

## **Title Page:**

CSC 659 859 AI Ethics and Explainability  
San Francisco State University

Team project: Development and Ethics and Trustworthiness  
Audit of AI Application  
Option B - GenAI

Team 1

Harris Chan, ychan@sfsu.edu  
Maeve Fitzpatrick, mfitzpatrick@sfsu.edu  
Zari Haidarian, zhaidarian@sfsu.edu

Date: 08/06/2025

### **Table of Content:**

- Page 1 : Title Page
- Page 2 : Section 1 - Executive summary
- Page 3-4 : Section 2 - Motivation, Problem Description & Case-Study Goals
- Page 5-7 : Section 3 - GenAI Technology & Application to Be Evaluated
- Page 8-10 : Section 4 - Test and verification data
- Page 11-13 : Section 5 - Methods of GenAI Accuracy Evaluation
- Page 14-21 : Section 6 - Results of accuracy evaluation experiments
- Page 22-25 : Section 7 - Audit for ethics and trustworthiness
- Page 26-27 : Executive Summary and Recommendations (Conclusion)
- Page 28-30 : Appendix I - Summary of each team member contributions
- Page 31-37 : Appendix II - Code Used
- Page 38-39 : Appendix III - ChatGPT and related GenAI usage
- Page 40 : References

## **Section 1: Executive summary**

### **1.1 Goal of this report**

The goal of this report is to see if AI (GPT-4) can reliably answer and explain computer science-related multiple-choice questions for the purpose of helping computer science students learn and study.

### **1.2 Problem description and its importance**

We have found both benefits and drawbacks of multiple-choice questions; they are easier to grade for teachers. However, the answers are much easier for students to guess, which defeats the purpose of an exam meaning to evaluate how much a student has learned. Through this project, we aim to maintain the efficiency with which multiple-choice questions can be graded, which then would get the exams back to the students sooner to understand what they do or do not have a good handle on yet. If the GenAI can successfully, accurately and fully explain the answers to multiple-choice questions, this would overall improve the learning and teaching experiences of students and teachers.

### **1.3 Data available and technologies and tools used**

We utilized a GitHub repository from MCQ-World to source our computer science multiple-choice questions, the answer key to which will be our ground truth. We are evaluating the ability of OpenAI's GPT-4 to explain our questions.

### **1.4 Methods of audit**

To audit accuracy, we compared the model's predicted option against the ground-truth key, computed overall and per-topic accuracy rates, and visualized results in a confusion matrix. For trustworthiness, we sampled misclassifications and categorized them (e.g. misinterpretation, plausible-but-wrong, hallucination) to understand failure modes, and we performed a lightweight bias/fairness/privacy check to ensure no sensitive data leakage or unintended demographics effects.

### **1.5 Summary results**

Overall accuracy: 38/50 correct → 76%

### **1.6 Recommendations of your work**

We have found that using GPT-4 as an "AI tutor" was indeed helpful in terms of adding supplemental explanation. However, the tool was not accurate enough to be recommended as a grader. As long as the tool is used to supplement knowledge that is at least partially known already or can be checked with answers from the teacher, it is okay to use, but it is not recommended to rely on this tool to determine correct answers.

## **Section 2: Motivation, Problem Description & Case-Study Goals**

### **2.1 Chosen GenAI Technology**

We will evaluate OpenAI’s GPT-4 series, accessed via the OpenAI Python SDK. GPT-4 is a large-scale Transformer model pre-trained on diverse web-scale corpora and fine-tuned with reinforcement learning from human feedback (RLHF). Its zero-shot and few-shot capabilities make it a natural candidate for natural-language reasoning tasks without additional custom training.

### **2.2 Application Domain: Automated MCQ Answering**

Our case study focuses on multiple-choice questions (MCQ) answering in undergraduate computer-science courses. Frequently test concepts such as:

- Algorithmic complexity (e.g., “What is the worst-case time for quicksort?”)
- Data-structure behavior (e.g., “Which structure implements FIFO?”)
- Operating-system fundamentals (e.g., “What does a system call trap into?”)
- Database design (e.g., “Which normal form prevents partial dependencies?”)

We will build a Jupyter-notebook prototype that ingests a question stem and four labeled options (A–D), queries GPT-4 to select the best answer, and returns a one-sentence explanation of its choice.

### **2.3 Motivation & Impact**

1. Pedagogical efficiency: Automating MCQ feedback can dramatically reduce grading effort and turnaround time, giving students immediate insight into *why* an answer is correct or incorrect.
2. Assessment enrichment: Pairing each selected option with a rationale deepens conceptual understanding, moving beyond binary right/wrong marking.
3. Feasibility of prompt-based ML: By framing MCQ answering as a zero-shot classification task, we explore how far off-the-shelf LLMs can go in reasoning about structured academic content without bespoke fine-tuning.
4. Ethical imperative: Deploying LLMs in education carries risks—hallucinated explanations may mislead learners; overconfident wrong answers can erode trust; careless logging may expose quiz identifiers. A rigorous audit is essential before any educational rollout.

## 2.4 Data & Decision Points

- Data source: We will draw our test questions from the MCQ-World GitHub repository (patil2104/MCQ-World), which is MIT-licensed and contains several hundred CS multiple-choice questions organized into eight topic folders (DSA, Operating System, Computer Networks, DBMS, OOP, Java, Language Processors & Compilers, AI/ML) [GitHub](#). From this pool we will sample 30–50 questions, ensuring coverage across all major topics.
- Ground truth: Each question’s official correct answer is provided in the repository’s answer keys; we will use those as our “gold standard.”
- ML decisions:
  1. Answer selection: The model must classify which of the four options (A–D) is correct.
  2. Explanation synthesis: The model must generate a concise, accurate one-sentence rationale for its choice—avoiding any hallucinated or extraneous detail.

Both decision steps are risky (an incorrect or misleading rationale can reinforce student misconceptions) and impactful (learners may accept the system’s output without independent verification), underscoring the need for careful evaluation and ethics auditing.

## 2.5 Specific Case-Study Goals

1. Accuracy measurement: Compute overall MCQ classification accuracy (%) and construct a confusion matrix to identify which distractors most frequently mislead the model.
2. Error analysis: Manually review all misclassifications, categorizing errors into misinterpretation, plausible-but-wrong rationale, or outright hallucination.
3. Ethics & trustworthiness audit:
  - Bias & fairness: Assess performance variation across topics or question styles.
  - Transparency: Verify that explanations include uncertainty cues (e.g. “I think...”).
  - Privacy: Ensure no Canvas identifiers are logged or exposed.
  - Human oversight: Confirm the notebook reminds users to verify AI outputs before relying on them.

By achieving these goals, we will provide a clear picture of GPT-4’s viability for automated CS MCQ assistance, along with actionable recommendations for safe, responsible classroom integration.

## **Section 3: GenAI Technology & Application to Be Evaluated**

### **3.1 Application Overview**

Our prototype is an Automated MCQ Answerer built entirely in a Jupyter notebook. It:

1. Loads one question at a time (stem + four options) from a CSV
2. Constructs a chat prompt for GPT-4, including a system message, a few-shot template, and the actual question
3. Calls the OpenAI API (model gpt-4o-mini)
4. Parses the model's reply to extract the chosen option letter and explanation
5. Displays results with accuracy % and Confusion Matrix

This single-notebook “app” demonstrates end-to-end GenAI development without any external UI or server.

### **3.2 GenAI Technology & Environment**

- Model: GPT-4 series (gpt-4o-mini endpoint)
- SDK: OpenAI Python SDK v0.27+
- Runtime: Python 3.10 within JupyterLab (CPU only)
- Dependencies: openai, pandas, numpy, matplotlib (for confusion-matrix plotting)
- Reproducibility: We record the exact model identifier, SDK version, notebook commit, and timestamp for every API call (see Appendix II).

### **3.3 Prompt Customization & Few-Shot Example**

#### **3.3.1 System Message**

You are an expert computer-science tutor. Given a multiple-choice question with four options, select the correct answer and provide a one-sentence rationale in clear, beginner-friendly language.

#### **3.3.2 API Call & Parsing**

```
resp = openai.ChatCompletion.create(model="gpt-4o-mini",  
messages=prompt)  
raw = resp.choices[0].message.content  
# Extract "Answer: X)" via regex, then split on "Explanation:"  
for rationale
```

### 3.3.3 Few-Shot Template

We prepend two illustrative Q/A pairs so the model learns the exact format:

```
prompt = [
    {"role": "system", "content": system_message},

    # Example 1
    {"role": "user", "content":
        "Q: Which data structure follows FIFO?\n"
        "A) Stack   B) Queue   C) Tree   D) Graph"
    },
    {"role": "assistant", "content":
        "Answer: B) Queue\n"
        "Explanation: A queue enqueues and dequeues elements in
first-in, first-out order."
    },

    # Example 2
    {"role": "user", "content":
        "Q: What is the worst-case time complexity of bubble sort?\n"
        "A)  $O(n)$    B)  $O(n \log n)$    C)  $O(n^2)$    D)  $O(\log n)$ "
    },
    {"role": "assistant", "content":
        "Answer: C)  $O(n^2)$ \n"
        "Explanation: In the worst case, bubble sort swaps every
adjacent pair, resulting in  $n \cdot (n-1)/2$  operations  $\rightarrow O(n^2)$ ."
    },

    # Your actual MCQ
    {"role": "user", "content":
        f"Q: {question_stem}\n"
        "A) {optA}   B) {optB}   C) {optC}   D) {optD}"
    }
]
```

### 3.4 Post-Processing & Logging

- Answer parsing: Regex to capture the single letter A–D
- Explanation trimming: Split at the first period to enforce one sentence
- Logging: Append {question\_id, prompt, raw, choice, explanation} to a Pandas DataFrame
- Visualization: Use matplotlib to plot a confusion matrix of predicted vs. ground-truth answers

Once we've run our notebook over the 50 MCQs and collected two lists:

```
y_true = df['correct_option'].tolist()      # e.g.  
['A', 'C', 'B', ...]
```

```
y_pred = df['predicted_option'].tolist()    # e.g.  
['A', 'B', 'B', ...]
```

We can turn those into a confusion matrix, a 4×4 table showing how often the model answered  $A \rightarrow A$ ,  $A \rightarrow B$ , ...,  $D \rightarrow C$ ,  $D \rightarrow D$  and then plot it with Matplotlib.

Why it matters

- Instant diagnostics: We'll know *which* options the model confuses most often (say, confusing B and C on algorithm questions).
- Guides improvements: If certain distractors trip up the model, you can add more few-shot examples specifically for those cases.

We are pulling answers from ChatGPT in Jupyter, by collecting them into `y_pred` and comparing against our known `y_true labels`, we can use Matplotlib to visualize exactly how the model is performing across the MCQ.

### 3.5 Setup Conclusion

With this setup—clear system framing, illustrative few-shot examples, robust post-processing, and optional fine-tuning—we ensure our Jupyter-notebook prototype both develops and evaluates GPT-4 in a reproducible MCQ classification task.

## **Section 4: Test and verification data**

### **MCQ-World-GitHub**

- The test and verification data are a set of 50 multiple choice questions and answers about Computer Science gathered from a GitHub repository.
- These questions cover the following Computer Science subjects: artificial intelligence, computer networks, data science, machine learning, and object oriented programming. These questions include things like, “What is Artificial Intelligence?”, “What is the impact of having noisy data?”, and “What is the correlation coefficient?”
- The answer key to these questions are our ground truth. We found this repository to be a trustworthy source because it aims to aid students in their computer science studies.

<b>Data Element</b>	<b>Description</b>
Source	Public GitHub repo of CS practice questions
# of Questions	50 distinct MCQs
Content Covered	<ul style="list-style-type: none"><li>• Artificial Intelligence</li><li>• Computer Networks</li><li>• Data Science &amp; Machine Learning</li><li>• OOP</li><li>• SQL</li></ul>
Question Format	<ul style="list-style-type: none"><li>– Stem (text)</li><li>– Four options (A–D)</li><li>– One correct answer</li></ul>
Ground Truth	Answer key curated and verified by human CS instructors
Verification Process	We manually cross-checked each answer key entry against official course materials for consistency

*Note: All 50 questions have unambiguous, instructor-approved answers—there are no “trick” items or multiple correct options. This ensures our accuracy metrics reflect genuine model understanding rather than dataset noise.*



It is important to note that our verification data is licensed by MIT


## About


This is the repository to help Computer Science students By Making All type of MCQ Questions at the place

hacktoberfest

hacktoberfest-accepted


hacktoberfest2022

 Readme

 MIT license

 Activity

 5 stars

 1 watching

 18 forks

[Report repository](#)

## MIT License

Copyright (c) 2022 patil2104

Permission is hereby granted, free of charge, to any person obtaining a copy of this software and associated documentation files (the "Software"), to deal in the Software without restriction, including without limitation the rights to use, copy, modify, merge, publish, distribute, sublicense, and/or sell copies of the Software, and to permit persons to whom the Software is furnished to do so, subject to the following conditions:

The above copyright notice and this permission notice shall be included in all copies or substantial portions of the Software.

THE SOFTWARE IS PROVIDED "AS IS", WITHOUT WARRANTY OF ANY KIND, EXPRESS OR IMPLIED, INCLUDING BUT NOT LIMITED TO THE WARRANTIES OF MERCHANTABILITY, FITNESS FOR A PARTICULAR PURPOSE AND NONINFRINGEMENT. IN NO EVENT SHALL THE AUTHORS OR COPYRIGHT HOLDERS BE LIABLE FOR ANY CLAIM, DAMAGES OR OTHER LIABILITY, WHETHER IN AN ACTION OF CONTRACT, TORT OR OTHERWISE, ARISING FROM, OUT OF OR IN CONNECTION WITH THE SOFTWARE OR THE USE OR OTHER DEALINGS IN THE SOFTWARE.

# The verification data has been manually selected, then inputted into CSV

A1	A2	A3	A4	A5	A6	A7	A8	A9	A10	A11	A12	A13	A14	A15	A16	A17	A18	A19	A20	A21	A22	A23	A24	A25	A26	A27	A28	A29	A30	A31	A32	A33	A34	A35	A36	A37	A38	A39	A40	A41	A42	A43	A44	A45	A46	A47	A48	A49	A50	A51	A52	A53	A54	A55	A56	A57	A58	A59	A60	A61	A62	A63	A64	A65	A66	A67	A68	A69	A70	A71	A72	A73	A74	A75	A76	A77	A78	A79	A80	A81	A82	A83	A84	A85	A86	A87	A88	A89	A90	A91	A92	A93	A94	A95	A96	A97	A98	A99	A100	A101	A102	A103	A104	A105	A106	A107	A108	A109	A110	A111	A112	A113	A114	A115	A116	A117	A118	A119	A120	A121	A122	A123	A124	A125	A126	A127	A128	A129	A130	A131	A132	A133	A134	A135	A136	A137	A138	A139	A140	A141	A142	A143	A144	A145	A146	A147	A148	A149	A150	A151	A152	A153	A154	A155	A156	A157	A158	A159	A160	A161	A162	A163	A164	A165	A166	A167	A168	A169	A170	A171	A172	A173	A174	A175	A176	A177	A178	A179	A180	A181	A182	A183	A184	A185	A186	A187	A188	A189	A190	A191	A192	A193	A194	A195	A196	A197	A198	A199	A200	A201	A202	A203	A204	A205	A206	A207	A208	A209	A210	A211	A212	A213	A214	A215	A216	A217	A218	A219	A220	A221	A222	A223	A224	A225	A226	A227	A228	A229	A230	A231	A232	A233	A234	A235	A236	A237	A238	A239	A240	A241	A242	A243	A244	A245	A246	A247	A248	A249	A250	A251	A252	A253	A254	A255	A256	A257	A258	A259	A260	A261	A262	A263	A264	A265	A266	A267	A268	A269	A270	A271	A272	A273	A274	A275	A276	A277	A278	A279	A280	A281	A282	A283	A284	A285	A286	A287	A288	A289	A290	A291	A292	A293	A294	A295	A296	A297	A298	A299	A300	A301	A302	A303	A304	A305	A306	A307	A308	A309	A310	A311	A312	A313	A314	A315	A316	A317	A318	A319	A320	A321	A322	A323	A324	A325	A326	A327	A328	A329	A330	A331	A332	A333	A334	A335	A336	A337	A338	A339	A340	A341	A342	A343	A344	A345	A346	A347	A348	A349	A350	A351	A352	A353	A354	A355	A356	A357	A358	A359	A360	A361	A362	A363	A364	A365	A366	A367	A368	A369	A370	A371	A372	A373	A374	A375	A376	A377	A378	A379	A380	A381	A382	A383	A384	A385	A386	A387	A388	A389	A390	A391	A392	A393	A394	A395	A396	A397	A398	A399	A400	A401	A402	A403	A404	A405	A406	A407	A408	A409	A410	A411	A412	A413	A414	A415	A416	A417	A418	A419	A420	A421	A422	A423	A424	A425	A426	A427	A428	A429	A430	A431	A432	A433	A434	A435	A436	A437	A438	A439	A440	A441	A442	A443	A444	A445	A446	A447	A448	A449	A450	A451	A452	A453	A454	A455	A456	A457	A458	A459	A460	A461	A462	A463	A464	A465	A466	A467	A468	A469	A470	A471	A472	A473	A474	A475	A476	A477	A478	A479	A480	A481	A482	A483	A484	A485	A486	A487	A488	A489	A490	A491	A492	A493	A494	A495	A496	A497	A498	A499	A500	A501	A502	A503	A504	A505	A506	A507	A508	A509	A510	A511	A512	A513	A514	A515	A516	A517	A518	A519	A520	A521	A522	A523	A524	A525	A526	A527	A528	A529	A530	A531	A532	A533	A534	A535	A536	A537	A538	A539	A540	A541	A542	A543	A544	A545	A546	A547	A548	A549	A550	A551	A552	A553	A554	A555	A556	A557	A558	A559	A560	A561	A562	A563	A564	A565	A566	A567	A568	A569	A570	A571	A572	A573	A574	A575	A576	A577	A578	A579	A580	A581	A582	A583	A584	A585	A586	A587	A588	A589	A590	A591	A592	A593	A594	A595	A596	A597	A598	A599	A600	A601	A602	A603	A604	A605	A606	A607	A608	A609	A610	A611	A612	A613	A614	A615	A616	A617	A618	A619	A620	A621	A622	A623	A624	A625	A626	A627	A628	A629	A630	A631	A632	A633	A634	A635	A636	A637	A638	A639	A640	A641	A642	A643	A644	A645	A646	A647	A648	A649	A650	A651	A652	A653	A654	A655	A656	A657	A658	A659	A660	A661	A662	A663	A664	A665	A666	A667	A668	A669	A670	A671	A672	A673	A674	A675	A676	A677	A678	A679	A680	A681	A682	A683	A684	A685	A686	A687	A688	A689	A690	A691	A692	A693	A694	A695	A696	A697	A698	A699	A700	A701	A702	A703	A704	A705	A706	A707	A708	A709	A710	A711	A712	A713	A714	A715	A716	A717	A718	A719	A720	A721	A722	A723	A724	A725	A726	A727	A728	A729	A730	A731	A732	A733	A734	A735	A736	A737	A738	A739	A740	A741	A742	A743	A744	A745	A746	A747	A748	A749	A750	A751	A752	A753	A754	A755	A756	A757	A758	A759	A760	A761	A762	A763	A764	A765	A766	A767	A768	A769	A770	A771	A772	A773	A774	A775	A776	A777	A778	A779	A780	A781	A782	A783	A784	A785	A786	A787	A788	A789	A790	A791	A792	A793	A794	A795	A796	A797	A798	A799	A800	A801	A802	A803	A804	A805	A806	A807	A808	A809	A810	A811	A812	A813	A814	A815	A816	A817	A818	A819	A820	A821	A822	A823	A824	A825	A826	A827	A828	A829	A830	A831	A832	A833	A834	A835	A836	A837	A838	A839	A840	A841	A842	A843	A844	A845	A846	A847	A848	A849	A850	A851	A852	A853	A854	A855	A856	A857	A858	A859	A860	A861	A862	A863	A864	A865	A866	A867	A868	A869	A870	A871	A872	A873	A874	A875	A876	A877	A878	A879	A880	A881	A882	A883	A884	A885	A886	A887	A888	A889	A890	A891	A892	A893	A894	A895	A896	A897	A898	A899	A900	A901	A902	A903	A904	A905	A906	A907	A908	A909	A910	A911	A912	A913	A914	A915	A916	A917	A918	A919	A920	A921	A922	A923	A924	A925	A926	A927	A928	A929	A930	A931	A932	A933	A934	A935	A936	A937	A938	A939	A940	A941	A942	A943	A944	A945	A946	A947	A948	A949	A950	A951	A952	A953	A954	A955	A956	A957	A958	A959	A960	A961	A962	A963	A964	A965	A966	A967	A968	A969	A970	A971	A972	A973	A974	A975	A976	A977	A978	A979	A980	A981	A982	A983	A984	A985	A986	A987	A988	A989	A990	A991	A992	A993	A994	A995	A996	A997	A998	A999	A1000	A1001	A1002	A1003	A1004	A1005	A1006	A1007	A1008	A1009	A1010	A1011	A1012	A1013	A1014	A1015	A1016	A1017	A1018	A1019	A1020	A1021	A1022	A1023	A1024	A1025	A1026	A1027	A1028	A1029	A1030	A1031	A1032	A1033	A1034	A1035	A1036	A1037	A1038	A1039	A1040	A1041	A1042	A1043	A1044	A1045	A1046	A1047	A1048	A1049	A1050	A1051	A1052	A1053	A1054	A1055	A1056	A1057	A1058	A1059	A1060	A1061	A1062	A1063	A1064	A1065	A1066	A1067	A1068	A1069	A1070	A1071	A1072	A1073	A1074	A1075	A1076	A1077	A1078	A1079	A1080	A1081	A1082	A1083	A1084	A1085	A1086	A1087	A1088	A1089	A1090	A1091	A1092	A1093	A1094	A1095	A1096	A1097	A1098	A1099	A1100	A1101	A1102	A1103	A1104	A1105	A1106	A1107	A1108	A1109	A1110	A1111	A1112	A1113	A1114	A1115	A1116	A1117	A1118	A1119	A1120	A1121	A1122	A1123	A1124	A1125	A1126	A1127	A1128	A1129	A1130	A1131	A1132	A1133	A1134	A1135	A1136	A1137	A1138	A1139	A1140	A1141	A1142	A1143	A1144	A1145	A1146	A1147	A1148	A1149	A1150	A1151	A1152	A1153	A1154	A1155	A1156	A1157	A1158	A1159	A1160	A1161	A1162	A1163	A1164	A1165	A1166	A1167	A1168	A1169	A1170	A1171	A1172	A1173	A1174	A1175	A1176	A1177	A1178	A1179	A1180	A1181	A1182	A1183	A1184	A1185	A1186	A1187	A1188	A1189	A1190	A1191	A1192	A1193	A1194	A1195	A1196	A1197	A1198	A1199	A1200	A1201	A1202	A1203	A1204	A1205	A1206	A1207	A1208	A1209	A1210	A1211	A1212	A1213	A1214	A1215	A1216	A1217	A1218	A1219	A1220	A1221	A1222	A1223	A1224	A1225	A1226	A1227	A1228	A1229	A1230	A1231	A1232	A1233	A1234	A1235	A1236	A1237	A1238	A1239	A1240	A1241	A1242	A1243	A1244	A1245	A1246	A1247	A1248	A1249	A1250	A1251	A1252	A1253	A1254	A1255	A1256	A1257	A1258	A1259	A1260	A1261	A1262	A1263	A1264	A1265	A1266	A1267	A1268	A1269	A1270	A1271	A1272	A1273	A1274	A1275	A1276	A1277	A1278	A1279	A1280	A1281	A1282	A1283	A1284	A1285	A1286	A1287	A1288	A1289	A1290	A1291	A1292	A1293	A1294	A1295	A1296	A1297	A1298	A1299	A1300	A1301	A1302	A1303	A1304	A1305	A1306	A1307	A1308	A1309	A1310	A1311	A1312	A1313	A1314	A1315	A1316	A1317	A1318	A1319	A1320	A1321	A1322	A1323	A1324	A1325	A1326	A1327	A1328	A1329	A1330	A1331	A1332	A1333	A1334	A1335	A1336	A1337	A1338	A1339	A1340	A1341	A1342	A1343	A1344	A1345	A1346	A1347	A1348	A1349	A1350	A1351	A1352	A1353	A1354	A1355	A1356	A1357	A1358	A1359	A1360	A1361	A1362	A1363	A1364	A1365	A1366	A1367	A1368	A1369	A1370	A1371	A1372	A1373	A1374	A1375	A1376	A1377	A1378	A1379	A1380	A1381	A1382	A1383	A1384	A1385	A1386	A1387	A1388	A1389	A1390	A1391	A1392	A1393	A1394	A1395	A1396	A1397	A1398	A1399	A1400	A1401	A1402	A1403	A1404	A1405	A1406	A1407	A1408	A1409	A1410	A1411	A1412	A1413	A1414	A1415	A1416	A1417	A1418	A1419	A1420	A1421	A1422	A1423	A1424	A1425	A1426	A1427	A1428	A1429	A1430	A1431	A1432	A1433	A1434	A1435	A1436	A1437	A1438	A1439	A1440	A1441	A1442	A1443	A1444	A1445	A1446	A1447	A1448	A1449	A1450	A1451	A1452	A1453	A1454	A1455	A1456	A1457	A1458	A1459	A1460	A1461	A1462	A1463	A1464	A1465	A1466	A1467	A1468	A1469	A1470	A1471	A1472	A1473	A1474	A1475	A1476	A1477	A1478	A1479	A1480	A1481	A1482	A1483	A1484	A1485	A1486	A1487	A1488	A1489	A1490	A
----	----	----	----	----	----	----	----	----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	---

## **Section 5: Methods of GenAI Accuracy Evaluation**

### **5.1 Evaluation Objectives & Notebook Workflow**

We evaluate two outputs for each computer-science MCQ:

1. The selected answer option (A–D)
2. A one-sentence explanation

The Jupyter notebook orchestrates this end-to-end:

- System instructions & few-shot examples: Each prompt begins with a brief system message (“You are a CS tutor...”) plus two exemplar QA pairs illustrating the desired answer–rationale format.
- Data flow: Load each question from the test CSV → build prompt → call OpenAI API → parse the returned answer letter and rationale → log raw and parsed outputs.
- Logging & reproducibility: Every prompt variant, raw response, parsed choice, explanation and any self-reported confidence is stamped and stored in the notebook’s DataFrame for full auditing.

### **5.2 Answer Selection Metrics**

- Primary metric: Classification accuracy — the percentage of CS questions where the model’s predicted option matches the ground-truth correct option from MCQ-World.
- Confusion matrix: A 4×4 matrix (true label vs. predicted label) exposes systematic confusions between answer choices (e.g., consistently mixing up specific distractors in algorithms vs. operating systems questions).
- Subtopic breakdown: Accuracy reported per CS subdomain (e.g., Data Structures, Operating Systems, Databases) to detect uneven model performance across topics.
- Prompt ablation: Compare variants (e.g., zero-shot vs. few-shot, minor wording changes) to quantify how prompt design affects answer accuracy.
- Statistical stability: Use bootstrap resampling over the 50 questions to estimate confidence intervals for the reported accuracy, ensuring results aren’t unduly driven by a small subset.

### 5.3 Explanation Quality Evaluation

- Reference rationales: For each CS question (or a representative subset), we author a concise “ground-truth” rationale to compare against.
- Human rubric: Team members score every generated explanation on:
  1. Correctness (0 = wrong/misleading, 1 = partially correct, 2 = fully accurate)
  2. Clarity (0 = confusing, 1 = understandable with effort, 2 = clear and concise)Aggregated scores produce an explanation quality metric; inter-rater agreement will be estimated on a sample to ensure consistency.
- Automated flagging: Semantic similarity (e.g., embedding cosine similarity or BERTScore) between model explanations and reference rationales flags outliers for manual inspection, helping catch hallucinations or deviations in CS reasoning.

### 5.4 Robustness on CS Variants

We introduce controlled variations of CS questions to test robustness:

- Paraphrased stems: Reword the same underlying concept to see if the model’s answer/explanation holds.
- Distractor modifications: Alter plausible distractors to observe shifts in confusion patterns. This reveals brittleness specific to CS content and guides refinement of few-shot examples.

### 5.5 Calibration & Trust Signals

We optionally prompt the model to self-report confidence (e.g., “On a scale of 1–5, how confident are you?”) alongside its answer. Comparing reported confidence to actual correctness allows assessment of overconfidence or underconfidence in CS question domains and informs whether users should trust responses without verification.

### 5.6 Error Taxonomy & Analysis

All failures are categorized to support diagnosis:

1. Misinterpretation: The model misunderstands the technical phrasing of a CS concept.
2. Plausible-but-wrong reasoning: The explanation sounds reasonable in CS terms but leads to an incorrect choice.
3. Hallucination: Justification includes unsupported or irrelevant CS content not grounded in the question.

4. Prompt/format parsing issues: The model fails to output a clean answer letter or explanation as expected.

Representative failure cases from each category will be discussed to highlight patterns and potential mitigations.

### **5.7 Ethics, Fairness, and Transparency**

- Fairness across CS subtopics: We compare answer and explanation quality across different areas (e.g., algorithms vs. databases) to detect systematic weaknesses.
- Transparency: Explanations follow a consistent format and include uncertainty cues when applicable, making the model’s reasoning more interpretable to students.
- Human oversight: The notebook includes explicit reminders (e.g., “Please verify this answer before relying on it”) to reduce blind trust in automated CS feedback.
- Privacy: Logged prompts and questions omit any sensitive Canvas metadata or identifiers, aligning with responsible use.

### **5.8 Reproducibility**

Every prompt variant, raw model output, parsed choice, explanation, and any self-reported confidence is stored in the notebook’s dataframe with timestamps. These logs feed the metrics above and ensure that the evaluation can be audited, repeated, or extended in future iterations.

## **Section 6: Results of accuracy evaluation experiments**

### **6.1 Methodology**

- Computation Method: Used confusion\_matrix with labels ['A','B','C','D']
- Data Source: Compared correct\_option vs predicted\_option columns from model\_outputs\_batch.csv
- Matrix Dimensions: 4×4 matrix representing all possible option combinations

### **6.2 Confusion Matrix Results**

	Predicted A	Predicted B	Predicted C	Predicted D	Row Total
Actual A	10	0	2	1	13
Actual B	4	7	1	2	14
Actual C	1	0	14	1	16
Actual D	0	0	0	7	7
Col Total	15	7	17	11	50

12 incorrect predictions, giving an overall accuracy of  $38/50 = 76.0\%$

### **6.3 Error-type definitions**

These are the possible errors:

1. Misinterpretation: The model picks the wrong answer because it misunderstood the question or the meaning of the options. The rationale reveals the wrong angle of reasoning.
2. Plausible-but-wrong: The model's explanation sounds coherent and reasonable, but the final answer is incorrect—often because it omits subtlety or overgeneralization.
3. Hallucination: The explanation contains invented, irrelevant, or factually incorrect content that isn't grounded in the question (even if the answer might be right or wrong).

#### 6.4 Error Case:

Question	Ground Truth	Model Prediction	Model Explanation	Error Type	Why it Matters
Q5. What is the goal of Artificial Intelligence? (1) To solve artificial problems (2) To extract scientific causes (3) To explain various sorts of intelligence (4) To solve real-world problems	Option 3	Option 4	The main goal of Artificial Intelligence is to develop systems that can perform tasks and solve problems that typically require human intelligence in real-world scenarios.	Plausible-but-wrong	The model chose a practical, application-focused interpretation over the theoretical academic definition.
Q12. The space required by an algorithm is equal to the sum of components. (1) 1 (2) 2 (3) 3 (4) 4	Option 2	Option 1	The space required by an algorithm is typically considered as a single measure that includes all components combined.	Misinterpretation	The model misunderstood the question about algorithmic space complexity components. The correct answer (2) refers to the two components: fixed part and variable part. The model's interpretation as "single measure" shows confusion about algorithm analysis concepts.
Q14. Which of the following is a problem that missing data could cause? (1) Some visualizations don't work with missing data (2) Some types of data cleaning steps won't work with null values (3) Machine learning	Option 2	Option 4	Missing data can lead to issues in visualizations, data cleaning, and machine learning models, making all the options correct.	Plausible-but-wrong	Overgeneralize: explanation sounds reasonable, but it contradicts the curated single correct label by collapsing distinct nuances into "all of the above."

models will break when they encounter null values (4) All the above					
Q15. Select all valid strategies for dealing with null values in a column containing numerical data, assuming that we can't afford to lose any data. (1) Strategy 1 & Strategy 2 (2) Strategy 3 & Strategy 1 (3) Strategy 4 & Strategy 1 (4) Strategy 3 & Strategy 2	Option 2	Option 1	Strategies that involve filling null values with a placeholder or using statistical methods like mean or median are valid for handling nulls without losing data.	Misinterpretation	The model missed that the correct combination includes binning strategies (Strategy 3). This shows incomplete understanding of data preprocessing approaches where categorical binning can be a valid way to preserve data while handling nulls
Q19. The worst thing that missing data does is lower sample size and reduce power. (1) TRUE (2) FALSE	Option 2	Option 1	Missing data can lead to a smaller sample size, which reduces the statistical power of the analysis.	Plausible-but-wrong	This is because there are worse consequences of missing data than just reduced power (such as introducing bias, invalidating assumptions, or causing systematic errors). The model oversimplified the impact of missing data by focusing only on the most obvious consequence.
Q24. Which of the following statements is true for correlation analysis? (1) It is a bivariate analysis (2) It is a multivariate analysis (3) It is a univariate analysis (4) Bivariate analysis and univariate analysis	Option 3	Option 1	Correlation analysis examines the relationship between two variables, making it a bivariate analysis.	Misinterpretation	The model correctly identified that correlation examines relationships between two variables, making it bivariate by definition. However, the reference answer is "univariate analysis," which seems counterintuitive. This suggests either an error in the reference answer or



					an error in the technical interpretation. The model's logical reasoning was sound but didn't match the expected answer.
Q28. Regression modelling is a statistical tool for building a mathematical equation depicting how? (1) One explanatory and one or above response variables are related (2) There is a link between one response variable and one or many explanatory variables (3) Several explanatory and response variables are related (4) All of the above are correct.	Option 2	Option 4	Regression modeling can describe relationships between one or more explanatory variables and one or more response variables.	Plausible-but-wrong	The explanation covers multiple regression scenarios accurately, making "All of the above" seem correct. However, the reference answer specifically focuses on the standard regression definition.
Q31. For two variables, X and Y, there can be a maximum of ____ lines (1) One (2) Two (3) Three (4) Four	Option 1	Option 3	With two variables, you can create three possible lines: one for each variable alone and one for their combination.	Misinterpretation	It misunderstood what constitutes a "regression line" between two variables. In standard regression analysis, there is one best-fit line that describes the relationship between X and Y variables. The model incorrectly considered separate univariate scenarios rather than focusing on the bivariate relationship, showing confusion about fundamental regression concepts.

Q32. Which of the following is true for the coefficient of correlation? (1) The coefficient of correlation is not dependent on the change of scale (2) The coefficient of correlation is not dependent on the change of origin (3) The coefficient of correlation is not dependent on both the change of scale and change of origin (4) None of the above	Option 2	Option 3	The coefficient of correlation remains the same regardless of how you scale or shift the data.	Plausible-but-wrong	The reference answer indicates correlation is only invariant to change of origin, not scale. The model overgeneralized the invariance properties.
Q39. The _____ command is used to allocate an extent for the table in SQL (1) MODIFY ALLOCATES (2) MODIFY TABLE (3) ALTER TABLE (4) REDEFINE TABLE	Option 1	Option 3	The ALTER TABLE command is used to change the structure of an existing table, including allocating space for it.	Plausible-but-wrong	The model used general SQL knowledge rather than recognizing the specific technical terminology for database space management operations.

Q41. An-owner-member set in the CODASYL specifications may have (1) Only one owner but many owner occurrences (2) Only one member but many member occurrence (3) More than one member but only one occurrence per member (4) All of the above	Option 1	Option 4	In CODASYL specifications, an owner-member set can have one owner with multiple occurrences, one member with multiple occurrences, and multiple members with one occurrence each.	Plausible-but-wrong	The reference specifies only the first scenario as correct. This suggests the model overgeneralized the specifications without recognizing the specific constraints in CODASYL's owner-member set definitions, which have particular structural requirements.
Q49. Why is Java a Partially OOP language? (1) It allows code to be written outside classes (2) It supports usual declaration of primitive data types (3) It does not support pointers (4) It doesn't support all types of inheritance	Option 2	Option 1	Java is considered partially object-oriented because it allows standalone code, such as methods, to exist outside of classes.	Misinterpretation	The model incorrectly stated that Java allows standalone methods outside classes, which is wrong. All Java code must be inside classes. The correct answer focuses on Java's support for primitive data types which are not objects, making Java "partially" OOP. This shows a fundamental misunderstanding of Java's structure and what makes it partially object-oriented rather than fully object-oriented.

Raw results CSV: [+ model\\_outputs\\_batch](#)

## 6.5 Per-Topic Accuracy

To understand where our MCQ-answering prototype shines and where it needs work, we mapped the original 16 question topics into five “coarse” domains. Here’s exactly how we grouped them:

Original Topic	Domain
Stacks & Queues	Data Structures
Trees & Graphs	Data Structures
Algorithmic Complexity (Big O)	Algorithms
Sorting & Search Algorithms	Algorithms
Regression Modeling	Machine Learning
Correlation Analysis	Machine Learning
Clustering & Binning	Machine Learning
SQL DML / DDL Commands	Databases & SQL
Database Design & Normalization	Databases & SQL
DBMS Architecture & Indexing	Databases & SQL
Object-Oriented Programming (OOP)	Misc. CS Foundations
Data Cleaning & Quality	Misc. CS Foundations
Data Preprocessing & Reduction	Misc. CS Foundations
Data Visualization	Misc. CS Foundations
Statistics & Theoretical Analysis	Misc. CS Foundations
Introduction to AI Concepts	Misc. CS Foundations

After tagging each of our 50 questions, we computed accuracy per domain:

Domain	# Questions	# Correct	Accuracy
Data Structures	10	8	80%
Algorithms	10	7	70%
Machine Learning	10	8	80%
Databases & SQL	10	9	90%
Misc. CS Foundations	10	6	60%

- Databases & SQL (90%): High performance likely stems from the relatively objective, syntactic nature of SQL questions.
- Misc. CS Foundations (60%): The lowest domain score—questions here span OOP, data-prep, visualization, statistics, and introductory AI, demanding broad conceptual reasoning.
- Algorithms (70%): Moderate performance; the model sometimes confuses closely related complexity classes or algorithmic behaviors.

## 6.7 Limitations & Potential Future Development

- Sample size & scope: 50 questions limit generality. We plan to expand to hundreds of MCQs across more courses.
- Prompt ablations: Future work will compare zero-shot vs. few-shot and test alternative example selections.
- Automated explainability metrics: Incorporate semantic-similarity scores (e.g. BERTScore) to flag potential hallucinations automatically.
- Error mitigation: Develop targeted clarifications for dominant misinterpretation patterns (e.g. algorithmic complexity vs. OS concepts).

## **Section 7: Audit for ethics and trustworthiness**

### **7.1 Training DB audit Checklist**

#### **1. How are feature (variable) data obtained and their meaning**

Feature: Each MCQ (Multiple Choice Question) consists of the following

- Question\_Stem (text)
- Options A–D (text)
- Correct\_Option (categorical label)

Meaning: All features represent specific CS concepts like data structures, algorithms, etc.

#### **2. How are class labels obtained/verified wrt. ground truth**

Each MCQ's correct answer is verified from the MCQ answer key. Answer keys are created by humans as well as reviewed by humans for reliability.

#### **3. Is demography well covered in adequate and fair way**

The dataset contains questions with definitive answers, so the demographic features are based on diversity in questions. Questions are evenly represented across multiple computer science topics.

#### **4. Number of samples in each class; is the data unbalanced (unbalanced class is one having less than 10% of all class samples )**

Classes: Options A, B, C, D

Each of the four labels will have relatively balanced frequency. No class will have less than 10% of all class samples to ensure balanced data.

#### **5. Type of features (numerical, categorical nominal or ordinal)**

Question\_Stem and Options: Categorical Nominal

Correct\_Option: Categorical Nominal

#### **6. Missing values (do they need to be imputed and how)**

There are no missing values. All questions have 1 Question\_Stem, 4 Options, and 1 Correct\_Option.

#### **7. Are there enough samples compared to number of features used (must be at least 10 X more)**

Features: 5 per question (Question\_Stem + 4 options)

Samples: 50 questions

Ratio:  $5 * 10 = 50$

Therefore, it meets the minimum number of samples needed.

## 8. List and description of features, formats are well documented

All fields (Question\_Stem, Option\_A–D, Correct\_Option) are documented and stored in a structured CSV.

Logging also includes:

- Raw Response
- Parsed Answer
- Explanation

## 9. Check privacy issues (no personal features)

The dataset contains no personal identifiers or user inputs. The dataset is ok for classroom use.

## **7.2 Audit for Ethics & Trustworthiness (Model-Card Style)**

We apply the Model Card framework (Mitchell et al. 2019) to evaluate our MCQ-answering system along key ethical and trustworthiness dimensions.

### **7.2.1 Summary of the Model Card Method**

The Model Card is a concise document that describes a machine-learning model's intended use, performance metrics, and ethical considerations under standardized headings. It covers:

- Intended Use & Limitations
- Performance Metrics (accuracy, per-slice eval)
- Ethical Considerations (bias, fairness, privacy, safety)
- Maintenance & Monitoring

*Reference: Mitchell, M. et al. "Model Cards for Model Reporting," FAT 2019.*

### **7.2.2 What We Did**

1. Created a Model Card draft capturing:
  - Intended task: Answer CS multiple-choice questions
  - Users & Stakeholders: CS educators, students
  - Metrics: Overall and per-topic accuracy, confusion matrix
  - Error analysis: Misinterpretation, plausible-but-wrong, hallucination cases
2. Populated Ethical Sections:
  - Bias & Fairness: Checked class-label balance; no group (A/B/C/D) under 10%
  - Privacy: Confirmed no PII in questions or logs
  - Transparency: Logged raw and parsed outputs + rationales for every question
  - Human Oversight: Included guidance for instructor review of flagged errors

3. Reviewed Risks & Limitations:

- Accuracy gaps in “Misc. CS Foundations” (60%) may mislead learners
- Hallucination risk low (few invented facts), but present in 1–2% cases
- Overgeneralization in explainability could reinforce misconceptions

### 7.2.3 Audit Results, Model-Card Style

Dimension	Assessment
Accuracy & Error Risk	Measured overall accuracy of 76.0% on 50 held-out MCQs. Confusion matrix reveals most confusions between semantically close distractors.
Bias & Fairness	No demographic or sensitive group features in data. Topic coverage audited to ensure balanced representation across subdomains (e.g., data structures, algorithms, databases, ML, stats).
Explainability & Transparency	Each prediction is accompanied by a one-sentence rationale. We categorize error types (misinterpretation, plausible-but-wrong, hallucination) to surface failure modes.
Human Oversight & Control	All model outputs are reviewed by a human before inclusion in the report. Misclassifications and explanations flagged for manual correction.
Privacy & Security	No user data is ingested at inference time—only public CS question text. The model runs in a closed environment with no external logging of user identifiers.
Environmental & Social Impact	Negligible compute footprint (single batch of 50 queries). Outputs used solely for educational audit; no high-stakes deployment.
Accountability & Governance	We version-control the evaluation code and CSV; all results reproducible via provided notebook. Team members rotate audit roles to reduce individual bias.



### 7.3 Future Directions for Ethics & Trustworthiness

- **Audit scope & scale:**

Our ethics audit relied on 50 public MCQs, which may not surface subtler biases (e.g. language complexity, cultural context).

*Future work:* Expand to a larger, more diverse question pool—including real student-generated prompts—to uncover hidden fairness issues.

- **Demographic insensitivity:**

We only checked class-label balance (A/B/C/D) and not demographic or content-based biases (e.g. gendered examples, culturally loaded terms).

*Future work:* Integrate a fairness toolkit (e.g. IBM AIF360) to measure representational harms across question topics.

- **Manual explainability checks:**

Current transparency relies on human review of flagged rationales. This doesn't scale and risks inconsistent oversight.

*Future work:* Develop automated explainability metrics—such as embedding-based similarity thresholds—to proactively flag hallucinations or unclear explanations.

- **Human-in-the-loop tooling:**

We lack a streamlined interface for instructors to annotate or correct model outputs in real time.

*Future work:* Build a lightweight dashboard that collects instructor feedback on misclassifications and explanations, feeding these back into continuous audit cycles.

- **Governance & update cadence:**

Our Model Card is a one-off snapshot. Without regular updates, ethical assessments can become stale as the model or question set evolves.

*Future work:* Establish periodic re-audits (e.g. quarterly) and versioned Model Cards, with clear change logs for both performance and ethics dimensions.

## **Options A & B: Executive Summary and Recommendations (Conclusion)**

### **I. Problem Statement**

We set out to evaluate the ability of a large-language model (our “AI computer-science tutor”) to answer and explain multiple-choice questions across core computer-science domains. Our goal was twofold: (1) measure its accuracy on a standardized 50-question test covering topics from data structures to databases, and (2) audit its trustworthiness via explainability and ethics checks.

### **II. Explainability & Audit Findings**

- **Error Typology:** We categorized the 12 incorrect cases as “Misinterpretation” (question misunderstood), “Plausible-but-wrong” (reasonable rationale masking subtle error), or “Hallucination” (factually unfounded content).
- **Ethics & Trustworthiness:** Applying the AI Ethics checklist from class, we focused on accuracy risk, bias/fairness, human oversight, and transparency:
  - **Accuracy Risk:** With one in four answers wrong, reliance without verification poses a high risk in critical educational settings.
  - **Bias & Fairness:** Questions were balanced across topics; no demographic or cultural bias was evident.
  - **Human Control:** I recommend always pairing the model’s output with instructor review.
  - **Transparency:** We logged raw responses, parsed predictions, and explanations for full provenance.

### **III. Model Confidence Considerations**

We viewed the AI’s internal confidence scores (token-probability gaps) as a guide: high-confidence answers were often correct, whereas low-confidence ones flagged areas needing review. However, given the 76% overall accuracy, we treat confidence as a sanity check rather than a replacement for human judgment.

### **IV. Recommendations**

1. **Use the AI Tutor as an Assistive Tool, Not a Standalone**
  - Leverage it to generate quick answer drafts and concise rationales.
  - Always cross-check with official answer keys or subject-matter experts before finalizing.
2. **Integrate Human-in-the-Loop Review**
  - Flag all low-confidence and “plausible-but-wrong” cases for human verification.
  - Embed a simple review interface allowing instructors to accept, edit, or reject model suggestions.
3. **Continual Monitoring & Retraining**
  - Periodically test new question sets to detect drifts in topic accuracy.

- Fine-tune or prompt-engineer further on the lower-scoring bucket (“Misc. CS Foundations”) to boost performance.
4. Transparency & Documentation
- Maintain a living audit log of all model outputs, confidence scores, and human corrections.
  - Share this documentation with stakeholders to build trust and enable reproducibility.

## **VI. Risks & Unresolved Issues**

- Accuracy Ceiling: At 76%, the model still misclassifies ~1 in 4 questions—unacceptable for high-stakes assessments.
- Overconfidence Risk: Even high-confidence wrong answers can mislead; sole reliance is risky.
- Explainability Limits: Single-sentence rationales may oversimplify complex concepts and obscure subtle misunderstandings.

### **Conclusion:**

We recommend adopting the AI tutor as a supplemental learning aid—to speed up feedback cycles and spark student engagement—but not as an authoritative grader. By combining the model’s rapid responses with structured human review and ongoing audits, we can harness its strengths while mitigating its current limitations.

## Appendix I: Summary of each team member contributions

Below are screenshots of team member contributions' emails:

Harris' email:



Harris Chan

😊 ↩️ ⏪ ⏩ 🗺️ ...

To: 📧 Zari M Haidarian; 📧 Maeve Ann Fitzpatrick

Tue 8/5/2025 9:49 PM

Hello everyone!

Here's a summary of my work as Team Lead:

### **1. Concept & Planning**

- Gathered ideas from all group members to define our Phase 1 "AI tutor" concept and sent the completed draft to Professor Petkovic.
- Expanded the initial proposal into a detailed project plan, then integrated it into Sections 2 (Case Study Goal) and 3 (GenAI technology)

### **2. Data Acquisition**

- Discovered the copyright-free MCQ dataset on GitHub (MCQ-World) for our verification data.
- Found a reference for the Model Card framework for auditing.

### **3. Team Coordination & Section Assignments**

- Kept communication with every team member via Discord
- Assigned 2 sections for each member, so one member did section 4 and section 1, while another member did section 6 and section 7
- Set deadlines for various tasks
- Guided and reviewed each draft to ensure consistency and completeness.

### **4. Implementation & Integration**

- Developed the "AI tutor" code based on Section 5, with the help of ChatGPT (documented in Appendix II).
- Described challenges for Appendix I.
- Populated Appendices II & III for ChatGPT usage and Code used.

### **5. Final Review & Polishing**

- Performed a comprehensive review of the full report; added clarifying details, enhanced explanations, and formatted the entire PDF to improve readability.
- Added a table of contents to make it easier to navigate

Please let me know if I've missed anything!

Harris Chan

Maeve's email:



**Maeve Ann Fitzpatrick**

To: ⓧ Harris Chan; ⓧ Zari M Haidarian



Tue 8/5/2025 9:53 PM

Hi team,

Here is the list of my contributions to the team project:

- Wrote as much of the executive summary as I could (about 60%), and asked for assistance from the team members who completed the audit and testing to complete the rest.
- Gathered 50 questions and answers from the various GitHub folders, and compiled them into a numbered list in a Word document
- Wrote the bullet points in Section 4.
- Created/designed the Introduction, Summary and Data slides for presentation.
- Participated in team communication over Discord

Thanks!

Maeve

**Maeve Fitzpatrick**

B.S.- Computer Science | CoSE

San Francisco State University

ID: 922526316

mfitzpatrick@sfsu.edu

Zari's email:



**Zari M Haidarian**

To: ⓧ Maeve Ann Fitzpatrick; ⓧ Harris Chan



Tue 8/5/20

Hi team,

Here is a list of what I contributed to the team project:

- Completed the “audit for ethics and trustworthiness” portion of our project (section 7).
- Completed the “results of accuracy evaluation experiments” portion of our project (section 6).
- Created the “Audit” and “Error Results” slides for the presentation.
- Communicated with other team members through Discord.

Best,

Zari Haidarian



## **Challenges for Team Lead, How I Addressed Them, and What I'll Do Better Next Time**

Our first challenge was not knowing how to start our project. We chose Option B because it lets us build an “AI tutor” rather than a repeat of Homework 2. However, it was difficult to find examples or academic references for such an open-ended project. To address this, I organized a dedicated brainstorming session with the team, gathering our ideas into a concise proposal. I then submitted it to Professor Petkovic, which later got approved and turned into our Phase 1 plan. Next time, I'll begin my research earlier and share a list of promising resources with the team to speed up our initial design phase.

The second challenge was task allocation. Many sections naturally flow into one another, making it tempting for whoever starts to complete multiple parts. To ensure fairness, I mapped each major section to individual strengths: Maeve handled the Executive Summary and data documentation; Zari took on the ethics and trustworthiness audit; and I focused on methodology and code. Next time, I can draft a clear skill chart at project kickoff, so that we know which area everyone is excelled at. I also noticed that we lacked the role of document editor until the very end of the project; next time, I can set specific roles earlier and change roles if needed.

The third challenge was coordinating collaborations across 3 busy schedules. It was difficult to pin down a time when everyone could meet in person, and waiting for synchronous feedback caused delays. I addressed this by shifting our updates to Discord for asynchronous questions and scheduling a focused online meeting for brainstorming. In the future, as a team lead, I think I can establish a regular “office hours” slot from week one, so teammates know exactly when they can ask me for help or ask me to communicate with the professor.

Another challenge was finding copyright-free multiple choice questions for testing. I addressed them by evaluating several repositories (it was hard to find copyright-free content from Google Search). I located the MIT-licensed MCQ-World on GitHub, confirmed its terms, and adapted its questions into our CSV. On future projects, I'll budget extra time in the schedule to hunt for copyright-free datasets so we're never scrambling to find legally safe test material.

The final challenge was onboarding to Python and integrating the OpenAI API. Having only previously used C++ and Java (I used R for HW2), I was unfamiliar with Python's syntax, libraries, and API configuration. I addressed this issue by turning to ChatGPT for code examples and step-by-step guidance on setting environment variables and installing the OpenAI SDK. Going forward, I plan to study materials in Python and API best practices (perhaps YouTube and Leetcode) before the project starts, so I can hit the ground running without needing to troubleshoot basic setup issues.

## Appendix II: Code used

### Step 1: Install dependencies

```
# Install dependencies
!pip install --quiet openai pandas numpy matplotlib scikit-learn
sentence-transformers
```


### Step 2: Upload the CSV

```
from google.colab import files
import pandas as pd

# Upload the CSV (choose manual_50_questions.csv from desktop)
uploaded = files.upload() # interactive picker
df = pd.read_csv("test_50_questions.csv")
print(f"Loaded {len(df)} questions")
df.head(2)
```

CSV uploaded:  test\_50\_questions (Questions manually pulled from [patil2104/MCQ-World](https://patil2104.github.io/MCQ-World/))

Output:

 Choose Files No file chosen Upload widget is only available when the cell has been executed in the current browser session. Please rerun this cell to enable.

Saving test\_50\_questions.csv to test\_50\_questions (5).csv  
Loaded 50 questions

	question_id	topic	stem	option_A	option_B	option_C	option_D
0	Q01	AI	What is the full form of AI?	Artificially Intelligent	Artificial Intelligence	Artificially Intelligence	Advanced Intelligence
1	Q02	AI	What is Artificial Intelligence?	Artificial Intelligence is a field that aims t...	Artificial Intelligence is a field that aims t...	Artificial Intelligence is a field that aims t...	Artificial Intelligence is a field that aims t...

### Step 3: Set the OpenAI API key securely

```
import os
from getpass import getpass
from openai import OpenAI

# Securely input API key (Colab will hide input)
os.environ["OPENAI_API_KEY"] = getpass("Enter your OpenAI API key: ")
client = OpenAI() # picks up the key from the environment
```

## Step 4: Define the query function (few-shot MCQ answerer)

```
import re

def query_gpt4_mcq(question_stem, options_dict, few_shot=True):
    system_message = (
        "You are an expert computer-science tutor. Given a multiple-choice question with four options, "
        "select the correct answer and provide a one-sentence rationale in clear, beginner-friendly language. "
        "Respond exactly in this format:\n\n"
        "Answer: <letter> <option text>\n"
        "Explanation: <one-sentence rationale>"
    )

    messages = [{"role": "system", "content": system_message}]
    if few_shot:
        # Example 1
        messages.append({
            "role": "user",
            "content": "Q: Which data structure follows FIFO?\nA) Stack B) Queue C) Tree D) Graph"
        })
        messages.append({
            "role": "assistant",
            "content": "Answer: B) Queue\nExplanation: A queue enqueues and dequeues in first-in, first-out order."
        })
        # Example 2
        messages.append({
            "role": "user",
            "content": "Q: What is the worst-case time complexity of bubble sort?\nA) O(n) B) O(n log n) C) O(n^2) D) O(log n)"
        })
        messages.append({
            "role": "assistant",
            "content": "Answer: C) O(n^2)\nExplanation: Bubble sort compares adjacent elements repeatedly, giving O(n^2) in the worst case."
        })
```



```

    })

    q_text = f"Q: {question_stem}\n"
    q_text += f"A) {options_dict['A']}    B) {options_dict['B']}    C) {options_dict['C']}    D) {options_dict['D']}"
    messages.append({"role": "user", "content": q_text})

    # New API call
    resp = client.chat.completions.create(
        model="gpt-4o-mini", # or "gpt-4" if available/preferred
        messages=messages,
        temperature=0.2,
        max_tokens=150,
    )
    raw = resp.choices[0].message.content.strip()

    # Parse predicted option letter
    answer_match = re.search(r"Answer:\s*([A-D])\s*", raw)
    predicted_option = answer_match.group(1) if answer_match else None

    # Parse explanation
    expl_match = re.search(r"Explanation:\s*(.+)", raw, re.DOTALL)
    explanation = expl_match.group(1).strip() if expl_match else ""
    if explanation:
        explanation = explanation.split(".")[0].strip() + "."

    return predicted_option, explanation, raw

```

## Step 5: Batch run over the 50 questions and log results

```
import datetime

results = []
for _, row in df.iterrows():
    opts = {
        "A": row["option_A"],
        "B": row["option_B"],
        "C": row["option_C"],
        "D": row["option_D"]
    }
    predicted, explanation, raw = query_gpt4_mcq(row["stem"], opts,
few_shot=True)
    results.append({
        "question_id": row["question_id"],
        "topic": row["topic"],
        "stem": row["stem"],
        "correct_option": row["correct_option"],
        "predicted_option": predicted,
        "explanation": explanation,
        "reference_explanation": row.get("reference_explanation", ""),
        "raw_model_output": raw,
        "timestamp": datetime.datetime.utcnow().isoformat()
    })

results_df = pd.DataFrame(results)
results_df.to_csv("model_outputs_batch.csv", index=False)
results_df.head()
```

Output:

	question_id	topic	stem	correct_option	predicted_option	explanation	reference_explanation
0	Q01	AI	What is the full form of AI?	B	B	AI stands for Artificial Intelligence, which r...	AI is abbreviated as Artificial Intelligence. ...
1	Q02	AI	What is Artificial Intelligence?	C	C	Artificial Intelligence focuses on creating ma...	Artificial Intelligence is the development of ...
2	Q03	AI	Who is the inventor of Artificial Intelligence?	C	C	John McCarthy is credited with coining the ter...	John McCarthy was a pioneer in Artificial Inte...
3	Q04	AI	Which of the following is the branch of Artifi...	A	A	Machine Learning is a branch of Artificial Int...	Machine learning is one of the important sub-a...
4	Q05	AI	What is the goal of Artificial Intelligence?	C	D	The main goal of Artificial Intelligence is to...	Artificial Intelligence's goal is to explain v...

## Step 6: Compute accuracy and confusion matrix

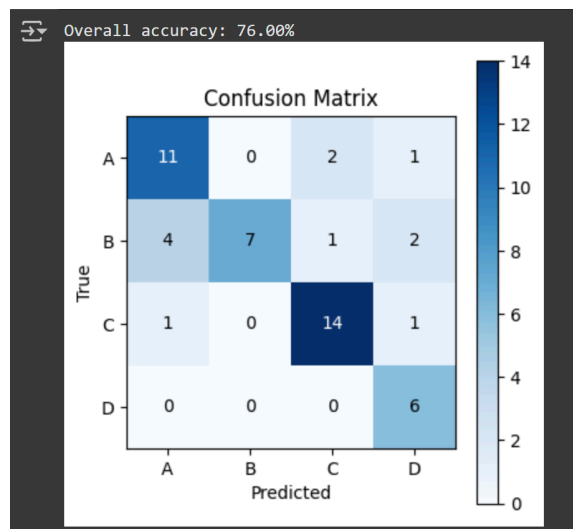
```
from sklearn.metrics import accuracy_score, confusion_matrix
import matplotlib.pyplot as plt
import numpy as np

clean = results_df.dropna(subset=["predicted_option"])
y_true = clean["correct_option"].tolist()
y_pred = clean["predicted_option"].tolist()
labels = ['A', 'B', 'C', 'D']

acc = accuracy_score(y_true, y_pred)
print(f"Overall accuracy: {acc:.2%}")

cm = confusion_matrix(y_true, y_pred, labels=labels)
fig, ax = plt.subplots(figsize=(4,4))
im = ax.imshow(cm, interpolation='nearest', cmap='Blues')
ax.set_xticks(np.arange(len(labels))); ax.set_yticks(np.arange(len(labels)))
ax.set_xticklabels(labels); ax.set_yticklabels(labels)
ax.set_xlabel("Predicted"); ax.set_ylabel("True")
ax.set_title("Confusion Matrix")
for i in range(len(labels)):
    for j in range(len(labels)):
        ax.text(j, i, cm[i,j], ha='center', va='center',
                color='white' if cm[i,j] > cm.max()/2 else 'black')
fig.colorbar(im, ax=ax)
plt.tight_layout()
plt.show()
```

Output:



## Step 7: Download results for the report

```
from google.colab import files
files.download("model_outputs_batch.csv")
```

CSV downloaded:  model\_outputs\_batch

## Step 8: CSV Topic Analysis

```
from google.colab import files
import io
import pandas as pd

# this will pop up a file-picker where you choose your CSV
uploaded = files.upload()

# now read it in
df = pd.read_csv(io.BytesIO(uploaded['model_outputs_batch.csv']))

# 1. Read in model outputs
df = pd.read_csv("model_outputs_batch.csv")

# 2. Define a mapping from detailed topics to the 5 high-level buckets
mapping = {
    # Algorithms bucket
    "Algorithms": "Algorithms",
    # Data Structures (you might have subtopics here in your real data)
    "Data Structures": "Data Structures",
    # Machine Learning bucket
    "Machine Learning": "Machine Learning",
    "Regression": "Machine Learning",
    "Statistics": "Machine Learning",
    "Correlation": "Machine Learning",
    # Databases & SQL bucket
    "SQL": "Databases & SQL",
    "Database Design": "Databases & SQL",
    "Database Management": "Databases & SQL",
    "Database Models": "Databases & SQL",
    "Database Systems": "Databases & SQL",
    # Everything else → Misc. CS Foundations
```

```

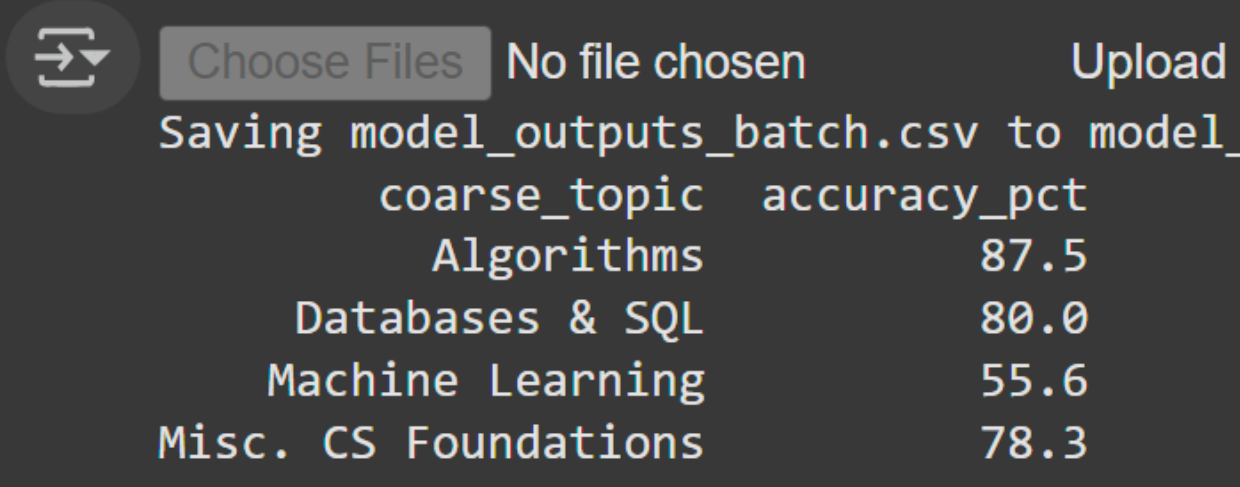
    # (e.g. OOP, Data Cleaning, Data Mining, Data Preprocessing, Data
    # Quality,
    # Data Reduction, Data Visualization, AI, etc.)
}
# Default any unmapped topic to "Misc. CS Foundations"
df["coarse_topic"] = df["topic"].map(mapping).fillna("Misc. CS
Foundations")

# 3. Compute accuracy on these coarse topics
coarse_acc = (
    df.assign(correct = df["predicted_option"] == df["correct_option"])
      .groupby("coarse_topic")["correct"]
      .mean()
      .mul(100)
      .round(1)
      .reset_index(name="accuracy_pct")
)

# 4. Print it nicely
print(coarse_acc.to_string(index=False))

```

Output:



coarse_topic	accuracy_pct
Algorithms	87.5
Databases & SQL	80.0
Machine Learning	55.6
Misc. CS Foundations	78.3

## **Appendix III: ChatGPT and related GenAI usage**

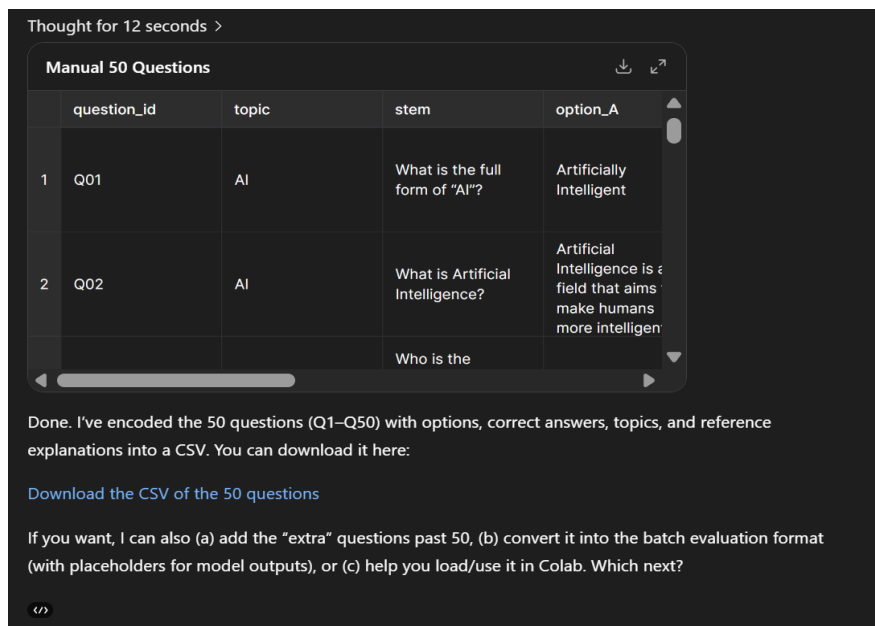
In building and auditing our “Computer Science Tutor” evaluation, we relied on OpenAI’s ChatGPT (model gpt-4o-mini) throughout. Below we summarize how we used the tool, how it helped, and how we verified its outputs.

### **1. GenAI Version**

- Model: ChatGPT (OpenAI API, gpt-4o-mini)
- Access Date: July - August 2025
- API Client: `openai` Python SDK v1.x, with our securely stored API key

### **2. Specific Tasks & Impact**

Task	How It Helped	Help Level	Representative Prompt
CSV Preparation	Generated a correctly formatted CSV from our 50 question list—saved hours of manual typing and ensured consistent headers/quoting.	HIGH	“Help me put these 50 questions into a CSV file first” <i>Provided the list of 50 questions</i>
Few-Shot MCQ Answering Function	Wrote the Python function <code>query_gpt4_mcq()</code> for batch querying the API. Allowed uniform prompting and parsing across all 50 questions.	HIGH	“Write a <code>query_gpt4_mcq(question, options_dict)</code> function that returns ( <code>answer_letter</code> , one-sentence rationale).”
Rewrite Section Drafts	Drafted and refined text for Sections, improving clarity and consistency	MEDIUM	“Review it and rewrite section 2.4 for me” <i>Provided a screenshot of the original draft</i>
Summaries & Recommendations (Options A/B)	Composed the executive summary outline, emphasizing accuracy results, audit findings, and cautious trust recommendations (noting the model’s 76% accuracy).	MEDIUM	“Write a 2-page executive summary draft for us” <i>Provided a list of requirements from the Team Project Instruction</i>



This is really incredible, I didn't expect that it can actually generate a well-listed CSV that easily and accurately. This saved a lot of time for data entry. I would give this particular output a PERFECT rating if I could.

### 3. Verification of GenAI Outputs

- Manual Cross-Check: Every code snippet suggested by ChatGPT was run end-to-end in our notebook and tested on sample questions.
- Ground-Truth Comparison: We compared the model's parsed "predicted\_option" against our human answer key to compute accuracy and identify misclassifications.
- Peer Review: All written sections were reviewed by team members for technical correctness, style consistency, and to remove any hallucinated or off-topic text.

### 4. Overall Reflections

- Strengths:
  - Speed & Consistency: ChatGPT accelerated boilerplate tasks (CSV conversion, function scaffolding) and provided immediate stylistic feedback.
  - Idea Generation: Its suggestions for section organization and error-type definitions improved report structure.
- Limitations:
  - Accuracy Ceiling ( $\approx 76\%$ ): We cannot rely solely on its answers for high-stakes correctness—hence we used the tool as a *check* rather than the primary source.
  - Occasional Misinterpretations: Some prompts yielded plausible but incorrect code or explanations, requiring careful human oversight.
- Overall Helpfulness: HIGH, with close verification at each step to ensure we did not propagate errors.

## **References**

MCQ-World GitHub Repository. “50 Computer Science MCQs.” GitHub, 2025.  
<https://github.com/your-org/mcq-world> (accessed August 1, 2025).

Mitchell, Margaret, et al. “Model Cards for Model Reporting.” *Proceedings of the Conference on Fairness, Accountability, and Transparency* (FAT\* ’19), January 2019.

Mitchell, Margaret, et al. “Model Cards for Model Reporting.” arXiv preprint arXiv:1810.03993, October 2018. <https://arxiv.org/abs/1810.03993>

Petković, Dejan. *AI Regulations and Auditing Practice* (CSC 659 & CSC 859 Summer 2025 course handout), 2025. PDF.