# Exploring the MIT Mathematics and EECS Curriculum Using Large Language Models

**Sarah J. Zhang**[*]
MIT
sjzhang@mit.edu

**Sam Florin**[*]
MIT
sflorin@mit.edu

**Ariel N. Lee**
Boston University
ariellee@bu.edu

**Eamon Niknafs**
Boston University
en@bu.edu

**Andrei Marginean**
MIT
atmargi@mit.edu

**Annie Wang**
MIT
annewang@mit.edu

**Keith Tyser**
Boston University
ktyser@bu.edu

**Zad Chin**
Harvard University
zadchin@college.harvard.edu

**Yann Hicke**
Cornell University
ylh8@cornell.edu

**Nikhil Singh**
MIT
nsingh1@mit.edu

**Madeleine Udell**
Stanford University
udell@stanford.edu

**Yoon Kim**
MIT
yoonkim@mit.edu

**Tonio Buonassisi**
MIT
buonassi@mit.edu

**Armando Solar-Lezama**
MIT
asolar@csail.mit.edu

**Iddo Drori**
MIT, Columbia University, Boston University
idrori@csail.mit.edu

## Abstract

We curate a comprehensive dataset of 4,550 questions and solutions from problem sets, midterm exams, and final exams across all MIT Mathematics and Electrical Engineering and Computer Science (EECS) courses required for obtaining a degree. We evaluate the ability of large language models to fulfill the graduation requirements for any MIT major in Mathematics and EECS. Our results demonstrate that GPT-3.5 successfully solves a third of the entire MIT curriculum, while GPT-4, with prompt engineering, achieves a perfect solve rate on a test set excluding questions based on images. We fine-tune an open-source large language model on this dataset. We employ GPT-4 to automatically grade model responses, providing a detailed performance breakdown by course, question, and answer type. By embedding questions in a low-dimensional space, we explore the relationships between questions, topics, and classes and discover which questions and classes are required for solving other questions and classes through few-shot learning. Our analysis offers valuable insights into course prerequisites and curriculum design, highlighting language models' potential for learning and improving Mathematics and EECS education.

---

[*]Equal contribution

# 1 Introduction

Large language models (LLMs) have demonstrated the ability to pass exams from individual college-level courses [17]. However, a systematic evaluation of this ability across entire curricula that students would be expected to go through to obtain their college degrees has yet to be explored due to the lack of a central repository of questions from actual exams and assignments at the curriculum level. In this paper, we introduce MITCOURSES, a curated dataset of 4,550 questions and their solutions spanning exams and assignments from all courses that form the curriculum for MIT's Mathematics and Electrical Engineering and Computer Science (EECS) majors. The dataset is seven times larger than our previously released dataset [7, 8] and covers all the course requirements for an undergraduate degree.

Using this dataset, we benchmark four state-of-the-art language models, GPT-4 [17], GPT-3.5, StableVicuna [22], and LLaMA [23], which vary in their sizes, their capabilities, and whether their weights are publicly available or not. We evaluate different LLM prompting techniques (few-shot [5], chain of thought [12], self-critique [11, 14, 19, 21, 9]) and document their effect on the model success rates. We propose a new prompting technique which we call *expert prompting*, where we ask the model to suggest named experts on a given question, then ask for the answer the named experts would have given and subsequently make a collective decision. Our experiments show expert prompting improves performance relative to prior prompt engineering techniques.

For models whose weights are publicly available, we also demonstrate that fine-tuning improves performance on our test set and other reasoning task benchmarks.

In addition to providing performance measures across various models, we demonstrate the application of our models for curriculum design for college majors. An essential aspect of curriculum design is determining the appropriate sequence of courses to ensure prerequisites are established effectively. Traditionally, this process is manual, relying on human input to identify key concepts and learning outcomes. However, this method is subjective, and coordinating input from various faculty members teaching different courses can be challenging. Instead, we use the embeddings for the questions of different courses to discover dependencies between courses. Given the challenges of accurate student evaluation in a world where large language models are readily available, we also propose the development of new meta-questions that focus on assessing the correctness and completeness of students' understanding rather than their ability to generate correct answers.

To prevent our dataset being incorporated as part of LLM training corpora, the dataset will not be made publicly available but will be made available to researchers upon request through a data use agreement (DUA). This is both to preserve the value of this dataset as a resource for benchmarking LLMs—which is lost if the dataset becomes part of the training set for future public models—as well as to respect the wishes of instructors who contributed to it and who want to preserve their ability to use problems from this dataset in the future. We release our open-source LLaMA models fine-tuned on our dataset, called MIT-LLM, in the supplemental material.

**Related Work.** Large language models such as GPT-4 demonstrate excellent results on standardized AP tests [17, 6]. Their performance across many reasoning tasks has been enhanced by few-shot learning and by providing intermediate reasoning steps and chain-of-thought (CoT) prompts [16, 25, 24, 26]. Iterative prompting techniques such as self-critique, self-refinement, self-feedback, self-reflection, and self-improvement [11, 14, 19, 21, 9] have further been shown to improve performance. While much of the early work on eliciting better reasoning capabilities of LLMs has focused on prompt engineering, their capabilities have recently expanded by combining LLMs with reinforcement learning. In this work, we write an optimization objective that accurately describes the problem, formulates these improvements, and provides an ablation study. LLMs are now routinely used for mathematical reasoning, solving university-level mathematics and computer science problem sets and final exams [7, 13, 20, 8, 2] at the human level. Our work extends this to entire undergraduate degrees.

Table 1: MIT Mathematics and EECS undergraduate degree questions by course. The table includes the level and number of questions and parts.

| ID | Course | Course Name | Level | Questions | Parts |
|----|--------|-------------|-------|-----------|-------|
| 1 | 6.100A | Intro to CS Programming in Python | 0 | 34 | 47 |
| 2 | 18.100B | Real Analysis | 1 | 60 | 66 |
| 3 | 18.102 | Intro to Functional Analysis | 2 | 68 | 104 |
| 4 | 18.C06 | Linear Algebra & Optimization | 1 | 77 | 195 |
| 5 | 6.1210 | Intro to Algorithms | 2 | 82 | 164 |
| 6 | 6.1220 | Design & Analysis of Algorithms | 3 | 44 | 158 |
| 7 | 6.3900 | Intro to Machine Learning | 2 | 114 | 619 |
| 8 | 18.303 | Linear Partial Differential Equations | 2 | 22 | 65 |
| 9 | 18.200 | Principles of Discrete Applied Math | 2 | 45 | 86 |
| 10 | 6.1800 | Computer Systems Engineering | 3 | 58 | 112 |
| 11 | 18.702 | Algebra II | 3 | 52 | 94 |
| 12 | 18.701 | Algebra I | 2 | 58 | 87 |
| 13 | 18.01 | Calculus I | 0 | 203 | 495 |
| 14 | 6.4110 | Rep., Inference, & Reasoning in AI | 3 | 54 | 324 |
| 15 | 6.1010 | Fundamentals of Programming | 1 | 22 | 31 |
| 16 | 18.704 | Seminar in Algebra | 3 | 16 | 25 |
| 17 | 6.4120 | Computational Cognitive Science | 3 | 10 | 67 |
| 18 | 6.1020 | Elements of Software Construction | 2 | 26 | 52 |
| 19 | 18.02 | Calculus II | 0 | 81 | 154 |
| 20 | 18.600 | Probability & Random Variables | 1 | 65 | 160 |
| 21 | 6.8611 | Quantitative Methods for NLP | 3 | 20 | 31 |
| 22 | 18.404 | Theory of Computation | 3 | 53 | 101 |
| 23 | 6.1910 | Computation Structures | 2 | 72 | 198 |
| 24 | 18.03 | Differential Equations | 1 | 66 | 160 |
| 25 | 6.2000 | Electrical Circuits | 2 | 27 | 97 |
| 26 | 18.300 | Principles of Continuum Applied Math | 2 | 43 | 90 |
| 27 | 6.3000 | Signal Processing | 2 | 55 | 258 |
| 28 | 6.2300 | Electromagnetic Waves & Applications | 2 | 37 | 142 |
| 29 | 6.3010 | Signals, Systems & Inference | 3 | 57 | 224 |
| 30 | 18.901 | Intro to Topology | 2 | 58 | 144 |
| Mean | | | | 55.97 | 151.67 |
| Total | | | | 1679 | 4550 |

# 2 Methods

## 2.1 Dataset

We collect and curate a comprehensive dataset of 4,550 questions and corresponding solutions from 30 MIT Mathematics and EECS courses required to graduate from the institute. This includes a broad range of core and elective courses, providing students with the foundational and specialized knowledge necessary to succeed in the field. The dataset spans eight MIT Mathematics and EECS undergraduate degree paths: (1) 6-1: Electrical Science and Engineering, (2) 6-2: Electrical Engineering and Computer Science, (3) 6-3: Computer Science and Engineering, (4) 6-4: Artificial Intelligence and Decision Making, (5) 18-1: General Mathematics, (6) 18-2: Applied Mathematics, (7) 18-3: Pure Mathematics, and (8) 18-C: Mathematics with Computer Science.

We use course materials from the past two years to construct this dataset. We download the PDF documents associated with each course, such as the syllabus, problem sets, midterm exams, and final exams, and manually curate the data. The breakdown of the level and number of questions and parts for each course is shown in Table 1. We ensure that the dataset contains no personally identifiable information (PII).

For each course, we curate all the questions and solutions from the problem sets, midterm exams, and final exams. We first extract questions and solutions for all parts of all types of questions from the PDF documents and use Mathpix [15] for an initial transcription. We then process the data by manually correcting each question and answer to ensure quality and correctness. The questions in the dataset are formatted by concatenating the question's set-up and context information first, followed by the part's context and question. For each question part, the task is categorized as an exercise, final exam, lab, midterm exam, project, or problem set. The question type can be image or text, and the solution type can be an expression, image, multiple choice, numerical, open, or programming. The breakdown of the number of question parts by each task type is shown in Table 2 and by each answer type is shown in Table 3.

Table 2: MIT Mathematics and EECS undergraduate degree questions parts by task type. The table includes the number of parts and GPT-3.5 solve rate for each task type.

| ID | Task Type | Parts | Solve Rate |
|---|---|---|---|
| 1 | Exercise | 198 | 0.73 |
| 2 | Problem Set | 2820 | 0.41 |
| 3 | Final Exam | 418 | 0.37 |
| 4 | Midterm Exam | 799 | 0.36 |
| 5 | Lab | 278 | 0.26 |
| 6 | Project | 37 | 0.10 |
| Mean | | 758.3 | 0.37 |
| Total | | 4550 | 0.36 |

Table 3: MIT Mathematics and EECS undergraduate degree questions parts by solution type. The table includes the number of parts and GPT-3.5 solve rate for each solution type.

| ID | Solution Type | Parts | Solve Rate |
|---|---|---|---|
| 1 | Programming | 234 | 0.51 |
| 2 | Multiple Choice | 710 | 0.33 |
| 3 | Numerical | 634 | 0.23 |
| 4 | Expression | 969 | 0.21 |
| 5 | Open | 1821 | 0.42 |
| 6 | Image | 182 | 0.06 |
| Mean | | 758.3 | 0.29 |
| Total | | 4550 | 0.36 |

## 2.2 Benchmark

Once the dataset is released publicly, it can be used for training LLMs, causing it to lose its value as a benchmark for evaluating these models. In addition, acquiring the necessary permissions to release the questions and answers from each course is challenging since instructors would like to reuse their class material. Therefore, instead of releasing the raw dataset, we fine-tune an open-source model using our dataset and then make the fine-tuned model available to the public. We divide the dataset into separate training and test sets to ensure a fair evaluation. We benchmark the open-source LLM by comparing its performance on the test set before and after fine-tuning on the training set. We then benchmark the performance of open-source and closed models on our data. This allows us to maintain the integrity of the dataset while still providing valuable resources to the research community.

## 2.3 Embedding Course Questions for Curriculum Design

We embed all course questions into 1536-dimensional vectors using OpenAI's *text-embedding-ada-002* [18]. The embedded text consists of the topic of the question concatenated with the question itself. This embedding performs well on sentence similarity, can process long course questions, and is computationally efficient, embedding all course questions in five minutes. We make curriculum decisions by analyzing the similarity between question embeddings from the same and different courses. For example, a dependency graph of prerequisite classes is inferred in which the graph nodes are classes. The directed edges between two nodes measure the ability to answer questions from one class by questions from another by few-shot learning.

## 2.4 Optimizing the Prompt

Let $W$ be the context of words, $Q$ be the question, $A$ be the LLM answer, and $S$ be the ground truth solution. Our dataset $D$ is a set of $\{Q, S\}$ pairs. Our goal is to find the context of words $W$ such that the answer by the LLM given the words $W$ and the question $Q$ is an answer $A$:

$$A = f(W, Q), \tag{1}$$

where $f$ is the LLM.

An optimization objective is to find the $W$ that maximizes:

$$\max_W g(f(W, Q), S), \tag{2}$$

where $g$ is a grading function that measures the grade of the answer $A$ compared with the solution $S$.

Our optimization objective is to find the $W$ that maximizes the expected grade over the dataset $D$:

$$\max_W \mathbb{E}_{(Q,S) \in D} g(f(W, Q), S), \tag{3}$$

4

where $\mathbb{E}$ denotes the expectation over the questions and solutions in the dataset.

The context $W$ is in a very high-dimensional discrete space. This work uses several prompt engineering heuristics to find a good $W$.

## 2.5 Heuristics

We propose several heuristics for $W$:

1. **Zero-shot Learning:** Without any data or example questions in the context, the LLM attempts to answer the question directly.

2. **Few-shot Learning:** The LLM is provided with a few example questions and answers in the context to guide its understanding of the task [5]. This can be represented as:

$$W = \{(Q_1, A_1), \ldots, (Q_n, A_n)\},$$

   where $n$ is the number of example question-answer pairs.

3. **Chain-of-Thought:** Uses prompt engineering to elicit a step by step answer [25].

4. **Tree-of-Thought:** Uses prompt engineering to generate a tree of answers and then searches this tree [26] using BFS or DFS, combining an LLM with classical search algorithms.

5. **Program Synthesis:** The LLM is prompted to write a program that solves the problem, and then the program is run in an interpreter [7].

6. **Critique:** An LLM generates a critique $C$ for an answer $A$, and the answer is iteratively refined using the critiques [11, 14, 19, 21, 9]. This process can be represented as follows:

$$Q, A \rightarrow C,$$
$$Q, A, C \rightarrow A',$$
$$Q, A, C, A' \rightarrow C',$$
$$\ldots$$

7. **Expert Prompting:** A novel contribution of this work is to use the LLM to identify experts $E$ in the field, generate answers as if the experts wrote them, and combine the experts' answers by collaborative decision-making. This process is represented by using a generic expert defined as the system role such as:

   $E = $ `You are an MIT Professor of Computer Science and Mathematics teaching Calculus I.`

   for questions from, e.g., Calculus I, or specific named experts generated by the LLM using the prompt:

   $P_3 = $ `Give an educated guess of the three experts most capable of solving this question.`

   The LLM then generates the names $E$ of multiple experts:

   $$P_3, Q \rightarrow E.$$

   Next, the LLM uses the named experts as the system role to generate an answer:

   $$Q^E \rightarrow A^E.$$

   where $Q^E$ is the question being asked with the system role being $E$, for example:
   System: You are $E$.
   User: Solve $Q$

8. **Fine-tuning:** The LLM is fine-tuned on a dataset $D$ to improve its performance on specific tasks.

## 2.6 Automatic Grading

Given the question $Q$, ground truth solution $S$, and LLM answer $A$, we use GPT-4 to automatically grade the answers:

$$Q, S, A \rightarrow G. \tag{4}$$

The grading is either binary (correct or incorrect) or scaled (between 0 and 5 inclusive). Automatic grading allows us to form a cascade of answers and prompts, accepting correct answers and transferring the remaining questions to the following heuristics, until achieving a perfect score. We perform automatic grading using the same generic expert for each class.

## 2.7 Meta-Questions

We develop new meta-questions about the correctness and completeness of GPT-4. Specifically, in addition to standard questions, our meta-questions consist of questions and their answers by GPT-4. The students are asked to identify whether the answer is correct or not for each question. If the answer is correct, then the students are asked to explain why. If the answer is incorrect, then the students are asked to write the correct answer and provide a complete explanation. By teaching students how to use LLMs in an educational setting, they develop the skills they need to navigate and critically evaluate the material. Through learning how to prompt LLMs, students can complete tasks more quickly and accurately, increasing their productivity and efficiency.

# 3 Results

Table 4: MIT Mathematics and EECS undergraduate degree questions by course. The table includes the number of questions and parts, a breakdown of GPT-3.5 solve rates by task and answer type, and the total GPT-3.5 solve rate for each course. The task types are denoted as follows: Exercise (E), Final Exam (FE), Lab (L), Midterm Exam (ME), Project (PR), and Problem Set (PS). The answer types are denoted as follows: Expression (E), Image (I), Multiple Choice (M), Numerical (N), Open (O), and Programming (P). Note that the last three courses marked with * do not have solutions, so the task type and answer type solve rates represent the proportions of each respective type.

| ID | Course | Task Type Solve Rate | Solution Type Solve Rate | Solve Rate |
|---|---|---|---|---|
| 1 | 6.100A | 1 (E), 1 (ME), 0.84 (PS) | 0.93 (P) | 0.93 |
| 2 | 18.100B | 0.90 (FE), 0.67 (ME), 0.79 (PS) | 0 (I), 0.80 (O) | 0.79 |
| 3 | 18.102 | 0.73 (FE), 0.71 (ME), 0.94 (PS) | 0.77 (O) | 0.77 |
| 4 | 18.C06 | 0.81 (FE), 0.76 (ME), 0.67 (PR), 0.71 (PS) | 0.72 (E), 0 (I), 1 (M), 0.84 (N), 0.78 (O), 0.71 (P) | 0.74 |
| 5 | 6.1210 | 0.78 (FE), 0.72 (ME), 0.66 (PS) | 0.21(E), 1 (I), 0.56 (M), 0 (N), 0.83 (O), 0.53 (P) | 0.72 |
| 6 | 6.1220 | 0.56 (FE), 0.51 (ME), 0.54 (PS) | 0.65 (E), 0.50 (M), 0.34 (N), 0.53 (O) | 0.53 |
| 7 | 6.3900 | 0.35 (E), 0.38 (FE), 0.68 (L), 0.57 (ME), 0.63 (PS) | 0.52 (E), 0.07 (I), 0.35 (M), 0.26 (N), 0.66 (O), 0.46 (P) | 0.51 |
| 8 | 18.303 | 0.08 (ME), 0.5 (PR), 0.71 (PS) | 0.37 (E), 1 (I), 1 (N), 0.47 (O) | 0.49 |
| 9 | 18.200 | 0.45 (ME), 0.50 (PS) | 0.07 (E), 0 (I), 0.28 (N), 0.67 (O) | 0.48 |
| 10 | 6.1800 | 0.51 (ME), 0 (PR), 0.59 (PS) | 0 (E), 0.54 (N), 0.25 (N), 0.07 (O) | 0.45 |
| 11 | 18.702 | 0.58 (ME), 0.36 (PS) | 0.23 (E), 0.49 (M), 0.76 (N), 0.43 (O) | 0.42 |
| 12 | 18.701 | 0.37 (ME), 0.46 (PS) | 0.39 (O) | 0.39 |
| 13 | 18.01 | 0.28 (FE), 0.37 (ME), 0.45 (PS) | 0.44 (E), 0.07 (I), 0.25 (M), 0.27 (N), 0.57 (O) | 0.36 |
| 14 | 6.4110 | 0.23 (FE), 0.40 (ME), 0.37 (PS) | 0.26 (E), 0.33 (M), 0.16 (N), 0.42 (O), 0.52 (P) | 0.34 |
| 15 | 6.1010 | 0.38 (L), 0.26 (ME) | 0.56 (E), 0.17 (M), 0 (N), 0.41 (O), 0.35 (P) | 0.32 |
| 16 | 18.704 | 0 (PR), 0.61 (PS) | 0.31 (O) | 0.31 |
| 17 | 6.4120 | 0 (PR), 0.19 (PS) | 0 (E), 0.08 (O) | 0.31 |
| 18 | 6.1020 | 0.27 (ME), 0 (PR), 0.35 (PS) | 0.75 (E), 0.19 (M), 0.18 (O), 0.32 (P) | 0.30 |
| 19 | 18.02 | 0.13 (FE), 0.38 (ME), 0.35 (PS) | 0.23 (E), 0.31 (I), 0 (M), 0.20 (N), 0.57 (O) | 0.29 |
| 20 | 18.600 | 0.38 (FE), 0.13 (ME), 0.41 (PS) | 0.22 (E), 0.21 (N), 0.62 (O) | 0.29 |
| 21 | 6.8611 | 0.02 (PR), 0.52 (PS) | 0 (I), 0.30 (O) | 0.28 |
| 22 | 18.404 | 0.31 (FE), 0.12 (ME), 0.31 (PS) | 0 (E), 0 (I), 0.46 (M), 0.27 (O) | 0.27 |
| 23 | 6.1910 | 0.40 (E), 0.13 (ME), 0.04 (L) | 0.17 (E), 0 (I), 0.55 (M), 0.17 (N), 0.10 (O), 0.06 (P) | 0.19 |
| 24 | 18.03 | 0.05 (FE), 0.14 (ME), 0.29 (PS) | 0.11 (E), 0.21 (I), 0.27 (N), 0.29 (O) | 0.16 |
| 25 | 6.2000 | 0.02 (FE), 0.04 (ME), 0.29 (PS) | 0.02 (E), 0 (I), 0.38 (N), 0.37 (O) | 0.11 |
| 26 | 18.300 | 0 (PR), 0.12 (PS) | 0.07 (E), 0 (I), 0 (N), 0.12 (O) | 0.08 |
| 27 | 6.3000 | 0.05 (FE), 0.01 (ME), 0.13 (PS) | 0.13 (E), 0.02 (I), 0 (M), 0.06 (N), 0.22 (O) | 0.06 |
| 28* | 6.2300 | 0.09 (ME), 0.91 (PS) | 0.18 (E), 0.10 (I), 0.41 (N), 0.31 (O) | N/A |
| 29* | 6.3010 | 0.09 (FE), 0.12(ME), 0.79 (PS) | 0.41 (E), 0.09 (I), 0.06 (M), 0.09 (N), 0.36 (O) | N/A |
| 30* | 18.901 | 0.18 (ME), 0.82 (PS) | 0.02 (E), 0.01 (M), 0.97 (O) | N/A |
| Mean | | | | 0.33 |

Figure 1: (A) Number of question parts by course, (B) Solve rate by course, (C) GPT-3.5 solve rate by task type, and (D) Solve rate by answer type.
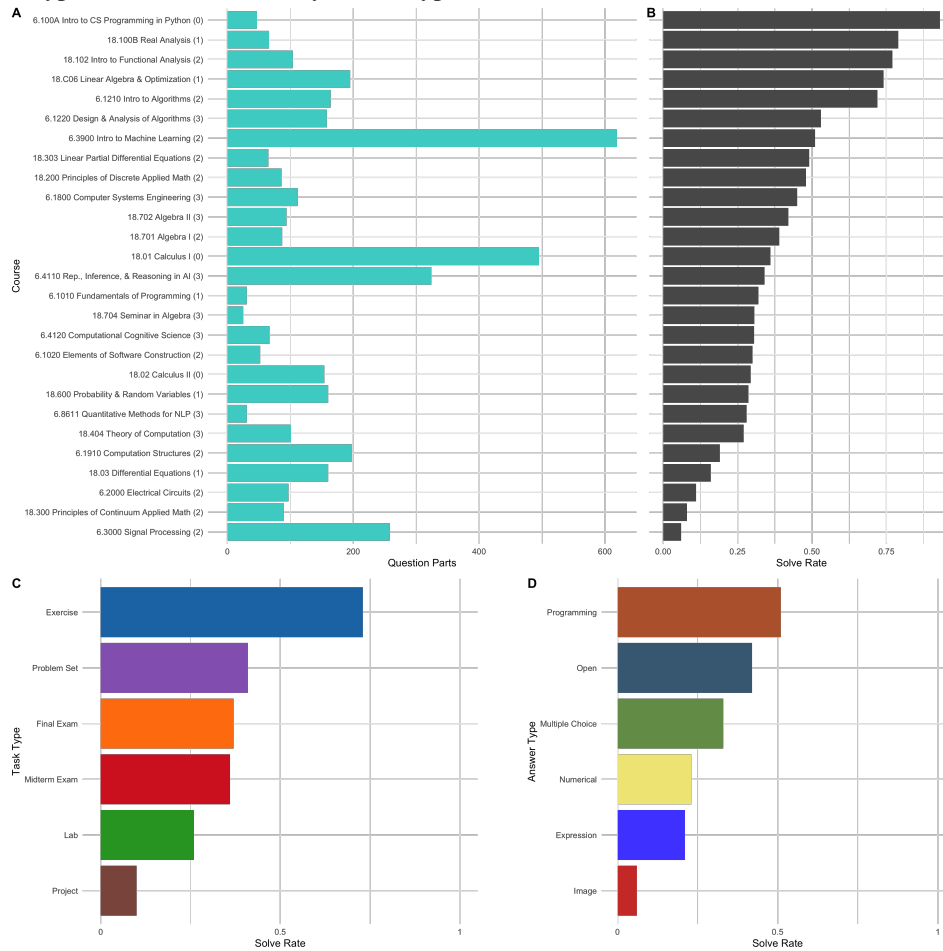


Table 4 shows the breakdown of the task type, solution type, and zero-shot GPT-3.5 solve rates for each course. Figure 1 shows course difficulty level by GPT-3.5 solve rate, task type, and answer type. The difficulty level of questions by task type from easy to hard is: exercises, problem sets, midterms, final exams, labs, and project. The difficulty of questions by answer type from easy to hard is: programming, open, multiple-choice, numerical, expression, and image.

Figure 2 shows a graph of questions, colored by course. The connectivity shows the nearest neighbor questions, in the low-dimensional embedding space, which serve as few-shot examples. We use the Force Atlas 2 algorithm [10] in Gephi [3] to produce this layout, which shows the clustering of similar courses and questions within them.

Figure 3 uses the Fruchterman Reingold layout algorithm [4], since there are much fewer nodes, one per course instead of one per question. Edges show prerequisite relationships between courses. We filter out cases where a later course in the curriculum supports a former course, and the edge weight shows the number of questions.

We apply few-shot, chain-of-thought, self-critique, and expert prompting as a cascade. Since grading is performed automatically, we apply each method to the questions that the previous methods do not solve perfectly. From the dataset of MIT questions, a test set of 288 questions is randomly selected amongst questions without images and with solutions. First, zero-shot answers to these questions are generated by the LLMs. These answers are then graded on a scale from 0 to 5 by GPT-4. Then, for all questions without a perfect score, few-shot answers from GPT-4 using the top 3 most similar questions under the embedding are generated. These answers are again graded by GPT-4 from 0 to

Figure 2: Graph of questions: nodes of the graph represent questions and edges represent the closest questions used for few shot-learning.



5. Then, few-shot and chain-of-thought are applied to questions not receiving a 5 out of 5 in the previous answers. Subsequently, self-critique is applied to the remaining questions. Finally, once grading on these answers is completed across all experts, all questions are solved correctly. A similar process is performed for the ReClor [27] validation set consisting of 500 questions.
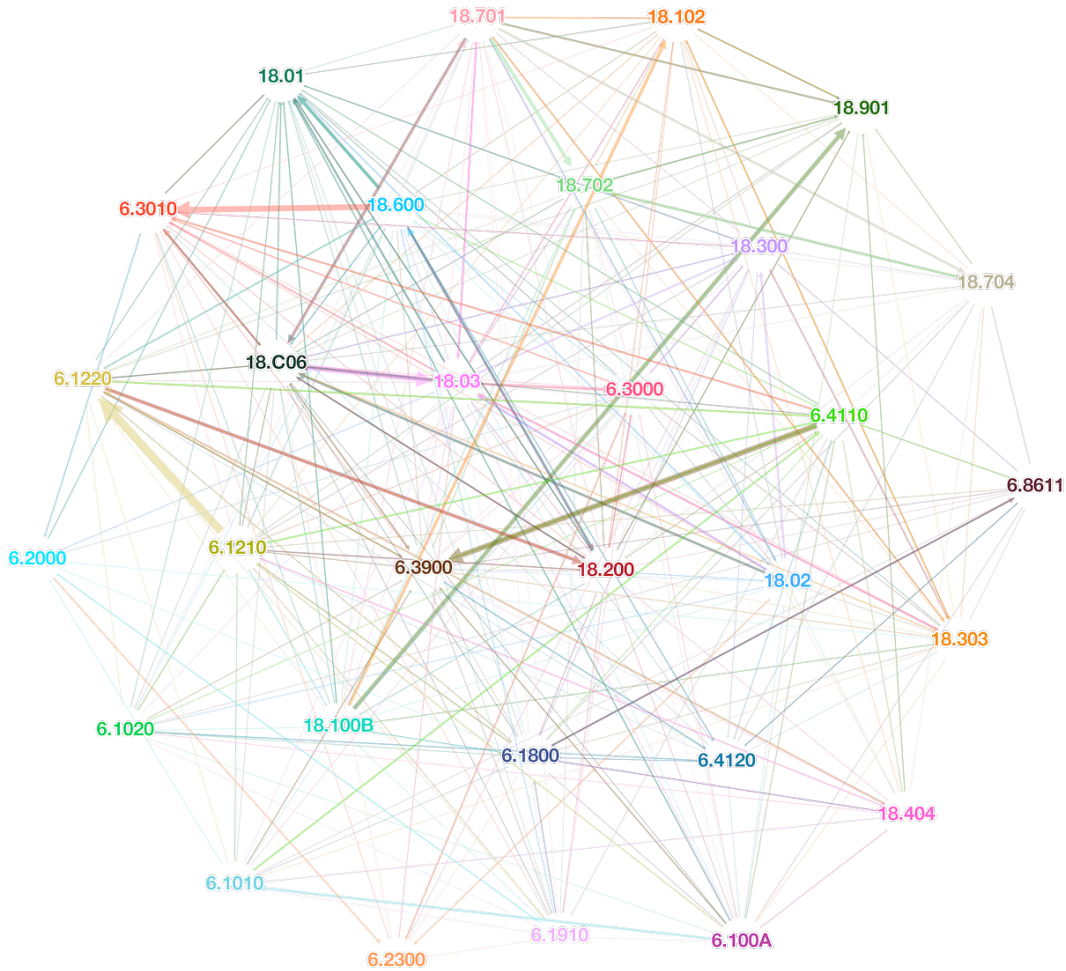
Table 5: Solve rates of LLMs on the MIT test set of 288 non-image questions of all types and the ReClor validation set of 500 multiple choice questions. Few-shot learning, chain of thought, self-critique, and experts improve performance.

| Model | MIT Test | ReClor Validation |
|---|---|---|
| StableVicuna-13B | 0.48 | 0.42 |
| LLaMA-30B | 0.35 | 0.15 |
| LLaMA-30B Fine-Tune MIT | 0.47 | 0.46 |
| LLaMA-65B | 0.39 | 0.65 |
| GPT-4 | 0.90 | 0.87 |
| GPT-4 + Few-Shot (FS) | 0.93 | N/A |
| GPT-4 + FS + CoT | 0.95 | N/A |
| GPT-4 + FS + CoT + Self-critique | 0.97 | N/A |
| GPT-4 + FS + CoT + Self-critique + Experts | 1 | N/A |

## 4 Conclusion

We compiled a comprehensive dataset containing all questions from the MIT Mathematics and EECS undergraduate curriculum, including problem sets, midterms, and final exams. We frame prompt engineering as an optimization problem and provide effective sample heuristics. Our evaluation demonstrates that GPT-4, combined with a system expert, few-shot learning, chain-of-thought, self-critique, and collaborative decision-making techniques, achieves a perfect solve rate on a randomly selected test set of these questions excluding image-based questions. We release a LLaMA model fine-tuned on the MIT curriculum, which improves performance on a benchmark for assessing logical reasoning abilities. By employing few-shot learning, or in-context learning, we determine which questions and courses should be prerequisites for other courses based on the data, thereby

8

Figure 3: Graph of classes: graph nodes represent classes and directed edges represent the data-informed prerequisites based on the questions used for few shot-learning.



identifying the foundational content necessary for advanced topics. Instead of prohibiting using LLMs in the classroom, we advocate for their integration by designing meta-questions that incorporate LLM-generated answers, requiring students to evaluate the completion and correctness of these responses.

Limitations of our work is that inference and automatic grading using GPT-4 is relatively slow, taking around a minute for each question, and has a limited context window of 8k tokens. Recent models use context windows of 100k tokens [1] and demonstrate the feasibility of windows of 500k tokens.

Using a dataset of problem sets, midterms, and final exam questions and answers, we identify course connections and propose an optimal sequence of prerequisites. This approach assists educators and academic administrators to evaluate, design, and enhance curricula. In the future, such tools may prove invaluable for introducing new course content, identifying gaps in a curriculum, pinpointing courses with weak connections to others, and reinforcing key concepts. This method could benefit residential and asynchronous learners, helping students make informed decisions about course selection and guiding professors in determining which content to teach.

# References

[1] Anthropic, 2023.

[2] Zhangir Azerbayev, Bartosz Piotrowski, Hailey Schoelkopf, Edward W. Ayers, Dragomir Radev, and Jeremy Avigad. ProofNet: Autoformalizing and formally proving undergraduate-level mathematics, 2023.

[3] Mathieu Bastian, Sebastien Heymann, and Mathieu Jacomy. Gephi: An open source software for exploring and manipulating networks. 2009.

[4] Lin Bi, Yue Wang, Jian-ping Zhao, Hui Qi, and Ying Zhang. Social network information visualization based on fruchterman reingold layout algorithm. In *2018 IEEE 3rd International Conference on Big Data Analysis (ICBDA)*, pages 270–273. IEEE, 2018.

[5] Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners. *CoRR*, abs/2005.14165, 2020.

[6] Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrke, Eric Horvitz, Ece Kamar, Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, et al. Sparks of artificial general intelligence: Early experiments with gpt-4. *arXiv preprint arXiv:2303.12712*, 2023.

[7] Iddo Drori, Sarah Zhang, Reece Shuttleworth, Leonard Tang, Albert Lu, Elizabeth Ke, Kevin Liu, Linda Chen, Sunny Tran, Newman Cheng, et al. A neural network solves, explains, and generates university math problems by program synthesis and few-shot learning at human level. *Proceedings of the National Academy of Sciences*, 119(32):e2123433119, 2022.

[8] Iddo Drori, Sarah J. Zhang, Reece Shuttleworth, Sarah Zhang, Zad Chin, Pedro Lantigua, Saisamrit Surbehera, Gregory Hunter, Derek Austin, Leonard Tang, Yann Hicke, Sage Simhon, Sathwik Karnik, Darnell Granberry, and Madeleine Udell. From human days to machine seconds: Automatically answering and generating machine learning final exams. 2023.

[9] Zhibin Gou, Zhihong Shao, Yeyun Gong, Yelong Shen, Yujiu Yang, Nan Duan, and Weizhu Chen. Critic: Large language models can self-correct with tool-interactive critiquing. *arXiv preprint arXiv:2305.11738*, 2023.

[10] Mathieu Jacomy, Tommaso Venturini, Sebastien Heymann, and Mathieu Bastian. ForceAtlas2, a continuous graph layout algorithm for handy network visualization designed for the gephi software. *PloS one*, 9(6):e98679, 2014.

[11] Geunwoo Kim, Pierre Baldi, and Stephen McAleer. Language models can solve computer tasks. *arXiv preprint arXiv:2303.17491*, 2023.

[12] Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. Large language models are zero-shot reasoners. *ArXiv e-prints*, 2022. `https://arxiv.org/abs/2205.11916` (accessed: June 9, 2022).

[13] Aitor Lewkowycz, Anders Andreassen, David Dohan, Ethan Dyer, Henryk Michalewski, Vinay Ramasesh, Ambrose Slone, Cem Anil, Imanol Schlag, Theo Gutman-Solo, et al. Solving quantitative reasoning problems with language models. *arXiv preprint arXiv:2206.14858*, 2022.

[14] Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegreffe, Uri Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, et al. Self-refine: Iterative refinement with self-feedback. *arXiv preprint arXiv:2303.17651*, 2023.

[15] Mathpix. Mathpix snip. `https://mathpix.com/`, 2023.

[16] Maxwell Nye, Anders Andreassen, Guy Gur-Ari, Henryk Witold Michalewski, Jacob Austin, David Bieber, David Martin Dohan, Aitor Lewkowycz, Maarten Paul Bosma, David Luan, Charles Sutton, and Augustus Odena. Show your work: Scratchpads for intermediate computation with language models. *arXiv:2112.00114*, 2021.

[17] OpenAI. Gpt-4 technical report. *arXiv*, 2023.

[18] OpenAI. OpenAI Embeddings. `https://beta.openai.com/docs/guides/embeddings`, 2023.

[19] Debjit Paul, Mete Ismayilzada, Maxime Peyrard, Beatriz Borges, Antoine Bosselut, Robert West, and Boi Faltings. Refiner: Reasoning feedback on intermediate representations. *arXiv preprint arXiv:2304.01904*, 2023.

[20] Sami Sarsa, Paul Denny, Arto Hellas, and Juho Leinonen. Automatic generation of programming exercises and code explanations using large language models. In *Proceedings of the 2022 ACM Conference on International Computing Education Research-Volume 1*, pages 27–43, 2022.

[21] Noah Shinn, Beck Labash, and Ashwin Gopinath. Reflexion: an autonomous agent with dynamic memory and self-reflection. *arXiv preprint arXiv:2303.11366*, 2023.

[22] The Vicuna Team. Vicuna: An Open-Source Chatbot Impressing GPT-4 with 90%* ChatGPT Quality. 2023. `https://lmsys.org/blog/2023-03-30-vicuna/`.

[23] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timo-thée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. LLaMA: Open and Efficient Foundation Language Models. *arXiv:2302.13971*, 2023.

[24] Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language models. In *Proceedings of ICLR*, 2023.

[25] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Ed Chi, Quoc Le, and Denny Zhou. Chain of thought prompting elicits reasoning in large language models. *arXiv preprint arXiv:2201.11903*, 2022.

[26] Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Thomas L Griffiths, Yuan Cao, and Karthik Narasimhan. Tree of thoughts: Deliberate problem solving with large language models. *arXiv preprint arXiv:2305.10601*, 2023.

[27] Weihao Yu, Zihang Jiang, Yanfei Dong, and Jiashi Feng. Reclor: A reading comprehension dataset requiring logical reasoning. *arXiv preprint arXiv:2002.04326*, 2020.

# Appendix

## 4.1 Code and Models

Links to the source code and fine-tuned models are available here: https://github.com/idrori/MITQ

## 4.2 MIT Mathematics and EECS Courses and Questions by Major

We collect and curate a comprehensive dataset of 4550 questions and corresponding answers from 30 MIT Math and EECS courses required to graduate from the institute. This includes a broad range of core and elective courses, providing students with the foundational and specialized knowledge necessary to succeed. The dataset spans eight MIT Math and EECS undergraduate degree paths: (1) 6-1: Electrical Science and Engineering, (2) 6-2: Electrical Engineering and Computer Science, (3) 6-3: Computer Science and Engineering, (4) 6-4: Artificial Intelligence and Decision Making, (5) 18-1: General Mathematics, (6) 18-2: Applied Mathematics, (7) 18-3: Pure Mathematics, and (8) 18-C: Mathematics with Computer Science. A detailed breakdown of each major is summarized in Tables 6-13.

Table 6: MIT EECS undergraduate degree questions for Electrical Science and Engineering. The table includes the number of units, questions, and parts for each course.

| ID | Number | Name | Units | Questions | Parts |
|----|--------|------|-------|-----------|-------|
| 1 | 6.100A | Intro to CS Programming in Python | 6 | 34 | 47 |
| 2 | 6.1010 | Fundamentals of Programming | 12 | 22 | 31 |
| 3 | 6.1210 | Intro to Algorithms | 12 | 82 | 164 |
| 4 | 6.1800 | Computer Systems Engineering | 12 | 58 | 112 |
| 5 | 6.1910 | Computation Structures | 12 | 72 | 198 |
| 6 | 6.2000 | Electrical Circuits | 12 | 27 | 97 |
| 7 | 6.2300 | Electromagnetics Waves & Applications | 12 | 37 | 142 |
| 8 | 6.3000 | Signal Processing | 12 | 56 | 258 |
| 9 | 6.3010 | Signals, Systems & Inference | 12 | 57 | 224 |
| 10 | 6.3900 | Intro to Machine Learning | 12 | 114 | 619 |
| 11 | 6.8611 | Quantitative Methods for NLP | 15 | 20 | 31 |
| 12 | 18.01 | Calculus I | 12 | 203 | 495 |
| 13 | 18.02 | Calculus II | 12 | 81 | 154 |
| 14 | 18.03 | Differential Equations | 12 | 66 | 160 |
| 15 | 18.C06 | Linear Algebra & Optimization | 12 | 77 | 195 |
| 16 | 18.200 | Principles of Discrete Applied Math | 15 | 45 | 86 |
| 17 | 18.404 | Theory of Computation | 12 | 53 | 101 |
| 18 | 18.600 | Probability & Random Variables | 12 | 65 | 160 |
| Mean | | | 12 | 64.94 | 181.89 |
| Total | | | 216 | 1169 | 3114 |

Table 7: MIT EECS undergraduate degree questions for Electrical Engineering and Computer Science. The table includes the number of units, questions, and parts for each course.

| ID | Number | Name | Units | Questions | Parts |
|---|---|---|---|---|---|
| 1 | 6.100A | Intro to CS Programming in Python | 6 | 34 | 47 |
| 2 | 6.1010 | Fundamentals of Programming | 12 | 22 | 31 |
| 3 | 6.1020 | Elements of Software Construction | 15 | 26 | 52 |
| 4 | 6.1210 | Intro to Algorithms | 12 | 82 | 164 |
| 5 | 6.1220 | Design & Analysis of Algorithms | 12 | 44 | 158 |
| 6 | 6.1800 | Computer Systems Engineering | 12 | 58 | 112 |
| 7 | 6.1910 | Computation Structures | 12 | 72 | 198 |
| 8 | 6.2000 | Electrical Circuits | 12 | 27 | 97 |
| 9 | 6.3000 | Signal Processing | 12 | 56 | 258 |
| 10 | 6.3900 | Intro to Machine Learning | 12 | 114 | 619 |
| 11 | 6.8611 | Quantitative Methods for NLP | 15 | 20 | 31 |
| 12 | 18.01 | Calculus I | 12 | 203 | 495 |
| 13 | 18.02 | Calculus II | 12 | 81 | 154 |
| 14 | 18.03 | Differential Equations | 12 | 66 | 160 |
| 15 | 18.C06 | Linear Algebra & Optimization | 12 | 77 | 195 |
| 16 | 18.200 | Principles of Discrete Applied Math | 15 | 45 | 86 |
| 17 | 18.404 | Theory of Computation | 12 | 53 | 101 |
| 18 | 18.600 | Probability & Random Variables | 12 | 65 | 160 |
| Mean | | | 12.17 | 63.61 | 173.22 |
| Total | | | 219 | 1145 | 3118 |

Table 8: MIT EECS undergraduate degree questions for Computer Science and Engineering. The table includes the number of units, questions, and parts for each course.

| ID | Number | Name | Units | Questions | Parts |
|---|---|---|---|---|---|
| 1 | 6.100A | Intro to CS Programming in Python | 6 | 34 | 47 |
| 2 | 6.1010 | Fundamentals of Programming | 12 | 22 | 31 |
| 3 | 6.1020 | Elements of Software Construction | 15 | 26 | 52 |
| 4 | 6.1210 | Intro to Algorithms | 12 | 82 | 164 |
| 5 | 6.1220 | Design & Analysis of Algorithms | 12 | 44 | 158 |
| 6 | 6.1800 | Computer Systems Engineering | 12 | 58 | 112 |
| 7 | 6.1910 | Computation Structures | 12 | 72 | 198 |
| 8 | 6.3000 | Signal Processing | 12 | 56 | 258 |
| 9 | 6.3900 | Intro to Machine Learning | 12 | 114 | 619 |
| 10 | 6.4110 | Rep., Inference, & Reasoning in AI | 12 | 54 | 324 |
| 11 | 6.8611 | Quantitative Methods for NLP | 15 | 20 | 31 |
| 12 | 18.01 | Calculus I | 12 | 203 | 495 |
| 13 | 18.02 | Calculus II | 12 | 81 | 154 |
| 14 | 18.03 | Differential Equations | 12 | 66 | 160 |
| 15 | 18.C06 | Linear Algebra & Optimization | 12 | 77 | 195 |
| 16 | 18.200 | Principles of Discrete Applied Math | 15 | 45 | 86 |
| 17 | 18.404 | Theory of Computation | 12 | 53 | 101 |
| 18 | 18.600 | Probability & Random Variables | 12 | 65 | 160 |
| Mean | | | 12.17 | 65.11 | 185.83 |
| Total | | | 219 | 1172 | 3345 |

Table 9: MIT EECS undergraduate degree questions for Artificial Intelligence and Decision Making. The table includes the number of units, questions, and parts for each course.

| ID | Number | Name | Units | Questions | Parts |
|----|--------|------|-------|-----------|-------|
| 1 | 6.100A | Intro to CS Programming in Python | 6 | 34 | 47 |
| 2 | 6.1010 | Fundamentals of Programming | 12 | 22 | 31 |
| 3 | 6.1020 | Elements of Software Construction | 15 | 26 | 52 |
| 4 | 6.1210 | Intro to Algorithms | 12 | 82 | 164 |
| 5 | 6.1220 | Design & Analysis of Algorithms | 12 | 44 | 158 |
| 6 | 6.1800 | Computer Systems Engineering | 12 | 58 | 112 |
| 7 | 6.1910 | Computation Structures | 12 | 72 | 198 |
| 8 | 6.3000 | Signal Processing | 12 | 56 | 258 |
| 9 | 6.3900 | Intro to Machine Learning | 12 | 114 | 619 |
| 10 | 6.4110 | Rep., Inference, & Reasoning in AI | 12 | 54 | 324 |
| 11 | 6.4120 | Computational Cognitive Science | 12 | 10 | 67 |
| 12 | 6.8611 | Quantitative Methods for NLP | 15 | 20 | 31 |
| 13 | 18.01 | Calculus I | 12 | 203 | 495 |
| 14 | 18.02 | Calculus II | 12 | 81 | 154 |
| 15 | 18.03 | Differential Equations | 12 | 66 | 160 |
| 16 | 18.C06 | Linear Algebra & Optimization | 12 | 77 | 195 |
| 17 | 18.200 | Principles of Discrete Applied Math | 15 | 45 | 86 |
| 18 | 18.404 | Theory of Computation | 12 | 53 | 101 |
| 19 | 18.600 | Probability & Random Variables | 12 | 65 | 160 |
| Mean | | | 12.16 | 62.21 | 179.58 |
| Total | | | 231 | 1182 | 3412 |

Table 10: MIT Math undergraduate degree questions for General Mathematics. The table includes the number of units, questions, and parts for each course.

| ID | Number | Name | Units | Questions | Parts |
|----|--------|------|-------|-----------|-------|
| 1 | 6.100A | Intro to CS Programming in Python | 6 | 34 | 47 |
| 2 | 6.1010 | Fundamentals of Programming | 12 | 22 | 31 |
| 3 | 6.1210 | Intro to Algorithms | 12 | 82 | 164 |
| 4 | 6.1220 | Design & Analysis of Algorithms | 12 | 44 | 158 |
| 5 | 18.01 | Calculus I | 12 | 203 | 495 |
| 6 | 18.02 | Calculus II | 12 | 81 | 154 |
| 7 | 18.03 | Differential Equations | 12 | 66 | 160 |
| 8 | 18.C06 | Linear Algebra & Optimization | 12 | 77 | 195 |
| 9 | 18.100B | Real Analysis | 12 | 60 | 66 |
| 10 | 18.200 | Principles of Discrete Applied Math | 15 | 45 | 86 |
| 11 | 18.404 | Theory of Computation | 12 | 53 | 101 |
| 12 | 18.600 | Probability & Random Variables | 12 | 65 | 160 |
| 13 | 18.701 | Algebra I | 12 | 58 | 87 |
| 14 | 18.704 | Seminar in Algebra | 12 | 16 | 25 |
| Mean | | | 11.79 | 64.71 | 137.79 |
| Total | | | 165 | 906 | 1929 |

Table 11: MIT Math undergraduate degree questions for Applied Mathematics. The table includes the number of units, questions, and parts for each course.

| ID | Number | Name | Units | Questions | Parts |
|----|--------|------|-------|-----------|-------|
| 1 | 6.100A | Intro to CS Programming in Python | 6 | 34 | 47 |
| 2 | 6.1010 | Fundamentals of Programming | 12 | 22 | 31 |
| 3 | 6.1210 | Intro to Algorithms | 12 | 82 | 164 |
| 4 | 6.1220 | Design & Analysis of Algorithms | 12 | 44 | 158 |
| 5 | 18.01 | Calculus I | 12 | 203 | 495 |
| 6 | 18.02 | Calculus II | 12 | 81 | 154 |
| 7 | 18.03 | Differential Equations | 12 | 66 | 160 |
| 8 | 18.C06 | Linear Algebra & Optimization | 12 | 77 | 195 |
| 9 | 18.100B | Real Analysis | 12 | 60 | 66 |
| 10 | 18.102 | Intro to Functional Analysis | 12 | 68 | 104 |
| 11 | 18.200 | Principles of Discrete Applied Math | 15 | 45 | 86 |
| 12 | 18.300 | Principles of Continuum Applied Math | 12 | 43 | 90 |
| 13 | 18.303 | Linear Partial Differential Equations | 12 | 22 | 65 |
| 14 | 18.404 | Theory of Computation | 12 | 53 | 101 |
| 15 | 18.600 | Probability & Random Variables | 12 | 65 | 160 |
| 16 | 18.701 | Algebra I | 12 | 58 | 87 |
| 17 | 18.704 | Seminar in Algebra | 12 | 16 | 25 |
| Mean | | | 11.82 | 61.12 | 128.71 |
| Total | | | 201 | 1039 | 2188 |

Table 12: MIT Math undergraduate degree questions for Pure Mathematics. The table includes the number of units, questions, and parts for each course.

| ID | Number | Name | Units | Questions | Parts |
|----|--------|------|-------|-----------|-------|
| 1 | 6.100A | Intro to CS Programming in Python | 6 | 34 | 47 |
| 2 | 6.1010 | Fundamentals of Programming | 12 | 22 | 31 |
| 3 | 6.1210 | Intro to Algorithms | 12 | 82 | 164 |
| 4 | 6.1220 | Design & Analysis of Algorithms | 12 | 44 | 158 |
| 5 | 18.01 | Calculus I | 12 | 203 | 495 |
| 6 | 18.02 | Calculus II | 12 | 81 | 154 |
| 7 | 18.03 | Differential Equations | 12 | 66 | 160 |
| 8 | 18.C06 | Linear Algebra & Optimization | 12 | 77 | 195 |
| 9 | 18.100B | Real Analysis | 12 | 60 | 66 |
| 10 | 18.102 | Intro to Functional Analysis | 12 | 68 | 104 |
| 11 | 18.200 | Principles of Discrete Applied Math | 15 | 45 | 86 |
| 12 | 18.600 | Probability & Random Variables | 12 | 65 | 160 |
| 13 | 18.701 | Algebra I | 12 | 58 | 87 |
| 14 | 18.702 | Algebra II | 12 | 52 | 94 |
| 15 | 18.704 | Seminar in Algebra | 12 | 16 | 25 |
| 16 | 18.901 | Intro to Topology | 12 | 58 | 144 |
| Mean | | | 11.81 | 64.44 | 135.63 |
| Total | | | 189 | 1031 | 2170 |

Table 13: MIT Math undergraduate degree questions for Mathematics with Computer Science. The table includes the number of units, questions, and parts for each course.

| ID | Number | Name | Units | Questions | Parts |
|----|--------|------|-------|-----------|-------|
| 1 | 6.100A | Intro to CS Programming in Python | 6 | 34 | 47 |
| 2 | 6.1010 | Fundamentals of Programming | 12 | 22 | 31 |
| 3 | 6.1020 | Elements of Software Construction | 15 | 26 | 52 |
| 4 | 6.1210 | Intro to Algorithms | 12 | 82 | 164 |
| 5 | 6.1220 | Design & Analysis of Algorithms | 12 | 44 | 158 |
| 6 | 6.3900 | Intro to Machine Learning | 12 | 114 | 619 |
| 7 | 18.01 | Calculus I | 12 | 203 | 495 |
| 8 | 18.02 | Calculus II | 12 | 81 | 154 |
| 9 | 18.03 | Differential Equations | 12 | 66 | 160 |
| 10 | 18.C06 | Linear Algebra & Optimization | 12 | 77 | 195 |
| 11 | 18.100B | Real Analysis | 12 | 60 | 66 |
| 12 | 18.200 | Principles of Discrete Applied Math | 15 | 45 | 86 |
| 13 | 18.404 | Theory of Computation | 12 | 53 | 101 |
| 14 | 18.600 | Probability & Random Variables | 12 | 65 | 160 |
| 15 | 18.701 | Algebra I | 12 | 58 | 87 |
| 16 | 18.704 | Seminar in Algebra | 12 | 16 | 25 |
| Mean | | | 12 | 65.38 | 162.5 |
| Total | | | 192 | 1046 | 2600 |

## 4.3  GPT-4 Solve Rates on Entire Dataset

Tables 14, 15, 16 shows the breakdown of the question type, solution type, and zero-shot GPT-4 solve rates for each course. Text questions are easier than image question. The difficulty of questions by solution type from easy to hard is: programming, open, image, expression, numerical, and multiple-choice.

Table 14: MIT Mathematics and EECS undergraduate degree questions parts by question type. The table includes the number of parts and GPT-4 solve rate for each question type.

| ID | Question Type | Parts | Solve Rate |
|----|---------------|-------|-----------|
| 1 | Text | 2946 | 0.90 |
| 2 | Image | 879 | 0.79 |
| Mean | | 1912.5 | 0.84 |
| Total | | 3825 | 0.88 |

Table 15: MIT Mathematics and EECS undergraduate degree questions parts by solution type. The table includes the number of parts and GPT-4 solve rate for each solution type.

| ID | Solution Type | Parts | Solve Rate |
|----|---------------|-------|-----------|
| 1 | Programming | 221 | 0.94 |
| 2 | Multiple Choice | 690 | 0.81 |
| 3 | Numerical | 545 | 0.86 |
| 4 | Expression | 822 | 0.87 |
| 5 | Open | 1400 | 0.91 |
| 6 | Image | 147 | 0.89 |
| Mean | | 637.5 | 0.88 |
| Total | | 3825 | 0.88 |

Table 16: MIT Mathematics and EECS undergraduate degree questions parts by course. The table includes the number of parts and GPT-4 solve rate for each Course.

| ID | Course | Parts | Solve Rate |
|----|--------|-------|-----------|
| 1 | 6.100A | 47 | 0.98 |
| 2 | 18.100B | 63 | 0.95 |
| 3 | 18.102 | 103 | 0.95 |
| 4 | 18.C06 | 194 | 0.93 |
| 5 | 6.1210 | 164 | 0.87 |
| 6 | 6.1220 | 158 | 0.87 |
| 7 | 6.3900 | 619 | 0.88 |
| 8 | 18.303 | 9 | 0.89 |
| 9 | 18.200 | 86 | 0.93 |
| 10 | 6.1800 | 105 | 0.84 |
| 11 | 18.702 | 64 | 0.92 |
| 12 | 18.701 | 60 | 0.92 |
| 13 | 18.01 | 495 | 0.91 |
| 14 | 6.4110 | 324 | 0.84 |
| 15 | 6.1010 | 19 | 0.83 |
| 16 | 18.704 | 24 | 0.96 |
| 17 | 6.1020 | 52 | 0.88 |
| 18 | 18.02 | 151 | 0.91 |
| 19 | 18.600 | 160 | 0.90 |
| 20 | 6.8611 | 31 | 0.89 |
| 21 | 18.404 | 101 | 0.90 |
| 22 | 6.1910 | 198 | 0.81 |
| 23 | 18.03 | 160 | 0.91 |
| 24 | 6.2000 | 97 | 0.79 |
| 25 | 18.300 | 82 | 0.84 |
| 26 | 6.3000 | 258 | 0.74 |
| 27 | 6.3010 | 1 | 1.00 |
| Mean | | 141.7 | 0.89 |
| Total | | 3825 | 0.88 |

## 4.4 MIT Mathematics and EECS Course Descriptions

Table 17: MIT EECS undergraduate degree courses. The table includes the description for each course.

| ID | Number | Name | Description |
|---|---|---|---|
| 1 | 6.100A | Intro to CS Programming in Python | Introduction to computer science and programming for students with little or no programming experience. Students develop skills to program and use computational techniques to solve problems. Topics include the notion of computation, Python, simple algorithms and data structures, testing and debugging, and algorithmic complexity. |
| 2 | 6.1010 | Fundamentals of Programming | Introduces fundamental concepts of programming. Designed to develop skills in applying basic methods from programming languages to abstract problems. Topics include programming and Python basics, computational concepts, software engineering, algorithmic techniques, data types, and recursion. Lab component consists of software design, construction, and implementation of design. |
| 3 | 6.1020 | Elements of Software Construction | Introduces fundamental principles and techniques of software development: how to write software that is safe from bugs, easy to understand, and ready for change. Topics include specifications and invariants; testing, test-case generation, and coverage; abstract data types and representation independence; design patterns for object-oriented programming; concurrent programming, including message passing and shared memory concurrency, and defending against races and deadlock; and functional programming with immutable data and higher-order functions. Includes weekly programming exercises and larger group programming projects. |
| 4 | 6.1210 | Intro to Algorithms | Introduction to mathematical modeling of computational problems, as well as common algorithms, algorithmic paradigms, and data structures used to solve these problems. Emphasizes the relationship between algorithms and programming, and introduces basic performance measures and analysis techniques for these problems. |
| 5 | 6.1220 | Design & Analysis of Algorithms | Techniques for the design and analysis of efficient algorithms, emphasizing methods useful in practice. Topics include sorting; search trees, heaps, and hashing; divide-and-conquer; dynamic programming; greedy algorithms; amortized analysis; graph algorithms; and shortest paths. Advanced topics may include network flow; computational geometry; number-theoretic algorithms; polynomial and matrix calculations; caching; and parallel computing. |
| 6 | 6.1800 | Computer Systems Engineering | Topics on the engineering of computer software and hardware systems: techniques for controlling complexity; strong modularity using client-server design, operating systems; performance, networks; naming; security and privacy; fault-tolerant systems, atomicity and coordination of concurrent activities, and recovery; impact of computer systems on society. Case studies of working systems and readings from the current literature provide comparisons and contrasts. Includes a single, semester-long design project. Students engage in extensive written communication exercises. |
| 7 | 6.1910 | Computation Structures | Provides an introduction to the design of digital systems and computer architecture. Emphasizes expressing all hardware designs in a high-level hardware language and synthesizing the designs. Topics include combinational and sequential circuits, instruction set abstraction for programmable hardware, single-cycle and pipelined processor implementations, multi-level memory hierarchies, virtual memory, exceptions and I/O, and parallel systems. |
| 8 | 6.2000 | Electrical Circuits | Fundamentals of linear systems, and abstraction modeling of multi-physics lumped and distributed systems using lumped electrical circuits. Linear networks involving independent and dependent sources, resistors, capacitors, and inductors. Extensions to include operational amplifiers and transducers. Dynamics of first- and second-order networks; analysis and design in the time and frequency domains; signal and energy processing applications. Design exercises. Weekly laboratory with microcontroller and transducers. |
| 9 | 6.2300 | Electromagnetic Waves & Applications | Analysis and design of modern applications that employ electromagnetic phenomena for signals and power transmission in RF, microwaves, optical and wireless communication systems. Fundamentals include dynamic solutions for Maxwell's equations; electromagnetic power and energy, waves in media, metallic and dielectric waveguides, radiation, and diffraction; resonance; filters; and acoustic analogs. Lab activities range from building to testing of devices and systems (e.g., antenna arrays, radars, dielectric waveguides). Students work in teams on self-proposed maker-style design projects with a focus on fostering creativity, teamwork, and debugging skills. |
| 10 | 6.3000 | Signal Processing | Fundamentals of signal processing, focusing on the use of Fourier methods to analyze and process signals such as sounds and images. Topics include Fourier series, Fourier transforms, the Discrete Fourier Transform, sampling, convolution, deconvolution, filtering, noise reduction, and compression. Applications draw broadly from areas of contemporary interest with emphasis on both analysis and design. |
| 11 | 6.3010 | Signals, Systems & Inference | Covers signals, systems and inference in communication, control and signal processing. Topics include input-output and state-space models of linear systems driven by deterministic and random signals; time- and transform-domain representations in discrete and continuous time; and group delay. State feedback and observers. Probabilistic models; stochastic processes, correlation functions, power spectra, spectral factorization. Least-mean square error estimation; Wiener filtering. Hypothesis testing; detection; matched filters. |
| 12 | 6.3900 | Intro to Machine Learning | Introduces principles, algorithms, and applications of machine learning from the point of view of modeling and prediction; formulation of learning problems; representation, over-fitting, generalization; clustering, classification, probabilistic modeling; and methods such as support vector machines, hidden Markov models, and neural networks. |
| 13 | 6.4110 | Rep., Inference, & Reasoning in AI | An introduction to representations and algorithms for artificial intelligence. Topics covered include: constraint satisfaction in discrete and continuous problems, logical representation and inference, Monte Carlo tree search, probabilistic graphical models and inference, planning in discrete and continuous deterministic and probabilistic models including MDPs and POMDPs. |
| 14 | 6.4120 | Computational Cognitive Science | Introduction to computational theories of human cognition. Focus on principles of inductive learning and inference, and the representation of knowledge. Computational frameworks covered include Bayesian and hierarchical Bayesian models; probabilistic graphical models; nonparametric statistical models and the Bayesian Occam's razor; sampling algorithms for approximate learning and inference; and probabilistic models defined over structured representations such as first-order logic, grammars, or relational schemas. Applications to understanding core aspects of cognition, such as concept learning and categorization, causal reasoning, theory formation, language acquisition, and social inference. |
| 15 | 6.8611 | Quantitative Methods for NLP | Introduces the study of human language from a computational perspective, including syntactic, semantic and discourse processing models. Emphasizes machine learning methods and algorithms. Uses these methods and models in applications such as syntactic parsing, information extraction, statistical machine translation, dialogue systems. Instruction and practice in oral and written communication provided. |

Table 18: MIT Math undergraduate degree courses. The table includes the description for each course.

| ID | Number | Name | Description |
|----|--------|------|-------------|
| 16 | 18.01 | Calculus I | Differentiation and integration of functions of one variable, with applications. Informal treatment of limits and continuity. Differentiation: definition, rules, application to graphing, rates, approximations, and extremum problems. Indefinite integration; separable first-order differential equations. Definite integral; fundamental theorem of calculus. Applications of integration to geometry and science. Elementary functions. Techniques of integration. Polar coordinates. L'Hopital's rule. Improper integrals. Infinite series: geometric, p-harmonic, simple comparison tests, power series for some elementary functions. |
| 17 | 18.02 | Calculus II | Calculus of several variables. Vector algebra in 3-space, determinants, matrices. Vector-valued functions of one variable, space motion. Scalar functions of several variables: partial differentiation, gradient, optimization techniques. Double integrals and line integrals in the plane; exact differentials and conservative fields; Green's theorem and applications, triple integrals, line and surface integrals in space, Divergence theorem, Stokes' theorem; applications. |
| 18 | 18.03 | Differential Equations | Study of differential equations, including modeling physical systems. Solution of first-order ODEs by analytical, graphical, and numerical methods. Linear ODEs with constant coefficients. Complex numbers and exponentials. Inhomogeneous equations: polynomial, sinusoidal, and exponential inputs. Oscillations, damping, resonance. Fourier series. Matrices, eigenvalues, eigenvectors, diagonalization. First order linear systems: normal modes, matrix exponentials, variation of parameters. Heat equation, wave equation. Nonlinear autonomous systems: critical point analysis, phase plane diagrams. |
| 19 | 18.C06 | Linear Algebra & Optimization | Introductory course in linear algebra and optimization, assuming no prior exposure to linear algebra and starting from the basics, including vectors, matrices, eigenvalues, singular values, and least squares. Covers the basics in optimization including convex optimization, linear/quadratic programming, gradient descent, and regularization, building on insights from linear algebra. Explores a variety of applications in science and engineering, where the tools developed give powerful ways to understand complex systems and also extract structure from data. |
| 20 | 18.100B | Real Analysis | Covers fundamentals of mathematical analysis: convergence of sequences and series, continuity, differentiability, Riemann integral, sequences and series of functions, uniformity, interchange of limit operations. Shows the utility of abstract concepts and teaches understanding and construction of proofs. Places more emphasis on point-set topology and n-space. |
| 21 | 18.102 | Intro to Functional Analysis | Normed spaces, completeness, functionals, Hahn-Banach theorem, duality, operators. Lebesgue measure, measurable functions, integrability, completeness of L-p spaces. Hilbert space. Compact, Hilbert-Schmidt and trace class operators. Spectral theorem. |
| 22 | 18.200 | Principles of Discrete Applied Math | Study of illustrative topics in discrete applied mathematics, including probability theory, information theory, coding theory, secret codes, generating functions, and linear programming. Instruction and practice in written communication provided. |
| 23 | 18.300 | Principles of Continuum Applied Math | Covers fundamental concepts in continuous applied mathematics. Applications from traffic flow, fluids, elasticity, granular flows, etc. Also covers continuum limit; conservation laws, quasi-equilibrium; kinematic waves; characteristics, simple waves, shocks; diffusion (linear and nonlinear); numerical solution of wave equations; finite differences, consistency, stability; discrete and fast Fourier transforms; spectral methods; transforms and series (Fourier, Laplace). Additional topics may include sonic booms, Mach cone, caustics, lattices, dispersion and group velocity. Uses MATLAB computing environment. |
| 24 | 18.303 | Linear Partial Differential Equations | Provides students with the basic analytical and computational tools of linear partial differential equations (PDEs) for practical applications in science and engineering, including heat/diffusion, wave, and Poisson equations. Analytics emphasize the viewpoint of linear algebra and the analogy with finite matrix problems. Studies operator adjoints and eigenproblems, series solutions, Green's functions, and separation of variables. Numerics focus on finite-difference and finite-element techniques to reduce PDEs to matrix problems, including stability and convergence analysis and implicit/explicit timestepping. |
| 25 | 18.404 | Theory of Computation | A more extensive and theoretical treatment of the material in 6.1400J/18.400J, emphasizing computability and computational complexity theory. Regular and context-free languages. Decidable and undecidable problems, reducibility, recursive function theory. Time and space measures on computation, completeness, hierarchy theorems, inherently complex problems, oracles, probabilistic computation, and interactive proof systems. |
| 26 | 18.600 | Probability & Random Variables | Probability spaces, random variables, distribution functions. Binomial, geometric, hypergeometric, Poisson distributions. Uniform, exponential, normal, gamma and beta distributions. Conditional probability, Bayes theorem, joint distributions. Chebyshev inequality, law of large numbers, and central limit theorem. |
| 27 | 18.701 | Algebra I | Focuses on group theory, geometry, and linear algebra. |
| 28 | 18.702 | Algebra II | Continuation of 18.701. Focuses on group representations, rings, ideals, fields, polynomial rings, modules, factorization, integers in quadratic number fields, field extensions, and Galois theory. |
| 29 | 18.704 | Seminar in Algebra | Topics vary from year to year. Students present and discuss the subject matter. Instruction and practice in written and oral communication provided. Some experience with proofs required. |
| 30 | 18.901 | Intro to Topology | Introduces topology, covering topics fundamental to modern analysis and geometry. Topological spaces and continuous functions, connectedness, compactness, separation axioms, covering spaces, and the fundamental group. |

## 4.5 Meta-Questions

Figure 4: Example of meta-question about correctness.

**Instructions**: This meta-question consists of a question and its answer by GPT4.
Identify whether the answer is correct.
If the answer is correct then explain why.
If the answer is incorrect then write the correct answer and explanation.

- - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -

**Question**: Consider an MDP with a set of possible states $\mathcal{S} = (0, 1, 2, 3)$ and a set of possible actions $\mathcal{A}$ = (b, c). The reward function for an action **a** in a given state **s**:

$$R(s, a) = \begin{cases} 1 & \text{if } s = 1 \\ 2 & \text{if } s = 3 \\ 0 & \text{otherwise} \end{cases}$$

You get the reward associated with a state on the step when you exit that state. Notice that the reward function is the same for each of the two possible actions. The state transition probability matrix for action c are given below. Rows of the transition matrix correspond to the input states, and columns correspond to the output states. $T[i, x, j]$ is the conditional probability $P(s_{t+1} = j | a_t = x, s_t = i)$.

$$T(s_t, \text{c}, s_{t+1}) = \begin{bmatrix} 0.0 & 0.1 & 0.9 & 0.0 \\ 0.9 & 0.1 & 0.0 & 0.0 \\ 0.0 & 0.0 & 0.1 & 0.9 \\ 0.9 & 0.0 & 0.0 & 0.1 \end{bmatrix}$$

We would like to find the value function associated with the policy that always chooses action **a** = c in a finite horizon case. What are the horizon 2 values of the states under this policy?

**GPT4 Answer**: To find the horizon 2 values of the states under the policy that always chooses action a = c, we apply the Bellman equation for finite horizon MDPs again. This time, we will use the horizon 1 values we previously calculated:

$$V^2(s) = R(s, c) + \sum_{s'} T(s, c, s')V^1(s')$$

For each state, we can calculate the value as follows:

$$V^2(0) = R(0, c) + \sum_{s'=0}^{3} T(0, c, s')V^1(s') = 0 + (0.0 \cdot 0 + 0.1 \cdot 1 + 0.9 \cdot 0 + 0.0 \cdot 2) = 0.1$$

$$V^2(1) = R(1, c) + \sum_{s'=0}^{3} T(1, c, s')V^1(s') = 1 + (0.9 \cdot 0 + 0.1 \cdot 1 + 0.0 \cdot 0 + 0.0 \cdot 2) = 1.1$$

$$V^2(2) = R(2, c) + \sum_{s'=0}^{3} T(2, c, s')V^1(s') = 0 + (0.0 \cdot 0 + 0.0 \cdot 1 + 0.1 \cdot 0 + 0.9 \cdot 2) = 1.8$$

$$V^2(3) = R(3, c) + \sum_{s'=0}^{3} T(3, c, s')V^1(s') = 2 + (0.9 \cdot 0 + 0.0 \cdot 1 + 0.0 \cdot 0 + 0.1 \cdot 2) = 2.2$$