

WIP: Using Machine Learning to Automate Coding of Student Explanations to Challenging Mechanics Concept Questions

Harpreet Auby

Harpreet Auby is a STEM Education MS and Chemical Engineering PhD student at Tufts University. He is a graduate research assistant working with Dr. Milo Koretsky within the Institute for Research on Learning and Instruction (IRLI). Harpreet received his BS in Chemical and Biomolecular Engineering at the University of Illinois at Urbana-Champaign. His current work focuses on machine learning applications in educational research and evaluation, learning assistants, and uptake of an online technology tool emphasizing concept-based learning called the Concept Warehouse. His broad research interests include engineering education, learning and sensemaking in STEM, and liberatory pedagogies in STEM Ed.

Namrata Shivagunde

Namrata Shivagunde is a phd student in computer science at UMass Lowell. She is working with Prof. Anna Rumshisky at Text Machine Lab. Her research is in application of deep learning techniques in natural language processing. Previously she did MS in Applied and Computational Mathematics from UMass Lowell.

Anna Rumshisky

Anna Rumshisky is an Associate Professor of Computer Science at the University of Massachusetts Lowell, where she leads the Text Machine Lab for NLP. She is also a Visiting Academic at Amazon Alexa AI. Her primary research area is machine learning for natural language processing, with a focus on deep learning techniques.

Milo Koretsky (McDonnell Family Bridge Professor)

Milo Koretsky (he/him/his) is the McDonnell Family Bridge Professor holding a joint appointment in Chemical and Biological Engineering and Education at Tufts University. He received his BS and MS degrees from UC San Diego and his PhD from UC Berkeley, all in chemical engineering. He is interested in integrating technology into effective educational practices and in promoting the use of higher-level cognitive and social skills in engineering problem solving.

WIP: Using Machine Learning to Automate Coding of Student Explanations to Challenging Mechanics Concept Questions

Introduction

This work-in-progress paper describes a collaborative effort between engineering education and machine learning researchers to automate analysis of written responses to conceptually challenging questions in mechanics. These qualitative questions are often used in large STEM classes to support active learning pedagogies; they require minimum calculations and focus on the application of underlying physical phenomena to various situations. Active learning pedagogies using this type of questions has been demonstrated to increase student achievement (Freeman et al., 2014; Hake, 1998) and engagement (Deslauriers, et al., 2011) of all students (Haak et al., 2011).

To emphasize reasoning and sense-making, we use the Concept Warehouse (Koretsky et al., 2014), an audience response system where students provide written justifications to concept questions. Written justifications better prepare students for discussions with peers and in the whole class and can also improve students' answer choices (Koretsky et al., 2016a, 2016b). In addition to their use as a tool to foster learning, written explanations can also provide valuable information to concurrently assess that learning (Koretsky and Magana, 2019). However, in practice, there has been limited deployment of written justifications with concept questions, in part, because they provide a daunting amount of information for instructors to process and for researchers to analyze.

In this study, we describe the initial evaluation of large pre-trained generative sequence-to-sequence language models (Raffel et al., 2019; Brown et al., 2020) to automate the laborious coding process of student written responses. Adaptation of machine learning algorithms in this context is challenging since each question targets specific concepts which elicit their own unique reasoning processes. This exploratory project seeks to utilize responses collected through the Concept Warehouse to identify viable strategies for adapting machine learning to support instructors and researchers in identifying salient aspects of student thinking and understanding with these conceptually challenging questions.

Machine Learning of Constructed Responses

Machine learning has been leveraged in a number of educational applications (Zhai et al., 2020b, Zhai et al., 2021a, Burstein et al., 2020, Burstein et al., 2021), including analyzing constructed responses (short text) and essays (long text), diagnostic reasoning (Schulz et al., 2019), and studying learning processes through simulation and educational games (Zhai et al 2020b). In supervised learning, the machine learning model is trained using a training set (coded data) and is evaluated on a test set (uncoded data). SVM, Naive-Bayes, Random Forest and Logistic Regression have been most commonly used for constructed-response assessments in STEM (Zhai et al., 2021, Zhai et al., 2020a, Mao et al, 2018, Yik et al, 2021, Jescovitch et al., 2021, Rosenberg, 2021). Many studies also applied ensemble techniques like bagging, boosting on various text classification machine learning models to study student responses (Bertolini et al., 2021, Zhai et al., 2020a). Several studies have also used neural network models (Jiang et al., 2020; Luan et al., 2021; Rosenberg, 2021). However, to our knowledge, only a few studies for

educational applications in general have leveraged *Transformer-based machine learning models* (Vaswani et al., 2017, Devlin et al., 2018, Raffel et al., 2019, Brown et al., 2020). And none of the work used these models for assessing constructed responses of STEM students.

Transformer models, the current state of the art in natural language processing (NLP), are attention-based multi-layer neural networks pre-trained on large amounts of free text. The pre-training process uses a language modeling objective, i.e. the model is asked to predict a word or token, given other words in contexts. Such models are then fine-tuned on specific language tasks, or are used out of the box for *in-context learning*, where the model is queried with a prompt and asked to generate some text along with its interpretation. These models typically use a large amount of trainable parameters, ranging between a hundred million and a hundred billion, with larger models capable of more sophisticated prompted in-context learning.

While earlier studies required two to ten human coders to annotate 50 to 1000 samples of the data (Haudek et al., 2021, Mao et al., 2018, Jescovitch et al., 2021, Maestrales et al., 2021), the current state-of-the-art NLP models that leverage transfer learning can require substantially fewer annotated samples for fine-tuning and only a few annotated examples for in-context learning. In the present work, we leverage these capabilities and investigate the amount of annotated data needed to automatically analyze students' constructive responses of complex conceptual questions.

Methods

Context and setting

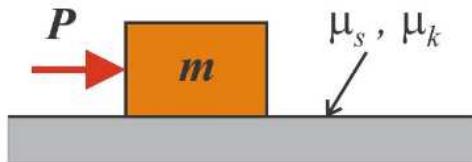
This study occurs in the context of a larger investigation which studies the propagation of the Concept Warehouse to mechanics courses in a diverse set of two- and four- year institutions (Koretsky et al., 2019; Nolen & Koretsky, 2020). Such service-oriented mechanics courses build foundational skills for upper-level engineering courses and develop students' problem-solving capabilities. The participants were consenting students from statics courses at different institutions. The instructors used the Concept Warehouse as part of their active learning course delivery. Eight questions were delivered by the participating instructors in Fall 2021. All questions are single correct-choice, qualitative conceptually challenging problems with little to no calculation involved. They all test the application of critical statics concepts to physical systems. The incorrect multiple-choice responses (distractors) are all conceptually significant, providing students the opportunity to carefully reason through the questions. The range correct from the eight questions (31% - 58.6%) indicates these questions are conceptually challenging for students. Thus, the associated written explanations are good candidates for machine learning analysis to reveal student reasoning.

Qualitative analysis

For the preliminary analysis reported here, one question, CW5703 - shown in Figure 1, was used for initial manual coding and machine learning coding. Using a combination of *a priori* and emergent approaches described in Creswell & Poth (2018), a coding scheme was developed to classify elements of student explanations and provide categories to train the machine learning algorithms. These elements were used to construct narratives of students' reasoning processes particular to each question or isomorphic question pair. The final code categories for question CW5703 are shown in Appendix A.1

Force $P = 10 \text{ N}$ is applied to the block of mass $m = 5 \text{ kg}$ on a horizontal rough surface with $\mu_s = 0.3$ and $\mu_k = 0.25$.

If $g = 9.81 \text{ m/s}^2$, what is the force of friction on the block?



- 10 N
- 12.26 N
- 14.7 N
- 45.1 N

Please explain your answer in the box below.



Please rate how confident you are with your answer.

substantially
unsure moderately
unsure neutral moderately
confident substantially
confident

-
-
-
-
-

Submit

Figure 1. Student view of a sample statics concept question (CW5703) used in this study. Students need to select an answer choice and justify their choice in writing.

Machine Learning

The machine learning task was formulated as a sequence labeling problem where the source is a student written explanation and the target is the human-coded response. INCEpTION (Klie et al., 2018) was used to translate manual coding to test spans in a tsv format for use in algorithm training. Transfer learning via fine-tuning and in-context learning techniques were used to respond to a prompt containing coding examples. Initial experimentation was carried out using Text-to-Text Transformer (T5) (Raffel et al., 2019) with fine-tuning employing Huggingface's transformer library and GPT3 (Brown et al., 2020) with in-context learning utilizing OpenAI GPT3 API. Experiments with T5 followed a prompt format from Raffel et al. (2019) with alterations that include an instruction sentence and prompt keywords. Examples of the instruction sentence and prompt keywords are shown in Appendix A2.1. In T5, every text processing task is reformatted as a text-to-text problem i.e. taking in a text as input and producing a new text as output. In our preliminary study, we reformatted our task into a text-to-text format and fine-tuned T5-base (220M parameters) and T5-large (770M parameters) with 20 to 240 human coded responses.

Experimentation with GPT-3 followed prompt instruction from Brown et al. (2020) and does not include fine-tuning. It involved the use of the GPT3-davinci-instruct (175 billion parameters) model which is the biggest model available that also works best with these kinds of instruction prompts. Two formats were used to support in-context learning and annotate several answers at once. The first format utilized an individual response and its annotation and the second format used group responses and group annotations. Examples of the prompt keywords are shown in Appendix A2.2.

Preliminary Results

After iteratively refining manual coding of student explanations, a narrative of students' reasoning processes was created that led to categorization according to three main, sequential cognitive processes: identification, comparison, and inference. Table A1 in Appendix A.1 shows the three processes, the 11 code categories that resulted, and some common language used to describe these concepts. Once all hand-coded spans were translated via INCEpTION, the machine learning models were run for both a training set and testing set. An example from the test set is shown in Table 1, with the results from various models. For this case, the models shown in Table 1 all identified the first two codes correctly, while the last two codes were only correctly identified by the T5 base models.

Table 1: Example of human-coded example (ground truth) with outputs from 3 out of 8 machine learning models run in initial experimentation. All codes are in red.

ground truth	gpt3-davinci-instruct	ft_large_ft_150	t5_base_ft_200
Max static friction is 14 N <Identification> <Static Friction Initial> so the box remains at rest <Inference> <Box Movement>. It will perfectly match the force of 10N <Inference> <Newton's Laws>	Max static friction is 14 N <Identification><Static Friction Initial>, so the box remains at rest <Inference><Box Movement>. It will perfectly match the force of 10N <Comparison><Static Friction to Applied Force>.	Max static friction is 14 N Identification> Static Friction Initial> so the box remains at rest. Inference> Box Movement> It will perfectly match the force of 10N. Inference> Newton's Laws>	Max static friction is 14 N Identification> Static Friction Initial> so the box remains at rest. Inference> Box Movement> It will perfectly match the force of 10N. Inference> Newton's Laws>

Overall, 290 responses were manually coded. To investigate which model performed the best, we evaluated model outputs for 50 held-out responses in comparison to the human-coded "ground truth." Table 2 shows results obtained with T5 and GPT3. We report precision, recall and F1 score for model-generated outputs. Precision is the percentage of correct model-generated codes, relative to the total number of model-generated codes. Recall is the percentage of human codes which the model was able to generate correctly. The F1 score is the harmonic mean of precision and recall. Ground truth is comprised of 175 human-assigned codes with 50 codes held-out for evaluation. T5-large fine-tuned on 150 samples performed best, with an F1 score of 0.73. T5-base fine-tuned on 240 codes had recall comparable to T5-large and second best F1 score, but had low precision. To gain further insight into the successes and failures of each model, we manually analyzed model-generated codes that did not match ground truth to determine what percentage of these codes in fact made sense. The breakdown between reasonable and meaningless model-generated codes is shown in columns 7 ("Misses but makes sense") and 8 ("Does not make sense") in Table 2. While GPT3 did not match as many ground truth responses as t5-large, in many cases, it generated meaningful responses. In fact, responses generated by

GPT3 turned out to be more creative and generated some new codes which were not present in the in-context examples, while T5 generated codes very similar or close to that of the fine-tuning dataset. Note that GPT3 also over-generated annotations, which explains the negative number of missed codes in the last column (“Codes missed”).

Table 2: Comparison of ground truth and model-generated responses. Best result is in bold.

Model	Correct codes	Total codes	Precision	Recall	F1	Misses but makes sense	Does not make sense	Codes missed
Ground truth	175							
t5-base-f20	0	0	0	0	0	0	0	175
t5-base-f50	40	49	0.82	0.23	0.36	2	7	126
t5-base-f100	60	90	0.67	0.34	0.45	14	16	85
t5-base-f150	80	92	0.87	0.46	0.60	7	5	83
t5-base-f200	93	126	0.74	0.54	0.62	19	14	49
t5-base-f240	105	133	0.79	0.60	0.68	14	14	42
t5-large-f150	107	118	0.91	0.61	0.73	6	5	57
gpt3-davinci-Instruct	89	189	0.47	0.51	0.49	52	48	-14

*We use t5-base-fXXX to indicate that t5-base was fine-tuned with XXX examples.

Implications

Our work shows promise for further application of machine learning in education. We seek to further characterize the feasibility of integrating machine learning tools into the Concept Warehouse to support instruction and research and to address the challenges faced during these preliminary experiments. Some of these goals include fine-tuning minimum data size, testing the ability to transfer to isomorphic questions, determining accuracy ranges of machine learning, and developing an automatic evaluation method for machine coded responses.

We envision that for instructors, such machine learning algorithms can enable processing of large amounts of data regarding student explanations to provide information on patterns, trends, and general ideas of student thinking that they could utilize in their instructional practices and pedagogical decision-making processes. For educational researchers, the machine learning algorithms could provide ways to determine the narrative of understanding students have in various institutional contexts at a scale not feasible with manual coding.

Acknowledgments

We acknowledge the support from the National Science Foundation (NSF) through grant DUE 2135190. Any opinions, findings, and conclusions or recommendations expressed are those of the authors and do not necessarily reflect the views of the NSF.

References

Bertolini, R., Finch, S. J., & Nehm, R. H. (2021). Testing the impact of novel assessment sources and machine learning methods on predictive outcome modeling in undergraduate biology. *Journal of Science Education and Technology*, 30(2), 193-209.

Brown, T.B., Mann, B., Ryder, N., Subbiah, M., Kaplan, J., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A. and Agarwal, S., (2020). Language models are few-shot learners. arXiv preprint arXiv:2005.14165.

Burstein, J., Horbach, A., Kochmar, K., Laarmann-Quante, R., Leacock, C., Madnani, Nitin., Pilan, I., Yannakoudakis, H., Zesch, T., Proceedings of the 16th Workshop on Innovative Use of NLP for Building Educational Applications, Association for Computational Linguistics, April 2021.

Burstein, J., Kochmar, K., Leacock, C., Madnani, Nitin., Pilan, I., Yannakoudakis, H., Zesch, T., Proceedings of the 15th Workshop on Innovative Use of NLP for Building Educational Applications, Association for Computational Linguistics, 2020.

Creswell, J. W., & Poth, C. N. (2018). Qualitative Inquiry & Research Design: Choosing Among Five Approaches (Vol. 4). SAGE.

Deslauriers, L., Schelew, E., & Wieman, C. (2011). Improved learning in a large-enrollment physics class. *science*, 332(6031), 862-864.

Devlin, J., Chang, M. W., Lee, K., & Toutanova, K. (2018). Bert: Pre-training of deep bidirectional transformers for language understanding. arXiv preprint arXiv:1810.04805.

Dollar, A., & Steif, P. (2004). Reinventing the teaching of Statics. 2004 Annual Conference Proceedings. <https://doi.org/10.18260/1-2--13940>

Freeman, S., Eddy, S. L., McDonough, M., Smith, M. K., Okoroafor, N., Jordt, H., & Wenderoth, M. P. (2014). Active learning increases student performance in science, engineering, and mathematics. *Proceedings of the National Academy of Sciences*, 111(23), 8410-8415.

Haak, D. C., HilleRisLambers, J., Pitre, E., & Freeman, S. (2011). Increased structure and active learning reduce the achievement gap in introductory biology. *Science*, 332(6034), 1213-1216.

Hake, R. R. (1998). Interactive-engagement versus traditional methods: A six-thousand-student survey of mechanics test data for introductory physics courses. *American journal of Physics*, 66(1), 64-74.

Haudek, K. C., & Zhai, X. (2021). Exploring the Effect of Assessment Construct Complexity on Machine Learning Scoring of Argumentation.

Jescovitch, L. N., Scott, E. E., Cerchiara, J. A., Merrill, J., Urban-Lurain, M., Doherty, J. H., & Haudek, K. C. (2020). Comparison of Machine Learning Performance Using Analytic and Holistic Coding Approaches Across Constructed Response Assessments Aligned to a Science Learning Progression. *Journal of Science Education and Technology*. doi:10.1007/s10956-020-09858-0.

Jiang, R., Gouvea, J., Hammer, D., Miller, E., & Aeron, S. (2020). Automatic coding of students' writing via Contrastive Representation Learning in the Wasserstein space. arXiv preprint arXiv:2011.13384.

Klie, J.-C., Bugert, M., Boullosa, B., Eckart de Castilho, R. and Gurevych, I. (2018): The INCEPTION Platform: Machine-Assisted and Knowledge-Oriented Interactive Annotation. In Proceedings of System Demonstrations of the 27th International Conference on Computational Linguistics (COLING 2018), Santa Fe, New Mexico, USA.

Koretsky, M. D., Brooks, B. J., & Higgins, A. Z. (2016a). Written justifications to multiple-choice concept questions during active learning in class. *International Journal of Science Education*, 38(11), 1747-1765.

Koretsky, M. D., Brooks, B. J., White, R. M., & Bowen, A. S. (2016b). Querying the questions: Student responses and reasoning in an active learning class. *Journal of Engineering Education*, 105(2), 219-244.

Koretsky, M. D., Falconer, J. L., Brooks, B. J., Gilbuena, D. M., Silverstein, D. L., Smith, C., & Miletic, M. (2014). The AiChE Concept Warehouse: A web-based tool to promote concept-based instruction. *Advances in Engineering Education*, 4(1), 7:1-27.

Koretsky, M. D., & Magana, A. J. (2019). Using Technology to Enhance Learning and Engagement in Engineering. *Advances in Engineering Education*.

Koretsky, M., Nolen, S., Self, B., Papadopoulos, C., Widmann, J., Prince, M., & Dal Bello, D. (2019). For Systematic Development of Conceptests for Active Learning. *EDULEARN Proceedings*, 1.

Luan, H., & Tsai, C. C. (2021). A review of using machine learning approaches for precision education. *Educational Technology & Society*, 24(1), 250-266.

Maestrales, S., Zhai, X., Touitou, I., Baker, Q., Schneider, B., & Krajcik, J. (2021). Using machine Learning to Score Multi-Dimensional Assessments of Chemistry and Physics. *Journal of Science Education and Technology*, 30(2), 239-254.

Mao, L., Liu, O. L., Roohr, K., Belur, V., Mulholland, M., Lee, H.-S., & Pallant, A. (2018). Validation of automated scoring for a formative assessment that employs scientific argumentation. *Educational Assessment*, 23 (2), 121–138.

Nolen, S. B., & Koretsky, M. D. (2020, June). WIP: An Ecosystems Metaphor for Propagation. In *ASEE annual conference proceedings*.

Plesha, M., Gray, G., & Costanzo, F. (n.d.). Problem solving in statics and dynamics: A proposal for a structured approach. 2005 Annual Conference Proceedings. <https://doi.org/10.18260/1-2--15371>

Raffel, C., Shazeer, N., Roberts, A., Lee, K., Narang, S., Matena, M., ... & Liu, P. J. (2019). Exploring the limits of transfer learning with a unified text-to-text transformer. arXiv preprint arXiv:1910.10683.

Rosenberg, J. M., & Krist, C. (2021). Combining machine learning and qualitative methods to elaborate students' ideas about the generality of their model-based explanations. *Journal of Science Education and Technology*, 30(2), 255-267.

Schulz, Claudia, Christian M. Meyer, and Iryna Gurevych. "Challenges in the automatic analysis of students' diagnostic reasoning." In *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 33, no. 01, pp. 6974-6981. 2019.

Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... & Polosukhin, I. (2017). Attention is all you need. In *Advances in neural information processing systems* (pp. 5998-6008).

Yik, B. J., Dood, A. J., de Arellano, D. C. R., Fields, K. B., & Raker, J. R. (2021). Development of a machine learning-based tool to evaluate correct Lewis acid–base model use in written responses to open-ended formative assessment items. *Chemistry Education Research and Practice*, 22(4), 866-885.

Zhai, X., Haudek, K. C., Stuhlsatz, M. A., & Wilson, C. (2020c). Evaluation of construct-irrelevant variance yielded by machine and human scoring of a science teacher PCK constructed response assessment. *Studies in Educational Evaluation*, 67, 100916.

Zhai, X., Krajcik, J., & Pellegrino, J. W. (2021a). On the validity of machine learning-based Next Generation Science Assessments: a validity inferential network. *Journal of Science Education and Technology*, 30(2), 298-312.

Zhai, X., Shi, L., & Nehm, R. H. (2021b). A meta-analysis of machine learning-based science assessments: factors impacting machine-human score agreements. *Journal of Science Education and Technology*, 30(3), 361-379.

Zhai, X., Yin, Y., Pellegrino, J. W., Haudek, K. C., & Shi, L. (2020a). Applying machine learning in science assessment: a systematic review. *Studies in Science Education*, 56(1), 111-151.

Zhai, X., Haudek, K., Shi, L., Nehm, R., & Urban-Lurain, M. (2020b). From substitution to redefinition: A framework of machine learning-based science assessment. *Journal of Research in Science Teaching*, 57(9), 1430-1459. doi:10.1002/tea.21658.

Appendix A: Manual and Machine Learning Coding Processes

A.1 Manual Coding

After several iterations of qualitative coding, the final reasoning categories and codes were developed. Table A1 also describes the definitions of these codes as well as some examples of common language students use to describe their reasoning.

Table A1. Final list of categories and conceptual codes for CW5703.

Reasoning Category	Code	Code Definition	Common Language
Identification	Friction (General)	student describes what it is qualitatively or quantitatively	<ul style="list-style-type: none"> - Maximum Static Friction - Kinetic Friction - Normal Force (N) - Formulae: $\mu_k N$, $\mu_s N$, etc. - Force - Friction
	Initial Assumptions Miscellaneous	any other assumptions identified by the student in the beginning	
	Kinetic Friction Initial	student describes what it is qualitatively or quantitatively	
	Normal Force	student describes what it is qualitatively or quantitatively	
	Static Friction Initial	student describes what it is qualitatively or quantitatively	
Comparison	Compare Kinetic Friction Force to Applied Friction Force	Student makes clear that they take the concept of kinetic frictional force and compare it to the applied force.	<ul style="list-style-type: none"> - force applied is less/more/higher/lower - force has not broken the static friction barrier - P is not large enough to overcome - P is less than the maximum force of static friction - maximum static friction force is larger/higher than 10 N/P
	Compare Static Friction Force to Applied Friction Force	Student makes clear that they take the concept of static frictional force and compare it to the applied force.	
	Solve for Own Coefficient of Friction	Student uses a coefficient they calculate and use the parameter for comparison	
Inference	Box Movement	Student explicitly states what will or will not happen to the box.	<ul style="list-style-type: none"> - the block is not in motion - the block won't slide - push back the same amount - the frictional force is only 10 N - equal and opposite - friction force and P must be equal - maintain static equilibrium
	Application of Newton's Laws	Student either quantifies the force on the box as 10 N or mentions the concepts of "static equilibrium" and needing the push/pull forces to be equal	
	Uncertainty		

A.2 Machine Learning

A.2.1 T5 Instruction Format and Prompts

The T5 model was used to support machine learning via fine-tuning and in-context learning. We largely used the input prompt formats described by Raffel et al. (2019). However, some alterations were made to better fit the nature of the student explanations. This included adding instruction keywords and prompt keywords to better have the model understand the task. This input format is shown below with the added source prompts in **blue** and manual coding in **red**.

Source : “*Given the question, annotate the answer. question: Force $P = 10 \text{ N}$ is applied to the block of mass $m = 5 \text{ kg}$ on a horizontal rough surface with $\mu_s = 0.3$ and $\mu_k = 0.25$. If $g = 9.81 \text{ m/s}^2$, what is the force of friction on the block?. answer: The maximum force of static friction ($.3 * 5 * 9.81$) is larger than the applied force in the x direction. This means that the force of static friction will be equal and opposite to the applied forces x component.”*

Target : “*The maximum force of static friction ($.3 * 5 * 9.81$) <Identification> <Static Friction Initial> is larger than the applied force in the x direction . <Comparison> <Static Friction to Applied Force> force of static friction will be equal and opposite to the applied forces x component . <Inference> <Newton's Laws> ”*

A.2.2 GPT-3 Instruction Format and Prompts

Inputs for the GPT-3 analysis were done in a different manner than T5. This was done to support both individual and group analysis of the text. The prompt and inputs to the algorithm are modeled below with instruction in **purple**, individual format in **orange**, and group format in **blue**. The prompt included four examples in individual format and four in group format.

“**Instructions**: Given the question and answers, annotate the span of the answers. The annotation should be wrapped within <> brackets . Each sentence can have a maximum of 3 annotations. **Question**: Force $P = 10 \text{ N}$ is applied to the block of mass $m = 5 \text{ kg}$ on a horizontal rough surface with $\mu_s = 0.3$ and $\mu_k = 0.25$. If $g = 9.81 \text{ m/s}^2$, what is the force of friction on the block?.

###

Individual

Answer: Because friction is only as much as is needed to keep the box at rest when using static friction.

Annotation : Because friction is only as much <Inference><Newton's Laws>, is needed to keep the box at rest <Inference><Box Movement>, when using static friction <Identification><Static Friction Initial>.

###

....

###

Group

Answer Text

1. **Answer**: Newton's third law states that for every action, there is a reaction. In this case, a force of 14.7 N is required to overcome static friction which it doesn't because there is only a 10 N force acting on it. So there is another 10 N frictional force reacting to the force P .

2. **Answer**: $F_s \text{ max}$ would be 14.715 N ($0.3 * 9.81 * 5$), which is greater than the applied force.

Therefore, the box will remain at rest, and the friction force would be equal to the applied force.

....

Annotations

1.Annotation: Newton's third law states that for every action , there is a reaction . <Inference> <Box Movement> a force of 14.7N <Identification> <Static Friction Initial> is required to overcome static friction which it doesn't because there is only a 10N force acting on it . <Comparison> <Static Friction to Applied Force> So there is another 10N frictional force reacting to the force P <Comparison> <Static Friction to Applied Force>.

*2.Annotation: $F_s \text{ max}$ would be 14.715 N ($0.3 * 9.81 * 5$) <Identification> <Static Friction Initial> which is greater than the applied force . <Comparison> <Static Friction to Applied Force> the box will remain at rest <Inference> <Box Movement> the friction force.*

....