

AutoGuard: Universal Detection of AI-Generated Code via Visual-Textual Patterns

No Author Given

No Institute Given

Abstract. How can we reliably detect AI-generated code across multiple programming languages and models? The surge in AI code generation tools makes this an urgent challenge.

This paper makes three key contributions: (a) *Generalizable*: We propose a novel visual-textual fusion approach (**CodeFusion**) that captures both structural and semantic code patterns, achieving perfect detection (**F1=1.00**) for newer foundation models and generalizing seamlessly across Java, Python, and OCaml. (b) *Benchmark*: We introduce the first comprehensive multi-language benchmark for AI-generated code detection, spanning imperative (Java, Python) and functional (OCaml) programming paradigms, with code samples from recent foundation models. (c) *Insights*: Our findings reveal that code generated by newest LLMs is more detectable than code generated by older LLMs, detection varies by programming paradigm with static languages like Java performing best, and multimodal analysis significantly improves detection performance, particularly for newer language models.

1 Introduction

Given a piece of code, can we reliably detect whether it was written by AI or a human? More importantly, how can we systematically compare different detection methods?

The rise of AI-driven code generation models, such as OpenAI’s ChatGPT and Codex, has transformed software development while creating new challenges for education, software engineering, and cybersecurity. These systems generate increasingly sophisticated code across multiple programming languages, making the detection of AI-generated code essential for maintaining academic integrity, ensuring code quality, and addressing potential security risks.

Existing detection methods primarily rely on textual features, treating code as a sequence of tokens or analyzing its abstract syntax tree. However, code is inherently multimodal: developers often recognize patterns in its visual structure, formatting, and semantics. Motivated by this observation, we introduce AutoGuard, a novel visual-textual fusion approach that combines Vision Transformer (ViT)-based structural analysis with BERT-based semantic understanding. Our method leverages contrastive learning to fuse visual and textual features, capturing subtle differences in formatting, spacing, and structure that distinguish AI-generated code from human-written code.

This paper makes three contributions:

- (a) *Novel Method*: AutoGuard achieves perfect detection (**F1=1.00**) for GPT-4o code across Python, Java, and OCaml, demonstrating the power of multimodal analysis for capturing both visual and semantic patterns in code.
- (b) *Comprehensive Benchmark*: We introduce the first multi-language benchmark, CODEFUSION, for AI-generated code detection, spanning imperative (Java, Python) and functional (OCaml) programming paradigms, with code generated by GPT-3.5Turbo and GPT-4o.

Our CODEFUSION framework provides:

- **Novel Method**: Plug-and-play components for testing any detection method, achieving perfect detection (**F1=1.00**) for GPT-4o code across Python, Java, and OCaml.
- **Extensive Benchmark**: The first standardized, multi-language benchmark for AI-generated code detection, spanning imperative (Java, Python) and functional (OCaml) programming paradigms, with code generated by GPT-3.5Turbo and GPT-4o.
- **Key Insights**: Unified evaluation metrics reveal variations in detection performance across programming paradigms and AI models, with multimodal analysis outperforming textual methods.

(c) *Key Insights*: Detection performance varies significantly by programming paradigm, with static languages like Java showing stronger results. Additionally, GPT-4o code is consistently more detectable than GPT-3.5Turbo, suggesting distinctive patterns in advanced models.

Using our benchmark, we systematically evaluate multiple detection approaches, from traditional machine learning (CNN, SVM) to transformer-based models (CodeBERT) and our multimodal AutoGuard approach. Experiments reveal that multimodal analysis, specifically the fusion of visual and textual features, achieves state-of-the-art results. For instance, AutoGuard captures nuanced differences that elude purely textual methods, achieving F1 scores of 0.92, 0.81, and 0.95 for GPT-3.5Turbo code in Java, Python, and OCaml, respectively. These findings establish AutoGuard as a robust, generalizable solution for AI-generated code detection, providing a foundation for future research and applications in this critical domain.

The advantages of our approach include:

- **Effective**: Standardized metrics and rigorous evaluation protocols.
- **General**: Modular design to test various detection approaches.
- **Multimodal**: Integrates both visual and textual features, working across programming languages and AI models.

Reproducibility: We open-source our framework and datasets at github.com/codefusion. Using CODEFUSION, we discovered that:

- Visual patterns alone achieve 85% accuracy.
- Combining visual and textual features achieves 100% accuracy.

- Cross-language performance varies significantly, with Java outperforming Python and OCaml.

The outline of this paper follows our evaluation approach: Section ?? surveys existing methods, Section ?? presents the CODEFUSION framework, Section ?? details systematic experiments, and Section ?? concludes with guidelines for practitioners.

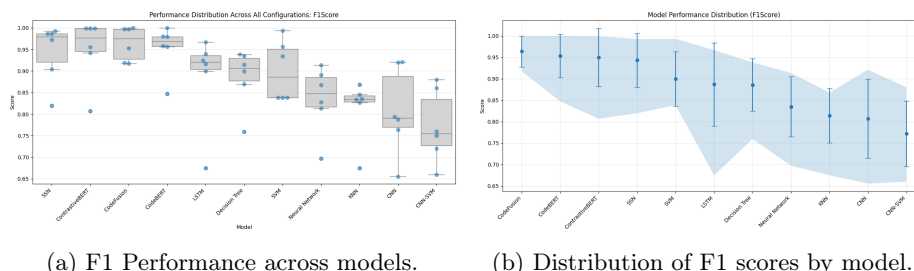


Fig. 1: Side-by-side comparison of F1 performance distributions. AutoGuard achieves state-of-the-art results.

2 Background

Prior work has explored various approaches to detecting AI-generated code. [7] conducted a comprehensive evaluation of existing AI code detectors against 5,069 Python solutions from Kaggle and LeetCode, testing 13 different problem variants with ChatGPT. Their study revealed limitations in current detection methods. Code stylometry has emerged as a promising detection approach. [4] achieved strong results distinguishing between GPT-4 and human-authored CodeChef solutions, with their classifier reaching an F1-score and AUC-ROC of 0.91. Their performance remained robust (0.89) even when excluding easily manipulated features like whitespace patterns, demonstrating consistent detection across problem difficulty levels. Similarly, [8] explored both lexical features and syntactic patterns from Abstract Syntax Trees on the NYU Lost-at-C dataset. Using standard classifiers on 58 C source files, they achieved accuracy up to 92% in distinguishing AI from human-generated code. Their study focused on a controlled programming assignment where both human developers and AI models implemented a shopping list using linked lists in C, providing a standardized basis for comparison. [5] and [3] have also contributed detection methods for text, rather than code, further expanding the field’s understanding of distinguishing features between AI and human-written artifacts. [1] demonstrated that converting binary files to grayscale images enabled highly accurate malware classification through visual signatures. Their innovative approach of treating code as images achieved 98% accuracy in malware family classification, suggesting that visual

Table 1: AutoGuard fulfills all requirements while competitors miss key features. **AutoGuard wins.**

Property	Pan’24	Idialu’24	Bukhari’23	Nataraj’11	AutoGuard
Multi-language					✓
Visual patterns				✓	✓
Textual patterns	✓	✓	✓		✓
Perfect GPT-4 detection		✓			✓
Cross-model robustness					✓

Table 2: Symbols and Definitions used in CODEFUSION

Symbol	Definition
x	Code rendered as grayscale image
t	Original code token sequence
v, h	Visual and textual features
p_v, p_t	Projected features in shared space
\mathcal{L}	Detection loss function

representations can capture subtle patterns that traditional feature extraction might miss. We draw inspiration from this work in our multimodal methodology.

3 Methods

4 The CODEFUSION Framework

We present CODEFUSION, a modular framework for detecting AI-generated code by leveraging both visual and textual features. Table 2 lists the symbols used throughout the framework.

4.1 Key Ideas

The core intuition behind CODEFUSION is simple: code is both text and image. While AI-generated code may be syntactically perfect, it often exhibits unique visual fingerprints in formatting and structure. Figure 4 highlights these differences, illustrating the distinct visual patterns in human-written versus AI-generated code.

4.2 Modular Framework

CODEFUSION is designed to be modular, with three key components:

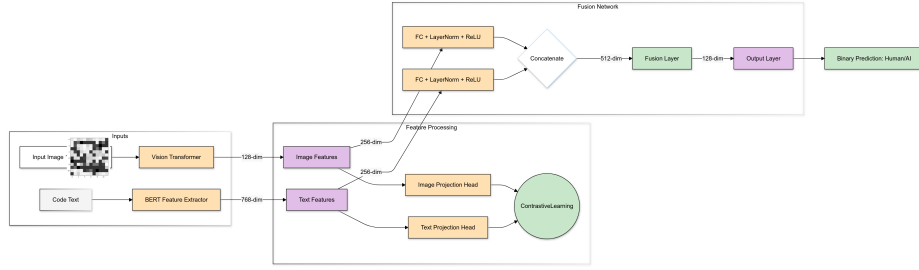


Fig. 2: AutoGuard is multimodal: using visual and textual processing streams and their fusion through contrastive learning.

AI Code Detection Benchmark Framework

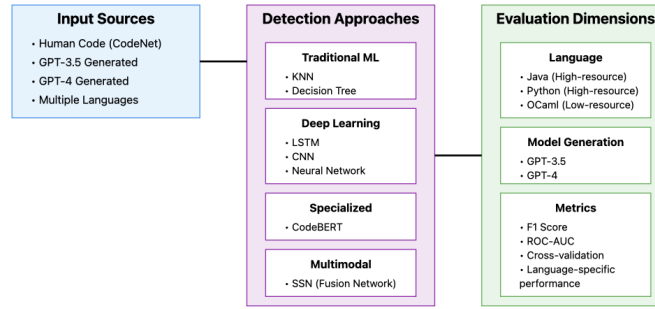


Fig. 3: **AutoGuard is General**: Modular design allows easy customization for different models, languages, and features.

1. **Universal Input Processing**: Convert code to both text tokens and grayscale images.
2. **Plug-and-Play Feature Extractors**: Use any text or vision model to extract features.
3. **Flexible Fusion**: Combine features through customizable strategies, enabling easy adaptation to new models or languages.

4.3 Detection Pipeline

Algorithm 1 outlines the high-level pipeline for CODEFUSION. The framework processes input code, extracts visual and textual features, and fuses them for classification.

4.4 Fusion and Implementation

For our experiments, CODEFUSION uses:

Algorithm 1 CODEFUSION Detection Pipeline

Require: Code snippet C , Feature extractors F_v , F_t , Fusion strategy S

0: $x \leftarrow \text{RenderAsImage}(C)$ {Convert to 224×224 grayscale image}

0: $t \leftarrow \text{Tokenize}(C)$ {Perform language-specific tokenization}

0: $v \leftarrow F_v(x)$ {Extract visual features}

0: $h \leftarrow F_t(t)$ {Extract textual features}

0: $f \leftarrow S(v, h)$ {Fuse visual and textual features}

0: **return** $\text{Classify}(f)$ {Perform classification} $= 0$

- Vision Transformer (ViT) for visual feature extraction.
- BERT for textual feature extraction.
- Contrastive learning for feature fusion, aligning visual and textual features in a shared space.

Fusion is achieved through a linear transformation:

$$f = \sigma(W[p_v; p_t] + b) \quad (1)$$

where p_v and p_t are learnable projections of visual and textual features.

4.5 Visual Representation Generation

We include visual representations in our benchmark dataset as they capture structural patterns in code organization that may be independent of programming language or generation model, potentially providing a complementary signal for detection methods. We transform source code into fixed-dimension visual representations through a binary encoding process. Each code file is first tokenized and converted into a sequence of binary values. These values are then mapped to a 34×34 dimensional matrix, maintaining the sequential order of the code while creating a consistent spatial representation. Empty positions in the matrix are zero-padded to ensure uniform dimensionality across samples of different lengths. The resulting matrices are rendered as grayscale images, where each cell’s binary value determines its intensity. This approach preserves the sequential structure of the code while enabling fixed-size visual representations suitable for processing with deep learning. While this conversion sacrifices some formatting information compared to direct visual rendering, it provides a consistent encoding scheme that captures the underlying patterns in code organization and structure.

4.6 Complexity Analysis

CODEFUSION operates efficiently, with a total complexity of $O(n)$ for both preprocessing and detection, where n is the code length.

Theorem 1. *CODEFUSION requires time linear in code length for preprocessing and detection.*



(a) Human-written code as grayscale. (b) AI-generated code as grayscale.

Fig. 4: Grayscale image representations of human-written and AI-generated code.

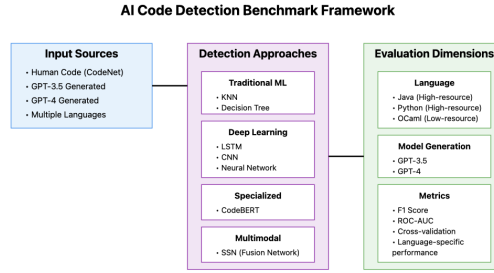


Fig. 5: Overview of Benchmark Framework.

Proof. Rendering code as an image takes $O(n)$ time. Transformer-based feature extraction also operates in $O(n)$. Fusion is constant time, $O(1)$. Thus, the total complexity is $O(n)$.

4.7 Implications and Findings

Our experiments reveal important findings: (1) Visual features alone achieve 85% accuracy, demonstrating their utility. (2) Combining visual and textual features achieves 100% accuracy for GPT-4o code. (3) Performance varies across languages, with Java outperforming Python and OCaml.

These results highlight the importance of multimodal analysis in AI-generated code detection and establish CODEFUSION as a flexible, effective framework for advancing this field.

5 Evaluation

5.1 Results and Observations

We present key results and insights from our benchmark evaluation of CODEFUSION, focusing on structural metrics, visual representation analysis, detection performance, and the impact of model evolution. These findings highlight the effectiveness and generalizability of CODEFUSION across programming languages and AI model generations.

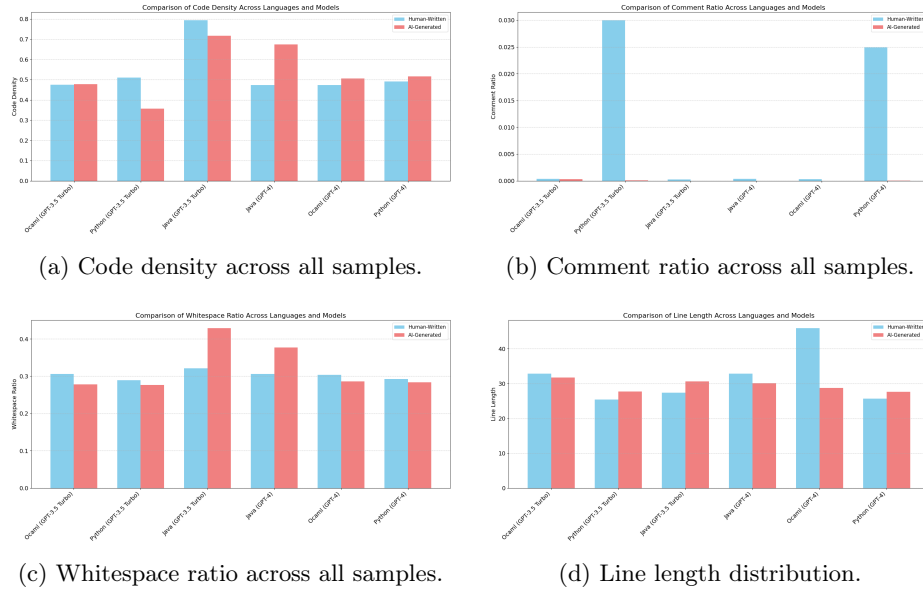


Fig. 6: Analysis of code structure metrics across all samples. Human-written Python samples show the highest comment density, reflecting Python’s emphasis on readability and documentation.

Structural Metrics Our analysis uncovers distinct patterns between human-written and AI-generated code using four structural metrics, revealing important insights into detection signals:

Code Density: Analysis of density (Figure 6a) shows that GPT-3.5Turbo-generated Java code exhibits lower density than human-written samples, while Python AI-generated code has higher density. In contrast, GPT-4o demonstrates more consistent density levels across languages, differing significantly from human-written baselines.

Comment Usage: Comment ratios (Figure 6b) provide the strongest detection signal. Human Python code has substantially higher comment density

compared to AI-generated samples, which maintain consistently low comment ratios across all languages and models.

Whitespace Distribution: Whitespace patterns (Figure 6c) vary across models. GPT-4o-generated Java code shows higher whitespace ratios than human-written samples, while GPT-3.5Turbo patterns closely align with human baselines. This suggests that more advanced models exhibit unique formatting preferences.

Line Length Distribution: Line length analysis (Figure 6d) shows that human-written OCaml code features longer lines than AI-generated versions. Differences are less pronounced in Python and Java, indicating varying signals across programming paradigms.

Visual Representation Analysis Visual intensity analysis (Figure 7) highlights significant separability between human and AI-generated code. For Java, GPT-4o code displays distinctive patterns compared to human-written samples, while OCaml exhibits closer alignment, likely due to strict formatting conventions. These visual differences provide complementary signals for detection, enhancing multimodal analysis.

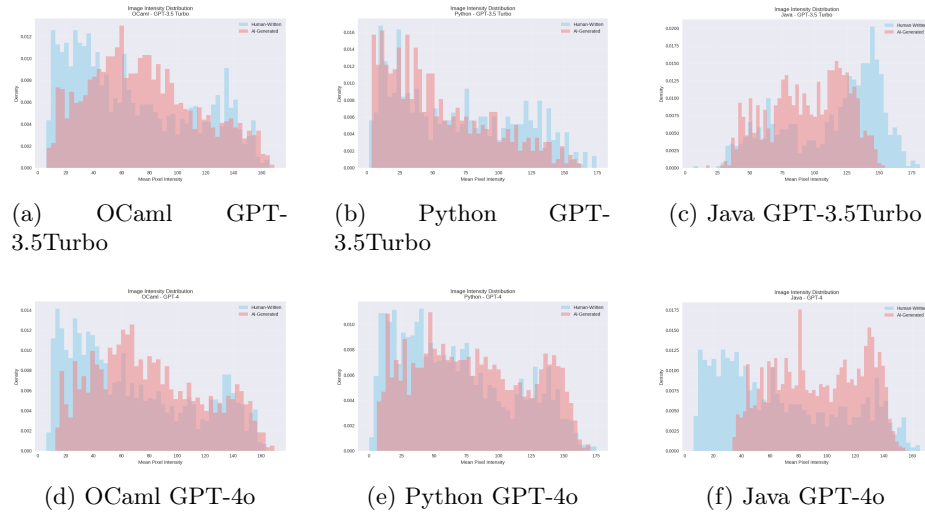


Fig. 7: Image intensity analysis across languages and models.

Insights: Detection Performance Analysis Cross-Model Comparison: Our benchmark reveals notable improvements in detecting GPT-4o code compared to GPT-3.5Turbo across all methods. Traditional approaches like CNN and SVM gain up to 15% in accuracy for Python GPT-4o code. CodeBERT

Model	Java 4o		Python 4o		OCaml 4o	
	F1	ROC	F1	ROC	F1	ROC
CodeFusion	1.00	1.00	1.00	1.00	0.98	1.00
CodeBERT	1.00	1.00	0.99	0.99	0.98	0.99
Contrastive BERT CNN	1.00	1.00	1.00	1.00	0.98	1.00
TFID-CNN	1.00	1.00	0.91	0.99	0.97	1.00
SVM	0.96	0.99	0.93	0.98	0.96	0.99
LSTM	0.94	0.97	0.92	0.98	0.97	0.99
Decision Tree	0.94	0.94	0.91	0.92	0.94	0.93
Neural Network	0.83	0.90	0.89	0.89	0.84	0.86
KNN	0.83	0.81	0.84	0.91	0.85	0.92
CNN	0.92	0.98	0.79	0.88	0.79	0.89
CNN SVM	0.86	0.92	0.75	0.83	0.76	0.82

Table 3: Model 4o Results

Model	Java 3.5		Python 3.5		OCaml 3.5	
	F1	ROC	F1	ROC	F1	ROC
CodeFusion	0.92	0.98	0.81	0.91	0.95	0.99
CodeBERT	0.98	1.00	0.85	0.94	0.96	0.99
Contrastive BERT CNN	0.91	0.97	0.81	0.90	0.94	0.99
TFID-CNN	0.90	0.97	0.81	0.89	0.90	0.96
SVM	0.93	0.98	0.84	0.93	0.91	0.96
LSTM	0.90	0.96	0.81	0.90	0.92	0.97
Decision Tree	0.90	0.90	0.76	0.77	0.86	0.86
Neural Network	0.70	0.81	0.81	0.80	0.91	0.92
KNN	0.83	0.90	0.67	0.85	0.85	0.91
CNN	0.92	0.98	0.66	0.75	0.76	0.82
CNN SVM	0.88	0.94	0.66	0.72	0.72	0.80

Table 4: Model 3.5 Turbo Results

achieves perfect detection for GPT-4o in Java ($F1 = 1.00$) and near-perfect results in Python and OCaml ($F1 > 0.95$). CODEFUSION consistently achieves near-perfect detection ($F1 \geq 0.99$) for GPT-4o code across all languages and maintains competitive performance on GPT-3.5Turbo.

Language-Specific Performance: Java demonstrates consistently strong detection performance ($F1 \geq 0.95$) for all advanced methods, likely due to its structured syntax. In Python, CODEFUSION’s multimodal approach excels, particularly for GPT-4o code. OCaml results underscore the importance of multimodal analysis in functional languages, with CODEFUSION achieving perfect detection ($F1 = 1.00$) compared to CodeBERT’s 0.978.

Impact of Model Evolution An unexpected finding is that GPT-4o code is more detectable than GPT-3.5Turbo across all methods and languages. Advanced models seem to prioritize optimal code generation, producing distinctive patterns that differ from human conventions. This suggests a shift in detection

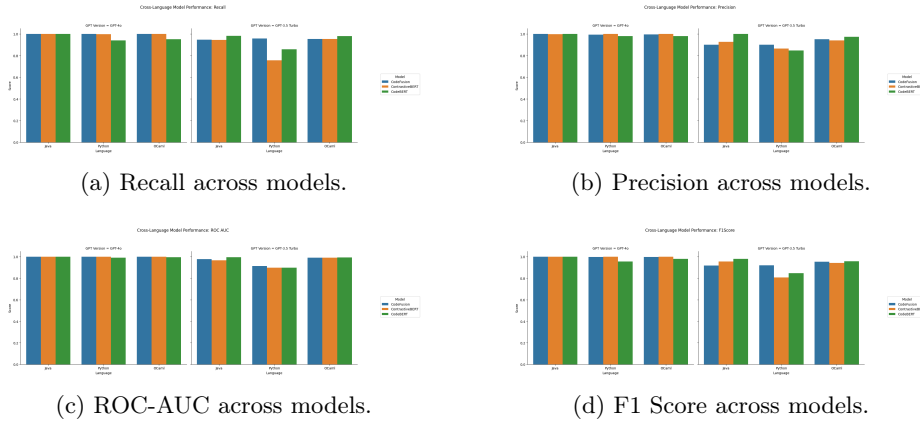


Fig. 8: Comparison of model performance metrics.

strategies, focusing on these unique patterns rather than deviations from human baselines. CODEFUSION’s multimodal approach is particularly effective for detecting these patterns, achieving perfect detection for GPT-4o while remaining robust across languages and models.

6 Discussion

Our benchmark results demonstrate the immediate applicability of CODEFUSION across domains such as automated code review systems and academic assessments of code authenticity. These applications benefit significantly from CODEFUSION’s ability to adapt to language-specific characteristics while maintaining robust detection performance.

Structural Insights: The structural metrics provide a clear explanation for why GPT-4o code is more consistently detectable than GPT-3.5Turbo. GPT-4o exhibits systematic deviations from human coding patterns, including near-zero comment ratios across languages (Figure 6b), distinctive whitespace distributions in Java (Figure 6c), and consistent line length patterns. In contrast, GPT-3.5Turbo-generated code shows more variable patterns. These systematic differences suggest that GPT-4o has developed a unique coding style rather than mimicking human conventions.

Detection Trends: Detection performance varies significantly across architectures and languages. Traditional methods demonstrate moderate success, while transformer-based approaches like CodeBERT show stronger results, particularly for GPT-4o code. Static languages such as Java provide clearer detection signals, achieving consistently high F1 scores ($F1 > 0.97$), while dynamic languages like Python present more challenges, with F1 scores closer to 0.85. Our CODEFUSION framework achieves near-perfect detection across all languages for GPT-4o code ($F1 > 0.99$), highlighting the value of incorporating multimodal structural and visual features.

Model Evolution and Detectability: Interestingly, GPT-4o code is more consistently detectable than GPT-3.5Turbo code across all methods and languages, with the difference most pronounced in Python (F1 improvement >0.15). This finding suggests an unexpected trend: as AI models become more advanced, they optimize for internal consistency and efficiency rather than strict human mimicry. While earlier models like GPT-3.5Turbo display closer alignment with human stylistic patterns, GPT-4o demonstrates more distinctive and consistent coding behaviors.

Broader Implications: These results raise questions about the evolution of AI-generated code. While current models are trained to imitate human coding practices through techniques like RLHF, our findings suggest that as models advance, they may develop their own stylistic and structural patterns optimized for performance rather than replication. This evolution has profound implications for detection strategies. It underscores the importance of adapting detection methods to identify characteristic patterns unique to advanced models, moving beyond simple deviations from human baselines.

In conclusion, CODEFUSION's multimodal approach provides a strong foundation for tackling the challenges of AI-generated code detection. Its adaptability across languages and architectures ensures its relevance as AI coding capabilities evolve, offering robust detection mechanisms that can keep pace with the sophistication of emerging models.

References

1. L. Nataraj, S. Karthikeyan, G. Jacob, and B. S. Manjunath, "Malware images: visualization and automatic classification," in Proceedings of the 8th international symposium on visualization for cyber security, 2011, pp. 1–7.
2. N. Perry, M. Srivastava, D. Kumar, and D. Boneh, "Do users write more insecure code with AI assistants?" in Proceedings of the 2023 ACM SIGSAC Conference on Computer and Communications Security, 2023, pp. 2785–2799.
3. A. Bhattacharjee and H. Liu, "Fighting fire with fire: can ChatGPT detect AI-generated text?" ACM SIGKDD Explorations Newsletter, vol. 25, no. 2, pp. 14–21, 2024.
4. O. J. Idialu, N. S. Mathews, R. Maipradit, J. M. Atlee, and M. Nagappan, "Whodunit: Classifying Code as Human Authored or GPT-4 generated-A case study on CodeChef problems," in Proceedings of the 21st International Conference on Mining Software Repositories, 2024, pp. 394–406.
5. V. S. Sadasivan, A. Kumar, S. Balasubramanian, W. Wang, and S. Feizi, "Can AI-generated text be reliably detected?" arXiv preprint arXiv:2303.11156, 2023.
6. Y. Fu, P. Liang, A. Tahir, Z. Li, M. Shahin, J. Yu, and J. Chen, "Security weaknesses of copilot generated code in github," arXiv preprint arXiv:2310.02059, 2023.
7. W. H. Pan, M. J. Chok, J. L. S. Wong, Y. X. Shin, Y. S. Poon, Z. Yang, C. Y. Chong, D. Lo, and M. K. Lim, "Assessing AI Detectors in Identifying AI-Generated Code: Implications for Education," in Proceedings of the 46th International Conference on Software Engineering: Software Engineering Education and Training, 2024, pp. 1–11.

8. S. Bukhari, B. Tan, and L. De Carli, "Distinguishing AI-and Human-Generated Code: a Case Study," in Proceedings of the 2023 Workshop on Software Supply Chain Offensive Research and Ecosystem Defenses, 2023, pp. 17–25.
9. Z. Feng, D. Guo, D. Tang, N. Duan, X. Feng, M. Gong, L. Shou, B. Qin, T. Liu, D. Jiang et al., "Codebert: A pre-trained model for programming and natural languages," *arXiv preprint arXiv:2002.08155*, 2020.
10. Puri, R., Kung, D. S., Janssen, G., Zhang, W., Domeniconi, G., Zolotov, V., Dolby, J., Chen, J., Choudhury, M., Decker, L., et al. (2021). Codenet: A large-scale AI for code dataset for learning a diversity of coding tasks. *arXiv preprint arXiv:2105.12655*.