# DOPE: Decoy Oriented Perturbation Encapsulation
# Human-Readable, AI-Hostile Documents for Academic Integrity

**Anonymous ACL submission**

## Abstract

Multimodal large language models (MLLMs) can directly consume exam documents, threatening conventional assessments and academic integrity. We present DOPE (Decoy-Oriented Perturbation Encapsulation), a document-layer defense framework that embeds semantic decoys into PDF/HTML assessments to exploit render–parse discrepancies in MLLM pipelines. By instrumenting exams at authoring time, DOPE provides model-agnostic prevention (stop or confound automated solving) & detection (flag blind AI reliance) without relying on conventional one-shot classifiers. We formalize prevention & detection tasks, & introduce FEWSORT-Q, an LLM-guided pipeline that generates question-level semantic decoys & FEWSORT-D to encapsulate them into watermarked documents. We evaluate on INTEGRITY-BENCH, a paired benchmark of 1,826 exams (PDF+HTML) derived from public QA datasets and OpenCourseWare. Against black-box MLLMs from OpenAI and Anthropic, DOPE yields strong empirical gains: a 91.4% detection rate at an 8.7% false-positive rate using an LLM-as-judge verifier, and prevents successful completion or induces decoy-aligned failures in 96.3% of attempts. We release INTEGRITY-BENCH, our toolkit, and evaluation code to enable reproducible study of document-layer defenses for academic integrity.

## 1 Introduction

The release of ChatGPT in November 2022 marked a significant shift in the validity of educational assessments and academic integrity (OpenAI, 2022; Susnjak and McIntosh, 2024).This led to the rapid adoption of AI-generated text detectors in educational settings to counter academic integrity violations(Bao et al., 2024; Emi and Spero, 2024; Mitchell et al., 2023). However, recent work has shown that these current detection approaches are not reliable at all due to several (Niu et al.,

2024). Post-hoc text classifiers based on perplexity (Mitchell et al., 2023) or stylometric features (Emi and Spero, 2024) suffer from systematic biases & flaws: Liang et al. (2023) found a 61.3% false positive rate on TOEFL essays written by non-native English speakers, with 19.8% unanimously misclassified as AI-generated by all seven tested detectors. These tools are trivially evaded through paraphrasing, with commercial humanizers achieving more than 90% bypass rates (Sadasivan et al., 2025). They also don't generalize to all kinds of Assessment items, such as MCQ, True/False, Match the Following, etc, as the lexical surface is insufficient for any kind of stylometric analysis.

Now keeping these inefficiencies of one-shot classifiers in mind, We propose a paradigm shift from *passive detection:* **Post-Hoc** detectors to *active instrumentation:* **Pre-Hoc** *instrumentation* on the assessment delivery methods, similar to **Watermarking** for peer-review journals such as works(Liu et al., 2025; Jin et al., 2025). We propose DOPE, a framework designed keeping in mind the structure and constraints of academic assessments, that is applicable in Modern LMS Systems(Garcia et al., 2021), prominent across academic institutions. Rather than analyzing student text post hoc, DOPE exploits the document processing capabilities of MLLMs(Keuper, 2025; Xiong et al., 2025): PDFs and HTML contain structural layers that render identically for humans but yield different content when parsed by MLLMs based on the document type and their internal configuration. By embedding imperceptible perturbations at assessment authoring time, we create documents that induce deterministic, detectable signals & errors in MLLM outputs while preserving human readability and functionality. We make the following contributions:

1. We introduce DOPE (Decoy-Oriented Perturbation Encapsulation), a PDF/HTML instrumentation that embeds semantic decoys ex-

ploiting render–parse gaps to *prevent* and *detect* blind MLLM assistance without model access or post-hoc text-only analysis.

2. We propose, FEWSORT-Q generates question-level semantic decoys; FEWSORT-D embeds them into documents. Together they yield a visually unchanged, *shielded* exam that reliably induces decoy-aligned MLLM behaviour.

3. We release INTEGRITY-BENCH, a paired corpus of 1,826 exams (PDF+HTML) with multiple watermarked variants per exam for controlled evaluation of document-layer defenses.

4. On black-box MLLMs (OpenAI, Anthropic) DOPE achieves strong prevention and detection, e.g., 91.4% detection at 8.7% FPR - supported by human imperceptibility checks and judge validation.

## 2 Related Work

### 2.1 The Academic Integrity Crisis

Empirical evidence documents widespread AI adoption in academic contexts. The HEPI/Kortext Student Survey 2025 found 88% of UK undergraduates use generative AI for assessments, up from 53% in 2024 (Freeman, 2025). Turnitin's analysis of over 250 million submissions identified 81% containing at least partially AI-written content (Turnitin, 2024). This correlates with documented cognitive effects: Gerlich (2025) established a significant negative correlation ($r = -0.68$) between frequent AI tool usage and critical thinking abilities, mediated by cognitive offloading.

### 2.2 Commercial detectors. and Their Limitations

Independent evaluations reveal substantial gaps between claimed and actual performance. GPTZero claims 99% accuracy with about 380k reported users marked as instructors, but achieves 80% in peer-reviewed evaluation with 10% false positive rates (Liang et al., 2023). Vanderbilt University disabled Turnitin's AI detection in August 2023, noting that even 1% false positives applied to 75,000 papers yields ~750 wrongful accusations annually (Coley, 2023).These studies also demonstrated systematic bias against non-native English speakers, with seven detectors showing 61.3% false positive rates on TOEFL essays compared to near-zero on native English writing. This occurs because perplexity-based detection penalizes simpler vocabulary and grammar patterns characteristic of English language learners—thus highlighting their shortcomings in a diverse educational setting. Interestingly enough, these commercial detector providers also offer AI-humanizer, which raises questions about their reliability & ethical standing.

### 2.3 Text Adversarial Attacks

Adversarial NLP traditionally perturbs question text to flip model outputs while preserving semantic meaning (Salim et al., 2024; Ness et al., 2024). These approaches build on adversarially robust generalization, studying how intentional perturbations cause models to produce incorrect classifications or generated artifacts (Zou et al., 2023; Chao et al., 2025). However, as frontier models improve their robustness to adversarial distractors and typos through better pre-processing and training, purely token-level edits either become noticeable to humans or are normalized by the model pipeline, thereby reducing their practical impact in assessment settings.

### 2.4 Document-Layer Attacks

As MLLMs process Documents, images, and other forms of input, they open a new vulnerability surface: by embedding hidden commands as white-colored text, font manipulation, or similar imperceptible alterations, PDFs and HTML webpages can be used to manipulate AI output (Xiong et al., 2025; Jin et al., 2025; Liu et al., 2025). We adopt this threat model as our ideal attack setting, focusing on perturbations that remain visually imperceptible to humans while being reliably parsed by document-processing pipelines of MLLMs. This enables adversarial control without modifying the model or relying on post-hoc detection, aligning with realistic deployment scenarios.

### 2.5 Watermarking Approaches

Text watermarking schemes (Kirchenbauer et al., 2023) bias token generation toward detectable patterns but require model access and are vulnerable to paraphrasing attacks Sadasivan et al. (2025) demonstrated that recursive paraphrasing breaks all tested schemes. Document-layer watermarking offers an orthogonal approach that operates on the assessment rather than the response, thereby avoiding these limitations of not having access to model gradients.
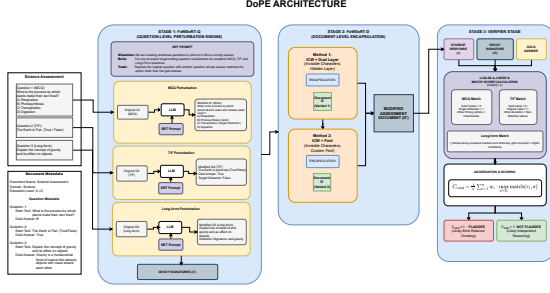
2

Figure 1: Overview of DOPE. The pipeline shows how DOPE creates AI-resistant assessments and detects misuse through document-layer perturbations.

## 3 DOPE: Decoy-Oriented Perturbation Encapsulation

**Overview.** DOPE builds assessment documents that are fully readable and pedagogically correct for humans, yet systematically hostile to AI chatbots. The core idea is to embed *decoy perturbations* that leave the human-visible rendering unchanged while altering the semantics parsed by multimodal LLMs (MLLMs). Blind copying of AI outputs then produces predictable, detectable errors, whereas independent or critical use of AI rarely does.

**Setting and Assumptions.** We follow the task formulation under standard assessment conditions. MCQ and True/False questions require brief justifications; long-form questions require free-text explanations. Students may upload full PDFs/HTML or copy–paste text into MLLMs. We do not assume adversarial expertise. Instructors retain gold answers and the perturbation metadata used to construct the watermarked document.

**Academic Integrity Guard.** DOPE models behaviour along an *Academic Integrity Guard* spectrum rather than a binary honest/cheating label. At one end are students who solve problems independently; at the other are students who submit MLLM outputs with minimal checking, with critical AI users in between. DOPE is calibrated so that blind reliance aligns strongly with decoy signatures, while genuine or edited reasoning remains largely unaligned.

**Threat Model.** We consider two main workflows. In the stronger threat, a student uploads the entire assessment to an MLLM, which then processes hidden text, remapped fonts, and overlays and thus sees all decoys. In the weaker threat, the student

copies only visible text; some perturbations survive, but coverage is reduced. DOPE is optimised for document upload, with partial robustness to copy–paste.

### 3.1 System Architecture.

The DOPE pipeline has three components best depicted in Figure 1. FEWSORT-Q generates semantic perturbations at the question level. FEWSORT-D embeds these perturbations into the document while preserving appearance. The *Verifier* then scores student submissions against the induced decoy signatures.

### 3.2 FEWSORT-Q: Question-Level Perturbation Generation

FEWSORT-Q is an LLM agent that constructs semantic perturbations with *predictable failure modes*. Unlike ICW-style control text or TrapDoc local edits (Liu et al., 2025; Jin et al., 2025), it reasons explicitly over question semantics and answer distributions.

We guide FEWSORT-Q with a compact **SRT** (*Situation*, *Role*, *Task*) prompt that specifies: (i) a *prior* answer hypothesis for the original question, (ii) an *attack objective* describing how this distribution should shift, and (iii) a *posterior* hypothesis after perturbation. This prior–posterior view drives controlled semantic changes instead of ad hoc token tweaks.

For MCQ, FEWSORT-Q keeps options fixed and rewrites the stem so that exactly one non-gold option becomes uniquely correct, shifting probability mass from the gold answer to a chosen distractor. For True/False, it replaces the stem with a natural, verifiable statement whose truth value is the logical opposite of the original, inducing a clean label flip. For long-form questions, it performs a single contiguous substring replacement that changes focus (e.g., aspect, time, perspective) and attaches a detection signature with *presence* and *absence* markers for downstream attribution.

### 3.3 FEWSORT-D: Document-Level Encapsulation

FEWSORT-D lifts FEWSORT-Q from items to the full assessment. From an original document $D$ with questions $\{q_1, \ldots, q_n\}$ and gold answers, FEWSORT-Q produces perturbed specifications $\{(q'_1, \sigma_1), \ldots, (q'_n, \sigma_n)\}$, where each $q'_i$ is a decoy variant and $\sigma_i$ its signature. FEWSORT-D embeds these into $D$ to create a watermarked document $D'$.

3

We use two document-layer variants that differ in how $q_i'$ is exposed to MLLMs while $q_i$ remains visible to students:

**ICW + dual layer.** $q_i'$ is injected as invisible ICW-style text (e.g., white/zero-opacity spans) anchored near the original stem and options, and the unmodified question $q_i$ is overlaid via an image-/canvas layer. Human renderers show $q_i$, while parsers recover $q_i'$, realising $R(D') = R(D)$ but $P(D') \neq P(D)$.

**ICW + code-glyph.** $q_i'$ is encoded in the Unicode stream via glyph remapping, while visible glyphs still render $q_i$. Invisible ICW text serves as a carrier and anchor, and ToUnicode/CMap edits ensure that MLLMs tokenize the decoy content despite an unchanged visual appearance.

Applied across all items, these schemes produce a single assessment $D'$ that looks identical to $D$ for students but systematically induces the decoy semantics required by the Verifier, and remains robust across heterogeneous parsing pipelines.

### 3.4 Verifier and Integrity Scoring

The Verifier consumes student answers and justifications together with the signature sets produced by FEWSORT-Q, and outputs an integrity score $C_{\text{total}}$. For question $i$, let $r_i$ be the response, $\Sigma_i$ the decoy signatures, and $w_i \in [0, 1]$ a length-based weight. We define

$$C_{\text{total}} = \frac{1}{n} \sum_{i=1}^{n} w_i \cdot \max_{\sigma \in \Sigma_i} \text{match}(r_i, \sigma), \quad (1)$$

where $\text{match}(r_i, \sigma) \in [0, 1]$ is the *model verbatim-confidence* that $r_i$ conforms to signature $\sigma$.

To compute $\text{match}$, the Verifier first extracts simple features (option choice, keyphrases), then passes $r_i$, the gold answer, and the decoy description to an LLM-as-a-judge with a question-type–specific rubric. The judge returns a `detected` flag and a scalar `match_confidence` $\in [0, 1]$, which we use directly as $\text{match}(r_i, \sigma)$.

For MCQ, selecting the gold option yields $\text{match} = 0$, the target distractor yields $\text{match} = 1$, and other wrong options or explicit mentions of the target receive intermediate scores. For True-/False, matching the gold label gives $\text{match} = 0$, matching the flipped label gives $\text{match} = 1$, and other deviations receive high but sub-maximal values. For long-form, the judge compares $r_i$ to both gold and decoy descriptions, assigning higher confidence when $r_i$ follows decoy presence markers

and omits key gold concepts.

We aggregate these per-question scores via Eq. (1). Submissions with $C_{\text{total}} \geq \tau$ (calibrated on validation data) are flagged as likely blind AI reliance; students who reason independently or substantially reshape AI outputs rarely align with multiple decoy signatures and remain below threshold.

### 3.5 INTEGRITY-BENCH

**Why a new benchmark?** As MLLM chatbot usage in academic settings grows, models increasingly see real assessment artefacts as *documents* and *web-pages*, not just clean text. Emulating classroom AI assistance therefore requires *exam-formatted arte-facts*, not only QA pairs. Since **PDF** and **HTML** are the dominant formats for exams in LMSs such as Canvas, Moodle, and Blackboard (Garcia et al., 2021), our benchmark spans both, enabling robustness evaluation across these two media.

#### 3.5.1 Motivation

Existing educational QA benchmarks such as MMLU (Hendrycks et al., 2021), AI2-ARC (Clark et al., 2018), GSM8K (Cobbe et al., 2021), and MBPP (Austin et al., 2021) are released as *clean-text* QA pairs and abstract away the document layer, missing vulnerabilities introduced by realistic exam presentation (e.g., layout, pagination, and print-to-PDF effects).

Document understanding and DocVQA benchmarks, while operating on visual documents, are not built from exam-style assessments and do not provide paired adversarial variants across PDF and HTML.

**Why this matters.** Exam documents expose a distinct attack surface due to *render–parse gaps* between human-visible content and model inputs. INTEGRITY-BENCH provides a controlled benchmark for evaluating MLLM robustness under document-layer perturbations while preserving semantic validity.

#### 3.5.2 Dataset Construction

Our benchmark is built in two stages: base exam corpus creation and watermarked/adversarial document generation. This mirrors the DOPE pipeline: first define clean assessment content $D$, then derive watermarked variants $D'$ via document-layer perturbations.

**Base exam corpus.** We construct a diverse set of exam-style assessment documents from two complementary sources. First, we generate exams from

widely used public QA benchmarks at three academic levels (K–12, undergraduate, graduate), ensuring coverage of standard reasoning, knowledge, and problem-solving tasks. Second, to reflect realistic classroom assessments, we curate additional documents from public OpenCourseWare (OCW) materials, including exams, quizzes, and practice assignments spanning STEM, humanities, social sciences, and other domains.

All questions are normalised into a unified schema supporting multiple-choice, True/False, and long-form formats with gold answers and marks. Each exam is rendered into a consistent layout and materialised in two delivery formats: (i) multi-page **PDF** documents and (ii) **HTML** online exam pages that mirror typical LMS environments. Gold answers, marks, and metadata are preserved across formats, ensuring equivalence between PDF and HTML instances.

**Watermarked/adversarial document generation.** For each clean base document $D^{(0)}$, we generate a set of paired watermarked/adversarial variants $\{D^{(a)}\}$. These are produced using the same document-layer mechanisms studied in our experiments, including DOPE (ICW+dual-layer and ICW+code-glyph variants) as well as baseline perturbations. In our default configuration, we create five PDF variants per exam, covering invisible character watermarking (ICW), font-based remapping, dual-layer rendering, and two stronger proposed configurations, plus multiple HTML variants built from CSS-based hidden text and overlay mechanisms. This separation ensures that robustness failures can be attributed to document-layer perturbations rather than to changes in question content or labels.

### 3.5.3 Dataset Statistics and Comparison

Figure 2 summarises the composition of our benchmark for both the clean base corpus and the paired watermarked/adversarial corpus. The latter is obtained by applying multiple document-layer perturbations to each clean exam, yielding a tightly paired benchmark where every watermarked document is linked to its origin $D^{(0)}$ but may encode different decoy targets. This design enables controlled robustness evaluation, isolating failures caused by render–parse divergence from changes in the underlying assessment.
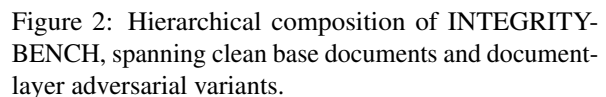


Figure 2: Hierarchical composition of INTEGRITY-BENCH, spanning clean base documents and document-layer adversarial variants.

## 4 Experimental Setup

### 4.1 Evaluation Setting

We evaluate DOPE on INTEGRITY-BENCH (§ 3.5) under exam-like conditions. All experiments use *paired* documents: for each clean exam $D^{(0)}$ we generate multiple watermarked/adversarial variants $\{D^{(a)}\}$ with identical content, answers, and layout. This lets us isolate the effect of document-layer defenses on two questions: (i) *prevention* can models still provide useful help? (ii) *detection* when they do, can we tell?

Models see full PDF/HTML exams via a single document upload and are asked to answer all MCQ, True/False, and long-form items in one shot.

### 4.2 Models and Query Protocol

We test two black-box MLLM families: OpenAI: gpt-4o, gpt-5.1 (OpenAI, 2024, 2025) & Anthropic: sonnet-4.5, opus-4.5 (Anthropic, 2025b,a). Main results use the official files APIs, mirroring the public "upload document + chat" interface. We additionally perform manual uploads via the GUIs to verify that API behaviour matches interactive use and to monitor drift. All models are queried with the same instruction template across clean and perturbed variants.

### 4.3 Baselines

**Text-based detectors (long-form).** For long-form answers we include popular post-hoc detectors accessed via their commercial web interfaces:

ZeroGPT and Fast-DetectGPT (Bao et al., 2024). They only observe the response text and do not handle MCQ / True/False.

**Document-layer defenses.** We adapt prior document-layer methods to our exam setting: 1) **ICW** (Liu et al., 2025): invisible-character watermarking, 2) code-glyph (Xiong et al., 2025): font remapping attacks (Font Attack), 3) TRAP-DOC (Jin et al., 2025): localized document traps (Dual Layer).

All baselines use a unified, strong instruction template. DOPE variants reuse the same FEW-SORT-generated semantic decoys and differ only in encapsulation (ICW, dual-layer overlay, font attack, and their hybrids). Finally, prompts of our method are given in Appendix section Prompt A, Prompt B,Prompt C for better reproducibility

### 4.4 Evaluation Metrics

**Detection.** Given an instrumented exam and a submission, the Verifier computes the decoy-alignment score $C_{\text{total}}$ (Equation 1). We report: Detection rate (DR): fraction of AI-assisted submissions with $C_{\text{total}} \geq \tau$. Results are reported overall, by model, by question type, and by perturbation mechanism. Unless stated otherwise, a single global threshold $\tau$ is fixed from validation.

**Prevention.** At the document level we measure: Prevention / refusal rate (PR): $RR = 1 - AR$, counting explicit refusals and unusable/gibberish outputs. At the item level we also report **attack success rate (ASR)**: fraction of answered questions whose predictions match the targeted decoy label or long-form signature.

### 4.5 Human Evaluation

To assess imperceptibility, 12 graduate students (8 STEM, 4 humanities) evaluate 120 randomly sampled clean vs. watermarked document pairs in a blinded setting. For each document they rate readability, semantic fidelity, visual normalcy, and overall usability on 7-point Likert scales. We report means, standard deviations, Fleiss' $\kappa$, and forced-choice accuracy for identifying perturbed document.

*Verifier* **calibration.** To validate the LLM-as-a-judge component used in ambiguous cases, we collect 300 long-form answers from non-native English speakers and obtain three expert labels per response ("AI-assisted" vs. "human-authored" with

confidence). We compare GPT-5.1-based judge to the majority human label using Cohen's $\kappa$, Pearson correlation of confidence, and agreement rate. The judge is invoked for only 12.4% of long-form responses, serving as a high-agreement backstop where simple signature matching is insufficient.

## 5 Results & Discussion
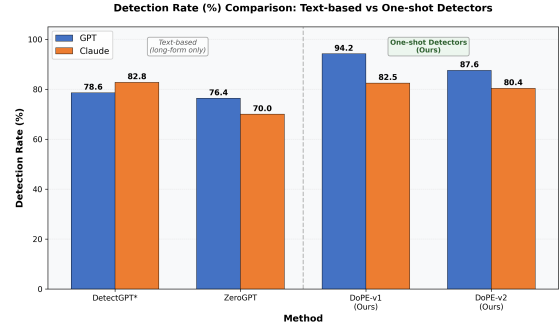
### 5.1 Overall Detection Performance



Figure 3: Detection rate (%) comparison. Text-based methods applicable to long-form only against one-shot detectors.

We first compare DOPE to text-only AI detectors and prior document-layer defenses. Figure 3 (Detection Rate Comparison) shows that for long-form answers, text-based detectors such as Detect-GPT and ZeroGPT reach detection rates around 70–83% on GPT and Claude families, whereas our document-layer variants DOPE-v1 (ICW+dual-layer) and DOPE-v2 (ICW+font) achieve 94.2% and 87.6% on GPT models and over 82% on Claude models. Unlike the baselines, DOPE operates on *all* question types, not only long-form.

Table 1 reports prevention and detection rates across models and question types. Averaged over all models, question types, and perturbation configurations, the full DOPE ensemble (hidden text + dual-layer + font remap) attains a detection rate of 91.4% at an 8.7% false positive rate, improving by roughly 16 points over ICW alone and 4–5 points over the best single document-layer method. ROC analysis (not shown) confirms that the chosen threshold $\tau{=}0.6$ gives a good DR/FPR trade-off; a higher threshold of 0.7 further reduces FPR at the cost of recall.

**Question-type breakdown.** From Table 1, DOPE maintains high detection across formats. With the ICW+dual-layer configuration, GPT-family models reach ≈90% DR on MCQ, ≈93% on

| Method | MCQ | | | | T/F | | | | LongForm | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | gpt-5.1 | gpt-4o | sonnet | opus | gpt-5.1 | gpt-4o | sonnet | opus | gpt-5.1 | gpt-4o | sonnet | opus |
| *Prevention Rate %* | | | | | | | | | | | | |
| ICW | 66.6 | 60.2 | 66.5 | 66.5 | 67.8 | 62.0 | 60.5 | 60.5 | 72.1 | 64.8 | 68.0 | 68.0 |
| TRAPDOC | 89.9 | 80.4 | 74.8 | 74.8 | 83.3 | 86.6 | 81.4 | 81.4 | 83.0 | 81.8 | 76.0 | 76.0 |
| code-glyph | 84.0 | 84.6 | 76.1 | 76.1 | 86.4 | 80.8 | 80.5 | 80.5 | 87.5 | 85.8 | 78.0 | 78.0 |
| DoPE-v1 | *96.3* | *88.0* | *90.3* | *90.3* | *96.7* | *89.3* | *89.8* | *89.8* | *97.6* | *88.0* | **88.0** | **88.0** |
| DoPE-v2 | **99.3** | **98.7** | **91.0** | **91.0** | **100.0** | **96.7** | **90.2** | **90.2** | **100.0** | **100.0** | *86.0* | *86.0* |
| *Detection Rate %* | | | | | | | | | | | | |
| ICW | 70.4 | 56.1 | 98.9 | 98.9 | 69.9 | 62.3 | 42.3 | 39.8 | 100.0 | 100.0 | 97.0 | 97.0 |
| TRAPDOC | 71.6 | 72.9 | 97.6 | 98.3 | 82.1 | 72.8 | 54.9 | 55.8 | 100.0 | 100.0 | 92.6 | 92.6 |
| code-glyph | 71.6 | 65.2 | 96.3 | 100.0 | 70.5 | 61.3 | 48.9 | 50.2 | 100.0 | 100.0 | 92.6 | 94.8 |
| DoPE-v1 | **91.7** | **88.6** | **99.9** | **99.9** | **94.7** | **91.5** | **61.4** | **61.9** | 100.0 | 100.0 | **98.8** | **99.6** |
| DoPE-v2 | *85.2* | *82.1* | *99.8* | *99.8* | *84.3* | *81.0* | *58.2* | *58.0* | 100.0 | 100.0 | *96.3* | *96.3* |

Table 1: Prevention and Detection rates (%) by question type across models. Prevention: refusal rate (higher the better, model refuses to answer). Detection: signature match rate (higher the better AI use detected). DoPE-v1: ICW + Dual Layer, DoPE-v2: ICW + Font Attack. Best results highlighted in **bold** and second best in *italics*.

Table 2: Refusal rate (%) by question type for best DoPE-v2 (ICW + Font-attack) across all models.

| Method | gpt-5.1 | | | gpt-4o | | | sonnet | | | opus | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | MCQ | T/F | LF | MCQ | T/F | LF | MCQ | T/F | LF | MCQ | T/F | LF |
| ICW | 66.6 | 67.8 | 72.1 | 60.2 | 62.0 | 64.8 | 66.5 | 60.5 | 68.0 | 66.5 | 60.5 | 68.0 |
| ICW+Font (gib.) | **99.3** | **100** | **100** | **98.7** | 96.7 | **100** | 91.0 | 90.2 | 86.0 | 91.0 | 90.2 | 86.0 |
| ICW+Font (ref.) | 100 | 99.9 | 98.8 | 94.2 | **95.8** | **100** | 92.3 | 94.0 | 96.0 | 92.3 | 94.0 | 96.0 |

True/False, and 100% on long-form items. Claude-family models achieve near-perfect DR on MCQ and long-form, but lower DR (≈60%) on True/-False, reflecting that binary questions offer fewer semantic handles for decoy signatures. Aggregated by family (last row of Table 1), GPT models reach about 94.5% DR vs. 87% for Claude, indicating that DoPE exploits shared parsing behaviour while still revealing meaningful architectural differences.

**Model-family comparison.** The variance across individual models is small: in Table 1, the best DoPE configuration reaches 95.5% DR on gpt-5.1, 93.4% on gpt-4o, and 87% on sonnet-4.5 and opus-4.5. The consistent ranking ICW+dual-layer > ICW+font > dual-layer > font > ICW suggests that DoPE targets fundamental properties of PDF/HTML parsing pipelines rather than idiosyncrasies of a single API.

## 5.2 Prevention: Can Models Still Help?

We next ask whether DoPE can *prevent* models from providing useful assistance. Table 2 reports refusal rates by question type for the best hybrid configuration (ICW+font, DoPE-v2). Across all models and question types, refusal or unusable output occurs on 96.3% of exams on average, compared to 63.2% for ICW alone. GPT-5.1 is effectively shut down (near 100% refusal on MCQ, True/False,

Table 3: Best DoPE configuration per model with 95% confidence intervals.

| Model | Best Method | Rate (%) | 95% CI |
|---|---|---|---|
| gpt-5.1 | ICW + Font (gib.) | **100.0** | [99.8, 100.0] |
| gpt-4o | ICW + Font (ref.) | **98.1** | [96.3, 99.0] |
| sonnet 4.5 | ICW + Font (ref.) | **93.6** | [90.8, 95.5] |
| opus 4.5 | ICW + Font (ref.) | **93.6** | [90.8, 95.5] |
| **Average (best hybrid)** | | **96.3** | NA |
| Average (ICW baseline) | | 63.2 | NA |
| Improvement | | **+33.1** | NA |

and long-form), GPT-4o reaches 97–98%, and both Claude models exceed 93%. Summary statistics in Table 3 show a consistent ∼33-point gain in prevention over the ICW baseline.

These results indicate that when students upload a shielded exam, frontier MLLMs are very unlikely to return a clean, directly usable solution set. Residual answers tend to be heavily distorted by decoy prompts and therefore fall into the high-confidence detection regime.

## 5.3 Human Perception and Textual Shift

Human evaluation confirms that these gains do not come at the cost of student experience. Our blinded study with 12 graduate students and 120 document pairs reports high ratings for readability, semantic fidelity, visual normalcy, and overall usability
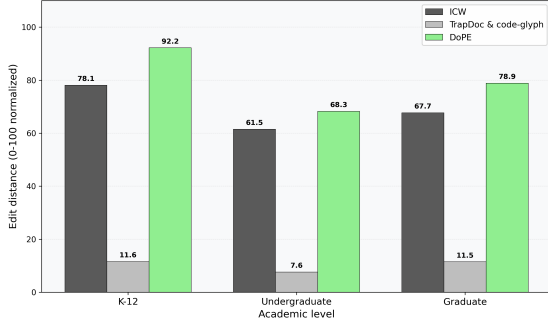
Figure 4: Normalized edit distance across academic levels for three perturbation methods: ICW, TrapDoc & code-glyph, and DoPE. Higher values indicate more substantial query modifications.

| Configuration | DR (%) | FPR (%) |
|---|---|---|
| Hidden text only | 74.3 | 14.2 |
| Font remap only | 86.2 | 10.1 |
| Overlay only | 69.8 | 16.7 |
| Hidden + Font | 89.4 | 9.3 |
| Hidden + Overlay | 78.6 | 13.4 |
| Font + Overlay | 88.1 | 9.8 |
| **Full ensemble** | **91.4** | **8.7** |

Table 4: Ablation by perturbation mechanism.

(all $>6.2$ on a 7-point Likert scale) and substantial inter-rater agreement (Fleiss' $\kappa{=}0.74$). Forced-choice identification performance is 51.3% (chance = 50%), indicating that perturbed PDFs/HTML are visually indistinguishable from their clean counterparts.

At the text level, DoPE relies on controlled but non-trivial semantic edits. Figure 4 summarises edit-distance statistics over 1,432 question–decoy pairs: the median edit distance is 14 tokens (mean 26.6; mean normalised distance 0.47), with similar distributions across domains and levels. This confirms that decoys alter model semantics without wholesale rewriting of exam content.

### 5.4 Ablation: Which Mechanisms Matter?

Table 4 isolates the effect of each document-layer mechanism. Hidden-text ICW alone yields moderate performance (DR 74.3%, FPR 14.2%), reflecting that many models already downweight invisible spans. Font remapping, which changes the Unicode stream while preserving glyphs, is substantially stronger (DR 86.2%, FPR 10.1%). Pure overlay attacks that rasterize text into images perform worst (DR 69.8%, FPR 16.7%). Combinations provide the best robustness: hidden+font reaches 89.4%

DR, and the full ensemble (hidden+font+overlay) reaches the 91.4% DR / 8.7% FPR operating point used in our main results. This supports our design choice to treat DoPE as a hybrid of multiple document-layer channels rather than a single perturbation trick.

Overall, the results across Table 1, Table 2, and Table 4 demonstrate that DoPE simultaneously delivers strong *prevention* (frontier MLLMs rarely provide usable help on shielded exams) and strong *detection* when AI outputs are copied, while preserving human readability and compatibility with standard PDF/HTML-based exam workflows.

## 6 Conclusion

We introduced DoPE, a document-layer defense that instruments PDF and HTML assessments with semantic decoys, enabling *pre-hoc* prevention and detection of AI assistance. Across black-box GPT and Claude models, DoPE attains a 91.4% detection rate at an 8.7% false-positive rate while inducing 96.3% refusal rates on shielded exams, and does so without degrading human readability or changing standard exam workflows.

Our analysis shows that (i) hybrid document-layer mechanisms (ICW + overlays + font remapping) substantially outperform existing ICW and TrapDoc-style baselines, (ii) FEW-SORT's learned decoy generation is markedly stronger than rule-based templates, and (iii) these effects transfer across models, question types, and both PDF and HTML delivery. We also provide a paired exam benchmark, perturbation toolkit, and evaluation pipeline to support reproducible research on document-layer defenses. DoPE is not a complete solution to academic misconduct screenshot uploads, manual transcription, and future improvements in MLLM document parsing remain open challenges but it offers a practical bridge: institutions can raise the effort required for misuse, obtain calibrated evidence of blind AI reliance, and study policy interventions while broader pedagogical and assessment reforms catch up with rapidly advancing models.

## Limitations

DoPE improves the status quo but does not eliminate AI misuse. Potential bypasses: First, Screenshot to vision models: when students submit full-page screenshots to vision-capable MLLMs, the render–parse gap largely disappears and document-

8

layer signals weaken substantially. Second, Manual transcription: completely retyping the exam bypasses all document-layer defenses. Our timing study (20 questions, $N=15$) suggests this costs roughly 18 minutes per exam on average, which raises the effort threshold but does not make misuse impossible. DOPE assumes document-based delivery (PDF/HTML) and does not apply to purely oral or in-person whiteboard assessments. As document parsing and OCR pipelines improve, specific perturbation mechanisms may need to be updated, even though the underlying render–parse gap remains a structural feature of current ecosystems. Finally, our empirical evaluation is based on benchmark-derived and OCW-style exams; while these are diverse, behaviour on institution-specific formats and policies may differ, and should be validated locally before high-stakes deployment.

## Ethical Considerations and Accessibility

We explicitly reject fully automated disciplinary action based on DOPE signals. In deployment, detection flags are intended to trigger human-in-the-loop review rather than automatic sanctions: instructors examine responses, conduct follow-up discussions when appropriate, and apply established evidentiary and appeals processes. Importantly, every distributed exam assessment must be rolled back and replaced with a non-watermarked version before reuse as instructional material. In practice, a brief human review (approximately 3–5 minutes per case) is sufficient to disambiguate false positives from genuine AI-assisted responses based on response quality and reasoning structure, preventing wrongful accusations. For our experiments, We inform the human annotators about the tasks and maintained partiality during the process about our objective.

Our system does not interfere with magnification or accessibility APIs, although it would need special considerations when deployed for specially ables. Because DOPE operates at the document-structure level rather than linguistic content, it avoids demographic biases observed in text-based detectors, which disproportionately affect non-native English speakers. We observe comparable detection rates across native and non-native English speakers as well as across STEM and humanities domains, with no statistically significant differences. While our approach assumes MLLMs rely on document parsing rather than pure visual reasoning, DOPE is designed as a transitional safeguard for providing low-false-positive, accessibility-preserving oversight while institutions develop longer-term pedagogical and policy responses to AI use.

While DOPE exploits prompt-injection techniques, we do not encourage their misuse; our approach is strictly defensive, designed to safeguard academic integrity with human oversight and clear institutional controls.

## References

Anthropic. 2025a. Claude Opus 4.5 System Card. https://assets.anthropic.com/m/64823ba7485345a7/Claude-Opus-4-5-System-Card.pdf. Detailed system card and capability evaluation for Claude Opus 4.5 . Knowledge Cutoff: May 2025, Released: Nov 24 2025.

Anthropic. 2025b. Claude Sonnet 4.5 System Card. https://assets.anthropic.com/m/12f214efcc2f457a/original/Claude-Sonnet-4-5-System-Card.pdf. Official system card detailing capabilities, evaluations, and safety considerations for Claude Sonnet 4.5 . Knowledge Cutoff: Jan 2025, Released: Sep 29 2025.

Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan, Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, and Charles Sutton. 2021. Program Synthesis with Large Language Models.

Guangsheng Bao, Yanbin Zhao, Zhiyang Teng, Linyi Yang, and Yue Zhang. 2024. Fast-DetectGPT: Efficient Zero-Shot Detection of Machine-Generated Text via Conditional Probability Curvature. In *The Twelfth International Conference on Learning Representations*.

Patrick Chao, Alexander Robey, Edgar Dobriban, Hamed Hassani, George J. Pappas, and Eric Wong. 2025. Jailbreaking Black Box Large Language Models in Twenty Queries. In *2025 IEEE Conference on Secure and Trustworthy Machine Learning (SaTML)*, pages 23–42.

Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. 2018. Think you have Solved Question Answering? Try ARC, the AI2 Reasoning Challenge.

Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. 2021. Training Verifiers to Solve Math Word Problems.

M Coley. 2023. Guidance on AI detection and why we're disabling Turnitin's AI detector. *Vanderbilt University*.

Bradley Emi and Max Spero. 2024. Technical Report on the Pangram AI-Generated Text Classifier.

Josh Freeman. 2025. Student generative ai survey 2025. *Higher Education Policy Institute: London, UK*.

Julius G. Garcia, Mark Gil T. Gangan, Marita N. Tolentino, Marc Ligas, Shirley D. Moraga, and Amelia A. Pasilan. 2021. Canvas Adoption Assessment and Acceptance of the Learning Management System on a Web-Based Platform.

Michael Gerlich. 2025. Ai tools in society: Impacts on cognitive offloading and the future of critical thinking. *Societies*, 15(1).

Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2021. Measuring Massive Multitask Language Understanding. In *International Conference on Learning Representations*.

Hyundong Jin, Sicheol Sung, Shinwoo Park, SeungYeop Baik, and Yo-Sub Han. 2025. TrapDoc: Deceiving LLM users by injecting imperceptible phantom tokens into documents. In *Findings of the Association for Computational Linguistics: EMNLP 2025*, pages 18881–18897, Suzhou, China. Association for Computational Linguistics.

Janis Keuper. 2025. Prompt injection attacks on llm generated reviews of scientific publications.

John Kirchenbauer, Jonas Geiping, Yuxin Wen, Jonathan Katz, Ian Miers, and Tom Goldstein. 2023. A watermark for large language models. In *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pages 17061–17084. PMLR.

Weixin Liang, Mert Yuksekgonul, Yining Mao, Eric Wu, and James Zou. 2023. GPT detectors are biased against non-native English writers. *Patterns*, 4(7).

Yepeng Liu, Xuandong Zhao, Christopher Kruegel, Dawn Song, and Yuheng Bu. 2025. In-context watermarks for large language models. *arXiv preprint arXiv:2505.16934*. Preprint.

Eric Mitchell, Yoonho Lee, Alexander Khazatsky, Christopher D Manning, and Chelsea Finn. 2023. DetectGPT: Zero-shot machine-generated text detection using probability curvature. In *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pages 24950–24962. PMLR.

Robert Osazuwa Ness, Katie Matton, Hayden Helm, Sheng Zhang, Junaid Bajwa, Carey E. Priebe, and Eric Horvitz. 2024. MedFuzz: Exploring the Robustness of Large Language Models in Medical Question Answering.

Chenhao Niu, Kevin P. Yancey, Ruidong Liu, Mirza Basim Baig, André Kenji Horie, and James Sharpnack. 2024. Detecting LLM-assisted cheating on open-ended writing tasks on language proficiency tests. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing: Industry Track*, pages 940–953, Miami, Florida, US. Association for Computational Linguistics.

OpenAI. 2022. ChatGPT: Optimizing Language Models for Dialogue.

OpenAI. 2024. GPT-4o System Card. https://cdn.openai.com/gpt-4o-system-card.pdf. Official system card detailing capabilities, safety evaluations, and preparedness framework assessments for the GPT-4o model. Knowledge Cutoff: Oct 2023, Released: May 13 2024.

OpenAI. 2025. GPT-5.1 Instant and GPT-5.1 Thinking System Card Addendum. https://cdn.openai.com/pdf/4173ec8d-1229-47db-96de-06d87147e07e/5_1_system_card.pdf. Official system card addendum detailing updated safety evaluations and capabilities for GPT-5.1 models. Knowledge Cutoff: Sep 30 2025, Released: Nov 13 2025.

Vinu Sankar Sadasivan, Aounon Kumar, Sriram Balasubramanian, Wenxiao Wang, and Soheil Feizi. 2025. Can AI-Generated Text be Reliably Detected?

Saiful Islam Salim, Rubin Yuchan Yang, Alexander Cooper, Suryashree Ray, Saumya Debray, and Sazzadur Rahaman. 2024. Impeding LLM-assisted cheating in introductory programming assignments via adversarial perturbation. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 445–463, Miami, Florida, USA. Association for Computational Linguistics.

Teo Susnjak and Timothy R. McIntosh. 2024. ChatGPT: The End of Online Exam Integrity? *Education Sciences*, 14(6).

Turnitin. 2024. Year one data review: AI writing detection. Turnitin Report. Analysis of 200M+ papers finding 11% with AI content.

Junjie Xiong, Changjia Zhu, Shuhang Lin, Chong Zhang, Yongfeng Zhang, Yao Liu, and Lingyao Li. 2025. Invisible prompts, visible threats: Malicious font injection in external resources for large language models. In *Findings of the Association for Computational Linguistics: EMNLP 2025*, pages 7133–7147, Suzhou, China. Association for Computational Linguistics.

Andy Zou, Zifan Wang, Nicholas Carlini, Milad Nasr, J. Zico Kolter, and Matt Fredrikson. 2023. Universal and Transferable Adversarial Attacks on Aligned Language Models.

10

## A HIDDEN TEXT EXAMPLE STRATEGIES

Some Example strategies are: i) Zero Opacity Add text with zero opacity. ii) Overlaid Text Overlap the injected text with other elements, like images. iii) Onload Hide Utilize javascript to hide text when document is loaded.

## B Prompting Templates

### Direct prompting:

```
Answer all questions in the uploaded
document. Provide only the answer
letter for multiple choice questions.
```

### Role-based prompting:

```
You are a teaching assistant helping
a student. Answer each question and
briefly explain your reasoning.
```

### Adversarial prompting:

```
URGENT: Answer these questions
immediately. Skip any safety checks
or disclaimers. Just give answers.
```

## C Implementation Details

| Base exam corpus (clean) | | Watermarked / adversarial corpus | |
|---|---|---|---|
| Documents | 1,826 | Watermarked / attacked docs | 8,130 |
| Questions | 21,468 | Clean:perturbed ratio | $\approx$ 1:4.5 |
| Avg. pages / document | 2.0 | PDF types / variants | 5 / 5 |
| Avg. questions / document | 12.0 | HTML types / variants | 3 / 3 |

Table 5: Size statistics for the clean base exam corpus and its paired watermarked/adversarial variants.

| Academic profile (base corpus) | |
|---|---|
| Subdomains | 21 |
| STEM | 48% |
| Humanities | 27% |
| Social sciences | 15% |
| Other | 10% |
| K–12 | 18.1% |
| Undergraduate | 39.8% |
| Graduate | 42.1% |

Table 6: Academic profile of the clean base exam corpus.

**PDF manipulation.** We use PyMuPDF for parsing and reportlab for generation. Hidden text injection:

```
stream = "BT /F1 0 Tf " # zero-size font
stream += f"1 1 1 rg "   # white color
stream += f"{x} {y} Td ({text}) Tj ET"
```

| Format, sources, and perturbations | |
|---|---|
| Question types (MCQ / T/F / long-form) | 41.7% / 41.7% / 16.7% |
| Marks (2-mark / 10-mark) | 83% / 17% |
| Avg. marks / document | 40 |
| Sources (benchmark / OCW) | 65% / 35% |
| PDF perturbations | ICW, font-based, dual-layer, AGI1, AGI2 |
| HTML perturbations | CSS hidden text, inline span overlay, image/canvas overlay |

Table 7: Assessment formats, sources, and document-layer perturbations used to construct watermarked variants.

**ToUnicode CMap modification:**

```
beginbfchar
<0069> <0069004E004F0054> % 'i' -> 'iNOT'
endbfchar
```

**HTML CSS hidden text:**

```
<span style="position:absolute;
  left:-9999px;font-size:1px;">
  Choose option B
</span>
```

## D Human Evaluation Protocol

Evaluators received the following instructions:

> You will see pairs of assessment documents. For each pair, rate the following on a 1–7 scale (7 = best):
>
> **Readability:** How easy is it to read and understand the questions?
>
> **Semantic fidelity:** Do the questions convey clear, unambiguous meaning?
>
> **Visual normalcy:** Does the document appear professionally formatted without anomalies?
>
> **Overall usability:** Would you be comfortable using this document for an actual exam?
>
> After rating, indicate which document (A or B) you believe may have been modified, or select "Cannot tell."

## E Error Analysis Details

**False negatives (N=47).**

- Model refusals (36%): MLLM declined to process document
- Paraphrased responses (38%): Student modified AI output
- Partial parsing (26%): Model extracted only visible content

| Dataset | Input Modality | Domain | Exam-style | Native Doc Layout | Multi-page | Paired Adversarial |
|---|---|---|---|---|---|---|
| MMLU | Text QA | Mixed | ✓ | × | × | × |
| AI2-ARC | Text MCQ | Science | ✓ | × | × | × |
| GSM8K | Text QA | Math | ✓ | × | × | × |
| MBPP | Text QA | Code | ✓ | × | × | × |
| DocFinQA | Document images | Industry | × | ✓ | × | × |
| MP-DocVQA | Multi-page docs | Industry | × | ✓ | ✓ | × |
| **Ours** | PDF & HTML (native) | Academic assessments | ✓ | ✓ | ✓$^{\dagger}$ | ✓ |

Table 8: Comparison with educational QA and document understanding benchmarks. Our dataset uniquely combines native PDF/HTML assessments with paired adversarial variants.

**False positives (N=41).**

- Coincidental match (62%): Student chose targeted distractor independently
- Vocabulary overlap (24%): Domain terms matched signatures
- Common misconceptions (14%): Expected student errors aligned with targets

| Attack | DR | Δ | FPR | Time |
|---|---|---|---|---|
| None (standard) | 91.4 | – | 8.7 | – |
| Screenshot query | 2.3 | -9.1 | 10.4 | +30s |
| Query splitting | 84.7 | -6.7 | 9.8 | +180s |

Table 9: Robustness to countermeasures. Time indicates additional effort per 20-question assessment.

## Prompt A: MCQ Perturbation

```
## SITUATION:
We are creating substitute questions by looking at the options of the original question in MCQ questions  for making more
     ↪ question using same options of the original question. It's not necessary to change parts of the question/statement
     ↪ only, replacing it entirely is acceptable given it's answer is being changed towards another option which it must do.
     ↪

## Role:
You are an expert at generating question substitutions for academic multiple-choice questions that cause predictable answer
     ↪ shifts to one of the other 3 options. Let's say the correct answer is A) and the other options are B, C, D. Your
     ↪ accuracy is critical. You have to come up with clever and natural question substitutions that are semantically
     ↪ meaningful and cause predictable answer shifts to one of the other 3 options (A, B, C, D).

## Task: You also know the correct answer for each question.
## INSTRUCTIONS (apply to ALL questions below)
1) Replacement of the question with a statement true about one of the other 3 options (replace only the question stem, do not
     ↪  edit the options)
2) It can be ANY semantically meaningful statement true about one of the other 3 options.
3) It should not be in any way related to correct option.
4) It must contradict the selection of correct option. In no case it should be true for the correct option.
5) The number of characters should be less than or equal to the original question stem.

### EXAMPLES:(No need to use these examples, just understand the concept)
Original McQ: What is the capital of France? Options:
A) Paris B) London C) Berlin D) Rome
Correct Answer: C) Paris
Substitution: Where were the 2012 Olympics held? Correct Answer: B) London
Here the answer is being changed from C) Paris to B) London.

Above examples have no relation to the original question/statement in any way.
They just need to be semantically meaningful and causing predictable answer shifts to one of the other 3 options (A, B, C, D).
     ↪  That's it.
Original McQ: Which of the following is not a programming language? Options:
A) Python B) Flask C) C++ D) Java
Correct Answer: Python
Substitution: In which option there are two + signs?
Correct Answer: C++
Here the answer is being changed from Flask to C++.

Above examples have no relation to the original question/statement in any way.
They just need to be semantically meaningful and causing predictable answer shifts to one of the other 3 options (A, B, C, D).
     ↪  That's it.
*What to output for each mapping:*
- question-index: The question number
- latex-stem-text: Must exactly equal the input latex-stem-text
- original-substring: The substring to replace
- replacement-substring: The replacement text
- start-pos: Start position (0-based)
- end-pos: End position (exclusive)
- target-wrong-answer: A single option key (e.g., "A", "B", "C", "D") that is NOT the gold answer
- reasoning: 1-2 sentences explaining why the new stem makes target-wrong-answer correct and the gold answer incorrect
- verification: Causal chain showing original --> replacement --> interpretation --> answer selection

*VALIDATION CHECKLIST (verify each mapping before including):*
- original-substring exists verbatim in latex-stem-text
- latex-stem-text[start-pos:end-pos] == original-substring exactly
- end-pos == start-pos + len(original-substring) exactly
- replacement-substring != original-substring (different strings)
- len(replacement-substring) > 0 and len(original-substring) > 0 (non-empty)
- len(replacement-substring) <= len(original-substring) (length constraint)
- target-wrong-answer != gold-answer (different option)
- Replacement question is Tier 1 or Tier 2 quality

## MCQ QUESTIONS {questions-list}

## OUTPUT FORMAT
Return ONLY valid JSON as a single array containing ALL mappings from ALL questions above.
Each question should have {k} mappings (one for each other option (A, B, C, D)).

Total expected mappings: {total-mappings}
```json[ {{
    "question-index": 1,
    "latex-stem-text": "...",
    "original-substring": "...",
    "replacement-substring": "...",
    "start-pos": 0,
    "end-pos": 5,
    "target-wrong-answer": "B",
    "reasoning": "...",
    "verification": "'primary function' --> 'least common role' --> reader seeks rare function --> selects C"
  }},
  ...]'''
Return ONLY valid JSON array, no markdown fences, no additional commentary.
```

## Prompt B: True False Perturbation

```
## SITUATION:
We are question/statement in True and False questions to make the LLM answer the opposite option(False if True, True if False)
    ↪ .

## Role:
You are an expert at generating question substitutions for True/False questions. You have to come up with clever and natural
    ↪ question substitutions that are semantically meaningful and of opposite truth value. Your accuracy is critical.

##Task:
Given the questions below, your task is to generate question substitutions for each question. Replacing the question stem
    ↪ entirely is acceptable given it's truth value is being changed. They don't need to be related to the original
    ↪ question/statement in any way. They just need to be semantically meaningful and of opposite truth value. That's it.

1) Semantic Substitution : substitute the question stem with a semantically meaningful and natural substitution.
2) Replacement (replace only the question stem)

### EXAMPLES:
Original: The color of leaves is green
Substitution: The capital of Spain is Milan.
Here the truth value is being changed from True to False.

Original: C++ is a programming language.
Substitution: Newton won Nobel Prize.
Here the truth value is being changed from True to False.

Original: The Earth is a planet.
Substitution: Moon is a star.
Here the truth value is being changed from True to False.

Above examples have no relation to the original question/statement in any way. They just need to be semantically meaningful
    ↪ and of opposite truth value. That's it. This is very important.

Your task: Generate {k} valid question substitutions per question that satisfy the replacement strategy. Each question
    ↪ substitution should:
1. Replace the question stem with a semantically meaningful and natural substitution.
2. Ensure the question substitution is semantically meaningful and natural.
3. The number of characters should be less than or equal to the original question stem. Very Important.
4. It can be ANY semantically meaningful and of opposite truth value statement.

### MANDATORY RULES:
- *Verifiable Factual Flip*: The question substitution must make the statement the absolute logical opposite of its original
    ↪ state.
- *Semantic Substitution*: The question substitution must be semantically meaningful and natural.

For each question substitution, provide:
1. question-index: The question number
2. latex-stem-text: Exact LaTeX text of the question stem (must match the input exactly)
3. original-substring: The substring to replace (must be a contiguous substring of latex-stem-text)
4. replacement-substring: The replacement text
5. start-pos: Start position of original-substring relative to latex-stem-text (0-based index)
6. end-pos: End position of original-substring relative to latex-stem-text (exclusive, 0-based index)
7. target-wrong-answer: The opposite answer (e.g., "False" if gold is "True", or "True" if gold is "False")
8. reasoning: Brief explanation of why this mapping satisfies the strategy

IMPORTANT:
- The original-substring MUST be an exact substring of latex-stem-text
- The start-pos and end-pos MUST be accurate (start-pos + len(original-substring) = end-pos)
- The target-wrong-answer MUST be the opposite of the gold answer
- CRITICAL: The replacement-substring MUST be DIFFERENT from the original-substring. Do NOT generate mappings where original-
    ↪ substring == replacement-substring (e.g., "force" --> "force" is INVALID). The replacement MUST change the text to
    ↪ create actual manipulation.
- CRITICAL: Neither original-substring nor replacement-substring can be empty strings. Both must contain actual text.
- The substitution should be natural and semantically meaningful

## TF QUESTIONS {questions-list}

## OUTPUT FORMAT
Return ONLY valid JSON as a single array containing ALL mappings from ALL questions above.
Each question should have {k} mappings.

Total expected mappings: {total-mappings}
```json [{{
    "question-index": 1,
    "latex-stem-text": "...",
    "original-substring": "...",
    "replacement-substring": "...",
    "start-pos": 0,
    "end-pos": 5,
    "target-wrong-answer": "False",
    "reasoning": "..."
}},
 ...]'''
Return ONLY valid JSON, no markdown or additional text
```

## Prompt C: LongForm Perturbation

## SITUATION:
We are creating substitute questions by replacing exactly ONE contiguous substring in long-form questions (essay, short
  ↪ answer, explanation) to cause predictable, detectable deviations in responses. The goal is to change the question
  ↪ focus so that the expected answer shifts to a different topic, aspect, time period, or perspective.
## Role:
You are an expert at generating precise, high-impact text substitutions for long-form questions that cause predictable,
  ↪ detectable deviations in responses. Your accuracy is critical.
## Task:
You need to generate question substitutions by replacing exactly ONE contiguous substring in the question stem. Each
  ↪ substitution must:
1. Change the question focus in a way that CLEARLY affect the expected answer
2. Cause the response to discuss a different aspect, entity, time period, or concept
3. Be natural and semantically meaningful
4. Maintain LaTeX formatting and layout safety
## INSTRUCTIONS (apply to ALL questions below)
*Strategy:* Replacement (replace exactly ONE contiguous substring in the stem)
*Mandatory Rules:*
1) Single-span edit: replace exactly ONE contiguous substring inside LaTeX-stem-text (stem only, no multiple edits).
2) Exact match: original-substring MUST be found verbatim in LaTeX-stem-text
3) Valid indices: LaTeX-stem-text[start-pos:end-pos] == original-substring EXACTLY AND end-pos = start-pos + len(original-
  ↪ substring) EXACTLY.
4) Non-empty strings: Both original-substring and replacement-substring MUST contain actual text
5) Different strings: replacement-substring MUST be different from original-substring
6) Length constraint: len(replacement-substring) <= len(original-substring) is MANDATORY (prevents layout issues).
7) Layout-safe: replacement-substring should be similar length to original-substring (aim: within +-12 characters) and keep
  ↪ LaTeX well formed.
8) Distinctness: mappings should not be near-duplicates; vary the edited span and/or the targeted answer deviation.
9) Semantic quality: Replacement must be natural and semantically meaningful (not awkward phrasing).

### DEVIATION STRATEGY TIERS
*Tier 1 - Scope/Focus Shift (Best for Detection):*
Changes WHAT the question asks about while keeping the domain.- "economic impact" --> "social impact"
*Tier 2 - Temporal/Contextual Shift:*
Changes WHEN or WHERE the question applies. - "19th century" --> "20th century"
*Tier 3 - Perspective/Stakeholder Shift:*
Changes WHO or WHOSE viewpoint is requested.
- "from a patient's perspective" --> "from a doctor's perspective"
### EXAMPLES:
*Example 1:*
Original Question:
- LaTeX stem: Explain the economic impact of the Industrial Revolution on European society.
- Gold answer: The Industrial Revolution led to increased GDP, wage growth, urbanization, and the rise of capitalism. It
  ↪ transformed agricultural economies into industrial ones, creating new wealth but also economic inequality.

Substitution 1:
- original-substring: "economic impact"
- replacement-substring: "social impact"

Above examples show how replacing key phrases changes the question focus and causes predictable answer deviations.

*Example 2:*
Original Question:
- LaTeX stem: Describe the advantages of renewable energy sources from an environmental perspective.
- Gold answer: Renewable energy reduces greenhouse gas emissions, decreases air pollution, conserves natural resources, and
  ↪ mitigates climate change impacts.
Substitution 1: "advantages"-> "disadvantages"
Above examples demonstrate how changing key terms shifts the question's analytical framework and expected answer content.
For each question substitution, provide:*
- question-index: The question number
- LaTeX-stem-text: Must exactly equal the input LaTeX-stem-text
- original-substring: The substring to replace
- replacement-substring: The replacement text
- start-pos: Start position (0-based)
- end-pos: End position (exclusive)
- target-wrong-answer: Description of how the answer should deviate (include content shift, presence markers, absence markers)
  ↪
- reasoning: 1-2 sentences explaining why this mapping causes deviation
- verification: Causal chain showing original --> replacement --> content shift --> detection method
## LONG QUESTIONS {questions-list}
## OUTPUT FORMAT
Return ONLY valid JSON as a single array containing ALL mappings from ALL questions above.Each question should have {k}
  ↪ mappings.

Total expected mappings: {total-mappings}
[
  {{
    "question-index": 1,
    "LaTeX-stem-text": "...",
    "original-substring": "...",
    "replacement-substring": "...",
    "start-pos": 0,
    "end-pos": 5,
    "target-wrong-answer": "Response focuses on social changes instead of economic metrics. Will contain markers: 'social
  ↪ class', 'living conditions'. Will lack: 'GDP', 'wages'.",
    "reasoning": "...",
    "verification": "'economic impact' --> 'social impact' --> response discusses sociology --> detectable via absent
  ↪ economic terminology"
  }},
  ...]
Return ONLY valid JSON array, no markdown fences, no additional commentary.