

FINAL PROJECT REPORT: DEVELOPING AN AI-ASSISTED GRADING SYSTEM
USING LARGE LANGUAGE MODELS

by

Andrei Modiga

A FINAL PROJECT REPORT

Presented to the Faculty of
The School of Computing at the Southern Adventist University
In Partial Fulfilment of Requirements
For the Degree of Master of Science

Major: Computer Science

Under the Supervision of Scot Anderson, Ph.D.

Collegedale, Tennessee

August, 2025

FINAL PROJECT REPORT: DEVELOPING AN AI-ASSISTED GRADING SYSTEM USING LARGE LANGUAGE MODELS

Andrei Modiga, M.S.

Southern Adventist University, 2025

Adviser: Scot Anderson, Ph.D.

We present a grading system that accelerates evaluation of open-ended student work across scanned and digital workflows. The system crops answer regions from PDFs, assigns submissions via OCR on identity regions only, and groups answers by visual semantics using a vision LLM. Instructors review and edit groups, apply rubric items once per group, and export grades from an on-screen table. The solution integrates Ghostscript rasterization, PdfPig page orchestration, SkiaSharp region extraction, Tesseract identity OCR, and GPT-4o Vision for grouping [1, 2, 3, 4, 5]. We detail the architecture, token-budgeted batching strategy, and persistence design, then describe testing results for grouping quality, time-on-task, and usability. The approach avoids brittle handwriting OCR while preserving instructor control, fairness, and auditability.

Contents

Contents	v
List of Figures	ix
List of Tables	xi
1 Introduction and Motivation	1
1.1 Problem Statement	1
1.2 Specific Project Goals/Requirements	1
1.3 Motivation and Benefits	2
1.4 Contributions	3
1.5 Assumptions and Scope	3
1.6 Report Organization	4
2 Background and Context	5
2.1 AI in Education	5
2.2 Automated Grading of Short Answers and Essays	5
2.3 LLMs and GPT-4/4o for Assessment	6
2.4 Bias, Fairness, and Student Perceptions	6
2.5 Similar Implementations	6

2.5.1	Commercial paper-exam graders (e.g., Gradescope, Crowdmark)	7
2.5.2	OMR/MCQ scanning (e.g., Akindi, ZipGrade)	7
2.5.3	LMS graders (e.g., Canvas SpeedGrader)	7
2.5.4	Autograding/algorithmic assessment (e.g., PrairieLearn, Möbius, CodeRunner/CodeGrade)	8
2.5.5	Positioning	8
2.6	Research Conclusions	8
3	Project Solution and Approach	9
3.1	Overview	9
3.2	High-Level Architecture	9
3.3	Sequence of Operations	11
3.4	Instructor-Facing UI Snapshots	12
3.5	Region Extraction and File Lifecycle	16
3.6	Identity Assignment	17
3.7	Grouping Heuristics	17
3.8	Data Model	17
3.9	Token Budget, Cost & Rate Limiting	19
3.10	Integration with ASP.NET	19
3.11	Service Isolation & Network Security	19
3.12	Student-Facing Assignment UI	20
3.13	Rubrics and Grading UX	20
3.14	Security & Privacy	21
3.15	Threat Model & Data Handling	21
3.16	Limitations & Risks	24
4	Testing and Evidence	25

4.1	Objectives	25
4.2	Test Data	25
4.3	Tests & Evidence	26
4.3.1	Basic Flow	26
4.3.2	Vision-Based Semantic Recognition (No Answer OCR)	26
4.3.3	Rubric Propagation & Clipboard Interoperability	27
4.4	UI Verification Checklist (Feature Presence)	27
4.5	UX Design Evaluation (Questionnaire)	31
4.6	Model Grouping Evidence (10-Student Run)	32
5	Results	33
5.1	UI Verification	33
5.2	Vision-Based Semantic Grouping	33
5.2.1	Runtime and Cost	34
5.2.2	Workload Reduction: Instructor Actions and Example	34
5.3	UX Design Outcomes	37
5.4	Clipboard Interoperability	38
5.5	Summary	39
6	Conclusion	41
6.1	Summary of Problem and Goals	41
6.2	Evaluation Summary	41
6.3	Final Outcomes and Deliverables	42
6.4	Lessons Learned and Future Work	42
A	Configuration and Deployment	43
A.1	Prerequisites	43

A.2 FastAPI Service: <code>.env</code>	43
A.3 Web App: <code>appsettings.json</code>	44
A.4 Ghostscript Location	46
A.5 Tesseract on macOS (dev builds)	47
A.6 Build and Run (Web App)	47
A.7 Operational Notes	48
Bibliography	49

List of Figures

3.1	High-level components and data flow.	10
3.2	Sequence for an auto-grouping job.	13
3.3	Assignment creation form. For free-form, authors enter questions and points only; for filled-form, authors also define identity/answer regions against a template.	14
3.4	Course assignments overview with status indicators and quick actions.	14
3.5	Region cropping widget used in two contexts: (i) instructor verification for filled-form scans and (ii) student free-form “mark your answers” flow.	15
3.6	Auto-grouping page with proposed clusters, edit tools, and rubric-first grading.	16
3.7	JSON Model in abbreviated form.	18
3.8	Student assignment page: layout indicator (filled vs. free-form), PDF upload (free-form only), download of submitted file, and grade display.	21
4.1	End-to-end processing and rubric application.	28
4.2	Vision-led semantic grouping from images only: coherent clusters (left, middle) and outliers (right).	29
4.3	Clipboard-based grade transfer using <i>Copy Table</i>	30

List of Tables

3.1	Key table: GroupingResults.	18
3.2	Abbreviated threat model and mitigations	23
4.1	UI checklist for instructor workflow (presence, not preference).	27
5.1	Clustering using GPT-4o vs. GPT-4o-mini.	35
5.2	Instructor manual moves by question (20 questions \times 10 students; lower is better).	35
5.3	Moves statistics per answer.	36
5.4	Per-item medians and IQR (1–5; NA excluded). Item labels abridge the list in §4.	37

Chapter 1

Introduction and Motivation

1.1 Problem Statement

Grading open-ended student work (handwritten or typed) is time-consuming, repetitive, and error-prone under deadline pressure. In large classes, feedback latency diminishes learning value. Typical bottlenecks include: (i) organizing mixed-format submissions (bulk scans vs. individual PDFs), (ii) locating answer regions for consistent review, and (iii) repeatedly applying identical rubric deductions to similar mistakes. Handwriting OCR is brittle, often forcing manual review even for simple cases.

1.2 Specific Project Goals/Requirements

This project delivers a practical, instructor-in-the-loop grading system that provides the several features.

Bulk Scanned *Filled-form* PDFs are split by known page counts. Ghostscript rasterizes pages; PdfPig validates page counts; SkiaSharp crops defined regions

to PNGs. We call these regions *crops*. Crops are used in several different ways throughout this project. The Tesseract OCR library provides the ability to extract name or ID from a designated identity region which the system uses to auto-assign submissions to students. Unresolved items fall back to a manual pick-list. No OCR is performed on answers; grouping uses GPT-4o Vision directly on the cropped images. Identity text crops are matched to the class roster.

Students directly upload free-form assignments which are automatically tied to the uploader’s account and therefore need no identity extraction. In these assignments, students identify crops that contain their answers. But beyond that, submissions are handled similarly to filled-form submissions.

In both cases, the *Editable AI Assistance* proposes semantic groups using a Vision LLM. Instructors can merge/split groups, move answers, and apply rubric items per-group.

The final feature is *Traceability*: All actions and groupings are persisted with timestamps for auditability. Grades are summarized in an on-screen table for review and/or copy.

1.3 Motivation and Benefits

The motivation for adoption by instructors is time-saving. If instructors do not see significant time savings, they are unlikely to adopt a new technology regardless of the other benefits. Nevertheless, additional benefits often factor into decisions once teachers realize time-saving benefits. These benefits include:

- *Faster feedback*: Instructors grade clusters of similar answers once, reducing turnaround time.
- *Consistency*: Per-group rubric application reduces drift across similar answers

and sessions.

- *Lower cognitive load*: The system automates extraction, grouping suggestions, and grade totals; instructors focus on assessment.
- *Reduced brittleness*: Avoiding handwriting OCR on answers eliminates a major failure mode.
- *Privacy-aware*: Only identity crops contain PII; answer crops are devoid of names or IDs.

1.4 Contributions

To achieve these benefits, we provide the following contributions:

- An *OCR-minimal*, vision-first pipeline that uses OCR only for identity assignment.
- A *token-budgeted batching* strategy (downscale + tiling) to bound latency/cost at class scale.
- A *traceable data model and APIs* with idempotent re-uploads and stable crop filenames.
- An *instructor-in-the-loop* UX with editable groups, explicit review status, and rubric-first grading.

1.5 Assumptions and Scope

We target short-answer problems with recognizable visual structure (boxes/lines). Free-form essays are supported via uploaded PDFs but are not auto-scored; the system focuses on grouping to speed human grading. We assume class rosters are available, that instructors can define crops once per assignment for filled-form

assignments, and that students can define crops for free-form assignments.

The scope includes the AI grouping assistance, and workflow to grade both filled-form and free-form assignments.

1.6 Report Organization

[Chapter 2](#) reviews related work and context. [Chapter 3](#) details the system, [Chapter 4](#) presents the evaluation plan, [Chapter 5](#) reports results, and the final chapter concludes with future work.

Chapter 2

Background and Context

2.1 AI in Education

AI has long promised efficiency gains and personalization in education, from adaptive tutoring to analytics that help instructors intervene earlier. Reviews highlight benefits such as individualized practice, faster feedback, and administrative automation when deployed with appropriate oversight [6, 7, 8, 9]. Ensuring teachers can review and revise AI-generated grades—and retain final say to confirm fairness—is essential for trust and real learning gains.

2.2 Automated Grading of Short Answers and Essays

Pre-LLM systems typically relied on feature engineering, keyword overlap, or supervised models trained on labeled answers. These reduce load but struggle with paraphrase and reasoning variance [10, 11]. Clustering similar answers to grade in batches is a recurring theme: once clusters form, instructors can assign rubrics at the group level.

2.3 LLMs and GPT-4/4o for Assessment

Recent work investigates LLMs for grading and feedback across STEM and writing tasks. Studies report promising alignment with human graders for mathematical reasoning and physics solutions when prompts focus evaluation criteria and preserve human oversight [12, 13, 14]. LLMs can also aid instructional design and rubric drafting [15]. We leverage GPT-4o Vision for grouping by meaning from images, bypassing handwriting OCR.

2.4 Bias, Fairness, and Student Perceptions

Bias can propagate into grading unless monitored and mitigated [16]. Student acceptance depends on clear processes and fairness [17]. Effective formative feedback principles—timely, specific, actionable—remain central whether drafted by AI or humans [18].

2.5 Similar Implementations

To situate this project, we survey adjacent tools and how our system compares. In short, we have not found public documentation of a production system that groups handwritten short answers directly from images using a general-purpose vision LLM without first using OCR on the answer content. Existing offerings fall into four families:

2.5.1 Commercial paper-exam graders (e.g., Gradescope, Crowdmark)

Gradescope [19] supports fixed-template paper exams with region-based workflows and “Answer Groups” that let instructors grade clusters of similar responses at once. However, the grouping method is not publicly documented and is presented at a high level as similarity-based. Crowdmark [20] provides strong scanning workflows (QR-coded booklets, automated student matching via OCR on cover pages) and can auto-grade multiple choice, but does not claim semantic grouping of open-ended answers. *Similarity*: our work also supports fixed templates, grouping, and rubrics-first grading. *Difference*: we group from *images only* (no handwriting OCR of answers), use a token-budgeted vision-LLM pipeline to bound cost/latency, archive prompts, and model versions for auditability.

2.5.2 OMR/MCQ scanning (e.g., Akindi, ZipGrade)

Akindi [21] and ZipGrade [22] excel at high-throughput multiple-choice grading from bubble sheets (including mobile scanning) and logistics like sheet sorting. *Similarity*: we likewise handle identity intake for large cohorts. *Difference*: OMR tools target selected-response scoring, not clustering of free-form handwritten work.

2.5.3 LMS graders (e.g., Canvas SpeedGrader)

Platform-native graders such as Canvas SpeedGrader [23] offer annotation and rubric workflows for uploaded files. *Similarity*: we present rubric-based grading and feedback at scale. *Difference*: LMS graders do not automatically cluster semantically similar answers for batch grading.

2.5.4 Autograding/algorithmic assessment (e.g., PrairieLearn, Möbius, CodeRunner/CodeGrade)

Systems like PrairieLearn [24], Möbius [25], CodeRunner [26], and CodeGrade [27] autograde parameterized or code questions (randomized variants, unit tests, CAS checks) with excellent coverage in constrained domains. *Similarity*: automation reduces repetitive grader effort. *Difference*: their strength is *automatic scoring* of structured responses; they do not aim to *group* heterogeneous, handwritten short answers for a human-in-the-loop rubric pass.

2.5.5 Positioning

Our system is closest in spirit to the “answer grouping” idea in commercial paper-graders, but our distinctives are:

1. *image-only* grouping of handwritten content to avoid brittle OCR,
2. a *cost/latency-bounded* vision-LLM pipeline (downscale + tiling + batching),
and
3. an explicitly *auditable, instructor-in-the-loop* workflow (merge/split/move, neutral labels, review status), integrated end-to-end with our site.

2.6 Research Conclusions

(1) Group-based grading is an effective accelerator; (2) LLMs help when auditable and editable; (3) OCR of handwriting is fragile—visual grouping bypasses failure modes; (4) fairness and oversight practices must be designed-in.

Chapter 3

Project Solution and Approach

3.1 Overview

The system comprises an ASP.NET Razor Pages app (instructor workflow), a Python FastAPI microservice (AI grouping), a MySQL database (persistence), and a file store (crops/exports). Tools include Ghostscript (rasterization), Pdfig (PDF orchestration), SkiaSharp (region extraction), and Tesseract (identity OCR). GPT-4o Vision provides semantic grouping over cropped answer images.

3.2 High-Level Architecture

At a glance, the platform is organized into three dashed *zones* that separate concerns and scaling boundaries (Figure 3.1): the *Application Server* (Razor Pages web app and extraction worker), the *AI Service Layer* (FastAPI microservice and the GPT-4o Vision endpoint), and *Data & Storage* (MySQL and the file store). Instructors interact via a browser over HTTPS with the Razor Pages app, which handles authentication, rubric management, and the grading UI. The app invokes an

extraction worker that orchestrates Ghostscript, PdFPig, and SkiaSharp to render pages and cut out per-answer crops; Tesseract OCR is used narrowly to read identifying fields (e.g., student ID) and is deliberately decoupled from semantic grouping. Crops and downstream exports are written to the file store while the web app serves these images directly to the UI.

Auto-grouping requests are queued from the web app to a Python FastAPI service (POST /autogroup). The service packages each question's crops, converts them to compact JPEGs, and sends them with prompt instructions to GPT-4o Vision. Proposed clusters are normalized server-side (stable UUIDs, neutral names, small-group collapse) and persisted to MySQL. The UI polls a job status endpoint and, when complete, renders groups for human review, edits, and grading. This arrangement keeps the web tier largely stateless, allows the extraction worker and AI service to scale independently, and localizes persistent state to MySQL and the file store. Figure 3.1 depicts the overall topology and data flow.

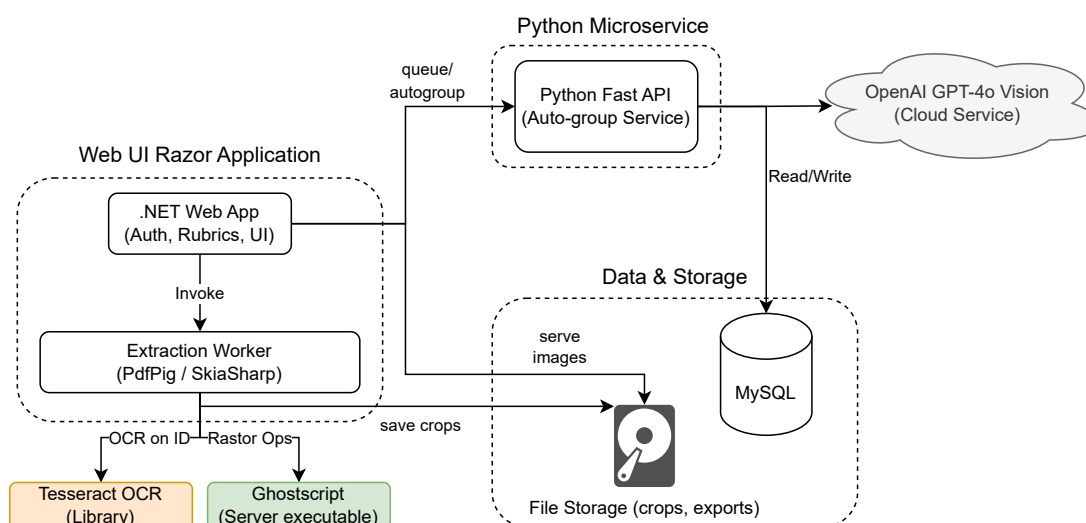


Figure 3.1: High-level components and data flow.

3.3 Sequence of Operations

An instructor initiates grouping from the web UI by selecting a question and confirming its configuration (maximum points, rubric, and the set of per-answer image paths). The browser submits this intent to the web app, which records a new job and enqueues a request to the FastAPI microservice at `/autogroup`. In parallel or beforehand (depending on question state), the extraction worker renders the relevant PDF pages with Ghostscript, enumerates answer regions via PdfPig, and uses SkiaSharp to crop each region to PNG. Lightweight OCR with Tesseract runs only on identifying fields (e.g., cover-page name/ID boxes) to support later reconciliation; the content of answers themselves is not OCR'd for grouping. All crops are written to the file store under stable, human-inspectable paths, and the web app exposes read-only URLs so the UI can preview exactly what will be grouped.

Upon receiving an `/autogroup` task, the FastAPI service prepares the model payload. Each PNG crop is converted to JPEG at a target quality of 50 to reduce bandwidth and context size while preserving legibility for short answers. The service computes a token budget using a 50% downscale heuristic that caps the number of pixels sent per image; if the source crop exceeds that budget, it is downsampled while maintaining aspect ratio. The prompt attaches each image with `detail:auto` and instructs the model to propose clusters of semantically similar answers. The prompt also requests that low-confidence or outlier responses be flagged for an *Ungrouped* bucket.

GPT-4o Vision returns candidate clusters that the service further *shapes* before persistence. Each cluster receives a stable UUID so subsequent UI edits (rename, merge, split) can be tracked independently of order. Neutral descriptions are

standardized (e.g., using a canonical exemplar from the cluster) and very small clusters are folded into *Ungrouped* based on a configurable minimum size. Optionally, if embeddings are available, near-duplicate microclusters are merged with a conservative similarity threshold to avoid over-fragmentation. The shaped result—including cluster membership, labels, and an audit of any collapsed/ungrouped items—is written to MySQL as one row per question with a foreign key to the job.

While the job runs, the UI polls `/status/{job_id}` with backoff to avoid excessive load. When the job is `Complete`, the browser fetches the grouped payload and renders cluster cards with thumbnails backed by the file store. Instructors may rename clusters, reassign individual answers, or merge/split clusters as needed; those edits are persisted incrementally. When grading begins, the UI ties rubric criteria and maximum points to each cluster, enabling a single grading action to fan out to all member answers. Final scores are written through to the `GroupingResults` and `GroupScore` tables. The complete message choreography for this workflow is shown in the sequence diagram in Figure 3.2.

3.4 Instructor-Facing UI Snapshots

This section highlights the key instructor workflows with inline references to the corresponding UI figures: assignment creation (Figure 3.3), the course assignments overview (Figure 3.4), the shared region-cropping tool (Figure 3.5), and the auto-grouping interface (Figure 3.6).

Assignment creation. Instructors configure the submission layout (*filled-form* vs. *free-form*), directions, and dates. As shown in Figure 3.3, the authoring adapts to the chosen layout:

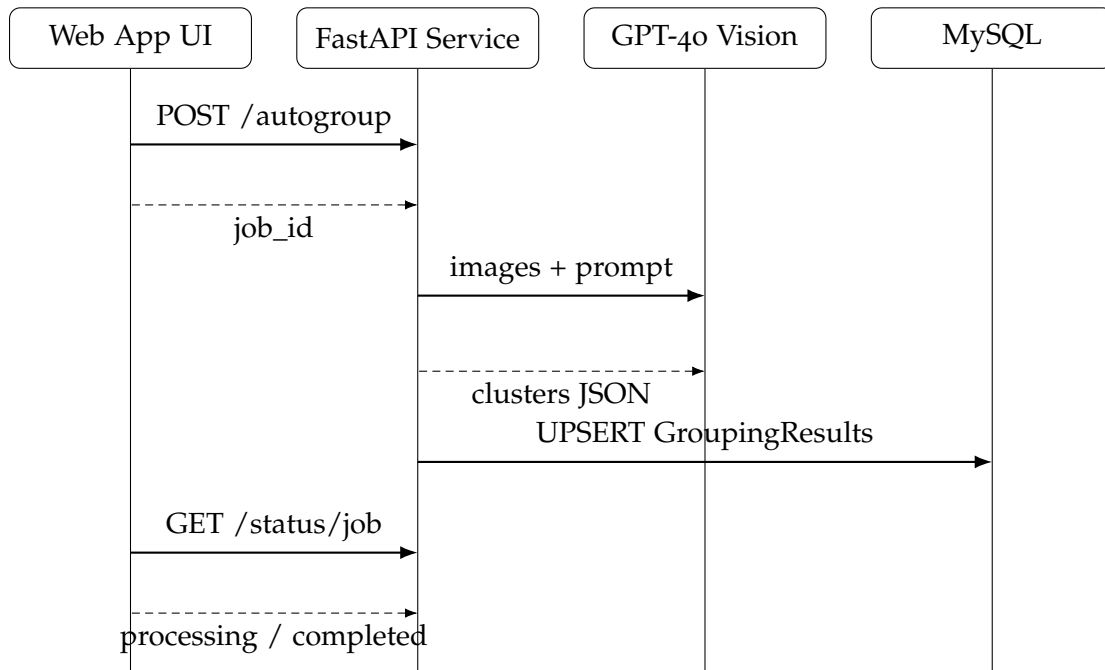


Figure 3.2: Sequence for an auto-grouping job.

- *Free-form*: the instructor enters *only* the question labels (e.g., Q1, Q2a) and max points per question. No region editor is shown here because students will upload a PDF and *define their own answer regions* in the next step (see also the cropping widget in Figure 3.5).
- *Filled-form*: in addition to questions/points, the instructor binds an exam template and draws the *identity* and *answer* regions once. Those regions are then used to crop all submissions automatically (the same widget used here is illustrated in Figure 3.5).

Course assignments page. The course view summarizes assignment state (open/-closed, submissions, grading progress) and links to grouping/grading. Figure 3.4 shows status indicators and quick actions that guide instructors into grouping (Figure 3.6) or back to the region editor (Figure 3.5) if setup needs adjustment.



The form is titled "Add Assignment" and contains the following fields and options:

- Assignment Name:** A text input field.
- Directions:** A large text area for instructions.
- Allow Late Submission:** A checkbox.
- Due Date:** A date and time picker showing "08/10/2025, 12:30 PM".
- Assignment Type:** A dropdown menu with "Auto grouping" selected.
- Layout:** A dropdown menu with "Filled-Form (scanned)" selected, "Free-Form (typed)", and "Blank Form PDF" as options.
- Blank Form PDF:** A file upload section with a "Choose File" button and "no file selected" text.
- Create Assignment:** A blue button to submit the form.
- Student List:** A link at the bottom to view the student list.

Figure 3.3: Assignment creation form. For free-form, authors enter questions and points only; for filled-form, authors also define identity/answer regions against a template.

CPTR-101
 Teacher: [REDACTED]
 Enrollment Code: BTHAN

[Add Assignment](#) [Student List](#)

Assignment Name	Assignment Type	Allow Late	Due Date	Cut-off Date	Actions
filled form	Auto grouping	No	2025-08-20 04:07 PM	N/A	See Groupings Edit Delete Grades Upload Scans
free form	Auto grouping	Yes	2025-08-19 07:43 PM	2025-08-27 07:43 PM	Auto Group Edit Delete Grades

Figure 3.4: Course assignments overview with status indicators and quick actions.

Region extraction (shared tool). The same cropping widget is used to verify identity/answer regions for filled-form scans and, in free-form flows, to let students mark their own answer areas. Figure 3.5 depicts the shared tool; the resulting crops flow directly into the grouping experience shown in Figure 3.6.

Home / Teacher Tools / Course Management / Define Regions

HW09 Chapter 3.5 Problems: 26, 27, 32 (Show your work!) 31pts

NOTE: Please keep the answers on the same page as the questions, this will help both you and I when you upload it to gradescope.

26. (4) Consider transferring an enormous file of L bytes from host A to host B. Assume an MSS of 536 bytes.

a. (2) What is the maximum value of L such that TCP sequence numbers are not exhausted (recall that TCP sequence number fields has 4 bytes)?

b. (2) For the L you obtain in (a), find how long it takes to transmit the file. Assume that a total of 66 bytes of transport, network, and data-link header are added to each segment before the resulting packet is sent out over a 155 Mbps link. Ignore flow control and congestion control so A can pump out the segments back to back and continuously.

Name: EXAM 09

Page 1 of 3

Identity 0

Identity

Question 26a 2

26

a

Answer

Question 26b 2

26

b

Answer

Save All

Figure 3.5: Region cropping widget used in two contexts: (i) instructor verification for filled-form scans and (ii) student free-form “mark your answers” flow.

Auto-grouping UI. After crops are generated, the grouping page proposes semantic clusters; instructors can merge/split groups, move items, and apply rubric items per group. Figure 3.6 shows the proposed clusters and edit tools; rubric-first grading executed here propagates scores to all members of a cluster.

Assignment: filled form

Select Student:
All Students

9/9 ready

Question 26a ●●● Question 26b ●●● Question 27a ●●● Question 27b ●●● Question 27c ●●● Question 27d ●●● Question 32a ●●●
Question 32b ●●● Question 32c ●●●

Question 26a (Max 2)

New Group (N) Assign to Group... Assign

Grouped Answers

Group: Group 1 ✓ AI Delete Group

Group Points: 2.0

Rubrics:

32 bits $\rightarrow 2^{32} = 4\text{ Gb}$

b. (2) For the L you obtain in (a), find how long it takes to transmit the file. Assume that a total of 66 bytes of transport, overhead, and data-link header are added to each segment before the smaller packet is sent out over a 150-Mbps link

Group: Group 1

Select

4 x 6.5 bits = 32 bits
2³² byte = 4.2949673 bytes = ~4GB

b. (2) For the L you obtain in (a), find how long it takes to transmit the file. Assume that a total of 66 bytes of transport, overhead, and data-link header are added to each segment before the smaller packet is sent out over a 150-Mbps link

Group: Group 1

Select

Group: Group 2 ✓ Manual Delete Group

Group Points: 1.0

Rubric

missing units (-1) ✎ ✖

Add New Rubric

Scheme:

Negative ▾

☐ Allow score < 0

☐ Allow score > max

Figure 3.6: Auto-grouping page with proposed clusters, edit tools, and rubric-first grading.

3.5 Region Extraction and File Lifecycle

For *filled-form* assignments, instructors define identity and answer regions once per assignment using the cropping widget (Figure 3.5); the worker then applies those regions to every scanned booklet and enforces stable filenames (e.g., Q27a.png) for reproducibility and idempotent re-uploads. For *free-form* assignments, students upload a PDF and mark their own answer regions in the same tool (Figure 3.5); the saved boxes are used to generate per-question crops that populate the grouping interface (Figure 3.6) for review and rubric-first grading. Debug images are suppressed outside development builds. Re-uploads replace prior crops and metadata to avoid stale data, and the updated crops reappear in the grouping view (Figure 3.6) so that any instructor edits remain aligned with the latest artifacts.

3.6 Identity Assignment

Filled-form batches use Tesseract [4] to extract roster identifiers from the pre-defined *identity* region (see Figure 3.5). Low-confidence or malformed results are flagged for manual resolution with a roster pick-list and a thumbnail preview. Free-form uploads are tied to the uploader’s account; therefore no OCR is needed. Identity OCR is used only for roster linkage—semantic grouping operates directly on cropped answer images (Section 3, Figures 3.1 and 3.2).

3.7 Grouping Heuristics

We instruct the model to produce *fewer, larger clusters*, to route unreadable/singletons into *Ungrouped*, and to emit neutral descriptions (“Group 1”, “Group 2”). Tiny groups under a threshold are collapsed into *Ungrouped*. Groups are sequentially re-numbered for clarity. All prompts and model versions are archived with the `job_id`.

3.8 Data Model

Auto-grouping results are persisted as a *single row per question*, providing a de-normalized snapshot that the UI can load quickly after the workflow in Figure 3.2 and during review in Figure 3.6. As summarized in Table 3.1, the JSON payload `GroupData` holds the proposed/edited clusters (neutral labels, membership, per-group points/flags), while `ScoringScheme` and the boolean guards (`AllowBelowZero`, `AllowAboveMax`) control rubric arithmetic. `CreatedAt/UpdatedAt` provide auditability. In practice, there is at most one record per (`AssignmentId`, `AssignmentQuestionId`); edits typically mutate `GroupData` and bump `UpdatedAt`.

Table 3.1: Key table: GroupingResults.

Column	Type	Notes
Id	INT (PK)	Surrogate primary key.
AssignmentId	INT (FK)	Links to the assignment entity.
AssignmentQuestionId	INT (FK)	Equals <code>template_region_id</code> sent by client.
GroupData	JSON	Array of groups with files, description, <code>is_correct</code> , <code>points</code> .
ScoringScheme	VARCHAR(32)	Default "Negative".
AllowBelowZero	TINYINT(1)	Boolean.
AllowAboveMax	TINYINT(1)	Boolean.
CreatedAt	DATETIME	UTC created timestamp.
UpdatedAt	DATETIME NULL	UTC last update; nullable.

Prompting & JSON schema. The service uses concise grouping rules (favor fewer, larger clusters; neutral labels; unreadable/singletons \rightarrow *Ungrouped*). Prompts and model/version are archived with each `job_id`. The model returns JSON shown in Figure 3.7 in abbreviated form.

```
{
  "groups": [
    {
      "group_id": "uuid",
      "files": [
        {
          "file_path": "...",
          "user_id": "...",
          "template_region_id": "...",
          "question_label": "...",
          "content_type": "image/jpeg",
          "user_name": "..."
        }
      ],
      "description": "Group 1",
      "is_correct": true | false | null,
      "points": 5.0
    }
  ]
}
```

Figure 3.7: JSON Model in abbreviated form.

3.9 Token Budget, Cost & Rate Limiting

For an image of width w and height h , the service estimates tokens after a 50% downscale: $w' = \lfloor w/2 \rfloor$, $h' = \lfloor h/2 \rfloor$. The number of 512×512 tiles is $T = \lceil w'/512 \rceil \cdot \lceil h'/512 \rceil$, and the cost estimate is $85 + 170 T$ tokens/image. Batches exceeding a threshold are split. On rate limits, the client retries up to 10 times with exponential backoff. For a representative 768×768 downsampled image ($T = 3$), the estimate is approximately 595 tokens/image.

3.10 Integration with ASP.NET

The web app calls `QueueAutoGroupAsync` (server) to POST to `/autogroup`, records a `GroupingJob`, and renders a progress UI that polls `/status/{job_id}`; when the background job completes, the page automatically refreshes into the results view. Subsequent instructor actions (save groups, apply/remove rubric items, scoring method toggles) are persisted in the relational database and mirrored in a `GroupScore` table; a background grade recalculation keeps submission grades in sync.

3.11 Service Isolation & Network Security

In production, the FastAPI auto-grouping service runs inside the same Docker Compose stack on a private bridge network and is not published to the host (container uses `expose` only—no ports mapping). As a result, `/autogroup` and related endpoints are reachable only from the web app over internal service DNS (e.g., `http://autogroup:8000/...`); the browser never talks to the service directly.

Compose-level network isolation and host firewall rules prevent external ingress to the service container.

3.12 Student-Facing Assignment UI

Students reach an assignment-specific page that adapts to the configured layout:

Filled-form (read-only). Students cannot upload Filled-form assignments. They can download the submitted booklet (when present) and view the grade once posted. Figure 3.8 shows a snippet of the page `Directions` and a simple status panel.

Free-form (student upload). Students upload a single PDF, then are routed to a `DefineAnswerRegions` step to mark answer boxes. When submissions are open (`DueDate/AllowLateWork/CutoffDate` enforced via `CanSubmit()`), they can resubmit; otherwise the page displays a “Resubmissions have closed” notice. *Note:* The region-marking widget for free-form uploads reuses the same cropping interface shown in Figure 3.5; we avoid duplicating the screenshot here.

3.13 Rubrics and Grading UX

While grading a group, instructors select and apply rubric items; zoom/pan is supported where available. Changing a rubric value propagates to all affected answers. The `Question` tab shows review status: each group displays a status badge (“Graded” or “Needs grading”). Instructors can either apply at least one rubric item or, if none are needed, click `Save All Groups for Question` to mark the question as reviewed.

[Home](#) / [Courses](#) / Grouping Assignment

Assignment: filled form

This is a filled-form assignment. You can view your grade and (if available) your submitted PDF.

Assignment Directions

test

Your Submission

[Download Submitted Assignment](#)

Grade

100.00%

Figure 3.8: Student assignment page: layout indicator (filled vs. free-form), PDF upload (free-form only), download of submitted file, and grade display.

3.14 Security & Privacy

We minimize student-data exposure by limiting OCR strictly to identity regions (therefore no OCR is needed on answers), sending only per-answer image crops—without student names—to the vision model for grouping, and confining operations to authenticated instructor actions with role-appropriate permissions. Only course roster identifiers and per-answer images are processed; no plaintext student content is transmitted for grouping. All reads and writes are logged for auditability, and the design follows least-privilege access control consistent with FERPA expectations [28]. To safeguard fairness, we monitor for potential bias by auditing cluster assignments across demographic-neutral cohorts and preserve human authority through full instructor override mechanisms for grouping and grading.

3.15 Threat Model & Data Handling

This section details how identities, images, and grouping results are protected throughout the workflow in Figure 3.2 and during review in Figure 3.6. For a

concise summary, see Table 3.2 at the end of this section.

Roles and access scope. *Students* may upload their own work and view only their own grades; they cannot view other students' submissions, crops, or groupings. *Instructors* (course owners/graders) can access assignments, submissions, identity crops for their course, grouping results, and grading tools; their access is scoped per course via ACLs. *Administrators* handle configuration and support tasks; by policy they do not browse student content unless temporarily granted course-scoped access for incident response.

Unauthorized access. *Risk.* A user without rights could read crops, grouping results, or grades. *Mitigation.* Role-based authorization (student/instructor/admin) plus per-course ACLs are enforced at the server for every endpoint that touches identity crops, answer crops, grouping JSON, or grade records. The UI never embeds direct file paths without an authorization check; URLs are resolved through the web app to ensure policy is applied. Audit logs provide a trail of who saw or changed what, enabling both deterrence and post-hoc review.

Model data exposure. *Risk.* Personally identifiable information (PII) might be sent to a third-party model. *Mitigation.* Only identity regions contain PII and are processed locally to assign roster IDs. Grouping uses per-answer crops that exclude names and form headers; no student plaintext is transmitted for clustering. Prompts avoid including course/roster metadata, and the shaped outputs (neutral group descriptions, UUIDs) intentionally carry no PII.

Prompt injection and model influence. *Risk.* Crafted markings within an answer image attempt to steer the model (e.g., "put all answers into one group"). *Mitigation.*

The service uses fixed, server-side prompts and normalizes model responses before persistence: clusters receive stable UUIDs and neutral labels; tiny or low-confidence clusters are folded into *Ungrouped*. Instructors retain full control to rename, merge, split, or reassign items, ensuring human oversight dominates over model suggestions.

Data retention and scope. *Risk.* Retaining artifacts longer than needed increases exposure. *Mitigation.* Identity crops and derived answer crops follow course-level retention settings; exports are versioned, and re-uploads purge and regenerate crops to avoid drift. Grouping JSON (GroupData) stores neutral labels and membership only; it excludes raw identity text. Retention and deletion actions are logged.

Table 3.2: Abbreviated threat model and mitigations

Threat	Risk	Mitigation
Unauthorized access	Disclosure of student data	Role-based auth; per-course ACLs; object-level checks; audit logs
Model data exposure	PII leakage to third party	PII limited to local identity OCR; answer crops exclude names; neutral outputs
Prompt injection	Manipulated grouping suggestions	Server-side prompts; normalize outputs; full instructor override
Data retention	Oversharing over time	Course retention settings; purge-on-reupload; logged deletions/exports

3.16 Limitations & Risks

- *Generalization*: Prompts tuned on one course may not transfer perfectly to other subjects; mitigated by neutral labels and editable groups.
- *Model drift*: Vision model updates can shift behavior; we pin model/version and archive prompts with run IDs.
- *Edge cases*: Faint pencil, skew, or multiple answers in one crop reduce grouping confidence; flagged to *Ungrouped*.
- *Human factors*: Instructor trust varies; we look to why/where groups changed and keep full override tools.

Chapter 4

Testing and Evidence

4.1 Objectives

We provide concise, visual evidence that the system:

1. processes submissions end-to-end,
2. groups answer images by semantic similarity (no answer OCR¹),
3. isolates outliers for review,
4. applies rubrics and propagates score changes consistently, and
5. supports clipboard-based grade transfer (*Copy Table*) into external systems.

4.2 Test Data

We used a small gold-standard set of 20 pages containing mixed typed and handwritten responses. Known-correct examples provide an interpretation of

¹No answer OCR means that text extracted from the image is not returned in the answer. In this case the answer is a grouping, so text extracted by the AI is not needed.

cluster coherence. The system does not rely on OCR results to group the answers.

4.3 Tests & Evidence

Four tests provide evidence for correctness in the areas of: Basic Flow (objective 1), Vision-based Semantic Recognition (objective 2-3), and Rubric propagation (objective 4). The final objective relies on the format of content copied to the clipboard and relies on a simple check to see if content is copied correctly.

4.3.1 Basic Flow

Basic flow tests check the process of Upload → process → results, to make sure that it matches the requirements for application flow.

Evidence: Figure 4.1 (a) shows a set of papers uploaded and processing. (b) shows the grouping results view, and (c) shows the graded group using the “Rubric Test”. Testing this under several scenarios demonstrates that the process correctly follows the required flow.

4.3.2 Vision-Based Semantic Recognition (No Answer OCR)

This tests if GPT-4o Vision groups answers *images* into semantically similar groups.

Evidence: Figure 4.2 shows three groups created by GPT-4o. In (a) we see identical answers being grouped together as expected. In (b), GPT-4o clustered answers that are semantically the same, with different amounts of work into the same Cluster. (c) shows diverse answers clustered into an “Ungrouped Answers” group as expected.

4.3.3 Rubric Propagation & Clipboard Interoperability

For rubric propagation, totals must update after a rubric edit and the rubric panel must expose a brief help tooltip clarifying positive vs. negative scoring semantics.

Evidence: Figure 4.2 (c) shows the results of a change in a rubric value assigned to group 1. Near the top, it clearly shows the correct points value of 0.0.

For clipboard Interoperability, the on-screen grade table must be able to be copied to the clipboard and pasted into an LMS/spreadsheet.

Figure 4.3 shows the results of a copy and paste operation that is correctly converted into XML. Note that real names have been redacted.

4.4 UI Verification Checklist (Feature Presence)

This checklist is completed once per build to verify that all expected UI elements and flows are present. Evaluators mark [X] or leave blank; notes capture anything missing or confusing.

Table 4.1: UI checklist for instructor workflow (presence, not preference).

Item	OK	Notes
Assignment creation (free-form + filled-form options)	<input type="checkbox"/>	
Template binding and identity/answer region editor (filled-form)	<input type="checkbox"/>	
Student upload flow; processing status/notifications	<input type="checkbox"/>	
Grouping view with clusters + <i>Ungrouped</i> bucket	<input type="checkbox"/>	
Cluster edit tools (rename, merge, split, reassign)	<input type="checkbox"/>	
Rubric panel; apply/remove rubric items per group	<input type="checkbox"/>	
Totals update immediately after rubric edit	<input type="checkbox"/>	
Copyable grade table (<i>Copy Table</i>)	<input type="checkbox"/>	
Pasted table preserves rows/columns in target app	<input type="checkbox"/>	
Audit-friendly identifiers (assignment, question, job id) visible	<input type="checkbox"/>	
Job status is legible during processing (progress, last update, error/hung states)	<input type="checkbox"/>	

Assignment: HW02

Select Student:
All Students

20/20 ready

Question 6a • Question 6b • Question 6c • Question 6d • Question 6e • Question 6f • Question 6g • Question 8a • Question 8b • Question 8c • Question 8d • Question 13a • Question 13b • Question 16 • Question 25a • Question 25b • Question 25c • Question 25d • Question 25e • Question 25f •

(a) Upload & processing

Question 13a (Max 1)

New Group (N) Assign to Group Assign

Ungrouped Answers

Group: Ungrouped AI

Group Points: 1.0

Rubrics:

NLR

Group: Ungrouped

Select

[Answer in bold]

Group: Ungrouped

Select

Rubric

Add New Rubric

Scheme:

Negative

Allow score < 0

Allow score > max

Grouped Answers

Group: Group 1 AI Delete Group

Group Points: 1.0

Rubrics:

Geometric sequence

$d_{\text{arg-queuing}} = (N-1) \frac{L}{2R}$

(b) Results view

Question 6a (Max 1)

New Group (N) Assign to Group Assign

Grouped Answers

Group: Group 1 AI Delete Group

Group Points: 0.0

Rubrics: Rubric Test (-1)

$d_{\text{prop}} = \frac{m}{s}$

Group: Group 1

Select

$d_{\text{prop}} = m / s$

Group: Group 1

Select

$d_{\text{prop}} = m/s$

Group: Group 1

Select

Rubric

Rubric Test (-1)

Add New Rubric

Scheme:

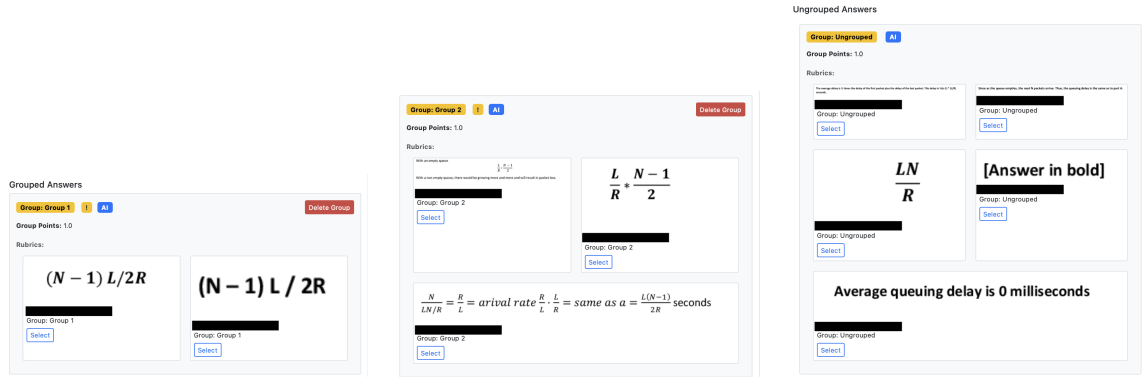
Negative

Allow score < 0

Allow score > max

(c) Rubric update (totals updated)

Figure 4.1: End-to-end processing and rubric application.



(a) Cluster A (same idea) (b) Cluster B (different idea) (c) Outliers (mixed/unclear)

Figure 4.2: Vision-led semantic grouping from images only: coherent clusters (left, middle) and outliers (right).

Grades for CPTR-328: Testing HW02

Email	Score	Feedback
[REDACTED]	95.0	
[REDACTED]	95.0	
[REDACTED]	100.0	
[REDACTED]	100.0	
[REDACTED]	100.0	
[REDACTED]	100.0	
[REDACTED]	100.0	
[REDACTED]	100.0	
[REDACTED]	95.0	
[REDACTED]	100.0	

Copy Grades

(a) Copyable grade table in the app

```

<table class="table">
  <tbody><tr>
    <th>Email</th>
    <th>Score</th>
    <th>Feedback</th>
  </tr>
  <tr>
    <td>[REDACTED]</td>
    <td>95.0</td>
    <td></td>
  </tr>
  <tr>
    <td>[REDACTED]</td>
    <td>95.0</td>
    <td></td>
  </tr>
  <tr>
    <td>[REDACTED]</td>
    <td>100.0</td>
    <td></td>
  </tr>
  <tr>
    <td>[REDACTED]</td>
    <td>100.0</td>
    <td></td>
  </tr>
  <tr>
    <td>[REDACTED]</td>
    <td>100.0</td>
    <td></td>
  </tr>
  <tr>
    <td>[REDACTED]</td>
    <td>100.0</td>
    <td></td>
  </tr>
  <tr>
    <td>[REDACTED]</td>
    <td>95.0</td>
    <td></td>
  </tr>
  <tr>
    <td>[REDACTED]</td>
    <td>100.0</td>
    <td></td>
  </tr>

```

(b) Pasted result in external tool

Figure 4.3: Clipboard-based grade transfer using *Copy Table*.

4.5 UX Design Evaluation (Questionnaire)

Participants include 3–6 instructors/TAs familiar with grading. The scope of the questionnaire covers Design/readability only, including layout, labels, hierarchy, contrast, discoverability — *not* speed or task timing.

We use the common likert scale for agreement: 1 - Strongly Disagree, 2 - Disagree, 3 - Neutral, 4 - Agree and 5 - Strongly Agree.

Likert (1–5) items

1. The layout feels clean and uncluttered.
2. Text is legible without zooming; font sizes are appropriate.
3. Color/contrast makes content easy to read (including low-light).
4. Visual hierarchy makes the primary actions obvious.
5. Labels and terminology are clear and unambiguous.
6. Icons and buttons are self-explanatory or clearly labeled.
7. Whitespace/spacing helps me scan and find things quickly.
8. Clusters are visually distinguishable from each other.
9. Thumbnails are large and crisp enough to compare answers.
10. Edit affordances (rename/merge/split/reassign) are easy to find.
11. The rubric panel is easy to locate and interpret.
12. Score changes and totals are clearly indicated.
13. Empty/error states are understandable and instructive.
14. Loading/progress indicators are clear and non-distracting.
15. Keyboard focus is visible and navigation feels predictable.
16. Color is not the only cue (labels/icons also indicate meaning).
17. Overall, the design feels easy to read and understand.

Free-response prompts. (1) What was hard to read, notice, or understand? (2) One UI element you would improve and how. (3) Accessibility notes: contrast, text size, color cues, keyboard focus, screen readers.

4.6 Model Grouping Evidence (10-Student Run)

Setup. Run one assignment with ~ 10 students. For a single question, create a *gold* clustering by manually assigning each answer crop to a semantic group (visual meaning).

Prediction. Record the model’s shaped output (cluster IDs, memberships, *Un-grouped*).

Teacher-facing evaluation. We quantify effort as the number of *manual moves* required to correct groupings to the gold (one move = reassigning a single tile). We additionally report ungrouped size and the presence of a dominant “correct” cluster when applicable. Formal clustering scores (e.g., B-Cubed, Pairwise F1) are out of scope for this study and omitted to keep the focus on instructor workload.

Chapter 5

Results

5.1 UI Verification

Table 4.1 was completed for the evaluated build. All required UI elements and flows are present (all boxes checked). These were visually verified by three different individuals.

5.2 Vision-Based Semantic Grouping

The system groups *answer crops* by visual meaning rather than OCR. Figure 4.2 shows two coherent clusters and an outlier set; each tile is a raw crop. As you can see, the crops represent not only semantically similar answers, but also the formatting of the answers as well. This has the added benefit of allowing a grader to see answers that are identical.

5.2.1 Runtime and Cost

On a 200-answer batch (20 questions, 10 students), end-to-end grouping/grading completed in approximately 10 minutes (≈ 20 answers/min; ≈ 3.0 s/answer).

Token use and cost (per run).

Model	Answers	Tokens	Cost (USD)	Tokens/answer
GPT-4o	200	112,227	\$0.79	561
GPT-4o-mini (compact)	200	112,220	\$0.03	561
GPT-4o-mini (verbose)	200	2,108,335	\$0.56	10,542

Response verbosity drives token variability for GPT-4o-mini; restricting outputs to IDs/memberships yields compact usage comparable to GPT-4o. As we can easily see, the cost to grade a single assignment of 20 questions for 10 students is insignificant to the cost of having a grader work through these papers even if all it saves is 7 minutes ($\$0.79 \div \$7.25/\text{hour} = .11$ hours or 6.6 minutes).

5.2.2 Workload Reduction: Instructor Actions and Example

We measure instructor effort as the number of *manual moves* required to correct the model’s grouping to a gold clustering (visual meaning). A *move* is a single reassign action of one answer tile. The following vignette is taken from an assignment in Networking #13-a. We analyze this example in terms of workload reduction across different versions of GPT.

Table 5.1 shows a typical question where GPT-4o produces a single coherent correct cluster plus one ungrouped set. However, GPT-4o-mini fragmented the same responses into multiple groups, where the groups differed from multiple

tests on the same sample set. The conclusion is clear, the application cannot rely on the cheaper “mini” version to reliably group answers.

Table 5.1: Clustering using GPT-4o vs. GPT-4o-mini.

Model	#Clusters	#Ungrp.	Largest correct cluster	Manual moves
GPT-4o	1	1	all correct answers	2
GPT-4o-mini	5	1	fragmented across clusters	<i>not reliably tallied (substantially higher)</i>

Table 5.2 shows the number of moves required for each question on the sample assignment. The results were reliably reproduced in subsequent tests.

Table 5.2: Instructor manual moves by question (20 questions \times 10 students; lower is better).

Q	Moves	Q	Moves
6a	4	6b	3
6c	2	6d	1
6e	2	6f	4
6g	2	8a	1
8b	1	8c	4
8d	4	13a	2
13b	5	16	2
25a	1	25b	1
25c	2	25d	2
25e	4	25f	2

Table 5.3 shows results across 200 answers requiring a total of 49 manual moves, i.e., *24.5 moves per 100 answers*. The per-question median was 2 moves with an interquartile range(IQR) of 1–4. This represents a conservative approach to grouping where the majority of moves are from the ungrouped set to another group.

Table 5.3: Moves statistics per answer.

Model	Moves / 100 answers	Median moves / Q	Notes
GPT-4o	24.5	2	49 moves over 200 answers
GPT-4o-mini	–	–	Substantially higher; not reliably tallied

Estimated time savings (conservative). Baseline page/PDF flipping at $t_{\text{page}}=8$ s per answer: $200 \times 8 \text{ s} = 26.7 \text{ min}$. With clustering:

$$T \approx (\text{clusters per batch}) \cdot t_{\text{group}} + (\text{moves}) \cdot t_{\text{move}},$$

using $t_{\text{group}}=2$ s and $t_{\text{move}}=5$ s. Assuming two clusters per question on average (correct + ungrouped) $\Rightarrow \sim 40$ group actions (≈ 80 s) and 49 moves (≈ 245 s), giving $T \approx 5.4$ min, an $\sim 80\%$ reduction vs baseline.

Note on GPT-4o-mini effort. During our audit, *GPT-4o-mini* required substantially more manual moves than we could reliably tally in-session due to fragmentation and frequent reassignments. We therefore report 4o-mini effort qualitatively as *higher* rather than providing a potentially misleading count.

On variability of student answers. Cluster counts naturally vary with how diverse student responses are for a given question (e.g., multiple valid solution paths or many non-answers increase the number of clusters/ungrouped). Even when clusters proliferate, presenting all answers for one question on a single screen still reduces handling time versus grading one student at a time or flipping through pages/PDFs, because instructors can scan, compare, and bulk-apply rubric actions across visually similar answers.

5.3 UX Design Outcomes

Participants included four subjects. One item (Q13) had a missing value; analysis uses available ratings per item.

Central tendency and dispersion: Across all items, ratings were ceilinged at 5: overall mean = 4.79, median = 5.0 ($Q_1 = 5.0$, $Q_3 = 5.0$, IQR = 0).

Per-item medians and IQR: Table 5.4 reports median and IQR per item (1–5), treating missing values as *NA*. Most items show perfect consensus (IQR = 0); localized dispersion appears most notably on Q14.

Table 5.4: Per-item medians and IQR (1–5; *NA* excluded). Item labels abridge the list in §4.

Q	Median	IQR	Item (abridged)
1	5.0	0.0	Layout clean/uncluttered
2	5.0	0.0	Text legible / sizes appropriate
3	5.0	0.0	Color/contrast readable
4	4.5	1.0	Visual hierarchy / primary actions
5	5.0	0.5	Labels/terminology clear
6	5.0	0.0	Icons/buttons clear
7	4.5	1.0	Whitespace/scanability
8	5.0	0.0	Clusters visually distinguishable
9	5.0	0.0	Thumbnails large/crisp
10	4.0	0.5	Edit affordances easy to find
11	5.0	0.5	Rubric panel easy to locate/interpret
12	5.0	0.0	Score changes/totals clear
13	5.0	0.5	Empty/error states understandable
14	4.5	1.5	Loading/progress indicators clear
15	5.0	0.5	Keyboard focus / navigation predictable
16	5.0	0.0	Color not sole cue
17	5.0	0.0	Overall easy to read/understand

Items needing attention. The lowest median was Q10 (*Edit affordances are easy to find*) with $\tilde{x} = 4.0$ (IQR = 0.5). The widest dispersion was Q14 (*Loading/progress*

indicators are clear) with $IQR = 1.5$, despite a high median (4.5). These indicate near-ceiling central tendency with pockets of uncertainty around (a) *where/how to rename/merge/split/reassign* and (b) *runtime feedback/telemetry* during processing.

Free-response themes (anonymized).

- **Scoring polarity help.** Missing or unclear help for negative vs. positive scoring semantics.
- **Job status clarity.** Ambiguity in whether jobs are running, complete, or hung when processing.
- **Selection feedback.** Highlighting the selected area in a different color was positively noted.
- **Demo coverage.** Some aspects were “not determinable from the video,” suggesting a documentation/demo gap.

Implications. Given the ceiling scores, changes should target *discoverability* and *status legibility*, not broad visual redesign: (1) add an inline *grading polarity* tooltip in the rubric panel; (2) strengthen signifiers and empty-state nudges for edit affordances (rename/merge/split/reassign); (3) improve job-status telemetry (progress %, “active/hung/complete” badges, timestamped last heartbeat); (4) update the demo/video to explicitly show these moments so raters can verify them without inference.

5.4 Clipboard Interoperability

Figure 4.3 shows the source grade table and pasted result.

5.5 Summary

Across these tests, we visually confirm end-to-end processing, vision-based semantic grouping with isolated outliers, consistent rubric propagation, and successful clipboard transfer of grades. The UI checklist verifies feature presence; the design-only UX questionnaire gauges readability and clarity; clustering metrics quantify grouping quality on a 10-student run. Crucially, the system reduces instructor workload: for this 200-answer batch, only 49 manual moves were required (median 2 per question), yielding an estimated $\sim 80\%$ reduction in handling time versus baseline page flipping.

Design-only ratings were ceilinged (median = 5, IQR = 0) with localized uncertainty around edit-tool discoverability (Q10) and runtime status indicators (Q14), which we address via stronger signifiers, inline scoring-polarity help, and clearer job-status telemetry.

Chapter 6

Conclusion

6.1 Summary of Problem and Goals

We addressed the effort and inconsistency of grading open-ended work by building an instructor-in-the-loop system that extracts regions, assigns identity via OCR, groups answers with a vision LLM, and enables rubric-first grading with auditability.

6.2 Evaluation Summary

Our visual tests demonstrated end-to-end processing, coherent semantic grouping (no answer OCR), consistent rubric propagation, reliable clipboard-based grade transfer using *Copy Table*, and a tangible reduction in instructor effort (49 moves over 200 answers; $\sim 80\%$ time savings under conservative assumptions).

6.3 Final Outcomes and Deliverables

We delivered the integrated web app, background services, a reproducible AI configuration (prompts and model settings), and a *Copy Table* flow that enables fast pasting of grades into external systems (e.g., LMSs or Sheets). Documentation includes an instructor guide and technical deployment notes.

6.4 Lessons Learned and Future Work

Implementing and using the system with instructors surfaced several practical insights about what most affected robustness and instructor experience, as well as concrete opportunities to extend the platform. Below we summarize key lessons and outline directions for future work.

- Visual pre-processing (downscale/tiling) mattered more for stability than minor prompt tuning.
- Instructors preferred neutral group names and an explicit *Ungrouped* bin.
- Future: direct CSV/Excel export in addition to *Copy Table*; domain-tuned prompts for math diagrams; adaptive thresholding for faint pencil; pre-clustering to cut LLM calls; regrade workflow.
- Future: prompting with answers to identify the correct group.
- Future: prompting with rubrics that the AI can assign to individual answers or groups.

Appendix A

Configuration and Deployment

This appendix captures the minimum configuration to run the system on a fresh machine. Replace placeholders (the ALL-CAPS bits) with values for your environment.

A.1 Prerequisites

- Windows 11 / Ubuntu 22.04 / macOS (developed & tested on all three).
- **.NET SDK** (web app).
- **Python 3.10+** (FastAPI grouping service).
- **MySQL/MariaDB**.
- **Ghostscript** (PDF rasterization) available on PATH or via GHOSTSCRIPT_EXE.
- **Tesseract OCR** (install language data eng at minimum).
- Vision LLM API access set via OPENAI_API_KEY.

A.2 FastAPI Service: .env

Create a .env file in the FastAPI project root:

```
# --- OpenAI / Vision LLM ---
OPENAI_API_KEY=YOUR_OPENAI_KEY

# --- Database used by the service ---
DB_HOST=YOUR_DB_HOST
DB_NAME=OICLearning
DB_USER=YOUR_DB_USER
DB_PASSWORD=YOUR_DB_PASSWORD

# --- Optional service bind (defaults shown) ---
HOST=0.0.0.0
PORT=8000
```

Start the service:

```
python -m venv .venv
# Windows: .venv\Scripts\activate
# macOS/Linux: source .venv/bin/activate
pip install -r requirements.txt
uvicorn app:app --host 0.0.0.0 --port 8000
```

A.3 Web App: appsettings.json

Use this structure for both `appsettings.json` and `appsettings.Development.json`. Only update the hostnames, passwords, folders, and API base URL; keep DB name and UID as shown.

```
{
  "ConnectionStrings": {
    "MySQLConnection":
```



```

        "Server=YOUR_DB_HOST;Database=OICLearning;Uid=www_oiclearning;Pwd=
        ↪ YOUR_DB_PASSWORD",
        "MySQLTestSite":
            "Server=YOUR_TEST_DB_HOST;Database=OICLearning;Uid=www_oiclearning;Pwd=
            ↪ YOUR_TEST_DB_PASSWORD"
    },

    "Logging": {
        "LogLevel": {
            "Default": "Information",
            "Microsoft.AspNetCore": "Warning"
        }
    },

    "AllowedHosts": "*",

    "RoleStrings": {
        "Teacher": ["Admin", "Teacher"],
        "Admin": ["Admin"],
        "Student": ["Student"]
    },

    "FileRepository": {
        "SubmissionFolder": "/path/to/submissions",
        "AutoGraderFolder": "/path/to/autograders"
    },

    "PythonApi": {
        "BaseUrl": "http://YOUR_FASTAPI_HOST:8000"
    }
}

```

Point the callers at PythonApi:BaseUrl

Update the two callers to use PythonApi:BaseUrl *without* a localhost fallback.

CourseService.cs (snippet)

```
var client = _httpClientFactory.CreateClient();
var baseUrl = _configuration.GetValue<string>("PythonApi:BaseUrl");
client.BaseAddress = new Uri(baseUrl);

var response = await client.PostAsJsonAsync("/autogroup", requestBody);
response.EnsureSuccessStatusCode();
```

GroupingService.cs (snippet)

```
var client = _httpClientFactory.CreateClient();
var baseUrl = _configuration.GetValue<string>("PythonApi:BaseUrl");
client.BaseAddress = new Uri(baseUrl);

// ...status polling / error handling continues...
```

A.4 Ghostscript Location

If Ghostscript is not on PATH, set GHOSTSCRIPT_EXE.

macOS (Homebrew):

```
export GHOSTSCRIPT_EXE=/opt/homebrew/bin/gs
```

Ubuntu/Debian:

```
export GHOSTSCRIPT_EXE=/usr/bin/gs
```

Windows (PowerShell):

```
$env:GHOSTSCRIPT_EXE="C:\Program Files\gs\gs10.03.0\bin\gswin64c.exe"
```

The extractor prefers the env var and falls back if needed:

```
var gsExe = Environment.GetEnvironmentVariable("GHOSTSCRIPT_EXE")  
    ?? "/opt/homebrew/bin/gs"; // adjust per OS
```

A.5 Tesseract on macOS (dev builds)

If you hit dylib resolution issues with Homebrew installs, use this post-build step:

```
<!-- OICLearning.csproj (snippet) -->  
<Target Name="link_deps" AfterTargets="AfterBuild">  
  <Exec Command="ln -sf /opt/homebrew/lib/libleptonica.dylib  
    $(OutDir)x64/libleptonica-1.82.0.dylib" />  
  <Exec Command="ln -sf /opt/homebrew/lib/libtesseract.dylib  
    $(OutDir)x64/libtesseract50.dylib" />  
</Target>
```

A.6 Build and Run (Web App)

1. Restore NuGet packages and build.
2. Apply EF Core migrations:

```
dotnet tool restore  
dotnet ef database update
```

3. Launch the app:

```
dotnet run
```

4. Ensure `FileRepository.SubmissionFolder` exists and is writable.

A.7 Operational Notes

- Re-uploads are idempotent; crops use stable names (e.g., Q27a.png).
- Only identity regions are OCR'd; answer crops go to the vision model.
- Groupings are editable; an *Ungrouped* bucket catches outliers.
- Grades appear in an on-screen table; CSV/Excel export is available.

Bibliography

- [1] “Ghostscript,” <https://ghostscript.com/>, Artifex Software, Inc., 2025, postScript/PDF interpreter.
- [2] “Uglytoad pdfpig,” <https://github.com/UglyToad/PdfPig>, UglyToad, 2025, .NET PDF reading library.
- [3] “SkiaSharp,” <https://github.com/mono/SkiaSharp>, .NET Foundation, 2025, .NET 2D graphics library (Skia bindings).
- [4] “Tesseract ocr,” <https://github.com/tesseract-ocr/tesseract>, Tesseract OCR Developers, 2025, open-source OCR engine.
- [5] “Openai gpt-4o and gpt-4o vision,” <https://platform.openai.com/docs/overview>, OpenAI, 2025, multimodal LLM used for grouping/vision.
- [6] V. Alto, *Modern Generative AI with ChatGPT and OpenAI Models: Leverage the capabilities of OpenAI’s LLM for productivity and innovation with GPT-3 and GPT-4*. Packt Publishing Ltd, 2023.
- [7] L. Chen, G. Chen, and X. Lin, “Artificial intelligence in education: A review,” *IEEE Access*, vol. 8, pp. 75 264–75 278, 2020.
- [8] R. Luckin, W. Holmes, M. Griffiths, and L. B. Forcier, “Intelligence unleashed: An argument for AI in education,” Pearson Education, Tech. Rep., 2016.

- [Online]. Available: <https://www.pearson.com/content/dam/one-dot-com/one-dot-com/global/Files/about-pearson/innovation/Intelligence-Unleashed-Publication.pdf>
- [9] O. Zawacki-Richter, V. I. Marín, M. Bond, and F. Gouverneur, “Systematic review of research on artificial intelligence applications in higher education – where are the educators?” *International Journal of Educational Technology in Higher Education*, vol. 16, no. 1, p. 39, 2019.
- [10] R. Weegar and P. Idestam-Almquist, “Reducing workload in short answer grading using machine learning,” *International Journal of Artificial Intelligence in Education*, vol. 32, pp. 611–643, 2022. [Online]. Available: <https://link.springer.com/article/10.1007/s40593-022-00322-1>
- [11] R. Könnecke and T. Zesch, “Automated scoring of content and style in short essays,” *Frontiers in Education*, vol. 5, p. 90, 2020.
- [12] T. Liu, J. Chatain, L. Kobel-Keller, G. Kortemeyer, T. Willwacher, and M. Sachan, “Ai-assisted automated short answer grading of handwritten university level mathematics exams,” *arXiv preprint arXiv:2308.11728*, 2023. [Online]. Available: <https://arxiv.org/abs/2308.11728>
- [13] G. Kortemeyer, “Toward AI grading of student problem solutions in introductory physics: A feasibility study,” *Physical Review Physics Education Research*, vol. 19, no. 2, p. 020163, 2023. [Online]. Available: <https://journals.aps.org/prper/abstract/10.1103/PhysRevPhysEducRes.19.020163>
- [14] G. Kortemeyer, J. Noh, and D. Onishchuk, “Grading assistance for a handwritten thermodynamics exam using artificial intelligence: An

- exploratory study,” *arXiv preprint arXiv:2306.17859*, 2023. [Online]. Available: <https://arxiv.org/abs/2306.17859>
- [15] B. D. Lund, D. Wang, K. Chao, and M. S. Gerber, “Chatgpt’s ability to provide timely, novel, and impactful ideas for research,” *arXiv preprint arXiv:2305.06566*, 2023. [Online]. Available: <https://arxiv.org/abs/2305.06566>
- [16] N. Mehrabi, F. Morstatter, N. Saxena, K. Lerman, and A. Galstyan, “A survey on bias and fairness in machine learning,” *ACM Computing Surveys*, vol. 54, no. 6, pp. 1–35, 2021.
- [17] C. C. Tossell, N. L. Tenhundfeld, A. Momen, K. Cooley, and E. J. de Visser, “Student perceptions of chatgpt use in a college essay assignment: Implications for learning, grading, and trust in artificial intelligence,” *IEEE Transactions on Education*, 2023. [Online]. Available: <https://ieeexplore.ieee.org/document/10120620>
- [18] D. J. Nicol and D. Macfarlane-Dick, “Formative assessment and self-regulated learning: A model and seven principles of good feedback practice,” *Studies in Higher Education*, vol. 31, no. 2, pp. 199–218, 2006.
- [19] “Gradescope,” <https://www.gradescope.com/>, Turnitin, LLC, 2025, online grading platform; fixed-template paper exams and Answer Groups.
- [20] “Crowdmark,” <https://crowdmark.com/>, Crowdmark Inc., 2025, collaborative grading and assessment for paper and online exams.
- [21] “Akindi,” <https://www.akindi.com/>, Akindi Inc., 2025, oMR bubble-sheet creation and scanning.

- [22] “Zipgrade,” <https://www.zipgrade.com/>, ZipGrade LLC, 2025, mobile multiple-choice scanning and analytics.
- [23] “Canvas speedgrader,” <https://www.instructure.com/canvas>, Instructure, Inc., 2025, IMS grading workflow with rubrics and annotations.
- [24] “Prairielearn,” <https://www.prairielearn.org/>, PrairieLearn, 2025, auto-graded, parameterized questions for STEM.
- [25] “Möbius,” <https://www.digitaled.com/mobius/>, DigitalEd, 2025, algorithmic assessment and math engine-based autograding.
- [26] “Coderunner,” <https://coderunner.org.nz/>, University of Canterbury, 2025, programming question autograder (often used via Moodle plugin).
- [27] “Codegrade,” <https://www.codegrade.com/>, CodeGrade B.V., 2025, code autograding and rubric workflows.
- [28] “The family educational rights and privacy act (ferpa),” Pub. L. No. 93-380 (1974), codified as 20 U.S.C. § 1232g and 34 CFR Part 99, 1974. [Online]. Available: <https://www.govinfo.gov/content/pkg/USCODE-2023-title20/pdf/USCODE-2023-title20-chap31-subchapIII-part4-sec1232g.pdf>