

All Our TeX Source Are Not Belong to You: PDF Obfuscation against arXiv’s TeX Source Policies

Zhongtang Luo
Purdue University
luo401@purdue.edu

Jianting Zhang
Purdue University
zhan4674@purdue.edu

Abstract—We present a black-box study of arXiv’s PDF classifier, which we nickname TeXRay, and show that arXiv’s policy of distinguishing TeX-generated PDFs from non-TeX PDFs is both practically fragile and conceptually unsound. Motivated by arXiv’s requirement that authors submit TeX source when TeX-generated PDFs are detected, we reverse-engineer TeXRay’s decision factors and find it relies primarily on a small set of surface features (PDF metadata fields and embedded font identifiers). Building on this understanding, we develop two complete bypasses: a tiny LaTeX macro that produces TeX-compiled PDFs that evade detection, and a 100-line Python post-processor that edits compiled PDFs to remove or alter the features TeXRay checks for. To show the limits of even a perfect detector, we implement a prototype obfuscation system that rewrites TeX-generated content streams to mimic PDFs produced by other toolchains (for example, Typst) while preserving near-identical visual output. These results demonstrate that distinguishing PDFs by their typesetting origin is inherently unreliable.

I. INTRODUCTION

arXiv is a curated research-sharing platform with over two million scholarly articles [1]. It is the de-facto online manuscript pre-print service for computer science, mathematics, physics, and many interdisciplinary communities. Since 1991, arXiv has offered a place for researchers to reliably share their work as it undergoes the process of peer-review, and for many researchers it is their primary source of literature [2].

The submission system of arXiv allows any authorized user to upload their manuscripts, which are then processed and hosted on arXiv’s website for free public access. arXiv has placed numerous policies on the submission system to ensure the quality and accessibility of the hosted manuscripts, including accepting only submissions from endorsed users, adopting a moderation process before hosting, and various formatting requirements [1].

Among them, arXiv has a standing policy to force authors to submit TeX source files when available, and runs a built-in TeX compiler to produce the final PDF hosted on arXiv [3]. It claims three reasons for not accepting PDF-only submissions when TeX files are available: (1) to process submissions into accessible formats such as HTML; (2) because source files are of higher archival value; and (3) source files offer greater insight (in TeX rendering techniques).

However, the policy is not without its debates. The built-in compiler has limited functionalities and has been complained as buggy [4]. In arXiv’s own accessibility study [5], 25% of the respondents cite requiring submitting TeX as a barrier

to submitting accessible papers. Moreover, such policy raises privacy concerns, as sensitive comments in the TeX source files may be accidentally released [6]. While arXiv claims that these concerns do not pose a significant problem [3], there are dedicated bots on X that scrape comments from TeX source files [7]. More recent criticism also concerns about potential data leaks to AI training datasets [8].

Moreover, not all papers use TeX/LaTeX as their typesetting system. Many fields not purely in STEM, such as computer science education, also feature a significant number of Word documents [9]. Furthermore, new technology stacks that focus on typesetting papers, such as Typst [10], are emerging as alternatives to TeX. Given that arXiv’s built-in compiler only supports pdflatex, one of TeX’s three main engines, as of now, and the vast amount of work needed to support every other typesetting system, as a compromise, arXiv currently allows authors to upload PDF files, as long as they are not generated from TeX source files.

To achieve this goal, arXiv runs a built-in checker, which we nickname TeXRay, to detect if any uploaded PDF file is generated from a TeX engine. Once detected, arXiv will block the upload and force the author to submit TeX source files instead. We give an example of such a TeXRay error message in Figure 1.

In this paper, we conduct a detailed black-box analysis of TeXRay that outlines its detection mechanisms. We discover that TeXRay is not a reliable detector of TeX-generated PDF files, and produce a 100-line Python script that completely bypasses its detection.

Furthermore, we argue that arXiv’s goal of distinguishing TeX-generated PDF files from other typesetting systems is *fundamentally unsound*. We stress that there is no intrinsic difference between TeX-generated PDF files and PDF files generated from other typesetting systems, such as Typst. To demonstrate this, we develop a prototype obfuscation system that transforms a TeX-generated PDF file to an obfuscated PDF file that looks like a Typst-generated PDF file by rewriting the content streams, while maintaining a visually near-identical output. This suggests that even if TeXRay were to be made more robust, it would still be possible to bypass it with more sophisticated obfuscation techniques.

Contribution. As our first contribution, we present a black-box analysis of TeXRay. We identify the key features that

test.pdf appears to have been produced by TeX

This file has been rejected as part your submission because it appears to be pdf generated from TeX/LaTeX source. For the reasons outlined at in the [Why TeX FAQ](#) we insist on submission of the TeX source rather than the processed version.

Our software includes an automatic TeX processing script that will produce PDF, PostScript and dvi from your TeX source. If our determination that your submission is TeX produced is incorrect, you should send e-mail with your submission number to [arXiv administrators](#).

Fig. 1. An example of TeXRay error message when uploading a PDF directly generated from TeX source files.

TeXRay uses for detection, including PDF metadata fields and embedded font information.

As our second contribution, we present numerous methods to bypass TeXRay and demonstrate its vulnerabilities. We show that we can either modify the LaTeX source files with a simple macro, or post-process the compiled PDF file with a Python script, to completely bypass TeXRay’s detection.

As our third contribution, we argue that arXiv’s goal of distinguishing TeX-generated PDF files from other typesetting systems is fundamentally unsound: we develop a prototype obfuscation system that transforms a TeX-generated PDF file to a Typst-generated PDF file.

Responsible Disclosure. We have reported our findings to arXiv on October 19, 2025. arXiv responded on October 21, assuring us that they did not consider our finding a security risk, as they actually “(do) not require TeX source in all situations,” and “there is always still a human in the loop (to perform moderation),” and allowed us to release what we have found. As such, we feel comfortable to release our findings.

II. TECHNICAL OVERVIEW

Analysis of TeXRay Mechanism. Based on our observations and experiments in Section IV, TeXRay appears to rely on two main detection mechanisms to identify TeX-generated PDF files:

- 1) **PDF Metadata Fields.** TeXRay inspects specific metadata fields in the PDF file to see if the PDF carries metadata from the TeX engine. Specifically, if the Producer or Creator field contains “tex” without being followed by the letter “t” (which forms “text”), TeXRay flags the PDF as TeX-generated.
- 2) **Embedded Font Information.** TeXRay analyzes the embedded font information within the PDF file, looking for font names and characteristics that are commonly associated with TeX-generated documents. Specifically, if the “BaseFont” field of any font contains “CMR” (Computer Modern Roman), TeXRay flags the PDF as TeX-generated.

Circumventing TeXRay. We show that our understanding of TeXRay’s detection mechanisms is complete by demonstrating two methods that can completely bypass TeXRay’s detection, as detailed in Section V:

- 1) **LaTeX Macro-based Bypass.** By incorporating a LaTeX macro that modifies the PDF metadata fields and uses non-TeX-specific fonts, we can compile TeX source files into PDF files that evade TeXRay’s detection.
- 2) **Post-processing-based Bypass.** By using a Python script to post-process the compiled PDF file, we can alter the relevant PDF metadata fields and embedded font information to avoid detection by TeXRay.

We give the full scripts in the respective sections.

TeX Translation. To argue that arXiv’s goal of distinguishing TeX-generated PDF files from other typesetting systems is fundamentally unsound and show that there are no intrinsic differences between PDF files from different engines, we develop a prototype obfuscation system that transforms a TeX-generated PDF file to a Typst-generated PDF file, as detailed in Section VI. We note that when the TeX source files are available, many modern toolchains (including LLMs) can directly convert TeX source files to other typesetting systems and then compile. Therefore, we focus on the scenario where only the TeX-generated PDF file is available. By modifying the content streams of the TeX-generated PDF file, we can produce an obfuscated PDF file that looks like a Typst-generated PDF file while maintaining a visually near-identical output. Figure 2 shows one such example from Section VI.

III. PRELIMINARIES

PDF File Structure. The PDF file format, specified as ISO 32000-2, mainly consists a list of objects that reference each other, loosely forming a tree structure. The Catalog object is the root of this tree, which points to other objects that define the structure of the PDF document, including the Pages object that points to the Page objects representing individual pages. Each Page object contains a Contents object and a Resources object. The Contents object holds the actual content stream of the page, including specific drawing commands for text, images, and graphics. The Resources object that defines the resources used on that page, such as fonts and images.

Additionally, the PDF file contains an Info object that holds metadata about the document, such as the author, title, and creation date.

It is worth noting that, since a document only uses a small subset of the glyphs in a font, PDF files often embed subsets

```

BT
/F31 9.96264 Tf
1 0 0 1 140.944 656.037 Tm [<00350049004A0054>-333<004A0
  ↪ 054>-333<0042>-333<0055004600540055>-333<00510042005
  ↪ 3004200480053004200510049000F>]TJ
1 0 0 1 303.509 96.112 Tm [<0012>]TJ
ET

1 0 0 -1 0 841.8898 cm
/d65gray cs
0 scn
/F31 9.96264 Tf
BT
1 0 0 -1 140.94 185.85 Tm
  ↪ [(\000\065\000\111\000\112\000\124\000\001\000\112\0
  ↪ 00\124\000\001\000\102\000\001\000\125\000\106\000\1
  ↪ 24\000\125\000\001\000\121\000\102\000\123\000\102\0
  ↪ 00\110\000\123\000\102\000\121\000\111\000\017)]
  ↪ TJ
1 0 0 -1 303.51 745.78 Tm [(\000\022)] TJ
ET

```

Fig. 2. An example of rewriting a TeX-generated content stream (top) to a Typst-style content stream (bottom) from Section VI.

of fonts to reduce file size. These embedded font subsets are typically given random prefixes in their names to avoid naming conflicts.

Figure 3 illustrates a simplified PDF structure generated from a LaTeX document. The PDF structure highlights the hierarchical relationships between various objects in the PDF file.

```

\documentclass{article}

\begin{document}

\section*{Test Document}
This is a test document to check the
  ↪ functionality of the LaTeX setup.
$a=b+c.$

\end{document}

Catalog (21 0 obj)
├─ Pages (8 0 obj)
│   └─ Page (2 0 obj)
│       └─ Contents (3 0 obj)
│           └─ "BT /F38 14.3462 Tf 133.768 707.125 Td
│               ↪ [(T)94(est)-375(Do)-31(cumen)31(t)]TJ/F28
│               ↪ 9.9626 Tf 0 -21.821 Td
│               ↪ [(This)-333(is)-334(a)-333(test)-333(do)-28(
│               ↪ cumen)28(t)-334(to)-333(c)28(hec)27(k)-333(t
│               ↪ he)-333(functionalit)28(y)-334(of)-333(the)-
│               ↪ 333(LaT)83(eX)-333(setup.)]TJ/F31 9.9626 Tf
│               ↪ 150.816 -21.918 Td [(a)]TJ/F28 9.9626 Tf
│               ↪ 8.033 0 Td [(=)]TJ/F31 9.9626 Tf 10.516 0 Td
│               ↪ [(b)]TJ/F28 9.9626 Tf 6.49 0 Td [(+)]TJ/F31
│               ↪ 9.9626 Tf 9.962 0 Td [(c:)]TJ/F28 9.9626 Tf
│               ↪ -16.452 -524.131 Td [(1)]TJ ET"
│           └─ Resources
│               └─ Font /F28 (5 0 obj, /MZCFWZ+CMR10)
│                   └─ Font /F31 (6 0 obj, /IOSIJS+CMMI10)
│                       └─ Font /F38 (4 0 obj, /LORVAG+CMBX12)
└─ Info (22 0 obj)

```

Fig. 3. A LaTeX document and its compiled PDF structure (simplified).

Preprint Repositories and arXiv.

are online platforms that allow researchers to share their manuscripts prior to formal peer review and publication in academic journals. These repositories facilitate rapid dissemination of research findings, enabling the scientific community to access and build upon new knowledge without the delays associated with traditional publishing processes.

arXiv is listed as one of the largest preprint repositories on Wikipedia [11] with >1,000,000 preprints for research manuscripts (“e-prints”) across physics, mathematics, computer science, and related disciplines. Established in 1991, it has become the primary venue for sharing preprints prior to journal publication, reaching over two million submissions by 2021 and currently receiving roughly 24,000 new papers per month. Although arXiv is not a peer-reviewed platform, each subject area is moderated to filter non-scientific or off-topic content. An endorsement mechanism, typically automatic for recognized academic authors, further helps maintain relevance within disciplines [1]. Many e-prints later appear in journals, yet arXiv itself has hosted landmark works that never underwent formal publication. A notable case is Grigori Perelman’s 2002 proof of the Poincaré conjecture, disseminated exclusively via arXiv but nonetheless acknowledged as one of the century’s most significant mathematical achievements [12].

arXiv has a list of policies governing submissions, among them a soft requirement to submit TeX files when available [1]:

Note: a PDF file created from a TeX/LaTeX file will typically be rejected, with exceptions granted on a case-by-case basis. There are good reasons why arXiv insists on TeX/LaTeX source if it is available. arXiv produces PDF automatically from all TeX submitted source.

On the other hand, other similar STEM preprint repositories, such as bioRxiv and Cryptology ePrint Archive, do not have such a requirement and ask for PDF-only submissions.

Typesetting Systems. Throughout the history of digital typesetting, multiple systems have been developed to facilitate the creation of high-quality documents. We introduce a few notable typesetting systems below.

TeX/LaTeX. TeX, created by Donald Knuth in 1978, is a typesetting system widely used in academia, particularly in mathematics, computer science, and physics. LaTeX, developed by Leslie Lamport in the 1980s, is a macro package built on top of TeX that simplifies document formatting and structuring. Most STEM field conferences and journals provide LaTeX templates for authors to use. While TeX supports compiling documents into various formats, including DVI and PostScript, PDF has become the dominant output format in recent years.

Word Processors. Microsoft Word, first released in 1983, has been the de facto standard word processing software since the 1990s. It provides a WYSIWYG (what you see is what you get) interface that allows users to create and format documents visually, without needing to understand the underlying markup

or code. Multiple interdisciplinary fields, such as human-computer interaction and computer science education, often accept Word-based submissions [9]. Word mainly uses a proprietary binary format (.doc) for older versions and an XML-based format (.docx) for newer versions. Multiple spinoff projects, such as LibreOffice Writer and WPS, also support the format. For export purposes, Word can generate PDF files directly from the application.

Markdown-based Systems. Modern typesetting systems often leverage lightweight markup languages like Markdown to simplify document creation. For instance, Typst [10], initially released in 2023, brands itself a new markup-based typesetting system for the sciences, and contains a few templates for academic papers. Typst uses its own .typst format for source files, and can export documents to PDF directly.

IV. ANALYSIS OF TEXRAY MECHANISM

We conduct our analysis of TeXRay by uploading PDF files with different features to see if arXiv rejects the file as being generated by TeX. Based on our empirical knowledge, we categorize the detection mechanism into two main components: the metadata analysis and the font analysis.

A. Metadata Analysis

Following ISO 32000-2, a PDF allows any number of key-value metadata fields, with a few common entries (title, author, etc.) standardized in the specification [13]. Moreover, it is well-known that TeX injects specific metadata into the PDF files it generates. We show an example of such metadata extracted from a PDF file generated by LuaLaTeX 2024 with pikepdf below.

```
"/Author": "",
"/CreationDate": "D:19791231160000-08'00'",
"/Creator": "LaTeX with hyperref",
"/Keywords": "",
"/ModDate": "D:19791231160000-08'00'",
"/PTeX.FullBanner": "This is LuaHBTeX,
Version 1.18.0 (TeX Live 2024/nixos.org)",
"/Producer": "LuaTeX-1.18.0",
"/Subject": "",
"/Title": "",
"/Trapped": "/False"
```

We observe that the metadata contains fields such as Creator, Producer, etc., which indicate the TeX engine used to generate the PDF. To understand which part of the metadata triggers the detection mechanism, we take a PDF file that otherwise passes the arXiv submission process and inject certain metadata fields. Then, we check if the injected file is rejected by arXiv. Table I summarizes our test results.

Based on the results, we hypothesize that if a PDF file contains “Creator” or “Producer” dictionary keys in its metadata, and the corresponding values of these keys contain “tex” without being followed by the letter “t” (as to form “text”), then the PDF file is detected by TeXRay.

Bypass. Given that the PDF metadata is freely editable both inside and outside the LaTeX system, there are many methods to bypass such detection mechanism. For instance, the following TeX command completely rewrites the PDF metadata.

TABLE I
RESULTS OF INJECTING DIFFERENT METADATA FIELDS INTO A PDF FILE AND CHECKING IF ARXIV’S TEXRAY DETECTS IT AS BEING GENERATED FROM TEX. ✓ INDICATES UNDETECTED, WHILE ✗ INDICATES DETECTED.

Field	Data	Undetected?
Creator	TeX	✗
Creator	ConTeXt	✓
Creator	texture	✓
Creator	plaintext	✓
Creator	pdfTeXt	✓
Creator	texttext	✓
Creator	texttex	✗
Creator	tex-text	✗
Creator	tex	✗
Creator	textex	✗
Creator	ConTeX	✗
Creator	TeX*	✗
Creator	*TeX	✗
Creator	LuaLaTeX	✗
Creator	pdfTeXs	✗
Creator	texas	✗
Creator	TeXas	✗
Creator	vertex	✗
Creator	paleocortex	✗
Producer	pdfTeX-1.40.26	✗
Producer	pdfConTeXt-1.40.26	✓
Producer	texture	✓
Producer	plaintext	✓
Producer	pdfTeXt	✓
Producer	texttext	✓
Producer	texttex	✗
Producer	tex-text	✗
Producer	pdfTeX-1.40.26	✗
Producer	pdfConTeX-1.40.26	✗
Producer	pdfTeX*-1.40.26	✗
Producer	pdf*TeX-1.40.26	✗
Producer	pdfLuaLaTeX-1.40.26	✗
Producer	TeX	✗
Producer	pdfTeXs	✗
Producer	texas	✗
Producer	TeXas	✗
Producer	vertex	✗
Producer	paleocortex	✗
Author	tex	✓
CreationDate	tex	✓
ModDate	tex	✓
PTeX.Fullbanner	tex	✓
Subject	tex	✓
Title	tex	✓
Trapped	tex	✓

```
\usepackage{hyperref}

% pdflatex
\ifdefined\pdfinfoomitdate
\pdfinfoomitdate=1
\fi
\ifdefined\pdfsuppressptexinfo
\pdfsuppressptexinfo=-1
\fi

% luatex
\ifdefined\pdfvariable
\pdfvariable suppressoptionalinfo 1023 \relax
\fi

\hypersetup{
  pdftitle = {},
```

name	type	encoding	emb	sub	uni	object	ID
HQWOEM+LMRoman12-Bold	CID Type 0C	Identity-H	yes	yes	yes	4	0
FSFNHQ+LMRoman10-Regular	CID Type 0C	Identity-H	yes	yes	yes	5	0
IOSIJS+CMMI10	Type 1	Builtin	yes	yes	no	6	0
RFLZJB+CMR10	Type 1	Builtin	yes	yes	no	7	0

Fig. 4. An example of embedded fonts in a LaTeX PDF file. LMRoman stands for Latin Modern Roman, and CM(MI/R) stands for Computer Modern (Math Italic/Regular).

```
pdfauthor = {},
pdfsubject = {},
pdfkeywords = {},
pdfcreator = {Not T3X},
pdfproducer = {Definitely not T3X}
}
```

Even when the TeX source is not directly available, plenty of tools allow editing the PDF metadata based on the PDF file alone. For instance, pdftk allows removing the metadata from any PDF file [8].

```
pdftk $PDFFILE dump_data | \
sed -e 's/\\(InfoValue:)\s.*\\1\\ /g' | \
pdftk $PDFFILE update_info - output clean-$PDFFILE
```

In our Python script, we used pikepdf to clear PDF metadata.

```
for key in list(pdf.docinfo.keys()):
    del pdf.docinfo[key]
```

B. Font Analysis

Being a portable format aimed for consistency, PDF allows embedding of fonts to ensure that the output is the same for every rendering software [13]. To reduce the size of the output, modern typesetting engines almost always include only a subset of the font that contains only the relevant glyphs used in the document, coded by “(random 6 capital letters) + (name of the font)”.

While Computer Modern, the default font used in LaTeX, is open-sourced and thus free to use in any typesetting software [14], its rarity still leads TeXRay to make the interesting decision that using it is the sure way to detect LaTeX submissions. We can observe the fonts embedded in a PDF file through pdffonts, a PDF font analyzer from the PDF rendering library Poppler [15]. Figure 4 shows the analysis of a sample PDF file generated by LaTeX.

Since PDF often only includes a different subset of glyphs, we infer that it is impractical for TeXRay to identify fonts based on the binary data, which can drastically change based on what glyphs are available. Instead, based on ISO 32000-2, each font carries a piece of metadata including FontName, FontFamily, and BaseFont, which contain the string information of the exact font being used in the PDF.

To isolate the relevant part that TeXRay uses for detection, we upload the same PDF to arXiv except that we include a font metadata field coming from Computer Roman, and record whether the PDF remains undetected from TeXRay. Based on our test results summarized in Table II, we conclude that TeXRay inspects the “BaseFont” field of each embedded font, and if any of them contains “CMR”, then the PDF file is detected as being generated from TeX. It is worth noting

that the font’s binary data itself, which consists of the font name, family name, and full name fields, does not trigger the detection.

TABLE II
RESULTS OF INJECTING DIFFERENT FONT METADATA FIELDS INTO A PDF FILE AND CHECKING IF ARXIV’S TEXRAY DETECTS IT AS BEING GENERATED FROM TEX. ✓ INDICATES UNDETECTED, WHILE ✗ INDICATES DETECTED.

Field	Data	Undetected?
Font’s FontName	QMHNZG+CMR10	✓
Font’s FamilyName	Computer Modern	✓
Font’s FullName	CMR10	✓
FontName	/QMHNZG+CMR10	✓
BaseFont	/QMHNZG+CMR10	✗
BaseFont	/QMHNZG+CMMI10	✓
BaseFont	/QMHNZG+CMBX10	✓

Bypass. Unfortunately for TeXRay, *Latin Modern*, the successor of Computer Modern with an identical look and improved support for Type 1 and Unicode, is included in TeX since 2003 and widely available [16]. It can be used in TeX through a simple package:

```
\usepackage{lmodern}
```

Furthermore, *New Computer Modern*, a font expanded from Latin Modern, is also publicly available and widely used in typesetting software such as Typst by default [10]. Therefore, TeXRay cannot simply mark the font as an indicator of the use of TeX. The font can be used in TeX through a simple package as well:

```
\usepackage[regular]{newcomputermodern}
```

Even without access to the TeX source and thus incapable of using new packages, we observe that every font residing in the PDF file can be extracted and replaced through utilities such as pikepdf. In our circumvention script, we use FontForge to edit the extracted font and make sure that the metadata remains undetected according to TeXRay.

V. CIRCUMVENTING TEXRAY

Based on our analysis of TeXRay above, we know that it is possible to bypass different parts of TeXRay’s detection with multiple simple techniques. In this section, we present two combined methods that can completely bypass TeXRay’s detection: a LaTeX macro that compiles with the PDF, and a Python post-processing script that modifies the relevant PDF data.

A. LaTeX Macro-based Bypass

Combining the LaTeX macro in both sections of Section IV gives us a macro that completely bypasses TeXRay’s detection. We give the full macro below:

```
\usepackage{hyperref}

% pdflatex
\ifdefined\pdfinfoomitdate
\pdfinfoomitdate=1
\fi
\ifdefined\pdfsuppressptexinfo
\pdfsuppressptexinfo=-1
\fi

% luatex
\ifdefined\pdfvariable
\pdfvariable suppressoptionalinfo 1023 \relax
\fi

\hypersetup{
  pdftitle = {},
  pdfauthor = {},
  pdfsubject = {},
  pdfkeywords = {},
  pdfcreator = {Not T3X},
  pdfproducer = {Definitely not T3X}
}

\usepackage[regular]{newcomputermodern}
```

B. Post-processing-based Bypass

In scenarios where modifying the LaTeX source is not possible, a post-processing script can be used to modify the PDF metadata and font information after compilation. We present a Python script that achieves this goal, available at <https://github.com/zhtluo/arxiv-scripts>. The script mainly fulfills the following two tasks:

- 1) It removes or alters the PDF metadata fields that TeXRay checks, such as Producer and Creator, to avoid detection.
- 2) It modifies the embedded font information to replace TeX-specific font identifiers with generic ones, preventing TeXRay from recognizing TeX-generated fonts.

VI. TEX TRANSLATION

When the LaTeX source code is available, translation to other typesetting systems like Typst is often straightforward through directly translating the source code to the other system. Several tools exist for this source-to-source conversion. For instance, pandoc offers robust document format conversion capabilities. The official Typst web application provides a built-in feature to import .tex files directly. Furthermore, recent advances in Natural Language Processing have enabled Large Language Models (LLMs) to perform direct translation from LaTeX to other systems with reasonable fidelity.

In this section, we focus on the scenario where the original .tex source is unavailable or otherwise cannot be easily translated, and only the final compiled PDF document is accessible. This paper focuses on this challenging context: transforming a PDF file to retain its (mostly) exact visual appearance while rewriting its underlying byte-level structure to be indistinguishable from a natively generated Typst document. This approach

bypasses the need for source code, operating directly on the rendered output.

A. Comparison of TeX and Typst PDF

We present a comparison between the PDF outputs generated by TeX and Typst for the same simple document in Figure 5.

Based on the comparison, we observe that the two PDF files share similar high-level structures, but differ significantly in their content streams. Therefore, we focus our effort on rewriting the content streams to match Typst’s style while preserving the visual appearance.

B. Rewriting Content Streams

As specified in Chapter 9 in the PDF standard, the text object in the content stream, surrounded by the BT (begin text) and ET (end text) operators, contains a sequence of operators that define how text is rendered on the page. We observe that there are a few differences regarding how TeX and Typst put text objects.

- 1) TeX puts the font resource operator (Tf) inside the text object, while Typst puts it outside.
- 2) Typst employs a concatenate matrix (cm) operator to essentially “flip” the Y axis for the text object, while TeX directly works with the original coordinate system.
- 3) Typst also employs additional fancy color space operators (cs and scn) to set the text color to black, while TeX relies on the default.

Similarly, within the text object, the two systems also differ in how they position text using the Tm (text matrix) operator and the TJ (show text with position adjustments) operator.

- 1) TeX rounds every coordinate in the Tm operator to three decimal places, while Typst uses up to two decimal places.
- 2) TeX uses the hex string format (e.g. <0035>) in the TJ operator to represent text, while Typst uses the literal string format with each character represented as a three-digit octal number (e.g. \000\001\000\002).
- 3) TeX uses position adjustments in the TJ operator to adjust spacing between characters, while Typst does not use any position adjustments and rely on using a space character.

We argue that these details are mostly superficial, and even if detectors like TeXRay were to inspect them, they could be easily modified to match Typst’s style without affecting the visual appearance. To demonstrate the idea, we focus on one aspect, the concatenate matrix, and explain how we rewrite it below.

Concatenate Matrix. A concatenate matrix specifies an affine map with parameters $[a\ b\ c\ d\ e\ f]$ interpreted as the 3×3 matrix

$$\begin{bmatrix} a & b & 0 \\ c & d & 0 \\ e & f & 1 \end{bmatrix}.$$

The last two parameters, e and f, are the x- and y-offsets (in PDF points): they shift the origin to (e, f). The linear part $L = \begin{bmatrix} a & b \\ c & d \end{bmatrix}$ encodes scale, rotation and shear.

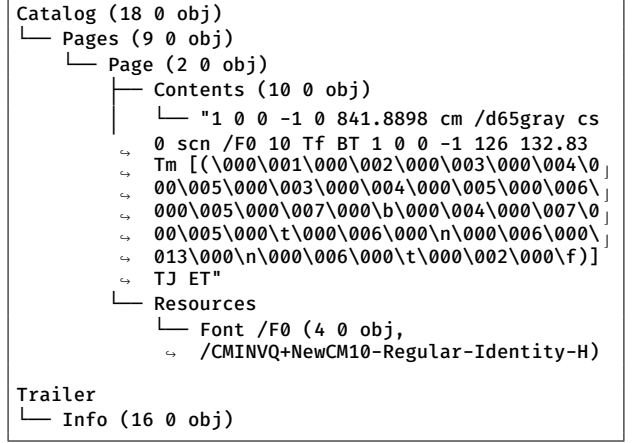
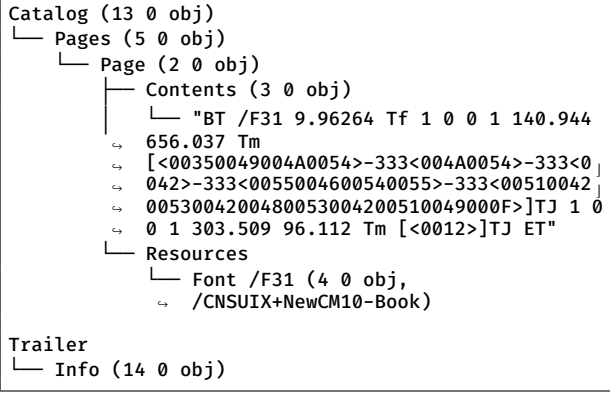


Fig. 5. Comparison of the PDF structures generated by TeX (left) and Typst (right) for the same simple document. We use the same font (New Computer Modern) in both cases to isolate structural differences.

Therefore, the concatenate matrix used by Typst

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & -1 & 0 \\ 0 & 841.8898 & 1 \end{bmatrix}$$

essentially encodes a “flip” of the Y axis that goes from the bottom-left origin (used by TeX) to the top-left origin (used by Typst). To match Typst’s style, we insert this matrix at the beginning of every text object and flip the Y coordinates in the Tm operators accordingly.

C. Demonstration of the Rewrite Process

Based on the observations above, we develop a 300-line Python script that rewrites the content streams of a TeX-generated PDF file to match Typst’s style while maintaining a visually near-identical output, available at <https://github.com/zhtluo/arxiv-scripts>. An example of rewriting a TeX-generated content stream to a Typst-style content stream is shown in Figure 6.

```

BT
/F31 9.96264 Tf
1 0 0 1 140.944 656.037 Tm [<00350049004A0054>-333<004A0
↵ 054>-333<0042>-333<0055004600540055>-333<00510042005
↵ 3004200480053004200510049000F>]TJ
1 0 0 1 303.509 96.112 Tm [<0012>]TJ
ET

1 0 0 -1 0 841.8898 cm
/d65gray cs
0 scn
/F31 9.96264 Tf
BT
1 0 0 -1 140.94 185.85 Tm
↵ [(\000\065\000\111\000\112\000\124\000\001\000\112\0
↵ 00\124\000\001\000\102\000\001\000\125\000\106\000\1
↵ 24\000\125\000\001\000\121\000\102\000\123\000\102\0
↵ 00\110\000\123\000\102\000\121\000\111\000\017)]
↵ TJ
1 0 0 -1 303.51 745.78 Tm [(\000\022)] TJ
ET

```

Fig. 6. An example of rewriting a TeX-generated content stream (top) to a Typst-style content stream (bottom).

We observe that the content stream output closely matches Typst’s style found in Figure 5, demonstrating the feasibility of our approach. Nevertheless, we do notice a few limitations in our current implementation.

- 1) Our current script only handles text-object-related content streams, and does not account for the full quirks of either TeX or Typst engines.
- 2) Our current script does not account for the font mapping differences between TeX and Typst. A more robust script would also rewrite the font resource objects to match Typst’s style.
- 3) Our current script replaces all spacing in TeX with space characters in Typst, which may lead to minor visual differences in certain cases.

We argue that most of these limitations can be addressed with more engineering effort, and do not affect the core idea. Thus, it remains a fact that TeX-generated PDF files can be transformed to be indistinguishable from Typst-generated PDF files.

VII. RELATED WORK

PDF Obfuscation and Related Attacks. PDF obfuscations have been studied under the context of spoofing AI agents, usually with some font modification schemes [17], [18]. The results, coupled with our idea in this paper, suggest that the binary data in PDF files is highly malleable, and thus cannot be reliably passed to automated systems.

Alternative Methods to Bypass TeXRay. Given that TeX is a very powerful typesetting language, as early as 2014, people have discovered that they can set up a shell TeX project that contains a single PDF file to be outputted directly, essentially bypassing the TeXRay checker altogether [4]. A sample shell TeX script is given in Figure 7.

To remedy the issue, arXiv appears to manually flag such submissions, and sends emails to require authors to submit the original TeX source files. One such email can be found online [19]. We include a copy of the email in Figure 8.

```

\documentclass{article}
\usepackage{pdfpages}

\begin{document}
\includepdf[pages=--]{the_real_article.pdf}
\end{document}

```

Fig. 7. A sample shell TeX script that outputs a PDF file directly.

Dear arXiv user,

Your submission appears to be a PDFLaTeX wrapper using pdfpages. This is an inappropriate submission, as it circumvents our TeX system. As a result, we have moved your submission to “Incomplete”.

Instead, please submit your TeX source. If there is a particular problem that you are encountering, please request assistance and include the specific error messages you receive, as well as your submit id.

Further submissions of this type may result in the loss of your submission privileges.

For more information about our TeX system and policies, see the following pages:

- <http://arxiv.org/help/submit>
- http://arxiv.org/help/submit_tex
- http://arxiv.org/help/submit_pdf
- <http://arxiv.org/help/faq/whytex>

Regards,
arXiv admin

Fig. 8. An example of email sent by arXiv to authors who upload shell TeX files.

We consider it a policy debate what kind of TeX source files are appropriate to upload. Given that this attack vector has little to do with uploading PDF files directly, we do not discuss it to details in this paper.

VIII. CONCLUSION

arXiv’s long-standing requirement to submit TeX source files when available aims to improve accessibility, archiving quality, and transparency of research manuscripts. However, our analysis shows that the enforcement mechanism behind this policy, TeXRay, is unreliable in practice and conceptually unsound in principle. Through black-box testing, we find that TeXRay bases its detection primarily on superficial PDF features such as metadata and font names, which can be trivially modified or obfuscated. We demonstrate both macro-level and post-processing-based bypasses that render TeXRay ineffective, and further argue that the very notion of distinguishing TeX-generated PDFs from those produced by other typesetting systems is ill-defined: identical visual results can always be obtained through obfuscation or translation.

These findings call into question the sustainability of enforcing source-based submission policies through file-level heuristics. Rather than maintaining fragile detection rules, we recommend that arXiv focus on improving moderation workflows, enhancing user awareness of source submission benefits, and supporting a broader ecosystem of modern typesetting systems. Our work highlights the importance of transparency and reproducibility in digital publishing infrastructure, and underscores the need for technically grounded, privacy-

conscious, and inclusive submission policies in open-access repositories.

ACKNOWLEDGMENT

We would like to thank the arXiv team for their prompt and understanding response to our responsible disclosure. Their cooperation allowed us to share our findings with the broader community in a timely manner. We would like to thank Deming Chu from Purdue University for his preliminary work on PDF obfuscation techniques.

REFERENCES

- [1] arXiv, “About arXiv,” <https://info.arxiv.org/about/index.html>, 2025, accessed: 2025-07-17.
- [2] C. B. Clement, M. Bierbaum, K. P. O’Keeffe, and A. A. Alemi, “On the use of arxiv as a dataset,” 2019. [Online]. Available: <https://arxiv.org/abs/1905.00075>
- [3] arXiv, “Why TeX?” <https://info.arxiv.org/help/faq/whytex.html>, 2025, accessed: 2025-07-17.
- [4] David Ketcheson, “How to upload LaTeX-generated PDF paper to arXiv without LaTeX sources,” TeX StackExchange Q&A, URL: <https://tex.stackexchange.com/questions/186068/how-to-upload-latex-generated-pdf-paper-to-arxiv-without-latex-sources>, accessed 2025-07-20.
- [5] arXiv, “Accessibility research report,” https://info.arxiv.org/about/accessibility_research_report.html, 2022, accessed: 2025-07-17.
- [6] J. Pont-Tuset, “arXiv LaTeX cleaner: safer and easier open source research papers,” <https://opensource.googleblog.com/2019/02/arxiv-latex-cleaner.html>, Feb. 2019, accessed: 2025-07-18.
- [7] Overheard on QuantPh (@QuantPhComments), “Overheard on QuantPh [X (Twitter) account],” <https://x.com/quantphcomments>, 2025, accessed: 2025-07-18.
- [8] Z. Luo, “Circumvent arxiv latex detection,” <https://zhtluo.com/misc/circumvent-arxiv-latex-detection.html>, 2023, accessed: 2025-07-18.
- [9] SIGCSE TS 2026 Organizing Committee, “SIGCSE TS 2026: Submission templates (papers track),” <https://sigcse2026.sigcse.org/track/sigcse-ts-2026-Papers#submission-templates>, 2025, accessed: 2025-07-18.
- [10] Typst GmbH, “Typst: Compose papers faster,” <https://typst.app/>, 2025, accessed: 2025-07-18.
- [11] Wikipedia contributors, “List of preprint repositories,” https://en.wikipedia.org/wiki/List_of_preprint_repositories, 2025, last edited 1 July 2025, retrieved 28 October 2025.
- [12] G. Perelman, “The entropy formula for the ricci flow and its geometric applications,” 2002. [Online]. Available: <https://arxiv.org/abs/math/0211159>
- [13] “Document management — portable document format — part 2: PDF 2.0,” International Organization for Standardization, Geneva, Switzerland, Standard ISO 32000-2:2020, 2020, second edition. [Online]. Available: <https://www.iso.org/standard/75839.html>
- [14] “Computer modern — the LaTeX font catalogue,” <https://tug.org/FontCatalogue/computermodern/>, TeX Users Group (TUG), 2021, last updated on 2021-01-19.
- [15] “Poppler: A PDF rendering library,” <https://poppler.freedesktop.org/>, freedesktop.org, 2025, version 25.08.0. Accessed 2025-08-25.
- [16] B. Jackowski and J. M. Nowacki, “Latin modern family of fonts (lm),” <https://www.ctan.org/tex-archive/fonts/lm/>, March 2021, version 2.005; GUST Font License (GFL).
- [17] J. Xiong, C. Zhu, S. Lin, C. Zhang, Y. Zhang, Y. Liu, and L. Li, “Invisible prompts, visible threats: Malicious font injection in external resources for large language models,” 2025. [Online]. Available: <https://arxiv.org/abs/2505.16957>
- [18] A. Creo, “Complete evasion, zero modification: Pdf attacks on ai text detection,” 2025. [Online]. Available: <https://arxiv.org/abs/2508.01887>
- [19] M. Sun, “How to bypass arXiv LaTeX-generated PDF detection in six lines,” <https://mssun.me/blog/how-to-bypass-arxiv-latex-generated-pdf-detection-in-six-lines.html>, Dec. 2016, updated: Sep 3, 2018; Accessed: 2025-07-20.