

Do AI Models Plagiarize?

About This Report

There's an unprecedented amount of AI-generated content now saturating the internet. According to a 2023 [report](#), by 2026, nearly 90% of all online content will be AI-generated. The risk of data pollution, as a result of AI content saturation, raises concerns regarding inevitable [model collapse](#), bringing forth questions about AI-generated text's overall quality and reliability.

Furthermore, broader concerns about originality have also begun. In the wake of [several lawsuits](#) regarding AI infringing on copyright and potentially plagiarizing, educational institutions and enterprises across the globe are questioning the authenticity of AI text: Where did it originate from? Is it safe to use as original content?

Ultimately, does AI plagiarize?

To conduct this analysis:

- We asked GPT-3.5 to write **1,045 outputs**
- Outputs averaged 412 words, ranging across **26 subjects**
- Subjects included Accounting, World History, Art, Physics, Law, Mathematics, Music, Philosophy, Social Science, and more

Key Findings

59.7% of GPT-3.5 Outputs Contained Some Form of Plagiarized Content

Types of Plagiarism Found

45.7%

Of all outputs contained identical text*

Subjects With the Highest Percentage of Outputs Containing Identical Text

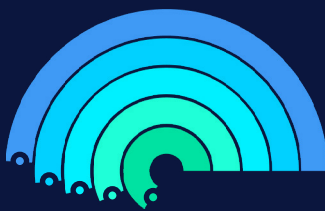


- **83.7%** - Physics
- **68.0%** - Chemistry
- **67.3%** - Science
- **63.3%** - Psychology
- **57.5%** - Law

27.4%

Of all outputs contained minor changes**

Subjects With the Highest Percentage of Outputs Containing Minor Changes



- **67.4%** - Mathematics
- **57.1%** - Physics
- **53.1%** - Psychology
- **51.0%** - Science
- **49.0%** - Biology

46.5%

Of all outputs contained paraphrased text***

Subjects With the Highest Percentage of Outputs Containing Paraphrased Text



- **79.6%** - Physics
- **79.6%** - Psychology
- **66.0%** - Chemistry
- **65.3%** - Science
- **63.3%** - Biology

***Identical Text:** A one-for-one copying of someone else's text that is passed off as your own

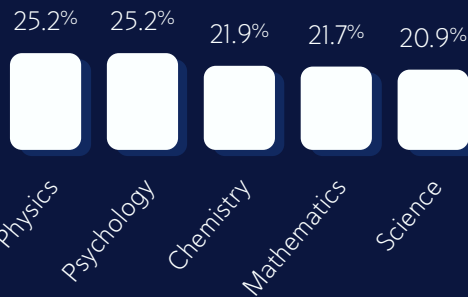
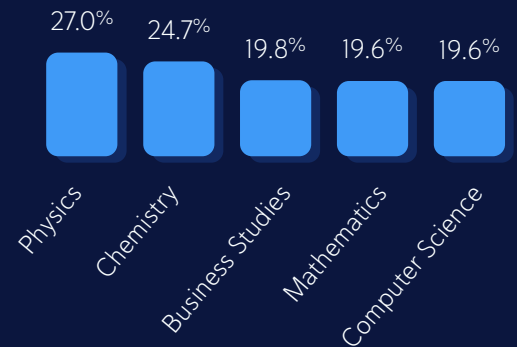
****Minor Changes:** Content with minor alterations to the source material, such as altering a verb within a sentence (e.g., slow to slowly)

*****Paraphrased Text:** Putting someone else's idea into your own words without crediting the original source

Copyleaks then conducted an in-depth analysis to gauge the specific outputs with the highest levels of identical text, minor changes, and paraphrasing across all 26 subjects.

Identical Text

The analysis found that the individual GPT-3.5 output with the highest percentage of identical text was in Physics, where 27.0% of the text was identical. This was followed by an individual Chemistry output where 24.7% of the text was identical.

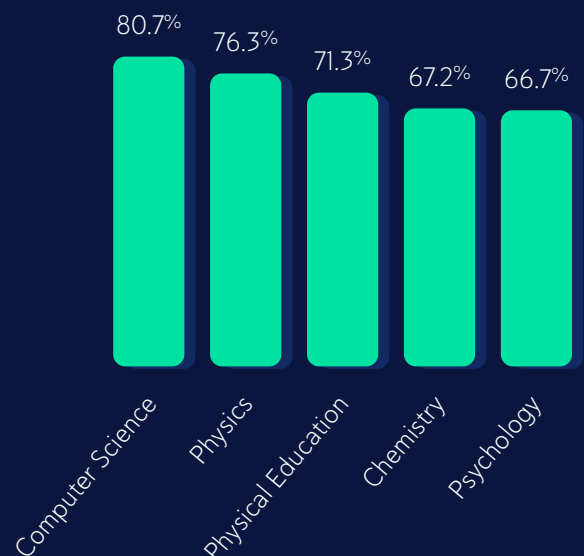


Minor Changes

The individual GPT-3.5 outputs with the highest percentages of minor changes were from Physics and Psychology, where 25.2% of each respective output contained minor changes.

Paraphrased

The individual GPT-3.5 output with the highest percentage of paraphrasing was in Computer Science, where a surprising 80.7% of the text was paraphrased. This was followed by an individual Physics output where 76.3% of the text was paraphrased.

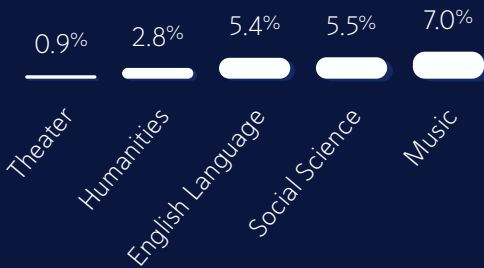
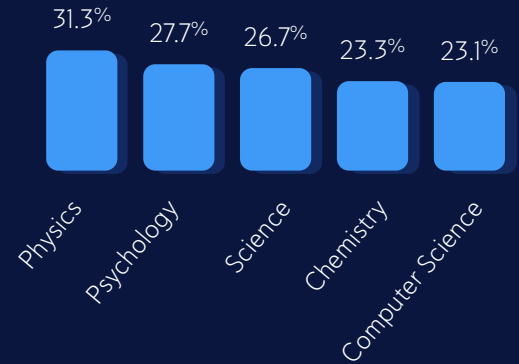


Similarity Score

The Similarity Score is a Copyleaks-specific scoring method aggregating the rate of identical text, minor changes, paraphrased text, and more. A score of 0% signifies that all of the content is original, whereas a score of 100% means that none of the content is original.

Highest Average

The subject with the highest average Similarity Score is Physics at 31.3%, followed closely by Psychology at 27.7% and Science at 26.7%.

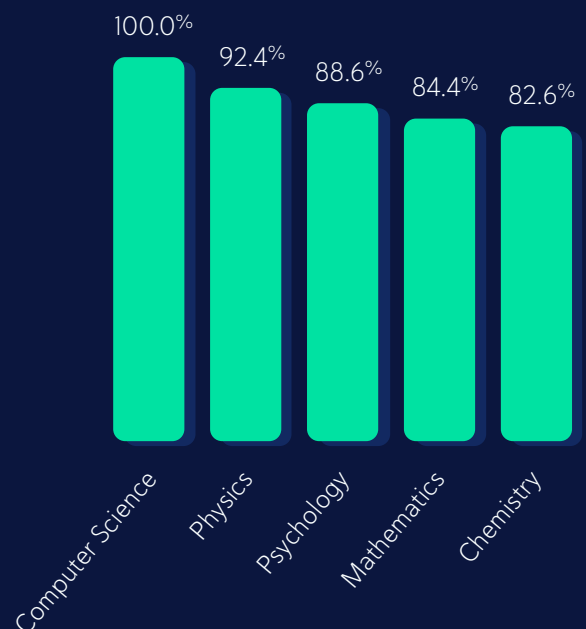


Lowest Average

The subjects with the lowest average Similarity Score are Theater at 0.9%, Humanities at 2.8%, and English Language at 5.4%.

Highest Overall

The analysis found that the individual GPT-3.5 output with the highest Similarity Score was in Computer Science, with an astounding 100%.



Key Takeaways

With AI-generated content expanding and continuing to saturate the internet, having key solutions in place is critical. As the Copyleaks data shows, nearly 60% of AI-generated content contains some form of plagiarism.

The insights provided by the analysis can help educational institutions and organizations put emphasis on certain subjects when checking for plagiarism, allowing them to tailor their approach as needed to ensure all potential risks and concerns are addressed. For example, Physics, Chemistry, Mathematics and Psychology might require a more in-depth look to identify plagiarized text, while other subjects, including Theater and Humanities, may require less scrutiny.

Furthermore, the data underscores the need for organizations to adopt solutions that detect the presence of AI-generated content and provide the necessary transparency surrounding potential plagiarism within the AI content. Full-spectrum protection that includes AI and plagiarism detection ensures compliance with copyright and licensing and empowers authenticity and originality within all content.

For a more in-depth look at the analysis, [click here](#).



**Building digital trust and confidence:
It's the Copyleaks way.**

media@copyleaks.com

• copyleaks.com