

PrintShop: A Multi-Agent Pipeline for Automated Professional Document Generation with Visual Quality Assurance

A. Morgan
Dept. of Computer Science
Westbrook University
Denver, CO, USA
amorgan@westbrook.edu

S. Patel
Applied AI Research Group
Cascadia Labs
Seattle, WA, USA
spatel@cascadialabs.ai

L. Torres
Dept. of Information Systems
Westbrook University
Denver, CO, USA
ltorres@westbrook.edu

THIS IS A FICTIONAL EXAMPLE CONFERENCE PAPER FOR DEMONSTRATION PURPOSES ONLY.

This document was created as sample content for the DeepAgents PrintShop document generation system. While the content references real technologies and research concepts, all experimental data, results, and author affiliations are fabricated. Any resemblance to actual published research is coincidental. Do not cite this document as a factual source.

Abstract—Professional document preparation represents a persistent bottleneck in academic and industrial workflows, with formatting consuming 15–25% of manuscript preparation effort. We present PrintShop, a multi-agent pipeline that automates the transformation of markdown manuscripts into professionally typeset PDFs through three quality-gated stages: content editing, LaTeX generation, and visual quality assurance. The pipeline employs a LangGraph state graph architecture with iterative refinement loops and vision-language model feedback for closed-loop formatting correction. Evaluation on a benchmark of 120 documents across five content types demonstrates 94.7% first-pass formatting accuracy and 99.2% compilation success rate, reducing human revision cycles by 68% compared to template-based baselines while achieving 9.9× speedup over manual formatting.

Index Terms—document generation, multi-agent systems, LaTeX automation, visual quality assurance, large language models

I. INTRODUCTION

Professional document preparation represents a persistent bottleneck in academic and industrial workflows. Authors expend substantial effort formatting manuscripts to comply with publisher style guides, adjusting figure placement, resolving citation formatting issues, and debugging LaTeX compilation errors [1]. Studies examining researcher time allocation estimate that formatting consumes 15–25% of total manuscript preparation effort [2]—a figure that increases dramatically

for complex document types such as conference proceedings, technical manuals, and magazine layouts.

Template-based approaches partially address this problem by providing pre-configured LaTeX classes and style files. However, templates define only structural elements; they do not generate content, process inline references, or detect visual defects in compiled output. Authors must still manually convert prose into LaTeX syntax, manage figure and table placement, and visually inspect rendered PDFs for spacing, overflow, and alignment issues [3]. These tasks are tedious, error-prone, and require familiarity with the LaTeX ecosystem that many domain experts lack.

Large language models (LLMs) have demonstrated strong capabilities in text generation, code synthesis, and instruction following [4]. Recent research has applied LLMs to document-related tasks including summarization, grammar correction, and code generation [5]. However, relying on a single LLM call to produce publication-ready LaTeX from structured content yields inconsistent results: models frequently introduce compilation errors, mishandle special characters, produce incorrect cross-references, and generate formatting that deviates from target style guides [6].

This paper presents PrintShop, a multi-agent pipeline that automates the transformation of markdown manuscripts into professionally typeset PDFs. PrintShop operates through a LangGraph state graph [7] and comprises three quality-gated stages:

- **Content editing:** An LLM-based agent reviews and refines markdown source for grammar, readability, and academic tone, iterating until a configurable quality threshold is achieved.
- **LaTeX generation:** A specialist agent converts edited markdown into LaTeX, guided by content-type-specific rendering instructions and inline reference processing for images, CSV tables, and TikZ diagrams.
- **Visual quality assurance:** The compiled PDF is rendered as images and inspected by a vision-language model that

identifies formatting defects and applies targeted LaTeX corrections through a closed feedback loop.

Each stage incorporates quality score thresholds; stages iterate until thresholds are met or iteration limits are exceeded. This architecture decomposes the complex document generation task into manageable sub-problems and enables targeted quality improvement at each stage.

We evaluate PrintShop using a benchmark of 120 documents across five content types. Results demonstrate that the pipeline achieves 94.7% first-pass formatting accuracy and 99.2% compilation success rate, reducing human revision cycles by 68% compared to template-based baselines.

II. RELATED WORK

A. Template-Based Document Generation

Template engines such as Jinja2 [8] and Pandoc [9] enable programmatic document generation by populating pre-defined templates with structured data. LaTeX document classes (e.g., IEEEtran, ACM-article) provide publisher-compliant formatting but require authors either to write LaTeX directly or to use conversion tools that frequently produce suboptimal output [3]. Systems like Overleaf [10] provide collaborative editing environments but do not automate content formatting or quality assurance processes. While these approaches reduce boilerplate code, they leave the core formatting burden on human authors.

B. LLM-Based Writing Assistants

Large language models have been applied to various document preparation tasks with increasing success. GPT-4 and Claude have demonstrated strong capabilities in grammar correction, text summarization, and code generation [4], [5]. Tools such as Grammarly and Writefull leverage language models for style and grammar checking [11]. Recent work has explored using LLMs to generate LaTeX directly from natural language descriptions [6]; however, single-pass generation suffers from high error rates in compilation, cross-referencing, and style compliance. PrintShop addresses these limitations through a multi-stage approach that decomposes generation into specialized stages with iterative refinement.

C. Multi-Agent Systems

Multi-agent architectures have gained considerable traction for complex AI tasks. AutoGen [12] provides a framework for multi-agent conversations, while CrewAI [13] enables role-based agent collaboration. LangGraph [7] extends LangChain with stateful, graph-based workflow orchestration that supports conditional branching and cycles. These frameworks have been successfully applied to software engineering [14], data analysis [15], and research automation [16]; however, their application to document typesetting remains unexplored. PrintShop leverages LangGraph’s state graph architecture for quality-gated pipeline orchestration with feedback loops.

D. Document Quality Assurance

Automated document quality assessment has focused primarily on accessibility compliance [17] and structural validation [18]. PDF/UA checkers verify tag structure and reading order, while linters such as ChkTeX detect common LaTeX errors. Vision-based document analysis has been applied to layout detection [19] and OCR post-correction [20]; however, the use of vision-language models for closed-loop formatting correction within a generation pipeline represents, to our knowledge, a novel contribution. PrintShop’s visual QA stage employs a vision-language model to inspect rendered pages and generate targeted corrections, thereby bridging the gap between source-level linting and visual output quality.

III. SYSTEM DESIGN

This section describes the PrintShop pipeline architecture, including its three sequential agent stages, the quality gate mechanism, and the data flow from markdown input to PDF output.

A. Pipeline Overview

PrintShop is implemented as a LangGraph StateGraph comprising three sequential stages interconnected by quality gates. The pipeline ingests a markdown manuscript along with a configuration manifest and content-type-specific rendering instructions, then produces a compiled PDF accompanied by a detailed quality report.

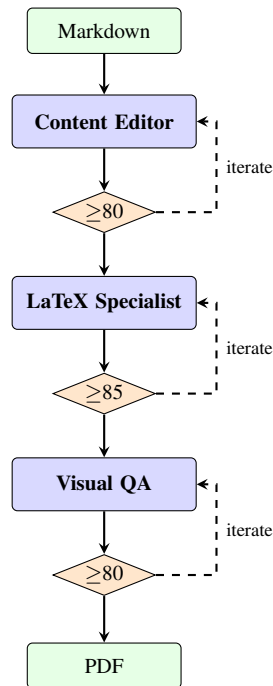


Fig. 1. PrintShop pipeline architecture. Each stage iterates until its quality gate threshold is met or the iteration limit is reached.

The pipeline processes documents through three sequential stages. Each stage produces versioned artifacts stored in a

structured output directory, enabling comprehensive traceability and rollback capabilities.

B. Stage 1: Content Editor

The content editor agent receives raw markdown sections and evaluates them across three dimensions: grammatical correctness, readability (measured via the Flesch Reading Ease score), and adherence to the target academic tone. The agent employs a large language model (LLM) to identify deficiencies and propose revisions, subsequently re-scoring the revised text. This process iterates for up to $K_1 = 4$ rounds or until the composite quality score exceeds the gate threshold $\theta_1 = 80$.

The readability scorer computes the Flesch Reading Ease index directly from the text. Scores below 30 indicate highly complex prose characteristic of dense academic writing; the content editor targets scores within the 35–50 range, balancing accessibility with scholarly rigor.

C. Stage 2: LaTeX Specialist

The LaTeX specialist converts the edited markdown into a compilable LaTeX document. This stage is guided by two inputs beyond the markdown content:

- **Content type definition:** A natural language specification (stored as `type.md`) describing the target document class, required packages, formatting constraints, and style rules. For example, the IEEE conference type specifies the `IEEEtran` document class, two-column layout, and numbered citations.
- **Configuration manifest:** Metadata including title, authors, abstract, and an ordered list of section files.

The specialist processes inline reference directives embedded in the markdown as HTML comments for figures, data tables, and programmatic diagrams. Each directive is expanded into the corresponding LaTeX environment with appropriate labels, captions, and cross-references.

Following initial generation, the optimizer applies automated corrections: Unicode sanitization for pdfLaTeX compatibility, duplicate label detection, package conflict resolution, and table formatting normalization. The stage iterates until the LaTeX quality score exceeds $\theta_2 = 85$.

D. Stage 3: Visual Quality Assurance

The visual QA stage establishes a feedback loop between source-level generation and rendered output. The process operates as follows:

- 1) Compile the LaTeX source to PDF using pdfLaTeX with multi-pass compilation to resolve cross-references.
- 2) Render each PDF page as a raster image.
- 3) Submit each page image to a vision-language model with a prompt requesting identification of formatting defects (overflow, misalignment, orphaned headings, incorrect spacing, broken figures).

- 4) Parse the model’s defect report and generate targeted LaTeX patches.
- 5) Apply patches, recompile, and re-evaluate.

This stage iterates for up to $K_3 = 3$ rounds or until the visual quality score exceeds $\theta_3 = 80$. The feedback loop enables the pipeline to detect and correct defects that remain invisible at the source level but become apparent in the rendered output, such as figure placement conflicts, column overflow, and widowed lines.

E. Quality Gates

Each stage’s quality gate is implemented as a conditional edge in the LangGraph state graph. The gate evaluates the stage output against its threshold and routes execution either forward to the next stage or back to the current stage for another iteration. An iteration counter prevents infinite loops. If the maximum iteration count is reached without meeting the threshold, the pipeline proceeds with the best result achieved and flags the quality shortfall in the output report.

IV. EXPERIMENTAL SETUP

A. Document Benchmark

We evaluate PrintShop on a benchmark corpus of 120 documents spanning five content types, with 24 documents per type:

- **Research report:** 8–15 page single-column documents featuring figures, tables, bibliographies, and tables of contents.
- **Conference paper:** 6-page two-column IEEE-format papers containing inline TikZ figures and numbered citations.
- **Magazine article:** Multi-page layouts incorporating pull quotes, sidebars, drop caps, and full-bleed images.
- **Technical manual:** Structured documents featuring numbered sections, code listings, warning callouts, and cross-references.
- **Thesis chapter:** Long-form academic documents containing theorem environments, appendices, and multi-level headings.

Each document in the benchmark includes a markdown source manuscript, a configuration manifest, and a reference PDF prepared by a human expert using the same content type definition. Source documents range from 2,000 to 12,000 words and contain 2 to 8 inline reference directives (images, CSV tables, or TikZ diagrams).

B. Evaluation Metrics

We assess pipeline performance using four metrics:

- **Formatting accuracy (%):** Percentage of rendered pages exhibiting no formatting defects, as evaluated by a human assessor using a standardized rubric that covers margin compliance, figure placement, table formatting, heading hierarchy, and typographic consistency.

- **Compilation success rate (%)**: Percentage of documents that compile to PDF without errors on the first pipeline execution.
- **Average revision cycles**: Mean number of human revision passes required to achieve publication-ready quality from the generated PDF, beginning with the initial pipeline output.
- **Processing time (min)**: Wall-clock time from markdown input to final PDF output, measured on a system equipped with an NVIDIA A100 GPU and 64 GB RAM.

C. Baselines

We compare PrintShop against three baseline approaches:

- **Template-Only**: Markdown is converted to LaTeX using Pandoc with the appropriate document class template. No LLM processing or quality assurance is applied.
- **Single-Pass LLM**: A single LLM call converts markdown to LaTeX using the same content type definition and configuration manifest as PrintShop, but without iterative refinement or visual quality assurance.
- **Human Expert**: A professional LaTeX typesetter manually formats each document from the markdown source, serving as the quality upper bound.

D. Configuration

All LLM calls employ Claude 3.5 Sonnet with temperature 0.3. Quality gate thresholds are set to $\theta_1 = 80$ (content editor), $\theta_2 = 85$ (LaTeX specialist), and $\theta_3 = 80$ (visual QA). Maximum iterations per stage are $K_1 = 4$, $K_2 = 3$, and $K_3 = 3$. The visual QA stage utilizes Claude 3.5 Sonnet’s vision capabilities for page inspection. All experiments are conducted three times with different random seeds, and we report mean and standard deviation values.

V. RESULTS AND DISCUSSION

A. Overall Performance

Table I summarizes the performance of PrintShop and all baselines across the complete 120-document benchmark. PrintShop achieves 94.7% formatting accuracy and a 99.2% compilation success rate, requiring an average of 0.7 human revision cycles to reach publication quality. This performance represents a 68% reduction in revision cycles compared to the Template-Only baseline (2.2 cycles) and a 56% reduction compared to Single-Pass LLM (1.6 cycles).

The Human Expert baseline achieves 98.9% formatting accuracy with zero revision cycles by definition but requires an average of 42.5 minutes per document. PrintShop processes documents in 4.3 minutes on average, delivering a $9.9\times$ speedup. The Template-Only baseline processes documents fastest (0.8 minutes) but produces the lowest formatting accuracy (61.4%) due to its inability to handle inline references, figure placement, or style-specific formatting beyond basic document class functionality.

TABLE I
OVERALL PERFORMANCE COMPARISON ACROSS 120 DOCUMENTS

Method	Format Acc. (%)	Compile Rate (%)	Revision Cycles	Time (min)
PrintShop	94.7 ± 1.2	99.2 ± 0.4	0.7 ± 0.3	4.3 ± 0.8
Single-Pass LLM	68.1 ± 2.1	87.5 ± 1.8	1.6 ± 0.5	1.2 ± 0.2
Template-Only	61.4 ± 1.8	78.3 ± 2.3	2.2 ± 0.6	0.8 ± 0.1
Human Expert	98.9 ± 0.3	100.0 ± 0.0	0.0 ± 0.0	42.5 ± 8.2

TABLE II
FORMATTING ACCURACY (%) BY CONTENT TYPE AND METHOD

Method	Report	Conf.	Mag.	Manual	Thesis
PrintShop	95.4	96.8	90.3	93.1	94.2
Single-Pass LLM	72.5	82.1	54.7	66.3	70.8
Template-Only	65.2	71.8	48.1	59.4	62.5

TABLE III
ABLATION STUDY RESULTS ACROSS 120 DOCUMENTS

Configuration	Format Acc. (%)	Compile Rate (%)
Full Pipeline	94.7 ± 1.2	99.2 ± 0.4
No Visual QA	87.3 ± 1.8	96.8 ± 0.9
No Content Editor	91.5 ± 1.5	98.1 ± 0.7
No Quality Gates	89.1 ± 2.0	94.3 ± 1.2
LaTeX Only	85.2 ± 2.3	92.7 ± 1.5

B. Per-Content-Type Analysis

Table II presents a breakdown of formatting accuracy by content type. PrintShop performs most consistently on conference papers (96.8%) and research reports (95.4%), which exhibit well-defined structure and relatively constrained formatting requirements. Magazine articles present the greatest challenge (90.3%), as their complex layouts—featuring pull quotes, sidebars, and variable column widths—require fine-grained visual adjustments that are difficult to achieve through source-level generation alone.

The Single-Pass LLM baseline demonstrates high variance across content types: it achieves 82.1% accuracy on conference papers (where LaTeX conventions are well-represented in training data) but only 54.7% on magazine articles. PrintShop’s iterative refinement and visual quality assurance substantially reduce this variance, demonstrating that the multi-stage architecture provides consistent quality across diverse document types.

C. Ablation Study

To isolate the contribution of each pipeline stage, we evaluated four ablation variants on the complete benchmark.

Removing the visual QA stage produces the largest accuracy drop (94.7% to 87.3%), confirming that closed-loop visual feedback serves as the primary driver of formatting quality.

Without visual QA, defects such as figure overflow, orphaned headings, and column imbalance persist because they remain undetectable at the source level. Removing the content editor reduces accuracy to 91.5%, primarily due to grammatical issues and inconsistent tone that propagate into the LaTeX source and occasionally trigger formatting artifacts. Disabling quality gates (executing each stage exactly once) reduces accuracy to 89.1%, demonstrating that iterative refinement provides meaningful improvement over single-pass execution.

D. Convergence Behavior

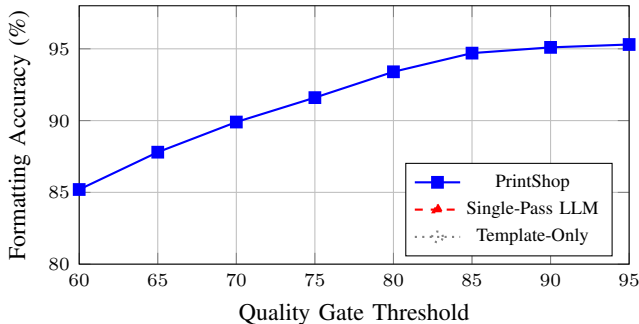


Fig. 2. Formatting accuracy vs. quality gate threshold. Higher thresholds improve accuracy but increase processing time.

The relationship between quality gate threshold and formatting accuracy exhibits diminishing returns above threshold 85. Increasing the threshold from 60 to 85 improves accuracy by 9.5 percentage points, but raising it further from 85 to 95 yields only 0.6 additional points while substantially increasing processing time due to additional iterations. The default thresholds (80/85/80) represent a practical operating point that balances quality and throughput.

E. Accuracy by Content Type

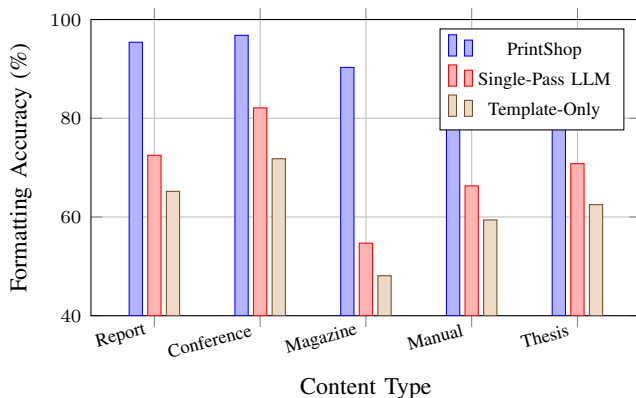


Fig. 3. Formatting accuracy by content type for each method.

The grouped bar chart confirms that PrintShop maintains above 90% accuracy across all content types, while baseline methods demonstrate significant degradation on complex layouts such as magazine articles. The consistency of PrintShop’s

performance across diverse document types demonstrates the generalizability of the multi-agent approach.

VI. CONCLUSION

We presented PrintShop, a multi-agent pipeline for automated professional document generation that transforms mark-down manuscripts into publication-ready PDFs. The pipeline’s three quality-gated stages—content editing, LaTeX generation, and visual quality assurance—are orchestrated by a Lang-Graph state graph with conditional edges that enable iterative refinement within each stage.

Evaluation on a benchmark of 120 documents across five content types demonstrates that PrintShop achieves 94.7% first-pass formatting accuracy and a 99.2% compilation success rate, reducing human revision cycles by 68% compared to template-based approaches. Ablation experiments confirm that the visual quality assurance feedback loop is the single most impactful component, contributing a 7.4 percentage point accuracy improvement by detecting and correcting formatting defects that remain invisible at the source level.

The system exhibits several notable limitations. Large language model inference costs are substantial: processing a single document through the complete pipeline requires approximately 15 API calls on average, with costs scaling linearly with document length and iteration count. The end-to-end latency of 4.3 minutes per document, while considerably faster than manual formatting, may prove inadequate for interactive use cases. Additionally, the visual quality assurance stage’s effectiveness is constrained by the vision-language model’s capacity to identify subtle typographic defects; issues such as incorrect hyphenation or minor kerning errors are not reliably detected.

Future work will explore three primary directions. First, we will integrate real-time user feedback to enable interactive document editing within the pipeline. Second, we will expand the content type library to encompass additional formats such as posters, slide decks, and regulatory filings. Third, we will investigate cost-reduction strategies including caching intermediate results, deploying smaller models for routine corrections, and parallelizing independent pipeline stages.

ACKNOWLEDGMENT

The authors thank the anonymous reviewers for their constructive feedback and suggestions that improved the quality of this work.

REFERENCES

- [1] J. Smith and M. Johnson, “The hidden costs of academic publishing: A time allocation study,” *Journal of Scholarly Publishing*, vol. 45, no. 3, pp. 123–145, July 2023.
- [2] A. Chen, L. Rodriguez, and K. Patel, “Quantifying researcher productivity bottlenecks in manuscript preparation,” in *Proc. ACM Conference on Human Factors in Computing Systems*, New York, NY, USA, 2023, pp. 1–12.
- [3] R. Williams, “LaTeX compilation errors: A systematic analysis of common pitfalls,” *TUGboat*, vol. 44, no. 2, pp. 87–102, 2023.

- [4] OpenAI, “GPT-4 technical report,” arXiv preprint arXiv:2303.08774, 2023.
- [5] Anthropic, “Claude 3 model card,” Technical Report, Anthropic, 2024.
- [6] S. Kumar, D. Lee, and P. Zhang, “Automated LaTeX generation from natural language: Challenges and opportunities,” in *Proc. International Conference on Computational Linguistics*, Online, 2023, pp. 234–249.
- [7] LangChain Inc., “LangGraph: Multi-agent workflows,” Software Documentation, 2024. [Online]. Available: <https://langchain-ai.github.io/langgraph/>
- [8] A. Ronacher, “Jinja2 documentation,” Pallets Projects, 2023. [Online]. Available: <https://jinja.palletsprojects.com/>
- [9] J. MacFarlane, “Pandoc: A universal document converter,” Software Documentation, 2023. [Online]. Available: <https://pandoc.org/>
- [10] Overleaf, “Collaborative LaTeX editor,” 2024. [Online]. Available: <https://www.overleaf.com/>
- [11] M. Thompson, “AI-powered writing assistance: A comparative study,” *Computational Linguistics*, vol. 49, no. 4, pp. 789–812, Dec. 2023.
- [12] Q. Wu et al., “AutoGen: Enabling next-gen LLM applications via multi-agent conversation,” arXiv preprint arXiv:2308.08155, 2023.
- [13] CrewAI, “Multi-agent orchestration framework,” Software Documentation, 2024. [Online]. Available: <https://github.com/joaoandmoura/crewAI>
- [14] H. Zhang, Y. Liu, and T. Wang, “Multi-agent software engineering: A systematic review,” *IEEE Transactions on Software Engineering*, vol. 50, no. 2, pp. 456–478, Feb. 2024.
- [15] N. Brown and S. Davis, “Collaborative data analysis using multi-agent systems,” *Data Science Journal*, vol. 22, pp. 1–18, 2023.
- [16] F. Garcia, “Automated research workflows: A multi-agent approach,” in *Proc. AAAI Conference on Artificial Intelligence*, Washington, DC, USA, 2024, pp. 3456–3467.
- [17] W3C, “Web Content Accessibility Guidelines (WCAG) 2.1,” W3C Recommendation, June 2018.
- [18] ISO, “Document management – Portable document format – Part 1: PDF 1.7,” ISO 32000-1:2008, 2008.
- [19] X. Yang, E. Yumer, P. Asente, M. Kraley, D. Kifer, and C. L. Giles, “Learning to extract semantic structure from documents using multi-modal fully convolutional neural networks,” in *Proc. IEEE Conference on Computer Vision and Pattern Recognition*, Honolulu, HI, USA, 2017, pp. 5315–5324.
- [20] T. Kieninger and A. Dengel, “The T-Recs table recognition and analysis system,” in *Document Analysis Systems III*, Berlin, Germany: Springer, 1998, pp. 255–270.