# Using Fine-Tuned LLMs to Grade Homework

Stephan Bohacek and Aishah Aseeri

Department of Electrical and Computer Engineering

bohacek@udel.edu

12th Annual CAE in Cybersecurity Community Symposium

# Outline

- Autograding scenario

- Rubrics

- Fine-tuning

- Demo – Walk Through www.guidance-rubrics.com

# Stephan Bohacek

- Associate Professor of Electrical and Computer Engineering at the University of Delaware
- Teach
    - Introduction to Cybersecurity: ~180 students
    - Introduction to Networking: ~60 students
- Continuous challenges with TAs and graders
    - Slow to grade
    - Incomplete or nonexistent comments
    - Incorrect grading
        - Experienced graders know that easy grading reduces the chances that the student will complain to the instructor

# Automatic Grading and Hint Generation Scenario

- Canvas Quizzes
  - Each student's answer is entered and assessed separately via API
- Submissions are automatically graded, and hints are automatically generated
- Students can resubmit as many times are they want, but grading takes at least 24 hours
  - Ideally, the delay reduces students' ability to "game" the system, but might cause students to lose interest
- If the student believes that the grading is wrong, they include the text "MANUAL REVIEW REQUESTED" and alerts are sent to instructors and graders

UNIVERSITY OF DELAWARE.

# Canvas Quiz



**Question 5** — 20 pts

What is the difference and similarity of MAC and a digital signature?

Edit  View  Insert  Format  Tools  Table

12pt  Paragraph  **B**  *I*  U  A

p

Region for student to enter their answer

Automatically generated grade

10 / 20 pts

**Question 5**

What is the difference and similarity of MAC and a digital signature?

Your Answer:

MAC proves that the message that was sent has not been changed. Digital signature proves that the message is from the person that has the private key. Both MAC and Digital signature encrypts a hash of a message.

Not quite. Try again. Your answer correctly addressed: MACs ensures the integrity of the message Your answer correctly addressed: Digital signature protect the integrity of the message. However, what type of encryption do digital signatures use? Symmetric key or asymmetric key encryption? Do MACs require a key? If so, is it a shared secret key, or asymmetric key? If you believe that the grade is incorrect, please resubmit the same answer but with the text 'MANUAL REVIEW REQUESTED' at the top and it will be reviewed manually. If you are stuck, you can see some hints by setting your answer to "Help" and nothing else.

Automatically generated comments

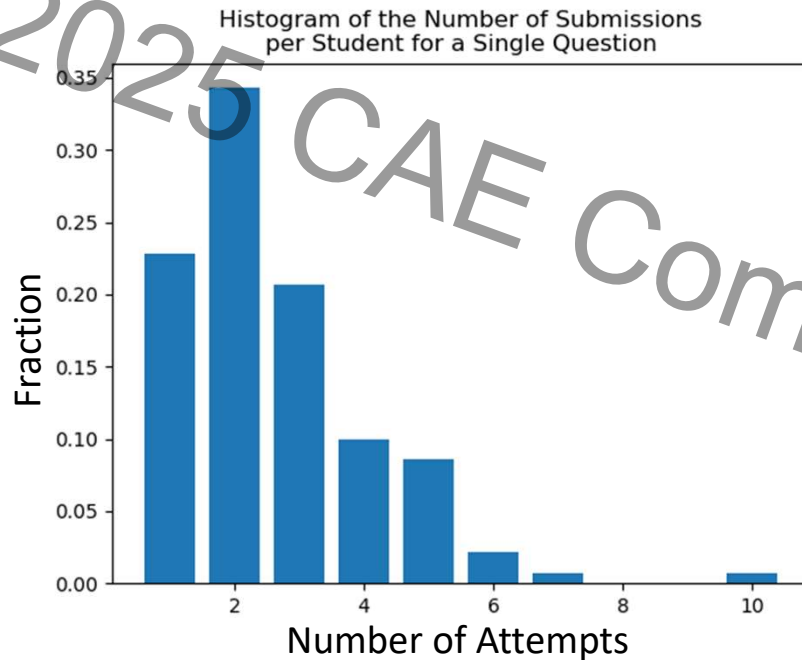UNIVERSITY OF DELAWARE.

# Guidance / Hints

- Rubrics to determine if the submission is correct

- Rubrics to determine if a particular hint is applicable

rubric:
- prompt: "Does the statement mention that George Washington was the United States of America's first president? Statement: <submission>"
- standard hint: Please review chapter 2 of the textbook.
- common errors:
  - prompt: "Does the statement mention that Thomas Jefferson was the United States of America's first president? Statement: <submission>"
  - hint: "While Thomas Jefferson was one of the founders, he was not the first president. Please review Section 2.3 of the textbook"

Targeted hint that includes
pointer to course material

# Resubmissions



Around 25,000 questions were graded in a single class with ~60 students

# How to get the LLM to grade

- Prompt v0
  - Below is a student's answer to the question "Who was the first president of the United States of America?" Grade the student's answer on a scale from 0 to 10. ….
    - This approach often fails.
    - The LLM is not given enough context into
      - What the student should know given what was presented in class
      - What a student at this level should know
      - The language (words and framework, e.g., risk) that the student should use
    - The LLM might not "know" the correct answer or be able to distinguish between a lay-persons answer and an exact answer
    - The LLM might get confused by simple issues such as definitions
    - The LLM might lack any understanding of the topic

UNIVERSITY OF DELAWARE.

# Rubric Components

- Question: Explain the significance of the Louisiana Purchase in American history and its impact on the expansion of the United States.
- Answer components:
    - Territorial Expansion: increased the nation's land area
    - Territorial Expansion: farming, settlement, and natural resource exploration
    - Economic Growth: agricultural potential helped spur economic growth and development
    - Strategic Importance: …
    - Exploration and Manifest Destiny: …
    - Diplomatic Relations: …
- Prompts:
    - Consider the statement below. Does this statement mention that the Louisiana Purchase greatly expanded the United States' land area? Statement: <SUBMISSION>
    - Consider the statement below. Does this statement mention that the Louisiana Purchase provided more space for farming, settlement, and natural resource exploration? Statement: <SUBMISSION>
    - Consider the statement below. Does this statement mention that the Louisiana Purchase greatly expanded agricultural potential and therefore helped spur economic growth and development? Statement: <SUBMISSION>
    - …

# Rubric Component Details

- Prompts should not be too specific
  - Prompt: Consider the statement below. Does this statement mention that the Louisiana Purchase greatly expanded agricultural potential and therefore helped spur economic growth and development. Statement: <SUBMISSION>
  - Does both economic growth and development need to be discussed? The prompt indicates that they both must be discussed
    - LLMs **can** be very literal
- Use Boolean expressions
  - OR
    - Consider the statement below. Does this statement mention that the Louisiana Purchase greatly expanded agricultural potential and therefore helped spur economic growth. Statement: <SUBMISSION>
    - Consider the statement below. Does this statement mention that the Louisiana Purchase greatly expanded agricultural potential and therefore helped spur economic development. Statement: <SUBMISSION>

# More Specific Prompts Tend to Give Better Results

- Instead of
  - OR
    - Consider the statement below. Does this statement mention that the Louisiana Purchase greatly expanded agricultural potential and therefore helped spur economic growth. Statement: <SUBMISSION>
    - Consider the statement below. Does this statement mention that the Louisiana Purchase greatly expanded agricultural potential and therefore helped spur economic development. Statement: <SUBMISSION>
- Use
  - AND
    - Consider the statement below. Does this statement mention that the Louisiana Purchase greatly expanded agricultural potential.  Statement: <SUBMISSION>
    - OR
      - Consider the statement below. Does this statement mention that the Louisiana Purchase helped spur economic growth. Statement: <SUBMISSION>
      - Consider the statement below. Does this statement mention that the Louisiana Purchase helped spur economic development. Statement: <SUBMISSION>

UNIVERSITY OF DELAWARE.

# Prompt Engineering is not enough – e.g., Diagram Grading



Reply: 'In the image, there are four lines or arrows between the client and the server named www.cnn.com . There are two lines that start at the client and end at www.cnn.com , and two lines that start at www.cnn.com and end at the client.'

Wrong: There are 8 lines. 4 start at the client and 4 start at www.cnn.com
Even with many different prompts, I was unable to get consistent interpretation of the diagram.

# Fine-tuning LLMs

**Huge Compute Cluster**

**General LLM**

**Huge Training Set**

**Small Compute Cluster**

**Specialized LLM**

**Small, Specialized Training Set**

# Fine-Tuning LLMs

- Leverages the huge expense spent to build general-purpose LLM
- With a high-quality fine-tuning dataset, the resulting LLM is well suited to answer questions aligned with the dataset's content.
- Most suitable for "yes" "no" types of questions
  - The training set provides inputs and desired outputs.
    - e.g., Input: "Is NYC the largest city in the world" Output: "No"
      - If the LLM answers "yes", then the weights are adjusted
    - e.g., Input: "Explain the taste of an orange." Output: "Oranges taste tangy"
      - If the LLM answers "Oranges have a tangy flavor", then the weights are adjusted, even though the answer is correct.
  - Consequently, fine-tuning is poorly suited for chain-of-thought prompts.
    - Without fine-tuning, chain-of-thought often results in the best performance.

UNIVERSITY OF DELAWARE

# Step to Fine-Tune for Autograding

www.guidance-rubrics.com provide a user interface to perform these steps

- Label submissions
- Select rubric elements to include in training set
  - Elements should be similar, e.g., graphical analysis of a particular type of plot
  - Fine-tuning for a single element is acceptable
- Build training set from one or more classes/sections
- Training LLM
- Make an updated rubric that uses the new LLM
- Compare results

# Work Flow

1. Make Rubric
2. Run grader
3. Label submissions
4. Explore errors
5. Fix errors
   1. Correct labels
   2. Correct prompts
   3. Fine-tune LLM
6. Go to 4

# Building Training Set (1)

# Building Training Set (2)

# Building Training Set (3)

# Update Rubric to Use Fine-tuned Model



Next, save and deploy new rubric

# Compare Rubrics

A diagram that shows communication between a client and two web servers, namely www.cnn.com and www.google.com.

**Prompt:** Does the image show exactly on...element id: ebcd38f0-fe9c-4ffe-828d-a35cda84f505　　∧

**Initial Prompt:** Does the image show exactly one a sequence diagram? Moreover, does this sequence diagram show the communication between a client and two servers, namely a server named www.cnn.com and a server named www.google.com and no other servers? Note that the servers must be labeled exactly as www.cnn.com and www.google.com. For example, labeling the servers cnn and google is incorrect and labeling a servers something like www.cnn.com/pic.jpg is also incorrect. Also, no other servers should be included in the diagram. On the other hand, the case of the text is not important, that is, WWW.CNN.COM and WWW.GOOGLE.COM is acceptable. Note that by sequence diagram, we mean that in this case the sequence diagram should contain three vertical lines. One vertical line should be labeled "www.cnn.com," one vertical line should be labeled "www.google.com", and the one vertical line can be named "client," however, labeling a vertical line is optional. Moreover, there should be diagonal lines between these vertical lines.

**Follow-up Prompt:** I don't understand. Please answer either "Yes, the image shows sequence diagram as described" or answer "No, the image does not show the sequence diagram as described.".

**number graded correctly:** 101
**number graded incorrectly:** 6

| metric | Version 21 | Version 39 |
|---|---|---|
| Number Correctly Graded | 76 | 101 |
| Number Incorrectly Graded | 31 | 6 |

# Examine Incorrectly Graded Submissions



This was mislabeled.
After fixing labeling, repeating the fine-tuning process might give better results.

# Examine Incorrectly Graded Submissions

2025 CAE Community Symposium

Guidance Rubrics

- Help
- Rubric Editor
- Student Submissions
- Create Submissions
- Manually Graded
- Student Progress
- Manage Questions
- Import Rubrics
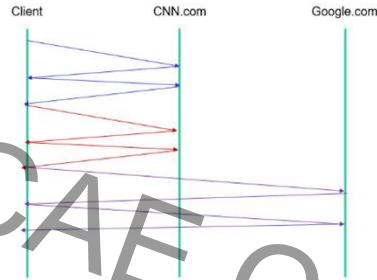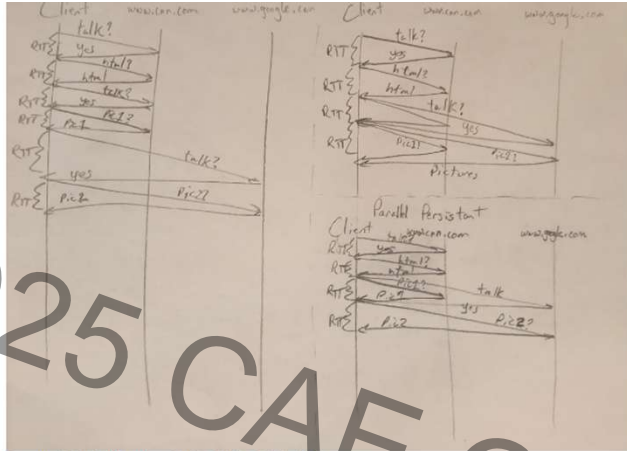- Manage Graders
- Status and Alerts
- Settings
- Fine-Tune
- Compare Rubrics
- Sign Out
- Sign In

Executed Rubric (Grades: 100%, 100%, 100%, 0%)

A diagram that shows communication between a client and two web servers, namely www.cnn.com and www.google.com. Grade:100%

Or (Grade: 1)

Grade: 1

Image Analysis: Does the image show exactly one a sequence diagram? Moreover, does this sequence diagram show the communication between a client and two servers, namely a server named www.cnn.com and a server named www.google.com and no other servers? Note that the servers must be labeled exactly as www.cnn.com and www.google.com. For example, labeling the servers cnn and google is incorrect and labeling a servers something like www.cnn.com/pic.jpg is also incorrect. Also, no other servers should be included in the diagram. On the other hand, the case of the text is not important, that is, WWW.CNN.COM and WWW.GOOGLE.COM is acceptable. Note that by sequence diagram, we mean that in this case the sequence diagram should contain three vertical lines. One vertical line should be labeled "www.cnn.com," one vertical line should be labeled "www.google.com", and the one vertical line can be named "client," however, labeling a vertical line is optional. Moreover, there should be diagonal lines between these vertical lines.

Desired reply: No, the image does not show the sequence diagram as described.

model: ft:gpt-4o-2024-08-06:personal:3429dadc3975450ab4deafa68293aad4-set3-v1:B0wFbJwl

Chat Messages

> user

v assistant

The image shows three sequence diagrams. Each sequence diagram illustrates communication between a client and two servers labeled as "www.cnn.com" and "www.google.com." There are no other servers included in the diagrams.
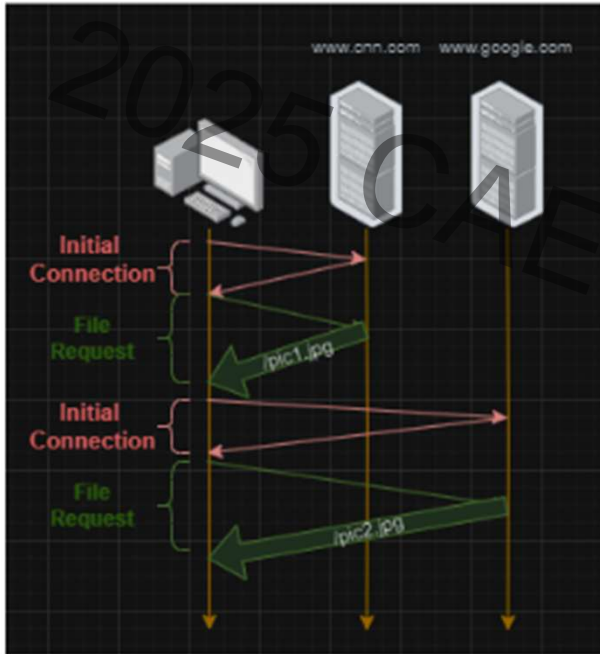
> user

v assistant

Yes, the image shows a sequence diagram as described.

> Common Errors

UNIVERSITY OF DELAWARE

# Examine Incorrectly Graded Submissions



## assistant

The image shows a sequence diagram with three vertical lines and diagonal lines between them, indicating communication. However, the servers are labeled "www.cnn.com" and "www.google.com," but these labels are not exact as required by the instructions. Thus, the diagram does not meet the criteria of having the vertical lines labeled exactly as "www.cnn.com" and "www.google.com."

Might be solvable with a different prompt

# Conclusions

- LLM enable autograding of open-ended questions

- Detailed rubrics are needed

- In some case, iteration is required to improve rubrics

- Iteration can include fine-tuning, prompt engineering, and rubric design

- Once rubrics are perfected, autograding is extremely useful

UNIVERSITY OF DELAWARE.