

---

# Human Preference Alignment in Financial Advice: A Generative AI Approach

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Financial advice is a highly regulated domain where unclear communication can  
2 cause consumer harm and regulatory breaches. Yet existing recommendation sys-  
3 tems often fail to adapt to individual risk preferences and comprehension levels.  
4 In this work, we investigate how generative AI can be used to improve both the  
5 clarity and personalisation of financial product communications. We first construct  
6 a benchmark of clarity by collecting human ratings of real financial product descrip-  
7 tions and construct a novel dataset with 25,000 synthetically generated variations.  
8 Using this dataset, we then explored two optimisation strategies for generative  
9 models: dynamic generation guided by classifier feedback, and an RLHF-style  
10 approach using the classifier as a reward model. Our findings show that clarity is  
11 shaped both by textual style and consumer profile, and that integrating preference  
12 signals significantly improves comprehensibility. This work contributes a bench-  
13 mark, models, and methods for aligning generative AI with human preferences in  
14 financial communication.

## 15 1 Introduction

16 The provision of financial advice has undergone a profound transformation in recent years, driven by  
17 the increasing demand for personalised solutions. Traditional recommendation systems in finance  
18 have primarily relied on statistical models, offering generic product suggestions that often fail to  
19 reflect the unique risk-reward preferences and financial literacy levels of individual clients. At the  
20 same time, the complexity of financial products and the opacity of communication have contributed  
21 to persistent challenges in consumer understanding and trust.

22 A substantial body of prior research has sought to evaluate how well consumers understand financial  
23 products by examining the readability of financial texts. For example, the study by Loughran and  
24 McDonald (2014a) analyzed financial disclosures using established readability metrics such as  
25 the Fog Index, the Flesch Reading Ease Score, and a measure inspired by the U.S. SEC’s plain  
26 English initiative (Loughran and McDonald, 2014b). Similarly, van Boom et al. (2016) assessed the  
27 comprehensibility of insurance contracts by applying multiple indicators, including the Common  
28 European Framework of Reference for Languages (CEFR) and the Dutch Flesch–Douma Reading  
29 Ease measure. In Burke and Fry (2019), online materials covering payday loans, personal loans, and  
30 credit card offers were evaluated through the Fog Index. While these approaches provide valuable  
31 insights into the readability of financial documents and online content, relatively little attention has  
32 been given to methods for generating texts that are not only measurable but also genuinely easier for  
33 consumers to understand.

34 Generative Artificial Intelligence (AI) offers new opportunities to address these limitations. Beyond  
35 assessing the readability of financial documents, generative models have the capacity to adapt financial  
36 communications to the linguistic and cognitive needs of diverse users. This capability is particularly

significant in a highly regulated domain such as financial advice, where unclear disclosures, or unsuitable recommendations can lead to substantial consumer harm and regulatory penalties.

In this paper, we focus on the clarity and effectiveness of financial promotions and client communications. We investigate how generative AI can be leveraged not only to generate personalised recommendations that respect risk-reward trade-offs, but also to produce explanations and disclosures that align with the preferences and comprehension levels of different client segments. Specifically, we propose and evaluate a model that aligns generative outputs with human preferences for clarity, transparency, and regulatory compliance in financial communication.

By bridging advances in generative modelling with the practical requirements of financial regulation and consumer protection, our work contributes to the growing body of research on trustworthy and human-centred AI in financial services.

2 Financial Clarity Benchmark

To establish a baseline of clarity in financial communications, we conducted a survey in which participants evaluated the comprehensibility of real product descriptions collected from online sources. Respondents were asked to rate clarity on a five-point Likert scale (1 = very unclear, 5 = very clear), the criterion being whether the information provided was sufficient to make an informed decision about the product.

To ensure representativeness, we focus on two of the most common financial products, credit cards and overdrafts, drawing material from 21 UK-based institutions offering these services. All texts were segmented into equal-length paragraphs by category (product description; risk disclosure; cost transparency; feature explanation). To normalize the ratings, we calculated the proportion of responses at each clarity level by dividing the number of ratings for that score by the total number of descriptions evaluated for specific financial institution. Our analysis revealed that low-clarity content (scores 1 and 2) made up to 44% of the evaluated text in some cases, suggesting that a substantial proportion of financial communications are not easily understandable to consumers.

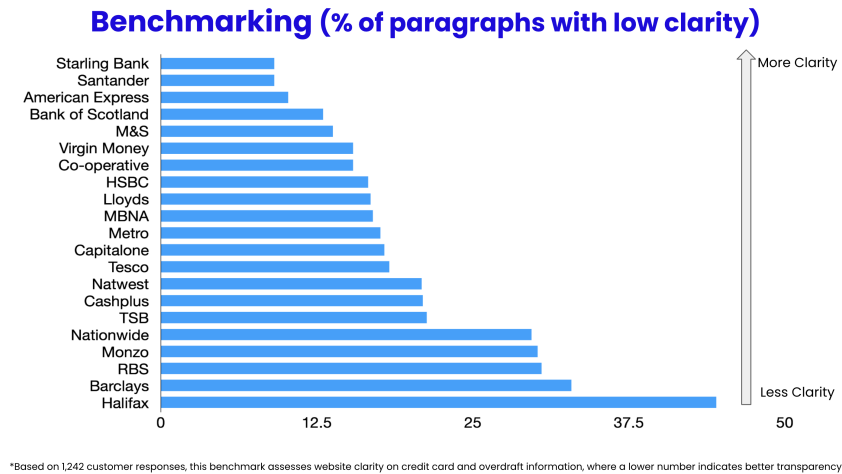


Figure 1: Distribution of clarity scores (1–5) across product descriptions from 21 UK financial institutions.

We further compared the share of low-clarity paragraphs with the number of complaints recorded for each institution in the most recent reporting period, as published by the Financial Conduct Authority. Our findings indicated a positive correlation between unclear communication and higher complaint volumes, supporting the hypothesis that insufficiently clear financial disclosures contribute to consumer dissatisfaction and regulatory risk.

This benchmark enabled us to capture patterns in how consumers perceive the clarity of financial communications across institutions and product types. Our preliminary analysis suggests that clarity

## Low Clarity Content is Positively Correlated with Higher Number of Complaints

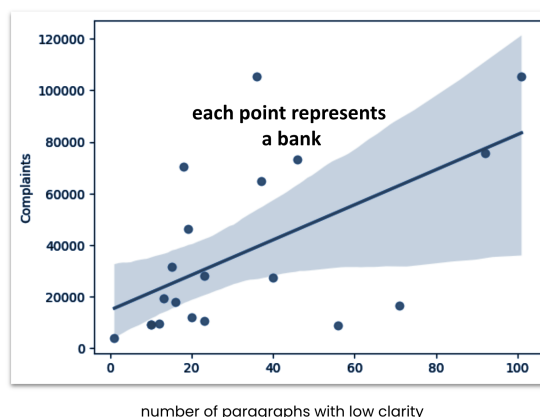


Figure 2: Positive correlation between the proportion of low-clarity content (scores 1–2) and the number of complaints reported to the Financial Ombudsman.

69 levels vary considerably across providers, highlighting systematic challenges in ensuring that financial  
70 promotions are both comprehensible and decision-useful.

### 71 3 Model training

72 For model training we synthetically generated 25,000 variations of financial product descriptions, each  
73 written in different styles. Using the same clarity evaluation framework as in our benchmark study,  
74 we asked survey participants to rank the clarity of these texts on a five-point scale. Alongside their  
75 ratings, we collected demographic information, including age, income, educational qualifications,  
76 employment status, marital status, native language, and number of financial dependents.

77 To ensure data quality, participants were required to provide a short explanation for each rating,  
78 justifying why they considered a description clear or unclear. Responses were then filtered and  
79 validated to remove inconsistent or low-effort answers.

80 We trained two types of BERT models (Devlin et al., 2019):

- 81 • A multi-class classification model to predict clarity scores based on the text alone.
- 82 • A multi-class classification model that incorporated both the text and consumer profile
- 83 attributes.

84 Our results demonstrated that incorporating consumer characteristics significantly improved predictive  
85 performance. Specifically:

- 86 • **BERT with demographic attributes** achieved an accuracy of 0.79.
- 87 • **BERT without attributes** achieved an accuracy of 0.71.

88 These findings indicate that clarity is not only a function of the text itself but also depends on the  
89 reader’s background and financial context, underscoring the importance of personalisation in financial  
90 communication.

### 91 4 Generative Model Training with Human Feedback

92 Generative models benefit substantially from **human feedback**, especially in domains such as finance  
93 where evaluation criteria like clarity and comprehensibility are inherently subjective. To improve the  
94 quality of generated financial communications, we explored two complementary methods.

## 95 4.1 Clarity Score with Dynamic Generation and Machine Feedback

96 In this approach, the generative model was conditioned on a predicted *clarity score*. A classifier  
97 trained on our benchmark dataset assessed each generated output and provided a clarity rating.  
98 The model iteratively refined its responses until the predicted clarity score exceeded a predefined  
99 threshold.

100 This method allowed for **scalable optimisation**, since it relied on machine-generated feedback loops  
101 rather than direct human annotation for every new sample. By continuously filtering generations  
102 through the classifier, the model was nudged toward producing texts that conformed to established  
103 clarity standards.

## 104 4.2 RLHF-style Training with a Reward Model

105 Our second approach was inspired by **Reinforcement Learning from Human Feedback (RLHF)**  
106 (Stiennon et al., 2020), commonly used to align large language models with user preferences. Instead  
107 of training a new reward model from pairwise comparisons, we repurposed our *clarity classifier* as a  
108 proxy reward function, since it had been trained on human-annotated clarity judgments.

109 The pipeline followed four steps:

- 110 1. **Reward Model:** The clarity classifier provided a reward signal based on human-labeled  
111 clarity scores.
- 112 2. **Policy Model:** A generative model produced candidate financial texts.
- 113 3. **Reward Signal:** The classifier evaluated each output and returned a clarity score.
- 114 4. **Optimisation:** The generative model was fine-tuned to maximise the reward, aligning  
115 generation with human preferences.

116 This RLHF-style approach offered **direct alignment with human judgments** while reducing annota-  
117 tion costs, as the reward model was trained once on human survey data and then reused for generative  
118 optimisation.

## 119 5 Conclusion

120 In this work, we investigated how generative AI can be leveraged to enhance the clarity and per-  
121 sonalisation of financial product communications. By constructing a benchmark of financial clarity,  
122 combining real-world product descriptions with 25,000 synthetic variations, and collecting human  
123 evaluations alongside demographic attributes, we demonstrated that clarity depends not only on  
124 textual features but also on consumer profiles.

125 Our experiments showed that models incorporating demographic information achieved substantially  
126 higher predictive accuracy than text-only baselines, underscoring the importance of tailoring commu-  
127 nication to user characteristics. Furthermore, we compared two generative optimisation strategies:  
128 dynamic generation guided by classifier-based feedback, and an RLHF-style approach using the  
129 classifier as a reward model. Both methods improved clarity, with the RLHF-style method offer-  
130 ing stronger alignment with human preferences, while the machine-feedback approach provided  
131 scalability.

132 These findings highlight the potential of generative AI to support clearer, more transparent, and  
133 user-aligned financial advice. At the same time, they raise important avenues for future research,  
134 including expanding the benchmark to additional product categories, integrating fairness constraints to  
135 ensure accessibility across diverse consumer groups, and exploring hybrid frameworks that combine  
136 machine and human feedback at scale. Ultimately, aligning financial communication with human  
137 preferences has the potential to increase trust, improve decision-making, and strengthen compliance  
138 in financial services.

## 139 References

- 140 Mary Burke and John Fry. How easy is it to understand consumer finance? *Economics Letters*, 177:  
141 1–4, 2019. doi: 10.1016/j.econlet.2019.01.012.

- 142 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of Deep  
143 Bidirectional Transformers for Language Understanding. In *Proceedings of the 2019 Conference of  
144 the North American Chapter of the Association for Computational Linguistics: Human Language  
145 Technologies*, pages 4171–4186, Minneapolis, Minnesota, 2019. Association for Computational  
146 Linguistics. doi: 10.18653/v1/N19-1423. URL <https://aclanthology.org/N19-1423>.
- 147 Tim Loughran and Bill McDonald. Measuring readability in financial disclosures. *the Journal of*  
148 *Finance*, 69(4):1643–1671, 2014a.
- 149 Tim Loughran and Bill McDonald. Regulation and financial disclosure: The impact of plain english.  
150 *Journal of Regulatory Economics*, 45(1):94–113, 2014b. doi: 10.1007/s11149-013-9222-4.
- 151 Nisan Stiennon, Long Ouyang, Jeffrey Wu, Daniel Ziegler, Ryan Lowe, Chelsea Voss,  
152 Alec Radford, Dario Amodei, and Paul Christiano. Learning to summarize with hu-  
153 man feedback. In *Advances in Neural Information Processing Systems*, volume 33,  
154 pages 3008–3021, 2020. URL [https://proceedings.neurips.cc/paper/2020/hash/  
155 1f89885d556929e98d3ef9b86448f951-Abstract.html](https://proceedings.neurips.cc/paper/2020/hash/1f89885d556929e98d3ef9b86448f951-Abstract.html).
- 156 Willem H. van Boom, Pieter Desmet, and Mark van Dam. “If It’s Easy to Read, It’s Easy to Claim”  
157 — The Effect of the Readability of Insurance Contracts on Consumer Expectations and Conflict  
158 Behaviour. *Journal of Consumer Policy*, 39(2):187–197, 2016. doi: 10.1007/s10603-016-9317-9.