

---

# A Dataset and Benchmark for Automatically Answering and Generating Machine Learning Final Exams

---

<b>Sarah Zhang</b> EECS MIT sazhang@mit.edu	<b>Reece Shuttleworth</b> EECS MIT rshuttle@mit.edu	<b>Derek Austin</b> CS Columbia University da2986@columbia.edu	<b>Yann Hicke</b> CS Cornell University ylh8@cornell.edu
--	--	---	---

<b>Leonard Tang</b> Mathematics Harvard University leonardtang@college.harvard.edu	<b>Sathwik Karnik</b> EECS MIT skarnik@mit.edu	<b>Darnell Granberry</b> EECS MIT darnellg@mit.edu
---	---	---

**Iddo Drori**  
EECS and CS  
MIT and Columbia University  
idrori@mit.edu, idrori@cs.columbia.edu

## Abstract

1 Can a machine learn machine learning? We propose to answer this question  
2 using the same criteria we use to answer a similar question: can a human learn  
3 machine learning? We automatically answer MIT final exams in Introduction to  
4 Machine Learning at a human level. The course is a large undergraduate class  
5 with around five hundred students each semester. Recently, program synthesis and  
6 few-shot learning solved university-level problem set questions in mathematics  
7 and STEM courses at a human level. In this work, we solve questions from final  
8 exams that differ from problem sets in several ways: the questions are longer,  
9 have multiple parts, are more complicated, and span a broader set of topics. We  
10 provide a new dataset and benchmark of questions from eight MIT Introduction  
11 to Machine Learning final exams between Fall 2017 and Spring 2022 and provide  
12 code for automatically answering these questions and generating new questions.  
13 We perform ablation studies comparing zero-shot learning with few-shot learning,  
14 chain-of-thought prompting, GPT-3 pre-trained on text and Codex fine-tuned on  
15 code on a range of machine learning topics and find that few-shot learning methods  
16 perform best. We make our data and code publicly available for the machine  
17 learning community.

## 18 1 Introduction

19 Can a machine learn machine learning? This work presents a new dataset of machine learning  
20 final exams with around 500 questions and a benchmark of baselines using transformers and their  
21 respective grade performance, demonstrating that the best baseline performs at a human level.



Figure 1: 500 MIT students taking the final exam in Introduction to Machine Learning in Fall 2021 (no student faces or identifying information appears in the image to maintain anonymity). The final exam is three hours long and is taken in an indoor track-and-field stadium that accommodates many students.

In university-level STEM courses, students complete assignments (including problem sets and labs) and exams throughout the course. Recent work has opened the door for a machine to solve course problem sets [4] using language models and few-shot learning. However, final exams remain challenging, and this work is the first to present a structured dataset of machine learning finals and a benchmark of baseline methods for answering them. Final exams are different from problem sets because they serve as a benchmark of cumulative understanding of material learned over a semester and evaluate the students' depth and breadth of expertise. Further, questions on final exams are longer, have multiple parts, span a broader set of topics, and are more complicated and nuanced.

All the above holds for MIT's Introduction to Machine Learning (Intro to ML), a rigorous undergraduate course with around 500 students each semester, as shown in Figure 1, making it one of MIT's largest undergraduate courses offered.

Intro to ML is a core class in MIT's computer science program. The course description<sup>1</sup> is:

“Introduction 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.”

The prerequisites for the course are Python Programming and Multivariate Calculus, with Introduction to Algorithms and Linear Algebra recommended. The class typically consists of weekly exercises, labs, quizzes, homework, a midterm, and a final exam. There were no final exams in Fall 2020 and Spring 2020 due to COVID-19.

Intro to ML final exams differ from problem sets (psets) in several ways, and the experience of solving each varies. First, Intro to ML finals are long, containing around nine questions with seven parts. Final exam questions are also multifaceted and multi-stepped: different parts of a single question require applying different concepts and problem-solving skills, and parts may build upon each other. While weekly problem sets focus on a single topic, finals span topics from the entire semester. Further, Intro to ML final questions are often story-based problems that may require mathematical modeling. Due to the time constraint of these exams, finals are also designed to test core understanding and application of course material over rote calculations. Thus, asking a machine to answer questions from finals allows for testing whether the model is able to learn a breadth and depth of topics beyond problem sets.

<sup>1</sup>MIT Registrars Office Catalog for Course 6: Electrical Engineering and Computer Science, 6.3900 (6.036) Introduction to Machine Learning <http://student.mit.edu/catalog/m6c.html>

53 In this work, we present a new dataset curated from the eight most recent final exams of MIT's  
 54 Intro to ML course, totaling nearly 500 questions spanning twelve topics in machine learning:  
 55 regression, classifiers, logistic regression, features, loss functions, neural networks, convolutional  
 56 neural networks (CNNs), Markov decision processes (MDPs), recurrent neural networks (RNNs),  
 57 reinforcement learning, clustering, and decision trees. Our dataset covers the finals given for the  
 58 semesters Fall 2017, Spring 2018, Fall 2018, Spring 2019, Fall 2019, Spring 2021, Fall 2021, and  
 59 Spring 2022. Due to the COVID-19 pandemic, there were no finals in the course during 2020.

Table 1: Spring 2021, Question 10.a: question part, solution, Codex zero-shot answer, program output, GPT-3 Zero-Shot answer, and GPT-3 Zero-Shot with CoT prompt answer.

Topic	Classifiers
Question Part	<p>Given a set of data <math>\mathcal{D}_{train} = \{(x_i, y_i)\}</math>, a weighted nearest neighbor regressor has the form</p> $h(x; \theta) = \frac{\sum_{(x_i, y_i) \in \mathcal{D}_{train}} f(x, x_i, \theta) y_i}{\sum_{(x_i, y_i) \in \mathcal{D}_{train}} f(x, x_i, \theta)}.$ <p>A typical choice for <math>f</math> is</p> $f(x, x', \theta) = e^{-\theta \ x - x'\ ^2}$ <p>where <math>\theta</math> is a scalar and <math>\ x - x'\ ^2 = \sum_{j=1}^d (x_j - x'_j)^2</math>. Assume our training data <math>\mathcal{D}_{train} = ((1, 1), (2, 2), (3, 6))</math>. What is <math>h(10, 0)</math>? That is, letting <math>\theta = 0</math>, what is our prediction for <math>x = 10</math>?</p>
Solution	3
Codex Zero-Shot Answer	<pre>import numpy as np  def f(x, x_i, theta):     return np.exp(-theta * np.linalg.norm(x - x_i)**2)  def h(x, theta):     D = np.array([(1, 1), (2, 2), (3, 6)])     numerator = 0     denominator = 0     for x_i, y_i in D:         numerator += f(x, x_i, theta) * y_i         denominator += f(x, x_i, theta)     return numerator / denominator  print(h(10, 0))</pre>
Program Output	3.0
GPT-3 Zero-Shot Answer	$h(10; 0) = \frac{f(10, 1, 0)y_1 + f(10, 2, 0)y_2 + f(10, 3, 0)y_3}{f(10, 1, 0) + f(10, 2, 0) + f(10, 3, 0)} = \frac{1 \cdot 1 + 1 \cdot 2 + 1 \cdot 6}{1 + 1 + 1} = \frac{9}{3} = 3$
GPT-3 ZS with CoT Answer	$f(10, (1, 1), 0) = e^{-0 \ 10 - (1, 1)\ ^2} = e^{-0 \ 9\ ^2} = e^{-0 \cdot 81} = 1$ $f(10, (2, 2), 0) = e^{-0 \ 10 - (2, 2)\ ^2} = e^{-0 \ 8\ ^2} = e^{-0 \cdot 64} = 1$ $f(10, (3, 6), 0) = e^{-0 \ 10 - (3, 6)\ ^2} = e^{-0 \ 7\ ^2} = e^{-0 \cdot 49} = 1$ $h(10, 0) = \frac{\sum_{(x_i, y_i) \in \mathcal{D}_{train}} f(x, x_i, \theta) y_i}{\sum_{(x_i, y_i) \in \mathcal{D}_{train}} f(x, x_i, \theta)} = \frac{1 \cdot 1 + 1 \cdot 2 + 1 \cdot 6}{1 + 1 + 1} = \frac{9}{3} = 3$

60 We prevent overfitting by including finals from 2021 and 2022 that are unavailable online. The Spring  
 61 2022 final exam was given after GPT-3 and Codex were last updated, verifying that the models are  
 62 not trained on this data. The final exam questions have many parts, each posing a new problem,  
 63 and each question in the dataset corresponds to one part. The questions in the finals are varied in

Semester	Mean Human Grade	Mean Human Non-Image Grade	Mean Machine Non-Image Grade
Spring 2021	75.84	80.77	62.09
Fall 2021	74.38	60.88	58.94
Spring 2022	69.07	70.82	68.86
Mean	73.10	70.82	63.29

Table 2: Mean human and machine grades on Introduction to Machine Learning final exams by semester. Spring 2021, Fall 2021, and Spring 2022 final exams were unavailable online when GPT-3 and Codex were trained, ensuring that the model is not overfitting. Non-image grades consider question parts that do not contain images that are required for solving the question.

64 topics and solution types. Solutions are multi-modal, with primarily open-ended questions and some  
65 true/false and multiple-choice questions on theory, math, and code implementations.

66 We make the dataset publicly available and welcome others to use it to aid in developing and assessing  
67 new language models and methods. Due to the diversity of Intro to ML final questions, our dataset  
68 uniquely assesses advanced problem-solving and reasoning skills in machine learning, math, and  
69 natural language processing. This dataset opens the door to achieving breakthroughs in machine  
70 learning performance in machine learning final exams.

71 In addition to the dataset, we present a benchmark using several baseline methods. We apply zero-shot  
72 and few-shot learning to GPT-3 and Codex, adding chain-of-thought prompting for GPT-3. We find  
73 that few-shot learning methods perform best. As shown in Table 2 the best performing methods  
74 pass the final exams, and their grade is comparable with human grades of MIT students on the same  
75 machine learning finals evaluated by the same human graders.

76 Following our previous work [4] that showed that questions generated by these models are indistin-  
77 guishable from human-written questions, we also generate new final exam questions.

## 78 1.1 Related Work

79 There is a conception that humans are generalists, whereas machines are specialists. However, large  
80 language models based on transformers such as GPT-3 [1], Gopher [7], and Palm [3], also called  
81 foundation models, are generalist learners. Specifically, in our setting, while humans care about the  
82 number of topics in an exam and therefore find finals more difficult than problem sets, foundation  
83 models effortlessly scale to many topics without re-training. Language models may be pre-trained on  
84 text and fine-tuned on specific datasets such as code [2] which allows generating programs from text.

85 There are several ways to improve the mathematical reasoning ability of language models: (1) using  
86 chain-of-thought (CoT) prompting [5, 11], (2) using the top-k ranking solutions [6] and merging  
87 them by voting [10] or least-to-most prompting [12], and (3) using program synthesis and few-shot  
88 learning to generate code that answers questions [4].

89 Solving university-level course questions is a relatively new endeavor. The first work to tackle  
90 university-level machine learning course problem set questions [9] used a transformer and GNN  
91 architecture and heavily relied on data augmentation. This resulted in overfitting and did not scale  
92 up to other types of questions or courses. Probability and statistics course problem-set questions  
93 have been answered [8] by probabilistic program synthesis with human performance. Problem-set  
94 questions from the core math courses [4] have been automatically solved using few-shot learning and  
95 program synthesis at a human level.

## 96 2 Dataset

97 We present a new dataset of 496 questions from the eight most recent final exams of MIT’s In-  
98 troduction to Machine Learning course. The dataset spans questions on the 12 machine learning

Semester	Questions	Parts
Fall 2017	10	61
Spring 2018	9	42
Fall 2018	9	59
Spring 2019	9	58
Fall 2019	8	61
Spring 2021	12	70
Fall 2021	8	86
Spring 2022	9	59
Mean	9.25	62
Total	74	496

Table 3: The number of questions and parts in the final for each semester of Introduction to Machine Learning. Spring 2020 and Fall 2020 did not have final exams due to COVID-19.

Topic	Questions	Parts
Regression	6	58
Classifiers	10	60
Logistic Regression	1	8
Features	3.5	21
Neural Networks	12.5	76
Loss Functions	2	14
CNNs	8	63
MDPs	9	75
RNNs	7	33
Reinforcement Learning	8	50
Clustering	1	8
Decision Trees	6	32
Mean	6.73	45.09
Total	74	496

Table 4: The number of questions and parts in the final for each topic of Introduction to Machine Learning.

99 topics covered in the course: (1) regression, (2) classifiers, (3) logistic regression, (4) features, (5)  
100 loss functions, (6) neural networks, (7) convolutional neural networks (CNNs), (8) Markov decision  
101 processes (MDPs), (9) recurrent neural networks (RNNs), (10) reinforcement learning, (11) clustering,  
102 and (12) decision trees. We make our data and code publicly available.<sup>2</sup>

103 The breakdown of questions, parts, points, and non-image points by each semester and topic are  
104 shown in Table 3 and Table 4. Each question in a final exam consists of multiple parts. Questions are  
105 written by providing set-up and context information first, followed by the question parts (which may  
106 come with additional information). Set-up and context information may contain (1) story elements  
107 (ex., character names, and motivations), (2) relevant definitions and equations, and (3) data points.  
108 We format questions in the dataset by concatenating the question context, any context or solutions  
109 from prior parts of the question required for answering the part, and the part’s context and question.  
110 We split the questions into their corresponding parts. Questions consist of English text, mathematical  
111 notation, and images. Mathematical notation is represented in the dataset by LaTeX and images by  
112 screenshots from pdfs files.

<sup>2</sup>Data and code: <https://github.com/idrori/mlfinalsQ>

113 The types of question answers are diverse. A few are multiple-choice or true/false questions. Most  
114 are open-ended, for which the evaluation requires modeling the problem, mathematical manipulation,  
115 or code writing. Many questions require providing an explanation and justification.

116 To curate the dataset, we transcribe questions into text. The six final exams of Fall 2017, Spring  
117 2018, Fall 2018, Spring 2019, Fall 2019, and Spring 2022 were available in pdf format, so we used  
118 mathpix.com to automatically extract transcriptions, including both text and LaTeX. The final exams  
119 of Spring 2021 and Fall 2021 were available in LaTeX. We extract questions and solutions for all  
120 parts of all types of questions, including those that rely on images.

121 We curated the five exams between 2017 and 2019 from publicly available pdf files. Spring 2020 and  
122 Fall 2020 do not have final exams due to COVID-19. The three exams between 2021 and 2022 were  
123 unavailable online; therefore, the model does not overfit their solutions. The aggregate average grades  
124 were available to the students and did not contain any personally identifiable information. Three  
125 duplicate questions were originally on the final exam of Fall 2017 (questions 1, 3, 6) and appear again  
126 in the final exam of Spring 2022.

## 127 **3 Benchmark**

### 128 **3.1 Baselines**

129 We provide a benchmark by comparing six baselines for answering the final exam questions: (1)  
130 GPT-3 with zero-shot learning, (2) GPT-3 with few-shot learning, (3) GPT-3 with zero-shot learning  
131 and chain-of-thought (CoT) prompting, (4) GPT-3 with few-shot learning and chain-of-thought (CoT)  
132 prompting, (5) Codex with zero-shot learning, and (6) Codex with few-shot learning.

133 Table 5 shows the prompt used for each approach. GPT-3 zero-shot uses the question as-is, whereas  
134 GPT-3 zero-shot with CoT uses the suffix “Let’s think step by step.” after the question to encourage  
135 multi-step output. Codex zero-shot uses the prefix “Write a program that answers” before the question  
136 within Python comments denoted by triple quotes “""" to encourage Codex to write code. GPT-3  
137 few-shot finds the closest questions in the embedding space, measured by cosine similarity, and  
138 uses these questions and their corresponding answers before the new question as examples in the  
139 prompt. Codex few-shot finds the closest questions in the embedding space also as measured by  
140 cosine similarity and uses these questions and their corresponding code as examples.

141 For students, a good study technique is to use previous final exams to review and practice for their  
142 upcoming final. We model this method by few-shot learning using the question–answer pairs (for  
143 GPT-3) or question–code (for Codex) with the closest question embeddings from previous finals. We  
144 implement this by considering all the exam questions, marking each question by its semester and  
145 year, and using only previous semesters’ questions for few-shot learning. The Fall 2017 and Spring  
146 2022 contain three questions that are the same. We handle these duplicate questions the same way  
147 humans do by allowing few-shot learning in Spring 2022 based on successful Fall 2017 zero-shot  
148 answers. In the same way, humans benefit by practicing for an exam by doing previous exams. Since  
149 it is the first final exam, we do not perform few-shot learning on the Fall 2017 exam.

### 150 **3.2 Grading**

151 We grade the answers and aim to keep all factors equal in grading human and machine answers.  
152 Human and machine answers are graded based on the number of points allocated to each question  
153 part, giving full, partial, or no credit for each answer. We approximate partial credit by assigning  
154 half-credit. The course staff graded the student final exams, which included graduate TAs and  
155 instructors. Two of the same graduate TAs and the instructor that graded the student answers also  
156 graded the machine answers. The grading instructions are the same for grading student answers as  
157 grading machine answers.

Method	Prompt
GPT-3 Zero-Shot	<question>
GPT-3 Few-Shot	Q: <similar question> A: <similar question's answer> Q: <question> A:
GPT-3 Zero-Shot with CoT	Q: <question> A: Let's think step by step.
GPT-3 Few-Shot with CoT	Q: <similar question> A: <similar question's answer> Q: <question> A: Let's think step by step.
Codex Zero-Shot	"" Write a program that answers the following question: <question> ""
Codex Few-Shot	"" Write a program that answers the following question: <similar question> "" <similar question's correct code> "" Write a program that answers the following question: <question> ""

Table 5: Input prompt for each of six baseline methods (1) GPT-3 Zero-Shot, (2) GPT-3 Few-Shot, (3) GPT-3 Zero-Shot with CoT, (4) GPT-3 Few-Shot with CoT, (5) Codex Zero-Shot, and (6) Codex Few-Shot.

### 3.3 Performance

Table 6 shows the machine grades by semester and Table 7 shows the machine grades by topic, excluding question parts that rely on images. We compare the average grade of GPT-3 with zero-shot (ZS), GPT-3 with ZS and chain-of-thought (CoT) prompting, GPT-3 with few-shot (FS) learning, GPT-3 with FS and CoT prompting, Codex with ZS, and Codex with FS. Fall 2017 is the first semester, so few-shot learning results based on previous semesters are unavailable (NA). Spring 2020 and Fall 2020 did not have final exams due to COVID-19. Spring 2021, Fall 2021, and Spring 2022 final exams were unavailable online when GPT-3 and Codex were trained, ensuring that the model is not overfitting content it has seen previously. The results consistently demonstrate that few-shot learning methods perform best across semesters and topics, as marked in bold.

### 3.4 Limitations

Our dataset consists of all question parts and their solutions, including images. However, our baseline methods do not handle questions that rely on an image containing the information required to solve the question since GPT-3 and Codex do not handle images. Tables 8 and 9 show the breakdown of the number of question parts and points of questions that do not rely on image information for answering the question. On average, 28.42% of the question parts, which make up 35.63% of the points in final exams, are questions that rely on image information. The points attributed to the non-image parts are tallied, recorded, and used to calculate non-image percentage grades.

### 3.5 Generating New Questions

We use the dataset of exam questions to generate new questions automatically. We use questions from our dataset as prompts to create new high-quality questions not present in our dataset. We create a list

Semester	GPT-3 ZS	GPT-3 FS	GPT-3 ZS CoT	GPT-3 FS CoT	Codex ZS	Codex FS
Fall 2017	<b>38.21</b>	NA	22.86	NA	21.43	NA
Spring 2018	44.35	60.48	38.71	<b>70.97</b>	32.26	67.74
Fall 2018	51.99	52.18	61.63	<b>64.17</b>	49.78	54.00
Spring 2019	43.45	54.23	41.07	<b>58.81</b>	15.54	41.55
Fall 2019	53.17	<b>62.70</b>	28.97	56.35	25.40	59.13
Spring 2021	44.33	55.81	53.45	60.21	33.62	<b>62.09</b>
Fall 2021	58.94	<b>58.94</b>	50.35	54.90	18.11	42.00
Spring 2022	42.78	<b>68.86</b>	32.03	53.48	51.01	65.46

Table 6: We benchmark different baselines for each semester, excluding question parts that rely on images. We compare the average grade of GPT-3 with zero-shot (ZS), GPT-3 with few-shot (FS) learning, GPT-3 with ZS, and chain-of-thought (CoT) prompting, GPT-3 with FS and CoT prompting, Codex with ZS, and Codex with FS. Fall 2017 is the first semester, so few-shot learning results based on previous semesters are unavailable (NA). Spring 2020 and Fall 2020 did not have final exams due to COVID-19. Spring 2021, Fall 2021, and Spring 2022 final exams were unavailable online when GPT-3 and Codex were trained, ensuring that the model is not overfitting. The highest-performing method is bolded for each semester.

Topic	GPT-3 ZS	GPT-3 FS	GPT-3 ZS CoT	GPT-3 FS CoT	Codex ZS	Codex FS
Regression	31.71	50.00	25.61	40.85	40.24	<b>50.00</b>
Classifiers	38.18	46.21	26.28	42.35	18.88	<b>53.74</b>
Logistic Regression	50.00	60.00	77.50	<b>77.50</b>	55.00	70.00
Features	58.65	75.96	53.85	77.31	68.85	<b>81.54</b>
Loss Functions	NA	NA	NA	NA	NA	NA
Neural Networks	48.34	60.23	44.54	<b>68.42</b>	37.82	63.45
CNNs	37.50	<b>53.58</b>	28.36	47.81	13.38	36.77
MDPs	49.19	52.01	46.03	<b>54.23</b>	24.38	38.03
RNNs	61.46	<b>71.88</b>	57.29	66.14	12.50	40.63
Reinforcement Learning	36.09	42.99	36.67	<b>50.11</b>	28.79	45.11
Clustering	100.00	100.00	100.00	<b>100.00</b>	50.00	50.00
Decision Trees	54.70	<b>71.80</b>	32.48	51.28	46.15	54.70

Table 7: We benchmark different baselines for each course topic, excluding question parts that rely on images. We compare the grade of GPT-3 with zero-shot (ZS), GPT-3 with few-shot (FS) learning, GPT-3 with zero-shot and chain-of-thought (CoT) prompting, GPT-3 with FS and CoT, Codex with zero-shot, and Codex with few-shot learning. The question parts on loss functions rely on image information and are therefore unavailable (marked NA). The highest-performing method is bolded for each semester.

Semester	Non-Image Parts / All Parts	Non-Image Points / All Points
Fall 2017	49 / 61	70 / 100
Spring 2018	27 / 42	62 / 100
Fall 2018	29 / 59	62 / 100
Spring 2019	41 / 58	70 / 100
Fall 2019	47 / 61	63 / 100
Spring 2021	50 / 70	61 / 100
Fall 2021	56 / 86	48 / 100
Spring 2022	46 / 59	69 / 100
Total	345 / 496	505 / 800

Table 8: The number of question parts that do not rely on images and the number of points that do not rely on images in Introduction to Machine Learning finals for each semester. Spring 2020 and Fall 2020 did not have final exams due to COVID-19.



Topic	Non-Image Parts / All Parts	Non-Image Points / All Points
Regression	33 / 58	41 / 67
Classifiers	52 / 60	101 / 123
Logistic Regression	8 / 8	10 / 10
Features	18 / 21	26 / 38
Loss Functions	0 / 14	0 / 18
Neural Networks	56 / 76	87 / 128
CNNs	52 / 63	57 / 76
MDPs	38 / 72	39 / 112
RNNs	27 / 33	48 / 67
Reinforcement Learning	42 / 50	55 / 79
Clustering	1 / 8	2 / 12
Decision Trees	18 / 32	39 / 72
Total	345 / 496	505 / 800

Table 9: The number of questions parts, and the number of points that do not rely on images, and number of points that do not rely on images in Introduction to Machine Learning finals for each topic of the course.

of various questions in our curated dataset and use the resulting list to prompt GPT-3 to create a new question. The supplementary material demonstrates the results of this process for each topic in the course. The Appendix consists of new generated questions and the closest question from our dataset as measured by the cosine similarity of the embedding of each question. These new questions are diverse and qualitatively similar to questions on previous MIT final exams. This provides an efficient way for course TAs and instructors to generate new final questions.

### 3.6 Implementation Details

We use the latest OpenAI GPT-3 and Codex models and do not re-train these very large language models. We fix all the hyperparameters of the models so that the answers are deterministic and reproducible. Specifically, we set both top P, which controls diversity, and sampling temperature, which controls randomness, to 0. The frequency and presence penalties are also set to 0, and we do not halt on any stop sequences. We allow diversity for generating new questions by setting the top P and temperature to 0.1. We run Codex with an upper bound of generating programs with 1024 tokens. We use the OpenAI text-davinci-002 and code-davinci-002 engines for generating text and programs. For few-shot-learning and question generation, we use the text-similarity-babbage-001 engine to embed the questions and find the closest questions in the dataset by cosine similarity. The running time for answering or generating each question part is a few seconds.

## 4 Conclusions

We present a dataset and benchmark for answering and generating university-level final exams in machine learning. The dataset is the first of its kind, and the benchmark compares state-of-the-art language model approaches. Machine performance and human performance are evaluated by the same graders and grading instructions. A comparison of baselines shows that few-shot learning methods perform best across semesters and topics. A limitation of our work is that our benchmark does not consider questions that rely on images for their solution. Potential negative societal impacts of the work may be applications that are used for helping answer finals during the exams. Potential positive societal impacts include improving students learning for final exams, helping course staff generate questions for finals, and comparing levels of difficulty of exams across semesters and schools. We hope this dataset and benchmark serve the machine learning community and advance the state-of-the-art in the field. Future work will extend the benchmark to handle questions that rely on images and extend the dataset to finals exams in many STEM courses and universities.

## References

- [1] Tom B Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. In *Proceedings of Advances in Neural Information Processing Systems (NeurIPS)*, volume 33, pages 1877–1901, 2020.
- [2] Mark Chen et al. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*, 2021.
- [3] Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. Palm: Scaling language modeling with pathways. *arXiv preprint arXiv:2204.02311*, 2022.
- [4] Iddo Drori, Sarah Zhang, Reece Shuttleworth, Leonard Tang, Albert Lu, Elizabeth Ke, Kevin Liu, Linda Chen, Sunny Tran, Newman Cheng, Roman Wang, Nikhil Singh, Taylor L. Patti, Jayson Lynch, Avi Shporer, Nakul Verma, Eugene Wu, and Gilbert Strang. A neural network solves, explains, and generates university math problems by program synthesis and few-shot learning at human level. *arXiv preprint arXiv:2112.15594*, 2022.
- [5] Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. Large language models are zero-shot reasoners. *arXiv preprint arXiv:2205.11916*, 2022.
- [6] Yujia Li, David Choi, Junyoung Chung, Nate Kushman, Julian Schrittwieser, Rémi Leblond, Tom Eccles, James Keeling, Felix Gimeno, Agustin Dal Lago, et al. Competition-level code generation with alphacode. *arXiv preprint arXiv:2203.07814*, 2022.
- [7] Jack W Rae, Sebastian Borgeaud, Trevor Cai, Katie Millican, Jordan Hoffmann, Francis Song, John Aslanides, Sarah Henderson, Roman Ring, Susannah Young, et al. Scaling language models: Methods, analysis & insights from training gopher. *arXiv preprint arXiv:2112.11446*, 2021.
- [8] Leonard Tang, Elizabeth Ke, Nikhil Singh, Bo Feng, Derek Austin, Nakul Verma, and Iddo Drori. Solving probability and statistics problems by probabilistic program synthesis at human level and predicting solvability. In *Proceedings of the International Conference on Artificial Intelligence in Education (AIED)*, 2022.
- [9] Sunny Tran, Pranav Krishna, Ishan Pakuwal, Prabhakar Kaffle, Nikhil Singh, Jayson Lynch, and Iddo Drori. Solving machine learning problems. In *Proceedings of the Asian Conference on Machine Learning (ACML)*, pages 470–485, 2021.
- [10] Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, and Denny Zhou. Self-consistency improves chain of thought reasoning in language models. *arXiv preprint arXiv:2203.11171*, 2022.
- [11] 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.
- [12] Denny Zhou, Nathanael Schärli, Le Hou, Jason Wei, Nathan Scales, Xuezhi Wang, Dale Schuurmans, Olivier Bousquet, Quoc Le, and Ed Chi. Least-to-most prompting enables complex reasoning in large language models. *arXiv preprint arXiv:2205.10625*, 2022.

## Checklist

The checklist follows the references. Please read the checklist guidelines carefully for information on how to answer these questions. For each question, change the default **[TODO]** to **[Yes]**, **[No]**, or **[N/A]**. You are strongly encouraged to include a **justification to your answer**, either by referencing the appropriate section of your paper or providing a brief inline description. For example:

- Did you include the license to the code and datasets? **[Yes]** See Section ??.
- Did you include the license to the code and datasets? **[No]** The code and the data are proprietary.
- Did you include the license to the code and datasets? **[N/A]**

Please do not modify the questions and only use the provided macros for your answers. Note that the Checklist section does not count towards the page limit. In your paper, please delete this instructions block and only keep the Checklist section heading above along with the questions/answers below.

### 1. For all authors...

- (a) Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope? **[Yes]**
- (b) Did you describe the limitations of your work? **[Yes]** See Section 3.4.
- (c) Did you discuss any potential negative societal impacts of your work? **[Yes]** See Section 4.
- (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? **[Yes]**

### 2. If you are including theoretical results...

- (a) Did you state the full set of assumptions of all theoretical results? **[N/A]**
- (b) Did you include complete proofs of all theoretical results? **[N/A]**

### 3. If you ran experiments (e.g. for benchmarks)...

- (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? **[Yes]** See our repository
- (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? **[Yes]** See Section 3.6.
- (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? **[N/A]**
- (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? **[Yes]** See Section 3.6.

### 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...

- (a) If your work uses existing assets, did you cite the creators? **[Yes]**
- (b) Did you mention the license of the assets? **[Yes]**
- (c) Did you include any new assets either in the supplemental material or as a URL? **[Yes]** See our repository
- (d) Did you discuss whether and how consent was obtained from people whose data you're using/curating? **[Yes]** See Section 2.
- (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? **[Yes]** See Section 2.

### 5. If you used crowdsourcing or conducted research with human subjects...

- (a) Did you include the full text of instructions given to participants and screenshots, if applicable? **[N/A]**

- 294 (b) Did you describe any potential participant risks, with links to Institutional Review  
295 Board (IRB) approvals, if applicable? [N/A]
- 296 (c) Did you include the estimated hourly wage paid to participants and the total amount  
297 spent on participant compensation? [N/A]