaged specifically during development beyond informing DACs | Does the Groundwater Sustainability Plan (GSP) document how Disadvantaged Community (DAC) stakeholders were given opportunities to engage in the GSP development process? Provide a one word answer to this question using the following Spectrum: No = no reference to DACs, no reference to stakeholder engagement during development; Somewhat= no specific reference to DAC engagement but lists DAC as a stakeholder, provides stakeholder engagement details; Yes= references DACs and how they were engaged specifically during development |

Our analytic process consisted of three primary steps:

1. Prompt Engineering, which consisted of altering the format of the general instructions to the LLM to consider when answering the questions
2. Retrieval-Augmented Generation (RAG) Tuning, which consisted of preprocessing the GSP content in a way that allowed our LLM to focus on the pertinent subsections of the document when it formulated its response; and
3. Fine Tuning, in which we honed and refined higher-order feature representations by training and validating a custom GPT on the human-generated GSP rubrics.

Methods

Large Language Models

The number and capabilities of foundational LLMs have increased substantially since the release of ChatGPT in 2022. At the time we started the research for this study, only OpenAI provided API access to the model with features we required like batch processing and fine tuning, so we focused on OpenAI's foundational models (GPT-3.5, GPT-4, GPT-4o). Subsequent studies have shown that OpenAI's LLMs have high performance on general tasks (Radford et al., 2018; Vaswani et al., 2017; Hadi et al., 2023) and the retrieval augmented generation tasks like the ones used in this study (Ke et al. 2025).

Vector Store

A vector store is a specialized database that can store and manage our text embedding vectors. Our text embedding vectors are numerical vectors that represent the semantic meaning and context for each of the GSP documents. Vector stores are optimized for handling high-dimensional data, enabling rapid querying and retrieval (LangChain Team, 2024). Because our text embeddings have high dimensionality, and

efficiency is a priority in this context, we choose to use a vector store to accelerate LLM response time and ensure future scalability.

Our study compared model performance between two candidate choices of vector stores: FAISS (Meta AI, 2024) and CHROMA (LangChain Team, 2024). We utilized the GPT-3.5 Turbo model to assess both vector stores' performance in terms of accuracy.

Prompt Engineering

Rewording the prompt in search of the optimal way to ask ChatGPT queries is known as prompt engineering. White et al. (2023) recently developed a catalog of prompt engineering techniques presented as patterns, offering solutions to common challenges encountered when interacting with LLMs. These prompt patterns serve as a method similar to software patterns, providing reusable solutions to common problems in output generation and interaction with LLMs within a specific context. To improve accuracy, we experimented with prompt engineering by refining both the prompts and the set of instructions given to the model.

We first experimented with adjusting the language used in the queries to encourage more critical responses and to prevent the model from being an excessively permissive judge. To refine the model's response and scrutiny, prompts were modified to emphasize "skeptical evaluation" and "critical assessment" in the instructions.

We also experimented with standardizing prompts to assess whether keeping question formats consistent and spelling out acronyms had an effect on performance. Given that human reviewers had only seen the original prompts, we had limited flexibility in changing their content and subject matter. However, many

of the 70 prompts contained inconsistent formatting and numerous acronyms. We sought to determine if using a consistent format and spelling out acronyms would enhance performance.

Retrieval-Augmented Generation

To optimize the RAG step, we conducted a series of experiments using FAISS as the vector store while experimenting with other model parameters. We first investigated the impact of different vector sizing strategies on model accuracy. Since the total text of the GSPs were too big for the LLM to evaluate all at once, we needed a method to send the most relevant chunks of text to the model. Vector sizing strategies are ways to break up long text documents into smaller chunks. We compared page-based vector sizing, where the vector size was adjusted to match the entire page length, with chunk-based vector sizing, which used a fixed chunk size of 1,000 words and a chunk overlap of 200 words. Additionally, we evaluated how varying the number of context vectors retrieved during inference affected the model's accuracy. Specifically, we tested the impact of retrieving the top 5, 10, and 15 most similar vectors to the input query to determine the optimal amount of contextual information to include in the model's response generation. Token size limits were met for values exceeding 15. Performance was compared across GPT-3.5 Turbo and GPT-4. Finally, we tested performance when including the appendices of the GSPs along with the main body of the report.

Problem Simplification

We hypothesized that the three-point spectrum (Yes/No/Somewhat) might be too nuanced for a large language model trained on GSPs in which the Somewhat category is ambiguous, subjective, and evaluator-dependent. To avoid this complication, we simplified the task by merging the "No" and "Somewhat" categories into a single "No" category, effectively creating a binary classification problem. For reference, 105 of the 350 answers in our 5 sample GSPs were initially labeled as "Somewhat."

Confidence Levels

After we converted our task into a binary classification problem, we altered our prompts so the GPT also provided soft predictions, in order to compare Receiver Operator Curves and Precision-Recall Curves across models. In addition, we hypothesized that forcing the model to consider its confidence level regarding a prediction could improve accuracy. The model instructions that we used are replicated below:

“You are a skeptical environmental scientist, tasked with answering questions about a section from a Groundwater Sustainability Plan (GSP) document. You are required to provide your confidence level along with the response. Format your response as 'X, Z' where 'X' is either 'Yes' or 'No' and 'Z' is your confidence level, which can be one of the following options: ["Extremely Confident, 100%", "Very Confident, 85%", "Fairly Confident, 75%", "Modest Confidence, 60%", "Random Guess, 50%"]. ' Answer the questions objectively, adhering to the provided spectrums. You should only state you are "Extremely Confident, 100%" if it is irrefutably true that your answer is correct according to the Groundwater Sustainability Plan.”

Fine-Tuned Models

Finally, we experimented with fine-tuning the GPT models used in GSP evaluation. LLM fine-tuning involves adapting a pre-trained large language model to a specific task or domain by training it further on a smaller, task-specific dataset, like our GSPs (OpenAI, 2024). This process leverages the model's existing knowledge while refining its performance on the desired application, enhancing its accuracy and relevance. We trained a fine-tuned GPT-3.5 model and a GPT-4o model (8 epochs, batch_size =1, learning rate =.5), using the 10 most similar pages of GSP text (determined via cosine similarity) as inputs. The fine-tuned GPT models were trained on the Butte and Santa Margarita GSP rubrics, using a combination of the “Relevant Text from GSP” rubric column and our engineered prompts as input. To assess the accuracy of the fine-tuned models we excluded the two GSPs that were used in fine tuning and added two more GSPs with human score responses to the standard questions (San Luis Obispo, Fillmore).

Scaling Up

After completing experiments on a sample of GSPs, we applied the best performing model to the full set of 65 GSPs using python code and OpenAI's API to automate the process. 2 of the GSP Rubrics (Tracy and Bedford) did not have a complete set of answers and we had technical formatting issues with 7 of the GSP PDF files (Los Molinos, South American, Santa Ynez Valley - Central, Santa Ynez Valley - Western, Pleasant Valley, Montecito, Carpinteria), thus we could not analyze them. 2 of the GSPs were used to fine-tune the model. This left 56 GSPs to score with the LLM and compare to the human review. Code for this pipeline is publicly available at: [https://github.com/ryanoh999/ChatGDE_Ryan].

Results

Vector Stores, RAG Parameters, and Model Versions

GPT 3.5 Turbo's "Yes/No/Somewhat" results indicated that FAISS outperformed CHROMA, achieving an accuracy of 52.83% compared to 33.96%. Thus, FAISS was selected as the preferred vector store for subsequent testing due to its superiority in handling and retrieving relevant information from the dataset. Including the appendices did not improve accuracy (see Supporting Information).

Initial testing of hyperparameter choices and models were evaluated on a test set of 255 questions from the 5 sample GSPs using "Yes/No/Somewhat" results (Fig. 2). Evaluating the GPT 3.5 Turbo model with different parameter choices showed that accuracy was maximized by using the report pages to split documents into text chunks instead of arbitrary chunks that were 1,000 characters long (Fig 2a). We then evaluated the number of chunks used to answer a question, and found that 10 chunks provided the best accuracy (Fig 2b). We also saw that depending on the GSP, GPT versions 4 or 4o outperformed the GPT-3.5 Turbo in terms of accuracy after optimizing hyperparameters (Fig. 2c).

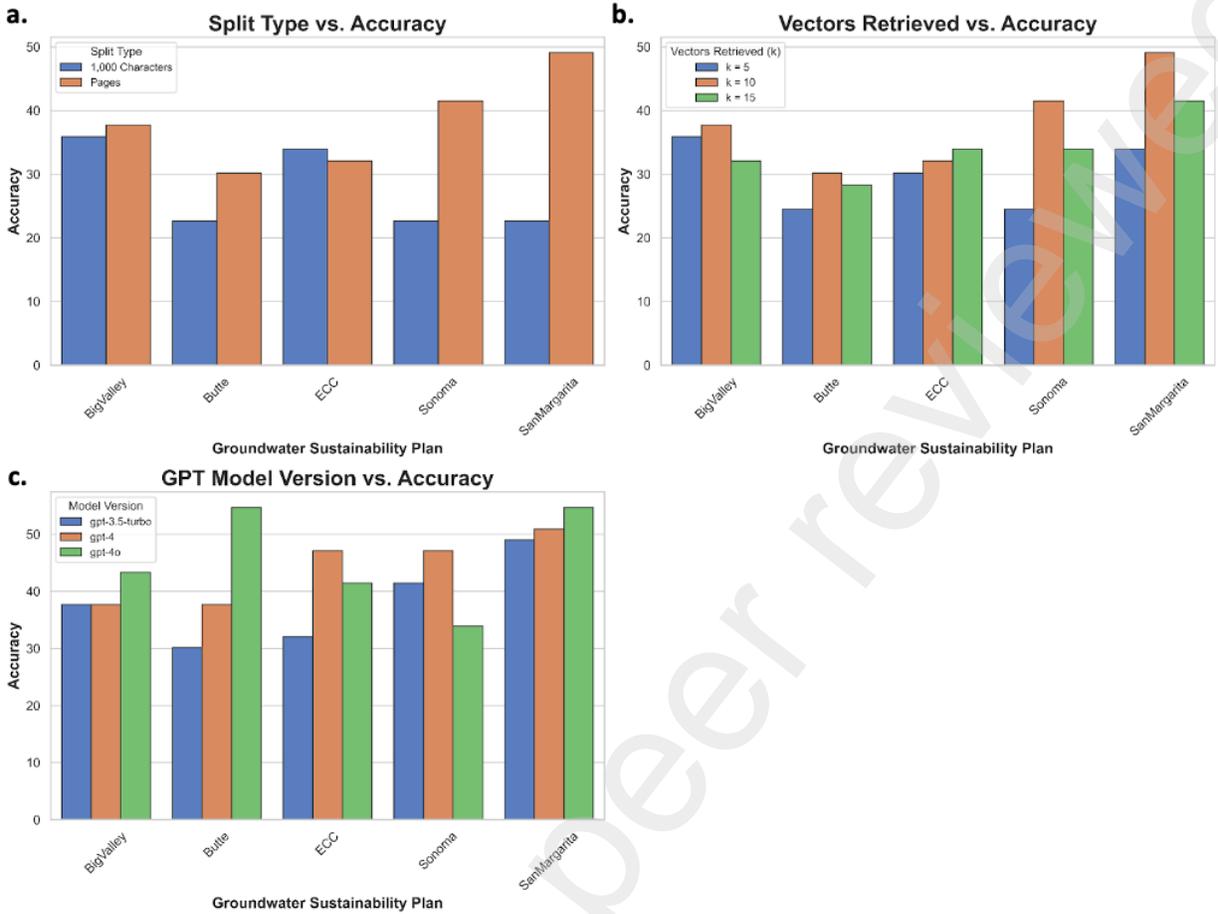


Figure 2: Accuracy(percent of questions answered correctly *100) comparison across model and hyperparameter choices for non-fine-tuned models. Panels a. and b. depict accuracy for experiments on split type to define vector chunk size and vector retrieval (chunk count) size using GPT-3.5 Turbo, while panel c. compares accuracy across model versions.

Pre-trained vs. Fine-Tuned Binary Classification

For the next set of experiments we simplified the classification problem from “Yes/No/Somewhat” results to "Yes/No" binary results by collapsing "No" and "Somewhat" answers together. We then compared the standard GPT-4o model with two fine-tuned models (Figure 3). The two of the original 5 sample GSPs that were used for fine tuning (Butte and Santa Margarita) were excluded from this accuracy assessment, and two new GSPs were added (Fillmore and San Luis Obispo). Both the GPT-3.5 and GPT-4o Fine-

Tuned models achieved an average accuracy of 72.9%, while the GPT-4o Base model only achieved an average accuracy of 62.7% using “Yes/No” results. This suggests that fine tuning had a positive impact on our results. This is because fine-tuning allows the model to learn domain-specific nuances and patterns that are not captured in the general training of the GPT models. If we assume that the human reviewers have a unified approach towards their evaluation, a fine-tuned model will be able to approximate these reviewers on unseen data.

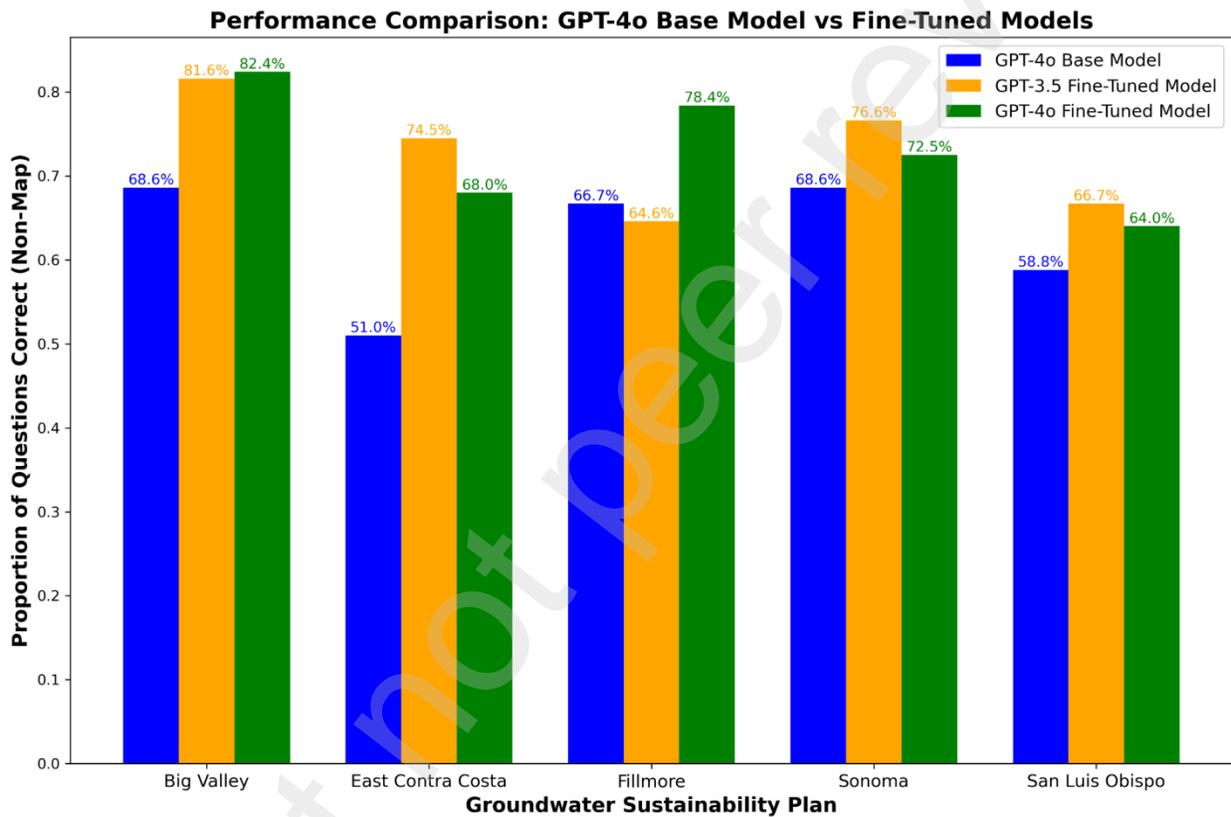


Figure 3: Performance Comparison:GPT-4o Base Model vs GPT-3.5 Fine-Tuned Model vs GPT-4o Fine-Tuned Model.

Table 3: Comparing Overall Model Accuracies with Bonferroni Correction across the 5 GSPs.

Significance level notation after applying Bonferroni correction: not significant (ns), 0.05 (*), 0.01 (**), 0.001(***)).

| Model Pair | Model 1 Accuracy | Model 2 Accuracy | Accuracy Difference | Z-Score | P-Value | Significance |
|---------------------------------|------------------|------------------|---------------------|---------|---------|--------------|
| 4o Base vs 3.5 Fine-Tuned | .627 | .729 | .102 | -2.465 | .0137 | * |
| 4o Base vs 4o Fine-Tuned | .627 | .729 | .102 | -2.465 | .0137 | * |
| 3.5 Fine-Tuned vs 4o Fine-Tuned | .729 | .729 | 0 | 0 | 1 | ns |

Figure 4a displays the Precision-Recall curves for three models: GPT-4o Base, GPT-3.5 Fine-Tuned, and GPT-4o Fine-Tuned. GPT-4o Fine-Tuned demonstrates the highest Area Under the Curve (AUC) at 0.67, followed by GPT-3.5 Fine-Tuned with an AUC of 0.60. GPT-4o Base has the lowest performance, with an AUC of 0.48. This graph shows how well each model balances precision and recall, with GPT-4o Fine-Tuned performing best in terms of having a high precision when also having a higher recall.

Figure 4b shows the ROC curves for the same models. Here, GPT-4o Fine-Tuned once again performs the best, achieving an AUC of 0.77, while GPT-3.5 Fine-Tuned has an AUC of 0.74. GPT-4o Base has the lowest AUC at 0.71. The ROC curve illustrates the trade-off between the true positive rate (sensitivity) and the false positive rate, with GPT-4o Fine-Tuned showing the best overall discrimination ability between positive and negative classes.

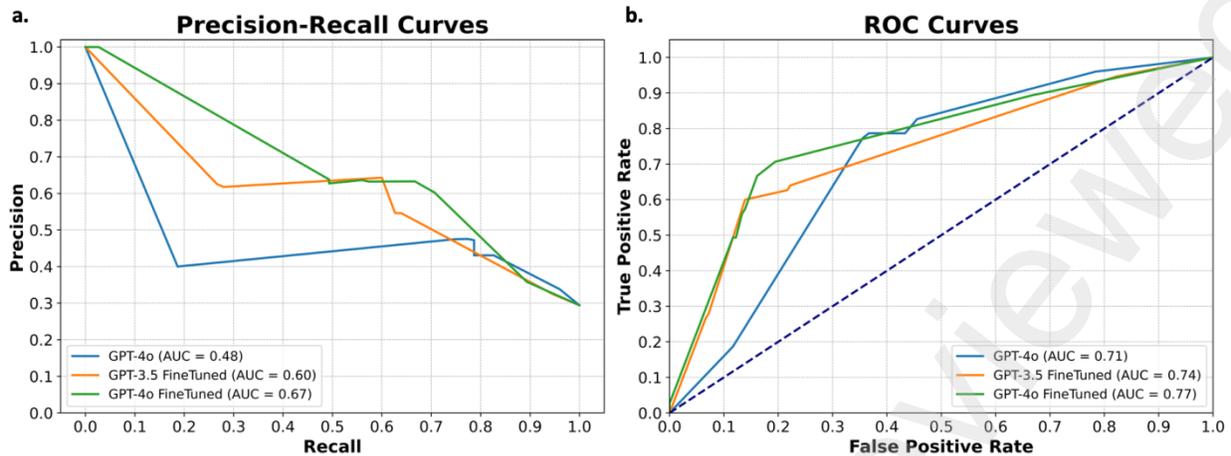


Figure 4: Precision-Recall Curves and ROC Curves for GPT-4o Base, GPT-3.5 Fine-Tuned and GPT-4o Fine-Tuned Models.

The GPT-4o Base model shows a relatively high number of 114 true negatives (114/255 or 44.71%), with 66 false positives (66/255 or 25.88%). For the positive cases, there are 59 true positives (59/255 or 23.14%) and 16 false negatives (16/255 or 6.27%), indicating a decent but moderate balance between correctly classified positive and negative responses.

The GPT-3.5 Fine-Tuned model improves on the false positives with a lower count of 39 (39/255 or 15.29%), but there are still 28 false negatives (28/255 or 10.98%). The true negative count is higher at 141 (141/255 or 55.29%), and the true positives are 47 (47/255 or 18.43%), showing a slightly better balance in classifying negative responses but a small drop in positive classification accuracy compared to the GPT-4o Base .

The GPT-4o Fine-Tuned model achieves the highest true negative count (155/255 or 60.78%) with the lowest number of false positives (25/255 or 9.80%). However, the false negatives are slightly higher than that of the other models at 32 (32/255 or 12.55%), while the true positives are 43 (43/255 or 16.86%).

This suggests that GPT-4o Fine-Tuned excels at identifying negative responses but has a slightly lower ability to capture positive responses accurately when compared to the other models.

Table 4: Average classification metrics for the 5 case study GSPs for Base GPT-4o, GPT-3.5 Fine-Tuned and GPT-4o Fine-Tuned Models. The best model according to each metric is bolded in each column.

| Model / Metric | Accuracy | AUROC | AUCPR | Precision | Recall | F1 |
|-----------------------|-----------------|--------------|--------------|------------------|---------------|-------------|
| GPT-4o Base | 62.7% | 0.71 | 0.48 | 0.47 | 0.79 | 0.59 |
| GPT-3.5 Fine-Tuned | 72.9% | 0.74 | 0.6 | 0.63 | 0.63 | 0.58 |
| GPT-4o Fine-Tuned | 72.9% | 0.77 | 0.67 | 0.63 | 0.57 | 0.60 |

Scalability

We further validated our findings by testing our best-performing model, GPT-4o Fine-Tuned, on the full cohort of 56 GSPs pre-evaluated by TNC to assess its scalability (Table S1). To summarize performance, we generated average performance metrics for the 54 GSPs that were not used to fine tune the model (Table 4). The model demonstrated consistent but imperfect performance, achieving an average AUC-ROC of 0.70 and an accuracy of 69%. While these metrics are below the 0.77 AUC-ROC and 72.9% accuracy recorded with our initial five-GSP subset, the sustained level of performance across the larger dataset strongly suggests our results are generalizable and not skewed by a small sample size.

Table 4: Average performance metrics for 54 GSPs for GPT-4o Fine-Tuned Model.

| Model / Metric | Accuracy | AUROC | AUCPR | Precision | Recall | F1 |
|-----------------------|-----------------|--------------|--------------|------------------|---------------|-----------|
| Fine-Tuned 4o | 69.0% | 0.70 | 0.49 | 0.50 | 0.67 | 0.57 |

Discussion

In this paper, we explored the application of LLMs for evaluating GSPs through a series of experiments. Specifically, we tested the impact of different context retrieval strategies, various prompt engineering

strategies, and the choice of different pre-trained models on model accuracy. Furthermore, we investigated the effect of model fine-tuning on accuracy, precision-recall, and ROC metrics.

Based on our results, we conclude that while an LLM can offer accurate results that are significantly better than random guessing, it is not currently able to replace a human reviewer. We found that modifying system instructions to emphasize skepticism influenced both the distribution and accuracy of our model's responses. Giving the model greater amounts of context as opposed to smaller amounts of context, and using page length as opposed to chunk size gave us better accuracy. Context-window limits will almost assuredly be eliminated as LLMs continue to develop and scale, so we can assume that accuracy will continue to improve.

Our fine-tuning results suggest that an LLM specifically fine-tuned on GSPs will outperform a pre-trained LLM that has not explicitly learned from this domain-specific data. Interestingly, updates to different versions of ChatGPT had a relatively insignificant impact on performance for this particular task. Our best-performing model, with regards to accuracy, AUCROC, and AUCPR, was a fine-tuned GPT-4o model. This model was tuned over 8 epochs with a batch size of 1 and a learning rate of 0.5, utilizing the 10 most similar pages of the text as RAG inputs. Crucially, when this model was applied to our entire cohort of 56 GSPs to evaluate scalability, it demonstrated strong generalization capabilities, with only a minor performance decrease of less than 10% in both AUC and Accuracy compared to its performance on our smaller experimental subset.

A major concern discovered during the process of evaluating the LLMs was that model output varied substantially even under consistent parameters. The stochastic nature of LLM responses is due to transformer architecture, developer updates and model variations (Su, W. (2025)). If researchers and practitioners are going to effectively implement LLMs as official evaluators, there is a clear need for model developers and the companies behind them to address and control for these inconsistencies and

provide more transparent version control in their APIs. Alternatively, one could implement a “Panel of Experts” approach to stabilize model performance. This process uses multiple unique LLMs (a “panel”) to gather predictions, and then makes a final prediction that is aggregated from the contributions of each LLM (Sourcery AI, 2024). Such an approach could further improve the robustness, consistency, and accuracy of our predictions.

Our results suggest the possibility of a hybrid approach to GSP review in which human reviewers use an LLM-based application to help quickly identify relevant passages and give them a second opinion on evaluation. The implementation of groundwater sustainability plans in California has multiple future engagement points, as annual reports and five-year assessments and plan re-evaluations are required for medium and high-priority groundwater subbasins. LLM-assisted plan review has the potential to increase The Nature Conservancy’s future engagement efficiency.

More broadly, this research explores the applicability of LLM-assisted content review. One of the main benefits of integrating LLMs is the time and cost savings associated with reviewing large volumes of text. For example, the time required for humans to manually review each GSP was approximately 8 hours, whereas the time required to review each GSP with the fine-tuned LLM review averaged just under 2 minutes and cost \$1.40 per plan. If LLMs can significantly reduce the duration of human review processes, organizations that review and comment publicly on policy documents (e.g., environmental, advocacy, and research and policy institutes) would be able to save considerable amounts on future time and resource requirements. Organizations may be able to increase their scale of content review, stakeholder education and engagement practices, and policy evaluation as a result of LLM-assisted review workflows.

LLMs could also serve as an effective science education tool for less-technically-resourced stakeholder groups. LLMs have been identified as effective education tools, providing personalized topic instruction,

adaptive feedback, and language translation capabilities (Bektik et al. 2024). Previous research has shown that increased representation of underrepresented stakeholder groups in environmental review processes improves outcomes for groundwater sustainability policy (Perrone, Rohde, Wagner et al. 2023). However, groundwater sustainability plans are often inaccessible to less-technically resourced stakeholder groups as a function of both length (GSPs typically exceed 200 pages), and technical subject matter (GSP subject matter assumes a background in the foundational concepts of groundwater hydrology).

Additional work remains in terms of lifting evaluation metrics, but this research establishes the general framework and proof of concept for the utility of LLMs in evaluating GSPs and more generally, in the field of environmental reviews. Future research should explore ensemble approaches (panels of multiple LLM judges), reinforcement learning from human feedback, and models outside of the OpenAI ecosystem.

Conclusion

Our study demonstrates that LLMs can support the content review process for environmental advocacy, reducing time and monetary costs and achieving reasonable agreement with human reviewers. While not a total replacement for experts, these tools can support human capacity in addressing pressing international environmental policy challenges. With ongoing model improvements, better context retrieval, and careful and relevant fine-tuning, LLM-assisted reviews hold promise for enabling non-profits and agencies to do more with limited resources, ultimately enabling more efficient and sustainable natural resource management.

Acknowledgements

The authors would like to thank Alicia Canales for help with the cartography to generate the locator map.

References

- Bommarito, M, et al. (2022) GPT Takes the Bar Exam. <https://arxiv.org/pdf/2212.14402.pdf>
- Bektik, D., Edwards, C., Whitelock, D., & Antonaci, A. (2024). Use of LLM tools within higher education: Report 1.
- Biswas, S. S. (2023). Role of Chat GPT in Public Health. *Annals of Biomedical Engineering*, 51(5), 868–869.
- Dumas, L. (2017). Implementing SGMA—An Update on California’s Foray into Groundwater Regulation. *World Environmental and Water Resources Congress 2017*.
- Eamus, D., & Froend, R. (2006). Groundwater-dependent ecosystems: the where, what and why of GDEs. *Australian Journal of Botany*, 54(2), 91-96.
- Frieder, S., Pinchetti, L., Griffiths, R. R., Salvatori, T., Lukasiewicz, T., Petersen, P., & Berner, J. (2024). Mathematical capabilities of chatgpt. *Advances in neural information processing systems*, 36.
- Goyal, T, Li, J., Durrett, G. (2023). News Summarization and Evaluation in the Era of GPT-3. arXiv: 2209.12356.
- Greif, L., Röckel, F., Kimmig, A., & Ovtcharova, J. (2025). A systematic review of current AI techniques used in the context of the SDGs. *International Journal of Environmental Research*, 19(1), 1.
- Hadi, M. U., et al. (2023). Large language models: a comprehensive survey of its applications, challenges, limitations, and future prospects. *Authorea Preprints*.
- Harrer, S., Howard, J., Rohde, M., et al. (2020). Groundwater Dependent Ecosystems under the Sustainable Groundwater Management Act. <https://www.scienceforconservation.org/assets/downloads/GDEsUnderSGMA.pdf>
- Howard, J. K., et al. (2023). Ecosystem services produced by groundwater dependent ecosystems: a framework and case study in California. *Frontiers in Water*, 5, 1115416.
- Ke, Y. H., Jin, L., Elangovan, K., Abdullah, H. R., Liu, N., Sia, A. T. H., ... & Ting, D. S. W. (2025). Retrieval augmented generation for 10 large language models and its generalizability in assessing medical fitness. *npj Digital Medicine*, 8(1), 187. <https://www.nature.com/articles/s41746-025-01519-z.pdf>.
- LangChain Team. (2024) Vector Stores. Retrieved from https://js.langchain.com/v0.1/docs/modules/data_connection/vectorstores/.
- Mello, R, et al. (2023) Education in the age of Generative AI: Context and Recent Developments. arXiv preprint arXiv:2309.12332.

Meta AI. (2024) Faiss. Retrieved from <https://ai.meta.com/tools/faiss/>.

OpenAI. (2024) Fine-tuning. Retrieved from <https://platform.openai.com/docs/guides/fine-tuning>.

Perrone, D., Rohde, M. M., Hammond Wagner, C., Anderson, R., Arthur, S., Atume, N., ... & Remson, E. J. (2023). Stakeholder integration predicts better outcomes from groundwater sustainability policy. *Nature Communications*, 14(1), 3793.

Radford, A., Narasimhan, K., Salimans, T., & Sutskever, I. (2018). Improving language understanding by generative pre-training. OpenAI. https://cdn.openai.com/research-covers/language-unsupervised/language_understanding_paper.pdf

Rohde, M. M., Biswas, T., Housman, I. W., Campbell, L. S., Klausmeyer, K. R., & Howard, J. K. (2021). A machine learning approach to predict groundwater levels in California reveals ecosystems at risk. *Frontiers in Earth Science*, 9, 784499.

Rohde, M. M., Albano, C. M., Huggins, X., Klausmeyer, K. R., Morton, C., Sharman, A., ... & Stella, J. C. (2024). Groundwater-dependent ecosystem map exposes global dryland protection needs. *Nature*, 632(8023), 101-107.

Sourcery AI. (2024) Better LLM Prompting using the Panel-of-Experts. Retrieved from <https://sourcery.ai/blog/panel-of-experts/>

Su, W. (2025). Do Large Language Models (Really) Need Statistical Foundations? arXiv preprint arXiv:2505.19145.

Thawkar, O, et al. (2023) Xraygpt: Chest radiographs summarization using medical vision-language models. arXiv preprint arXiv:2306.07971.

Ullah, F., Saqib, S., & Xiong, Y. C. (2024). Integrating artificial intelligence in biodiversity conservation: bridging classical and modern approaches. *Biodiversity and Conservation*, 1-21.

Vaswani, A, et al. (2017) Attention Is All You Need. arXiv preprint arXiv:1706.03762.

Wu, S, et al. (2023) Bloomberggpt: A large language model for finance. arXiv preprint arXiv:2303.17564.

White, J, et al. (2023) A Prompt Pattern Catalog to Enhance Prompt Engineering with ChatGPT. arXiv:2302.11382.

Wu, Y, et al. (2016). Google's Neural Machine Translation System: Bridging the Gap between Human and Machine Translation. arXiv preprint arXiv:1609.08144.

Yang, K, et al. (2020) Recent progress in the design and application of MOF-based materials for highly selective CO₂ capture. ACS Energy Letters, 5(11), 3490-3496.

Supporting Information

Groundwater Dependent Ecosystems in California:

Groundwater Dependent Ecosystems (GDEs) are natural ecological communities that rely on groundwater to meet their water needs (Eamus and Froend, 2006; Rohde et al, 2024). These ecosystems include several types of habitats, such as springs, wetlands, rivers, and certain types of forests, where the presence of groundwater is crucial for their survival. GDEs help maintain water quality by filtering pollutants and sediments, providing clean water for humans, plants, and animals (Howard, 2023). They support unique biodiversity, including species that are not found elsewhere, like the desert pupfish, which

can only survive in desert springs fed by groundwater (California Department of Fish and Wildlife, 2012). They also play a vital role in carbon storage, helping to mitigate climate change impacts (Yang et al. 2020). Additionally, GDEs support agricultural productivity by ensuring the availability of groundwater for irrigation. They are often sensitive to changes in groundwater levels and quality, making their management critical for environmental sustainability (Kreamer et al. 2015).

To facilitate sustainable groundwater management in California, The Nature Conservancy, the California Department of Fish and Wildlife, and the California Department of Water Resources collaboratively developed a comprehensive database of indicators of groundwater dependent ecosystems. This tool is designed for the identification and mapping of Groundwater Dependent Ecosystems (GDEs) across the state. The database provides an essential, detailed, and user-friendly map of GDEs, serving as a key resource for the effective implementation of the Sustainable Groundwater Management Act (SGMA) (Klausmeyer et al. 2018). This followed previous mapping efforts, which were done by independent volunteers (Howard et al. 2010).

Expanding on the mapping of Groundwater Dependent Ecosystems (GDEs) in California, Howard et al. (2023) explored their ecosystem functions and connections to various ecosystem services. This study found that GDEs significantly support pollination, with a third of pollinator-dependent crops within 1km of these ecosystems, enhancing potential agricultural yields. In large basins like the Central Valley, GDE's located between developed areas and bodies of water, are crucial for water quality regulation. Additionally, GDEs are key in climate regulation, storing 790 million tons of carbon dioxide, twice California's annual output Howard et al. (2023). These roles highlight GDEs' critical contribution to California's environmental and agricultural sustainability.

Building on the importance of GDEs, Rohde et al. (2021) used satellite remote sensing and machine learning to monitor groundwater levels from 1985 to 2019. Their findings revealed that 44% of GDEs

experienced significant long-term groundwater level declines, with declines intensifying in recent decades. It also found that groundwater declines are most prevalent in areas lacking sustainable groundwater management. This research underscores the urgency for robust groundwater management policies in California, highlighting the critical relationship between groundwater levels and management policies.

The Sustainable Groundwater Management Act and Groundwater Sustainability Plans

The Sustainable Groundwater Management Act (SGMA), enacted by the State of California in 2014, represents a transformative approach to groundwater management, targeting issues like over-extraction and depletion with lasting ecological and economic repercussions (Harrer et al., 2020). This legislation establishes a framework for the sustainable management of groundwater resources. Key to SGMA is the formation of Groundwater Sustainability Agencies (GSAs) in designated high and medium priority basins. These GSAs are charged with the critical task of formulating and executing Groundwater Sustainability Plans (GSPs). GSPs are essential for setting sustainability goals, monitoring and managing groundwater usage, and addressing substantial and unreasonable reductions in interconnected surface water, and reducing the risks of overdraft (Harrer et al. 2020).

Upon their adoption, GSPs undergo a comprehensive review process by the State, including a 60-day public comment period, and are classified as either adequate, conditionally adequate, or inadequate. This third classification could lead to State intervention if plans are deemed inadequate. The SGMA compliance journey is intricate and prolonged, necessitating both technical measures, such as establishing basin budgets and sustainability criteria, and consistent stakeholder engagement. GSPs must periodically be updated and refined every five years, taking into account the economic, social, and environmental impacts of the management strategies employed. Meeting interim milestones is crucial to circumvent

State control, with long-term monitoring being vital in tracking progress towards sustainability goals and facilitating adaptive management (Dumas Leslie, 2017).

For the five-year updates, constructing effective and inexpensive feedback mechanisms is essential. To improve these efforts, we sought to develop a large language model (LLM) designed to analyze a GSA's GSP and assess the plan according to a standardized rubric of evaluative questions. Development of such an LLM would enable a more efficient and insightful review process, ensuring that the GSAs receive comprehensive, data-driven recommendations to improve their groundwater management strategies, while avoiding prohibitive costs of consultant labor. This innovation could significantly enhance the adaptive management capabilities of GSAs, providing a dynamic and cost-effective tool to align with SGMA's objectives of sustainable groundwater management.

Large Language Models:

Large language models (LLMs) were developed to emulate human-like reading, writing, and communication skills through advanced natural language processing. These models are trained on vast datasets to predict word sequences. The evolution of LLMs began with rule-based natural language processing in the 1950s, utilizing hand-crafted linguistic rules (Chomsky 1956). The 1980s saw the advent of statistical models that used probabilistic methods for word sequence prediction (Rosenfeld, 2018). Google's Neural Machine Translation system in 2015 marked a milestone as the first major neural language model, which employs deep learning to analyze extensive textual data (Wu et al. 2016). The Transformer model, introduced in 2017, revolutionized LLMs with its ability to learn long-term dependencies and parallel training capabilities (Vaswani et al. 2017). OpenAI's introduction of GPT-1 in 2018, featuring a transformer-based architecture and 117 million parameters, was a significant progression (Radford et al. 2018). This was followed by the release of GPT-3 in 2020, which, trained on a massive dataset, set new benchmarks in generating coherent and natural-sounding text (Hadi et al. 2023).

OpenAI's ChatGPT, a generative AI model for natural language tasks, leverages advanced deep learning techniques. Trained on a vast array of internet-sourced texts like books, articles, and websites (Hadi et al. 2023), it excels in generating human-like, coherent responses. ChatGPT predicts word sequences to produce contextually relevant text. This capability stems from its training on complex language patterns, enabling it to engage in meaningful conversations and respond accurately to various queries.

GPT models, as cutting-edge tools in the workplace, are revolutionizing various industries with their advanced language processing capabilities. In the medical sector, LLMs are expected to play a significant role in decision-making and patient care. For example, XrayGPT is used to automate X-ray image analysis, allowing users to interact and inquire about these analyses via chat (Thawkar et al. 2023). In education, ChatGPT has been instrumental in offering personalized learning experiences. Platforms like Khan Academy are integrating ChatGPT to enhance tutoring (Mello et al. 2023). In the finance sector, models like BloombergGPT provide customer support and financial advice (Wu et al. 2023). In engineering, ChatGPT aids developers with code generation and debugging, showcasing the versatility and wide-ranging impact of these models in diverse fields. In Natural Language Processing research, ChatGPT has been proven to excel at text summarization; it performs extremely well in the benchmark domain of news summarization (Goyal et al. 2023). It has also taken, and passed, the Bar Exam (Bommarito et al. 2022).

Additional LLM Experiments

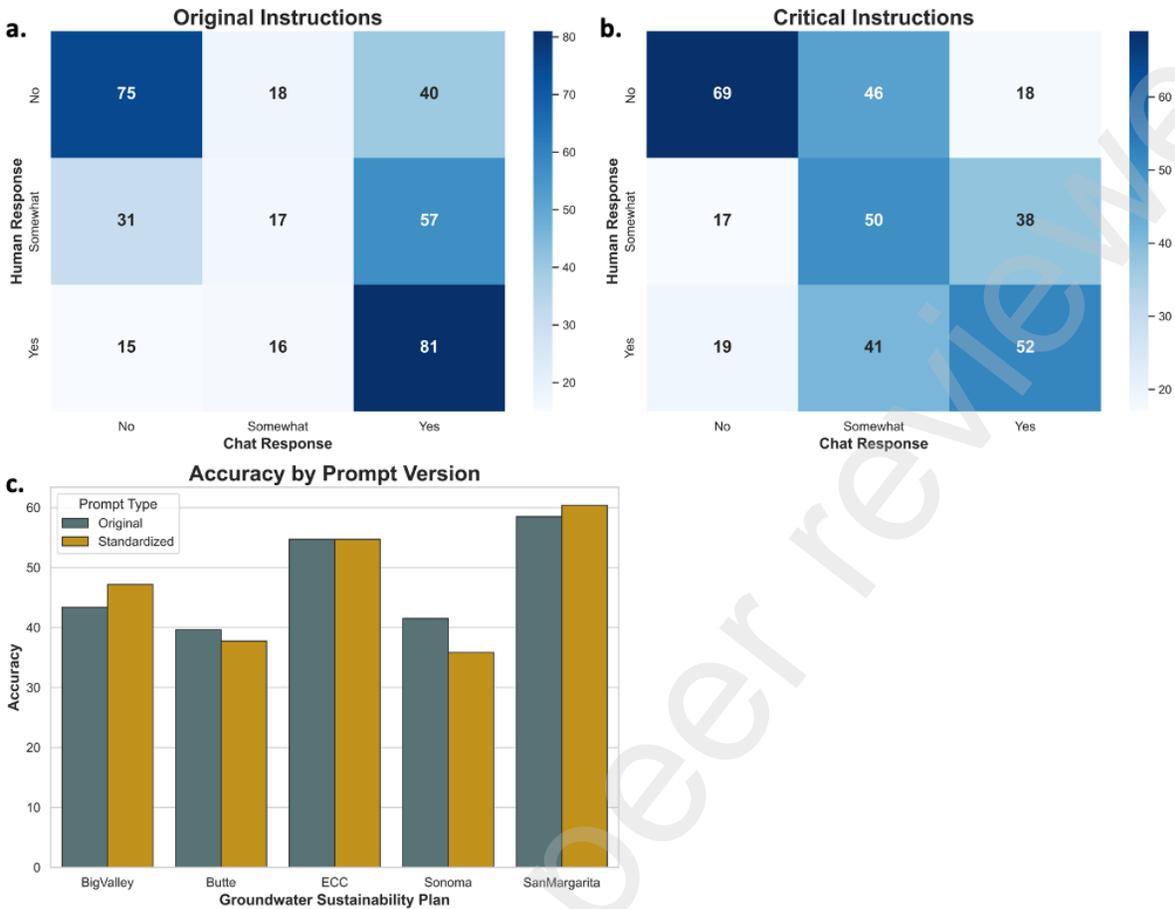


Figure S11: Performance comparison for different prompt engineering techniques.

Explicitly prompting the GPT 3.5 Turbo model to be a critical reviewer did not significantly influence model accuracy (67.8% original vs 67.1% critical).

We also ran a test that included the appendices for the 5 sample GSPs. The test with appendices answered 42% correct, whereas same model w/o appendices scored 51% correct. We found the model selected "Somewhat" 43% of the time when an appendix was included compared to 25% when the appendix was omitted. Humans selected "Somewhat" ~30% of the time on these sample GSPs. Thus, we concluded that adding the appendices did not improve accuracy and conducted further experiments without the appendices.

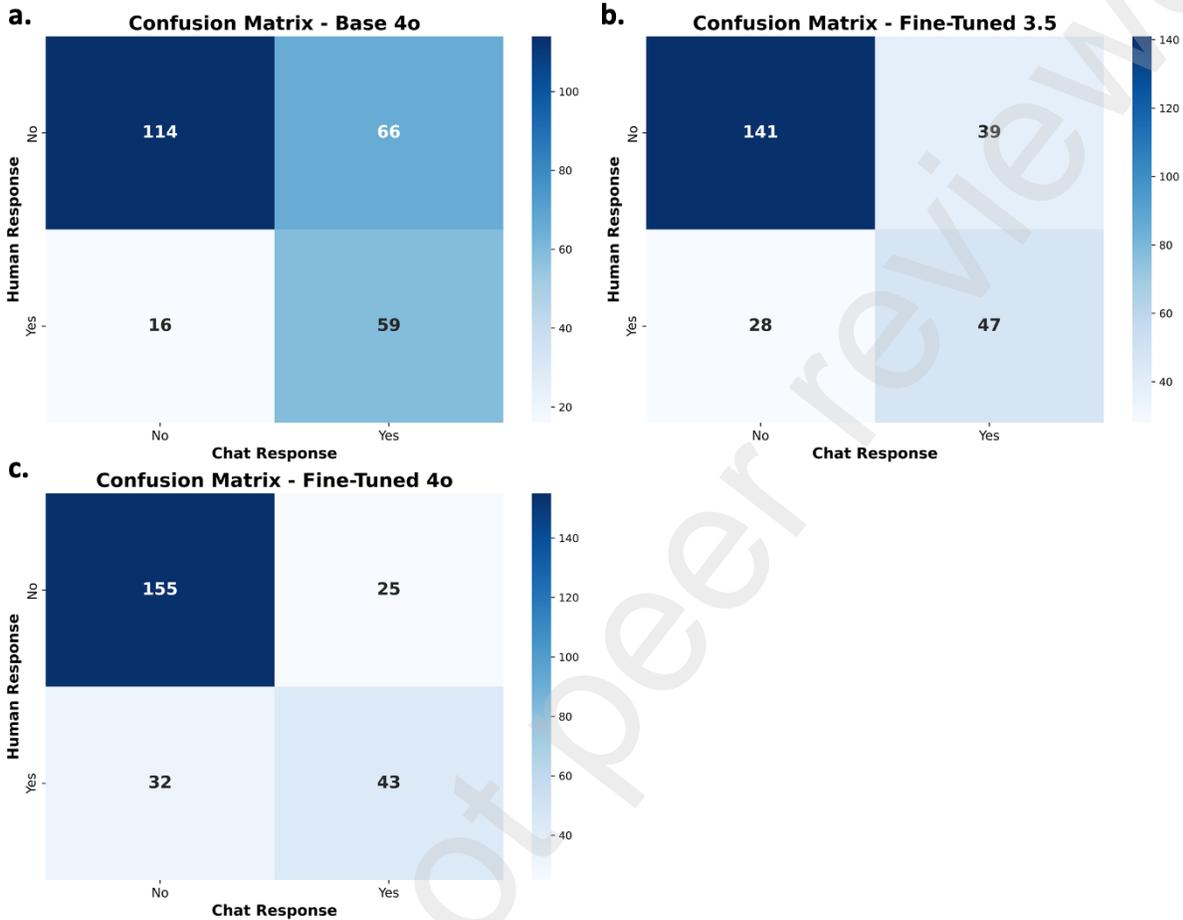


Figure SI2: Confusion Matrices for Base GPT-4o, GPT-3.5 Fine-Tuned and GPT-4o Fine-Tuned Models.

Figure SI2 displays raw confusion matrices for evaluating each of the three models (Base GPT-4o, GPT-3.5 Fine-Tuned and GPT-4o Fine-Tuned). Corresponding metrics and interpretations are found in Table 3 in the main text.

Table S1: AUC and Accuracy by GSP after applying the GPT-4o Fine-Tuned model

| GSP Number | GSP Name | AUC | Accuracy |
|------------|----------|-----|----------|
| | | | |

| | | | |
|----|-----------------------|--------|--------|
| 1 | BigValley | 0.8089 | 0.7800 |
| 2 | EelRiver | 0.6402 | 0.7600 |
| 3 | Yolo | 0.7561 | 0.7451 |
| 4 | SouthAmerican (*) | nan | nan |
| 5 | Colusa | 0.6611 | 0.6471 |
| 6 | NorthAmerica | 0.5750 | 0.5882 |
| 7 | Vina | 0.6548 | 0.6078 |
| 8 | LosMolinos (*) | nan | nan |
| 9 | Solano | 0.6482 | 0.5833 |
| 10 | Sutter | 0.6373 | 0.6667 |
| 11 | Butte (***) | 0.7891 | 0.6875 |
| 12 | Cosumnes | 0.7461 | 0.6863 |
| 13 | Tracy (**) | nan | nan |
| 14 | EastContraCosta | 0.6656 | 0.6889 |
| 15 | Fillmore | 0.7373 | 0.6735 |
| 16 | Piru | 0.5543 | 0.5745 |
| 17 | Mound | 0.7942 | 0.6939 |
| 18 | Shasta | 0.7839 | 0.7255 |
| 19 | SierraValley | 0.6363 | 0.6200 |
| 20 | ButteValley | 0.8333 | 0.7843 |
| 21 | Carpenteria (*) | nan | nan |
| 22 | SanGorgonioPass | 0.7459 | 0.7000 |
| 23 | SantaMonica | 0.8294 | 0.6875 |
| 24 | BedfordColdwater (**) | nan | nan |
| 25 | ElsinoreVally | 0.6625 | 0.8039 |
| 26 | NorthSanBenito | 0.7576 | 0.7000 |
| 27 | TuleLake | 0.5273 | 0.6458 |
| 28 | Montecito (*) | nan | nan |
| 29 | NapaValley | 0.7127 | 0.6600 |
| 30 | SonomaValley | 0.6896 | 0.7021 |
| 31 | OjaiValley | 0.7212 | 0.6905 |
| 32 | PetalumaValley | 0.6694 | 0.6383 |
| 33 | Anderson | 0.7526 | 0.7647 |
| 34 | Enterprise | 0.6814 | 0.7234 |
| 35 | Antelope | 0.6324 | 0.7400 |
| 36 | Corning | 0.7350 | 0.7400 |
| 37 | RedBluff | 0.6704 | 0.6667 |
| 38 | WyandotteCreek | 0.6035 | 0.5686 |
| 39 | EastSideAquifer | 0.6686 | 0.6383 |
| 40 | Langley | 0.7238 | 0.6596 |

| | | | |
|----|---|---------------|---------------|
| 41 | ForebayAquifer | 0.7319 | 0.6667 |
| 42 | Monterey | 0.5846 | 0.5745 |
| 43 | UpperValley | 0.7333 | 0.7447 |
| 44 | SanAntonioCreekValley | 0.7621 | 0.7447 |
| 45 | SanJacinto | 0.7621 | 0.7500 |
| 46 | Modesto | 0.8121 | 0.7959 |
| 47 | Turlock | 0.6314 | 0.6042 |
| 48 | PleasantValley (*) | nan | nan |
| 49 | WhiteWolf | 0.7961 | 0.7755 |
| 50 | SanLuisObispoValley | 0.7599 | 0.7200 |
| 51 | UpperSanLuisReyValley | 0.7321 | 0.7800 |
| 52 | SanPasqual | 0.6993 | 0.6667 |
| 53 | SantaClaraRiverValleyEast | 0.6745 | 0.7143 |
| 54 | EastBayPlain | 0.6542 | 0.6383 |
| 55 | SantaMargarita (***) | 0.6952 | 0.6800 |
| 56 | SantaRosaPlain | 0.6716 | 0.6800 |
| 57 | SantaYnezRiverValleyWestern (*) | nan | nan |
| 58 | SantaYnezRiverValleyCentral (*) | nan | nan |
| 59 | SantaYnezRiverValleyEastern | 0.6875 | 0.7000 |
| 60 | ScottRiverValley | 0.7375 | 0.6800 |
| 61 | Ukiah | 0.8222 | 0.6863 |
| 62 | Yucapia | 0.5702 | 0.6383 |
| 63 | Temescal | 0.6631 | 0.6809 |
| 64 | UpperVentura | 0.8005 | 0.7800 |
| 65 | BigValley(LakeCounty) | 0.6419 | 0.6600 |
| | <i>Average</i> | <i>0.7023</i> | <i>0.6893</i> |
| | <i>Average (excluding GSPs used in fine tuning)</i> | <i>0.7008</i> | <i>0.6895</i> |

Notes

*Technical issues with the PDF files for the GSP prevented the conversion of text to embeddings.

**Technical issues with the human coded rubrics prevented the accuracy assessment.

***Used for fine tuning the LLM.

References

Bommarito, M, et al. (2022) GPT Takes the Bar Exam. <https://arxiv.org/pdf/2212.14402.pdf>

California Department of Fish and Wildlife. (n.d.). Desert Pupfish (Cyprinodon macularis).<https://wildlife.ca.gov/Regions/6/Desert-Fishes/Desert-Pupfish>.

Chomsky, N. (1956). Three models for the description of language. *IRE Transactions on information theory*, 2(3), 113-124.

Dumas, L. (2017). Implementing SGMA—An Update on California’s Foray into Groundwater Regulation. World Environmental and Water Resources Congress 2017.

Eamus, D., & Froend, R. (2006). Groundwater-dependent ecosystems: the where, what and why of GDEs. *Australian Journal of Botany*, 54(2), 91-96.

Goyal, T, Li, J., Durrett, G. (2023). News Summarization and Evaluation in the Era of GPT-3. arXiv: 2209.12356.

Hadi, M. U., et al. (2023). Large language models: a comprehensive survey of its applications, challenges, limitations, and future prospects. Authorea Preprints.

Harrer, S., Howard, J., Rohde, M., et al. (2020). Groundwater Dependent Ecosystems under the Sustainable Groundwater Management Act.
<https://www.scienceforconservation.org/assets/downloads/GDEsUnderSGMA.pdf>.

Howard, J. & Merrifield, M. (2010). Mapping groundwater dependent ecosystems in California. *PLoS One*, 5(6), e11249.

Howard, J. K., et al. (2023). Ecosystem services produced by groundwater dependent ecosystems: a framework and case study in California. *Frontiers in Water*, 5, 1115416.

Klausmeyer, K., et al. (2018). Mapping indicators of groundwater dependent ecosystems in California: Methods report. San Francisco, California.

Kreamer, D. K., Stevens, L. E., & Ledbetter, J. D. (2015). Groundwater dependent ecosystems—Science, challenges, and policy directions. *Groundwater*, 205, 230.

Mello, R, et al. (2023) Education in the age of Generative AI: Context and Recent Developments. arXiv preprint arXiv:2309.12332.

Radford, A., Narasimhan, K., Salimans, T., & Sutskever, I. (2018). Improving language understanding by generative pre-training. OpenAI. https://cdn.openai.com/research-covers/language-unsupervised/language_understanding_paper.pdf

Rohde, M. M., et al. (2021) A machine learning approach to predict groundwater levels in California reveals ecosystems at risk. *Frontiers in Earth Science*, 9, 784499.

Rohde, M. M., Albano, C. M., Huggins, X., Klausmeyer, K. R., Morton, C., Sharman, A., ... & Stella, J. C. (2024). Groundwater-dependent ecosystem map exposes global dryland protection needs. *Nature*, 632(8023), 101-107.

Rosenfeld, Roni (2018). Two Decades of Statistical Language Modeling: Where Do We Go From Here?. Carnegie Mellon University. Journal contribution. <https://doi.org/10.1184/R1/6611138.v1>

Thawkar, O, et al. (2023) Xraygpt: Chest radiographs summarization using medical vision-language models. arXiv preprint arXiv:2306.07971.

Vaswani, A, et al. (2017) Attention Is All You Need. arXiv preprint arXiv:1706.03762.

Wu, S, et al. (2023) Bloomberggpt: A large language model for finance. arXiv preprint arXiv:2303.17564.

Wu, Y, et al. (2016). Google's Neural Machine Translation System: Bridging the Gap between Human and Machine Translation. arXiv preprint arXiv:1609.08144.

Yang, K, et al. (2020) Recent progress in the design and application of MOF-based materials for highly selective CO2 capture. *ACS Energy Letters*, 5(11), 3490-3496.

Title

Improving Conservation Efficiency: Accelerating Groundwater Sustainability Plan Reviews

Using Large Language Models

Author Names and Affiliations

Ryan Bernstein¹ 101 Howard St, San Francisco, CA 94105, United States of America

ryanoh999@gmail.com

Seneth Waterman¹ 101 Howard St, San Francisco, CA 94105, United States of America

seneth.waterman@gmail.com

Kirk R. Klausmeyer² 620 Davis St, San Francisco, CA 94111, United States of America

kklausmeyer@tnc.org

Nicholas Murphy² 620 Davis St, San Francisco, CA 94111, United States of America

nicholas.murphy@tnc.org

Melissa M. Rohde^{3,4} 1 Forestry Dr, Syracuse, NY 13210, United States of America

melissa@rohdeenvironmental.com

Cody Carroll¹ 101 Howard St, San Francisco, CA 94105, United States of America

cjcarroll@usfca.edu

¹ *University of San Francisco*

² *The Nature Conservancy of California, San Francisco, CA*

³ *Rohde Environmental Consulting, LLC, Seattle, WA*

⁴ *SUNY College of Environmental Science and Forestry, Syracuse, NY*

Corresponding Author

Ryan Bernstein