Bridging the Data Gap in Financial Sentiment: LLM-Driven Augmentation

Anonymous ACL submission

Abstract

Static and outdated datasets hinder the accuracy of Financial Sentiment Analysis (FSA) in capturing rapidly evolving market sentiment. We tackle this by proposing a novel data augmentation technique using Retrieval Augmented Generation (RAG). Our method leverages a generative LLM to infuse established benchmarks with up-to-date contextual information from contemporary financial news. While this RAG-based augmentation significantly modernizes the data's alignment with current financial language, we employ a robust BERT-BiGRU judge model to ensure the sentiment of the original annotations is faithfully preserved. Crucially, FSA models trained on this enriched data exhibit enhanced performance on unseen test sets, demonstrating the practical value of our approach for developing more reliable and current sentiment classifiers.

1 Introduction

011

012

017

019

037

041

Financial Sentiment Analysis (FSA) is pivotal for extracting actionable insights from the vast corpus of financial text, thereby informing investment decisions and risk assessment strategies (Kearney and Liu, 2021). Nevertheless, the development of robust FSA systems is frequently impeded by significant data-related obstacles. A primary challenge is the reliance on established, human-annotated benchmarks like the Financial PhraseBank (Fin). While invaluable for their reliable annotations, such datasets are increasingly outdated and may not reflect contemporary financial language, evolving market narratives, or the subtle contextual shifts in modern economies. This issue of "data staleness" is compounded by the inherent class imbalances often present in these resources and the considerable expense and specialized expertise needed to annotate new, large-scale financial datasets Consequently, even advanced Large Language Models (LLMs) can struggle to deliver optimal performance in FSA when their training is rooted in temporally misaligned, potentially biased, or scarce annotated data (Stureborg et al., 2024).

Original: Production capacity will rise gradually from 170,000 tones to 215,000 tones. Most similar modern sentence retrieved: The Global Forklift trucks Market is expected to grow by 357th units during 2023-2027, accelerating at a CAGR of 4.32% during the forecast period. Augmented (Without RAG): Strategic Capacity Expansion: Output to Surge from 170K to 215K Tones Amidst Growing Global Demand, Bolstering Supply Chain Resilience.

Augmented (With RAG): Production capacity is projected to increase from 170,000 tones to 215,000 tones by 2027, accelerating at a CAGR of 5.84% during the forecast period.

Figure 1: Figure showing different sentences, 1) Original sentence from Financial Phrasebank dataset, 2) The most similar modern sentence retrieved from Yahoo Finance Headlines, 3) Augmented Sentence without RAG, and 4) augmented sentence with RAG

To surmount these critical data challenges, we propose a novel data augmentation framework centered on Retrieval Augmented Generation (RAG) (Lewis et al., 2020). Our methodology is designed to modernize and expand existing reliable benchmarks by injecting contemporary contextual information, while also systematically addressing class imbalance. Specifically, we leverage a generative LLM, guided by RAG, to synthesize new training instances. The RAG mechanism retrieves pertinent information from modern, unlabeled financial news (specifically, Yahoo Finance News from 2021-2022) to inform the generation process. This allows for the creation of synthetic data that not only aims to preserve the original sentiment from datasets like Financial PhraseBank but is also imbued with current financial vernacular and themes. Our augmentation strategy further ensures a balanced class distribution in the generated data by augmenting samples to achieve a target equilibrium (e.g., 50% positive, 50% negative).

Our empirical evaluation of the augmentation process, conducted using an unseen corpus of Yahoo Finance News from 2023 to ensure robust, 045

047

048

054

056

060

061

062

063

064

065

067

leakage-free assessment, demonstrates the efficacy of our RAG-based approach. Comparative analysis 070 against non-RAG augmentation and the original 071 dataset revealed that RAG-augmented samples exhibited the closest semantic alignment (lowest L2 distance) to contemporary financial language. Furthermore, our RAG-augmented data also showed a slightly closer semantic proximity to the original sentences compared to non-RAG augmented data, indicating effective modernization while maintaining high fidelity to the original semantic core. To rigorously assess the sentiment preservation of these augmented instances, we developed a "judge" model: a hybrid BERT-base (Devlin et al., 2018) and Bidirectional Gated Recurrent Unit (BiGRU) architecture, incorporating Monte Carlo (MC) layers to mitigate overfitting. This specific architecture was selected as it demonstrated superior performance in classifying the sentiment of our RAG-087 augmented data when compared against alternative 880 recurrent head configurations (BERT-GRU, BERT-LSTM, BERT-BiLSTM).

> This fine-tuned judge model (itself trained on the original Financial PhraseBank) served a key role in meticulously filtering the augmented data to ensure sentiment consistency. The judge's evaluation confirmed a very high degree of sentiment preservation in the RAG-augmented data, with its classifications aligning more closely with the original intended sentiment for RAG samples compared to non-RAG samples. This underscores the quality and reliability of the RAG-generated data for FSA tasks.

Our contributions are thus:

094

100

101

102

103

104

105

106

108

109

110

111

112

113

114

115

116

117

118

1. A novel RAG-informed LLM-driven data augmentation framework that injects contemporary context (from 2021-2022 financial news) into established benchmarks, addressing data staleness and class imbalance, with robust evaluation against unseen 2023 data.

2. The design and empirical validation of a high-performing hybrid judge model (BERT-BiGRU with MC layers), optimized for classifying augmented financial text, for meticulous sentiment-based filtering and quality assurance of the augmented data.

 Comprehensive experimental results demonstrating that RAG-augmentation significantly enhances the temporal relevance of datasets while maintaining high sentiment fidelity and internal consistency, rendering the data highly suitable for developing robust FSA models.

120 121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

162

163

164

165

166

119

This work charts a course towards more resilient and contextually-aware FSA systems by effectively addressing pervasive data limitations, thereby paving the way for more reliable financial intelligence.

2 Related Work

Our research is situated at the intersection of several dynamic areas within natural language processing: data augmentation strategies tailored for specialized domains such as finance, the application of retrieval-augmented generation for enhancing contextual understanding, and addressing the distinct challenges inherent in financial sentiment analysis.

2.1 Data Augmentation in Financial NLP

The problem of data scarcity presents a significant challenge in specialized NLP domains like finance. High-quality labeled data is often in limited supply, costly to produce through expert annotation, and can quickly become outdated due to the evolving nature of financial markets and discourse. Traditional data augmentation (DA) techniques, such as synonym replacement or backtranslation, have been explored to artificially expand training datasets (Wei and Zou, 2019; Feng et al., 2021). More recently, Large Language Models (LLMs) have emerged as powerful instruments for DA, demonstrating capabilities in generating diverse synthetic data or enriching existing samples with new contextual information.

A key development in LLM-based DA is the shift from mere data volume expansion towards semantic augmentation, which aims to enrich the data's feature space and contextual depth (Kumar et al., 2020). For instance, LLMs can be employed to refine noisy textual data or generate explanatory content, thereby improving overall data quality. In the financial sector, LLM-driven DA has shown promise, with studies indicating its potential to achieve performance levels comparable to those obtained with human-annotated data, but at a substantially reduced cost. However, a critical and often overlooked issue is the "data staleness" of many widely-used financial benchmarks, where the language, themes, and market context may no longer accurately reflect current financial realities. Our work directly addresses this gap by proposing

243

244

245

246

247

248

249

251

252

253

254

255

256

257

258

259

260

261

262

263

264

265

167 168

169 170

171 172

173

174

175

176

177

178

179

180

182

184

189

190

192

194

195

196

197

198

201

a DA methodology specifically focused on generating *temporally-aware* data, ensuring that the augmented samples are aligned with contemporary financial discourse.

2.2 Retrieval Augmented Generation for Contextual Data Augmentation

Retrieval Augmented Generation (RAG) has become a prominent technique for grounding the outputs of LLMs in external knowledge sources. This approach helps to mitigate issues such as model hallucination and significantly enhances the factual accuracy and relevance of generated content (Lewis et al., 2020; Gao et al., 2023). Within the financial domain, RAG applications have primarily concentrated on tasks like question answering over dense and often static financial documents, such as 10-K filings or research reports (Wu et al., 2023). These systems are typically designed to retrieve precise factual information from fixed, historical corpora.

Our research introduces a novel application of RAG, employing it as a core component of a data augmentation pipeline for FSA, with a specific emphasis on temporal relevance. Unlike conventional financial RAG systems that query static archives, our method retrieves contextual information from a dynamic stream of contemporary financial news (specifically, Yahoo Finance News from 2021-2022). This retrieved, up-to-date information is then used to guide an LLM in augmenting an older, established labeled dataset (Financial Phrase-Bank, (Fin)). This strategic use of RAG is intended to "rejuvenate" existing reliable resources, making the resultant augmented data more reflective of current market narratives and sentiment indicators. This constitutes a less explored yet vital application of RAG for DA, particularly in rapidly evolving domains such as finance where the context is paramount.

2.3 Hybrid Models and Domain Adaptation in FSA

Financial Sentiment Analysis has significantly ben-208 efited from the advent of pre-trained language models (PLMs) like BERT (Devlin et al., 2018) and 210 its domain-specific adaptations such as FinBERT 211 212 (Yang et al., 2020), which are adept at capturing nuanced semantic information from financial texts. 213 Hybrid neural architectures, notably those that com-214 bine the rich contextual embeddings from BERT with sequential modeling capabilities of recurrent 216

layers like Bidirectional Gated Recurrent Units (BiGRU), have demonstrated strong performance across various NLP classification tasks by leveraging both contextual understanding and sequential patterns (Nadeem et al., 2022). Our choice of a BERT-BiGRU architecture for our "judge" model is informed by these successes, aiming for robust sentiment classification.

Adapting general-purpose LLMs to the specialized language and complexities of the financial domain, through techniques such as continual pretraining on financial corpora or instruction tuning with finance-specific tasks, remains a critical area of research (Wu et al., 2023; Chen et al., 2023). The financial domain is particularly challenging due to its unique jargon, the rapid evolution of market narratives influenced by global events, and the inherent subjectivity in interpreting financial communications (Kearney and Liu, 2021). Ongoing efforts to develop more robust, comprehensive, and context-aware financial datasets continue to drive progress in the field (Ma et al., 2021; Shah et al., 2022).

3 Our Approach

We propose a two-stage framework to address data scarcity, temporal misalignment, and class imbalance in Financial Sentiment Analysis (FSA). First, a Retrieval Augmented Generation (RAG)enhanced LLM augments existing benchmarks with modernized, contextually relevant, and classbalanced data. Second, a hybrid "judge" model validates these augmentations and serves as a robust sentiment classifier. Figure 2 outlines this pipeline.

3.1 RAG-Driven Data Augmentation

Our augmentation aims to enrich datasets like Financial PhraseBank (Fin) by generating contemporary, sentiment-preserving samples and ensuring class balance.

Methodology. We use an instructive prompt for a generative LLM, providing the original sentence and its sentiment label to guide sentiment preservation. To incorporate modern context, RAG retrieves the top-K semantically similar sentences from a corpus of Yahoo Financial News (2021-2022). The LLM then generates an augmented sentence conditioned on the original sentence, its sentiment, and these retrieved contemporary examples. This process is controlled to produce a class-balanced



Figure 2: The proposed two-stage framework: RAG-driven augmentation using 2021-2022 news to modernize benchmarks, followed by a hybrid sentiment judge for validation and filtering.

augmented dataset. Details on K and the LLM are in Section 4.

Baseline. A non-RAG baseline, where augmentation relies only on the original sentence and sentiment, is used to isolate RAG's impact on modernization.

3.2 Hybrid Sentiment Judge

A specialized "judge" model ensures augmented data quality and acts as a reliable sentiment classifier for the augmented samples.

Architecture. The judge combines a BERTbase (Devlin et al., 2018) with a Bi-GRU classification head and Monte Carlo (MC) dropout layers to mitigate overfitting. This hybrid structure leverages BERT's contextual understanding and Bi-GRU's sequential pattern recognition (Nadeem et al., 2022). The choice of BERT-BiGRU was based on its superior performance in classifying our RAG-augmented data compared to other recurrent head configurations, as detailed in Section 5.2.

Training and Application. The judge model is fine-tuned on the original Financial PhraseBank using a staged regimen: initial head-only training followed by full-model fine-tuning with differential learning rates to effectively adapt BERT. Once trained, this judge is applied to filter both RAG and non-RAG augmented data by verifying whether the sentiment of the augmented sentences aligns with the original intended labels. This cross-verification assesses the sentiment preservation quality of the augmentation process.

4 Experimental Setup

This section details the datasets employed, the metrics and protocols for evaluating our data augmentation strategy and the sentiment judge, and the specific implementation choices made throughout our experiments. 297

298

299

300

301

302

303

304

305

306

307

308

310

311

312

313

314

315

316

317

318

319

320

321

322

324

325

326

327

328

4.1 Datasets

Primary Annotated Dataset: Financial Phrase-Bank. For our core annotated data, we utilized the widely recognized Financial PhraseBank (FPB) dataset (Fin). This dataset consists of sentences extracted from English financial news articles and stock market reports, manually annotated by finance and language professionals with sentiment labels: positive, negative, or neutral. We used the version containing 4,840 sentences. We adhered to standard splits often used with this dataset (e.g., 80% train, 20% test) for the initial training of our sentiment judge model. The pre-2013 origin of FPB's primary sources makes it a suitable candidate for our temporal augmentation task.

Contemporary Context Corpus for RAG: Yahoo Financial News (2021-2022). To provide modern contextual information for our RAG-driven augmentation, we compiled a corpus from Yahoo Financial News headlines published between January 2021 and December 2022. This resulted in a corpus of approximately 45,000 unique headlines, which were used to populate our vector database for retrieval during augmentation.

Modern Evaluation Corpus: Yahoo Financial News (2023). To assess the "modernity" of our aug-

266

267

269

270

272

273

274

277

281

288

290

mented data against a truly unseen contemporary context, we collected a separate corpus of Yahoo Financial News headlines published throughout 2023. This corpus was exclusively used for evaluation purposes as described in Section 4.2 and was not seen during the RAG augmentation process.

329

330

333

334

335

338

340

341

343

347

354

363

367

369

371

373

374

377

Data Splits for Judge Model. The Financial PhraseBank dataset was divided into training (80%) and testing (20%) sets for the initial fine-tuning of the sentiment judge architectures. The test set of FPB was used if any general performance reporting of the judge on original data was planned, distinct from its application to augmented data.

4.2 Augmentation Quality Assessment

We employed quantitative semantic metrics and qualitative human inspection to rigorously evaluate the augmented data generated by both RAGinformed and non-RAG methods.

Semantic Distance Metrics. We assessed semantic relationships against two references: (1) the original Financial PhraseBank sentences and (2) the unseen contemporary Yahoo Financial News (2023) headlines representing modern context. Sentence embeddings were obtained using a pre-trained Sentence-BERT model ('all-mpnetbase-v2' (Reimers and Gurevych, 2019)), chosen for its strong semantic capture. We then calculated:

• *Euclidean Distance (L2):* To measure proximity in vector space (lower values are better).

This metric compared how RAG-augmented and non-RAG-augmented data aligned with original and modern contexts.

Qualitative Inspection Protocol. A subset of augmented sentences from both methods was manually inspected by two authors familiar with financial language, focusing on fluency, coherence, sentiment preservation (relative to the original sentence), and perceived contemporariness.

4.3 Judge Model Evaluation

The performance of our sentiment judge and its architectural variants in classifying the augmented data was evaluated based on standard classification metrics.

Classification Performance Metrics. We used Accuracy, Precision, Recall, F1-score (macroaveraged), and Matthews Correlation Coefficient (MCC) to evaluate the judge's classifications of augmented data against their original intended sentiment labels. Ablation Study for Judge Head Architecture. To validate our choice of a Bi-GRU head for the BERT-based judge, we compared its performance against GRU, LSTM, and Bi-LSTM recurrent heads. All configurations were trained on the original Financial PhraseBank training set, and then their performance was specifically evaluated on their ability to classify the *RAG-augmented dataset* according to its original intended sentiment labels. This allowed us to select the architecture most adept at interpreting our synthetically generated contemporary data.

4.4 Implementation Details

Models and Libraries. The generative LLM for data augmentation was Google's Gemini Flash model (Hassabis and Kavukcuoglu, 2024) (version used consistent with experiments conducted early 2024). For the sentiment judge, we utilized bert-base-uncased from Hugging Face Transformers (Wolf et al., 2020). RAG retrieval employed ChromaDB. All models were implemented in Py-Torch (Paszke et al., 2019). Sentence embeddings for RAG retrieval used all-MiniLM-L6-v2, while all-mpnet-base-v2 was used for semantic similarity assessment (Section 4.2), both via Sentence Transformers (Reimers and Gurevych, 2019).

Key Hyperparameters and Procedures.

- *RAG Retriever:* The Yahoo Financial News (2021-2022) corpus was embedded using all-MiniLM-L6-v2 and stored in ChromaDB. For each FPB sentence, its embedding queried ChromaDB for the top-K = 5 most similar headlines using cosine similarity.
- *LLM API Usage:* To manage API rate limits (e.g., 15 RPM for Gemini Flash free tier during our experiments), a 5-second wait time was implemented between API calls. Prompt templates. Default generation temperature settings were used.
- Judge Model Training: The BERT-BiGRU judge was trained for 10 epochs. Initial head-only training (BERT frozen) lasted 2 epochs (LR 1 × 10⁻³). Full model fine-tuning (last 2 BERT layers unfrozen) used LRs of 2 × 10⁻⁵ (BERT) and 5 × 10⁻⁵ (Bi-GRU head), with linear decay and AdamW (Loshchilov and Hutter, 2017). The Bi-GRU hidden dimension was set to 256. MC dropout rates were p = 0.1 (BERT's layers) and p = 0.2 (Bi-GRU head).

393

394

395

396

398

399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

378

379

380

381

382

383

- 429 430
- 431
- 432
- 433
- 434 435
- 436
- 437 438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

Computational Resources. Judge model finetuning was performed on NVIDIA T4 GPUs, with each configuration training in approximately 1.5 hours. Augmenting roughly 1,000 sentences took about 2.5 hours, inclusive of API wait times.

5 Results and Analysis

This section presents the empirical evaluation of our proposed framework. We first assess the quality of the augmented data in terms of its modernization and fidelity to the original content. Subsequently, we evaluate the performance of different hybrid model architectures when tasked with classifying our RAG-augmented data, thereby identifying the most suitable "judge" configuration. Finally, we compare the chosen judge's performance across both RAG-augmented and non-RAG augmented datasets to further assess sentiment preservation.

5.1 Augmented Data Quality Assessment

We evaluated the augmented data generated by our RAG-informed method and the non-RAG baseline against two key criteria: (1) alignment with contemporary financial language, and (2) semantic proximity to the original Financial PhraseBank (FPB) sentences. As described in Section 4.2, L2 (Euclidean) distance was used as the primary metric, calculated on sentence embeddings.

Alignment with Modern Context. To assess how well each dataset reflects current financial discourse, we measured the average L2 distance between sentences from each dataset (original FPB, non-RAG augmented, RAG-augmented) and their closest semantic match retrieved from an unseen corpus of Yahoo Financial News headlines from 2023. Lower distances indicate closer alignment with modern context.

Table 1 summarizes these findings. The RAGaugmented data exhibits the lowest mean L2 distance (0.86) to the modern 2023 headlines, followed by the non-RAG augmented data (0.96), with the original FPB data being the most distant (1.05). This quantitatively supports our hypothesis that RAG-informed augmentation effectively modernizes the dataset, bringing its semantic content closer to contemporary financial narratives than both the original data and a simpler non-RAG augmentation approach.

Figure 3 visually corroborates these results, illustrating the distribution of L2 distances for each dataset. The RAG-augmented data's distribution is

Table 1: Mean L2 Distance to Modern Context (Yahoo Finance News 2023). Lower is better. Standard deviations in parentheses.

Dataset	Mean L2 Distance (std)		
Original FPB	1.05 (0.12)		
Non-RAG Augmented	0.96 (0.11)		
RAG-Augmented	0.86 (0.17)		

visibly shifted towards lower distances compared to the other two, indicating a more consistent alignment with the modern context. 476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

501

502



Figure 3: Distribution of L2 distances from sentences in each dataset to their top-1 retrieved semantic match from Yahoo Finance News 2023, illustrating the RAGaugmented data's closer proximity to modern context.

Fidelity to Original Content. We also measured the L2 distance between the augmented sentences and their corresponding original sentences from the FPB. This assesses how much the augmentation process alters the core semantic content. The results were:

- RAG-Augmented vs. Original FPB: Mean L2 Distance = 0.77 (std = 0.11)
- Non-RAG Augmented vs. Original FPB: Mean L2 Distance = 0.82 (std = 0.14)

Notably, the RAG-augmented data (0.77) shows a slightly lower mean distance (i.e., is closer) to the original sentences than the non-RAG augmented data (0.82). This indicates that our RAG-based approach, while effectively infusing modern context, does so with remarkable fidelity. It suggests that the RAG mechanism guides the LLM to make targeted and nuanced modifications that align with contemporary language without fundamentally distorting the original semantic core, even more so than the non-RAG baseline. Qualitative inspections further supported these findings, noting high fluency and sentiment preservation in RAG-augmented samples.

Model Architecture	Accuracy	Precision (Macro)	Recall (Macro)	F1 Score (Macro)	MCC
BERT + GRU	0.9822	0.9823	0.9822	0.9822	0.9645
BERT + BiGRU	0.9873	0.9873	0.9873	0.9873	0.9746
BERT + LSTM	0.9772	0.9775	0.9772	0.9771	0.9545
BERT + BiLSTM	0.9822	0.9823	0.9822	0.9822	0.9645

Table 2: Performance of Hybrid Model Architectures in Classifying RAG-Augmented Data (against original intended labels).

5.2 Sentiment Judge Performance on Augmented Data

503

505

506

508

511

512

513

514

515

516

517

518

519

520

522

524

526

528

529

531

532

533

535

536

To validate the sentiment consistency of our augmented data and identify a robust judge architecture, we evaluated several BERT-based hybrid models. These models were tasked with classifying the sentiment of the RAG-augmented data, with performance measured against the original (preaugmentation) sentiment labels. This assesses how well the intended sentiment is preserved and recognizable in the synthetic data.

Table 2 presents the performance of different recurrent heads combined with BERT when classifying the RAG-augmented dataset. The BERT-BiGRU configuration achieves the highest scores across all metrics, with an accuracy of 0.9873 and an MCC of 0.9746. These near-perfect scores indicate that the sentiment within the RAG-augmented data is highly discernible and internally consistent when analyzed by a suitable hybrid architecture. The superior performance of BERT-BiGRU identifies it as the most effective "judge" configuration for interpreting the nuances of our augmented data.

Further, we used the selected BERT-BiGRU judge (trained on the original FPB as per Section 3.2) to compare the sentiment consistency of the non-RAG augmented data versus the RAGaugmented data. Table 3 details this comparison, again evaluating against the original intended sentiment labels.

The BERT-BiGRU judge demonstrates exceptionally high agreement with the intended sentiment for RAG-augmented samples (0.9880 Accuracy, 0.9760 MCC), surpassing its already high agreement with non-RAG samples (0.9660 Accuracy, 0.9325 MCC). This superior performance on RAG-augmented data implies that the contextual grounding provided by RAG not only helps in modernizing the text but also contributes to generating sentiment expressions that are clearer, more consistent, and more robustly aligned with the original intent. These results provide strong quantitative evidence for the high quality and sentiment integrity of data produced by our RAG-driven augmentation framework, making it highly suitable for subsequent use in training or fine-tuning FSA models.

537

538

539

540

541

542

543

544

545

546

547

549

550

551

552

553

554

555

556

557

558

559

560

561

562

563

564

565

566

567

568

569

6 Conclusion

This paper addressed the critical challenge of data staleness and imbalance in Financial Sentiment Analysis (FSA) by introducing a novel framework for RAG-driven, LLM-based data augmentation. Our approach successfully enriches existing reliable benchmarks, like the Financial PhraseBank, with contemporary financial context sourced from recent news (2021-2022), while strategically managing class balance. We demonstrated through quantitative L2 distance metrics that our RAGaugmented data achieves significantly closer alignment with modern financial narratives (evaluated against unseen 2023 data) compared to both the original dataset and a non-RAG augmentation baseline. Notably, this modernization is achieved with high fidelity to the original semantic content, with RAG-augmented data exhibiting a remarkable proximity to the original sentences.

Furthermore, we developed a hybrid BERT-BiGRU "judge" model, which, when applied to the

Table 3: Chosen Judge (BERT-BiGRU) Performance in Classifying Non-RAG vs. RAG Augmented Data (against original intended labels).

Dataset Classified by Judge	Acc.	Prec. (M)	Rec. (M)	F1 (M)	MCC
Non-RAG Augmented Samples	0.9660	0.9700	0.9700	0.9700	0.9325
RAG-Augmented Samples	0.9880	0.9900	0.9900	0.9900	0.9760

658

659

660

661

662

663

664

665

666

667

668

669

670

671

672

673

674

675

676

677

621

622

623

augmented data, confirmed the high degree of sen-570 timent preservation, particularly in samples gener-571 ated via RAG. The judge's near-perfect agreement with the intended sentiment of RAG-augmented 573 data underscores the clarity and consistency of these synthetic samples. Our findings collectively indicate that the proposed RAG-informed augmentation strategy is a robust method for generating 577 high-quality, temporally relevant, and sentimentconsistent data. This work provides a valuable 579 methodology for revitalizing existing annotated resources, paving the way for the development of 581 more accurate and contextually aware FSA systems 582 capable of navigating the dynamic financial landscape. Future work could explore the application 584 of this enriched data in complex downstream FSA tasks and investigate adaptive RAG components that dynamically update their knowledge sources.

References

588

590

591

592

594

595

598

599

601

604

606

611

612

613

614

615

616

617

618

619 620

- Zhi Chen, Ameya Kumar, Abhinandan Das, Linyong Ma, Mo Yu, and James Glass. 2023. FinGPT: Instruction tuning for financial sentiment analysis. *arXiv preprint arXiv:2310.04779*.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. arXiv preprint arXiv:1810.04805.
- Steven Y. Feng, Varun Gangal, Jason Wei, Sarath Chandar, Soroush Vosoughi, Teruko Mitamura, and Eduard Hovy. 2021. A survey on data augmentation for text classification. *arXiv preprint arXiv:2106.07158*.
- Yunfan Gao, Yunril Xiong, Xinyu Gao, Kangxiang Jia, Jinliu Pan, Yuxi Bi, Yi Dai, Jiawei Sun, Meng Wang, and Haofen Han. 2023. Retrieval-augmented generation for large language models: A survey. arXiv preprint arXiv:2312.10997.
- Demis Hassabis and Koray Kavukcuoglu. 2024. Introducing gemini 2.0: Our new ai model for the agentic era. Google DeepMind Blog.
- Colm Kearney and Sha Liu. 2021. Textual analysis in finance. *International Review of Financial Analysis*, 78:101833.
- Varun Kumar, Ashutosh Choudhary, and Mandar Jevalikar. 2020. Data augmentation using pretrained transformer models. *arXiv preprint arXiv:2003.02245*.
- Patrick Lewis, Ethan Pérez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Urvashi Khandelwal, Pontus Stenetorp, and Sebastian

Riedel. 2020. Retrieval-augmented generation for knowledge-intensive NLP tasks. In *Advances in Neural Information Processing Systems*, volume 33, pages 9459–9474. Curran Associates, Inc.

- Ilya Loshchilov and Frank Hutter. 2017. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*.
- Linyong Ma, Zhi Chen, Qian Gui, Zhenjie Yan, Mo Yu, and Paul Pu Liang. 2021. FinQA: A dataset of numerical reasoning over financial data. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 3697–3711. Association for Computational Linguistics.
- Muhammad Waseem Nadeem, Jamel Ali, Zaher Al Aghbari, and Masnida Mohd. 2022. A review on BERTbased hybrid models for text classification. *IEEE Access*, 10:65930–65953.
- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, and 2 others. 2019. Pytorch: An imperative style, high-performance deep learning library. In *Advances in Neural Information Processing Systems 32*, pages 8024–8035. Curran Associates, Inc.
- Nils Reimers and Iryna Gurevych. 2019. Sentence-BERT: Sentence embeddings using siamese BERTnetworks. In Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP), pages 3982–3992. Association for Computational Linguistics.
- Raghuraj Sailesh Shah, Ankur Anand, Srini Chodisetti, Ankit Gupta, Raj Sanjay Patel, Sameer Patel, and Manish Gupta. 2022. FinservNLP: A library of financial shared tasks and benchmarks. In *Proceedings* of the Third Workshop on Economics and Natural Language Processing, pages 144–150. Association for Computational Linguistics.
- Rickard Stureborg, Dimitris Alikaniotis, and Yoshi Suhara. 2024. Large language models are inconsistent and biased evaluators. *arXiv preprint arXiv:2405.01724*.
- Jason Wei and Kai Zou. 2019. EDA: Easy data augmentation techniques for boosting performance on text classification tasks. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 6382–6388. Association for Computational Linguistics.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz,

678	Joe Davison, Sam Shleifer, Patrick von Platen, Clara
679	Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven
680	Le Scao, Sylvain Gugger, and 3 others. 2020. Trans-
681	formers: State-of-the-art natural language processing.
682	In Proceedings of the 2020 Conference on Empirical
683	Methods in Natural Language Processing: System
684	Demonstrations, pages 38–45, Online. Association
685	for Computational Linguistics.
686	Shijie Wu, Ozan Irsoy, Steven Lu, Vadim Dabravolski,
687	Mark Dredze, Sebastian Gehrmann, Prabhanjan Kam-
699	badur David Rosenberg and Gideon Mann 2023

690

691

692

693

- i, **n**badur, David Rosenberg, and Gideon Mann. 2023. BloombergGPT: A large language model for finance. arXiv preprint arXiv:2303.17564.
- Yute Yang, Zhipeng Chen, Can Wang, Boyu Lu, Zhigang Liu, and Hongfeng Dong. 2020. FinBERT: A pretrained language model for financial communications. arXiv preprint arXiv:2006.08097.