



# Quant-GPT

Authors: Yuzhe Yang, Kangqi Yu, Junquan Peng

Student IDs: 121090684, 121090735, 120090556

Email: {yuzheyang, kangqiyu, junquanpeng}@link.cuhk.edu.cn

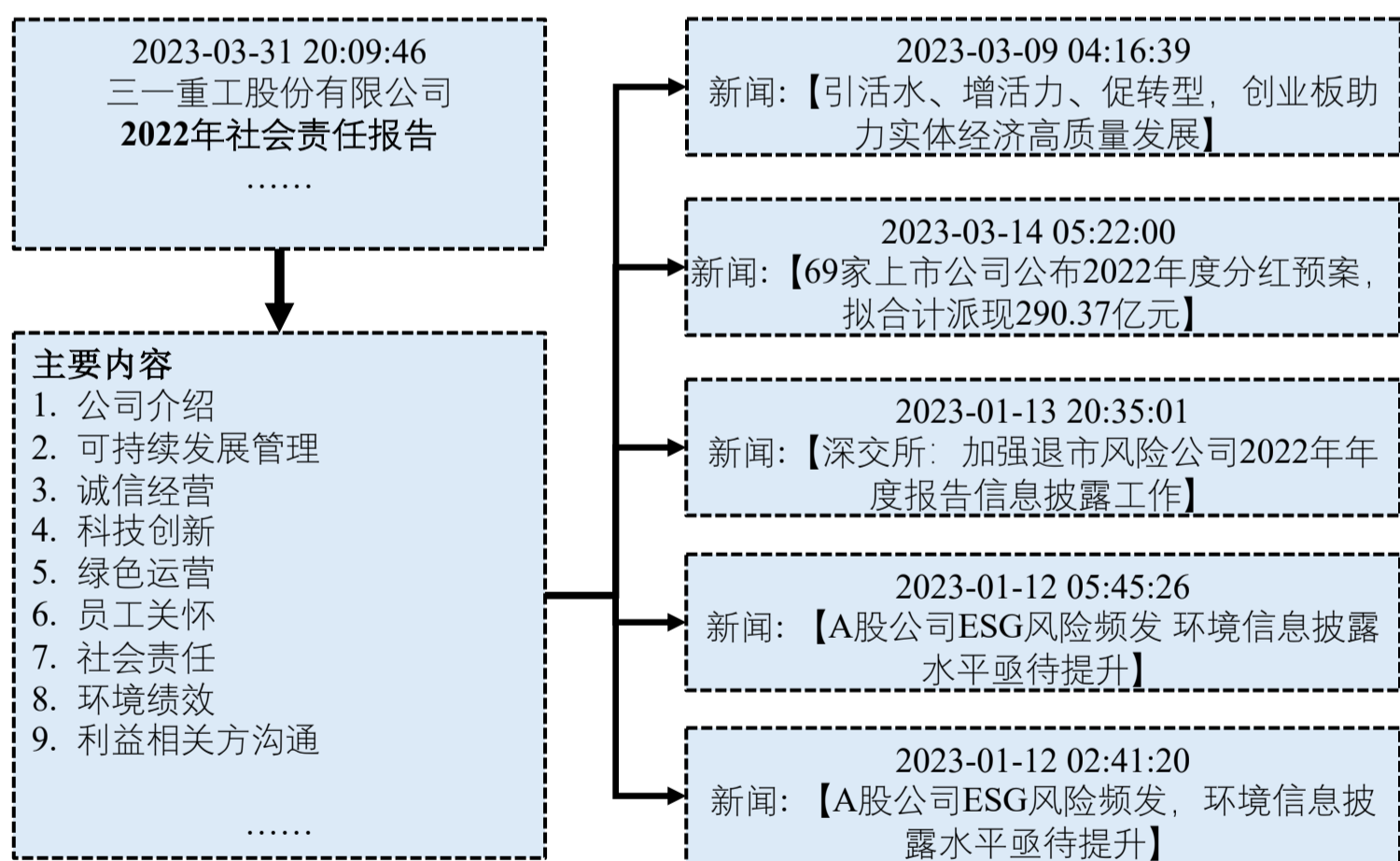
## Introduction

In this paper, we present Quant-GPT, a Large Language Model (LLM) for A-share market investment. The core recipe of Quant-GPT is to leverage both distilled sentiment analysis data from **ChatGPT** and real-world announcements from the **A shares market** in the supervised fine-tuning stage. This is not only because purely using **ChatGPT**-distilled data might cause "model collapse" and the weak causality between sentiment and expected return, but also because real-world data from **the A shares market** reflects the common expectation of all the investors.

To synergize the strengths of finance news, we introduce RAG (Retrieval-Augmented Generation) where a searching tool is designed to retrieve related news of company announcements, assisting Quant-GPT make more accurate judgments on the expected return of the announcement.

Experimental results (backtest metrics including annualized return, max drawdown, sharpe ratio) demonstrate that Quant-GPT achieves state-of-the-art results in investment decisions among open-source LLMs. It is worth noting that by using additional real-world data and RAG, the distilled language model (i.e., Quant-GPT) outperforms its teacher model (i.e., ChatGPT) in most cases. See our demo by scanning the upper QR code.

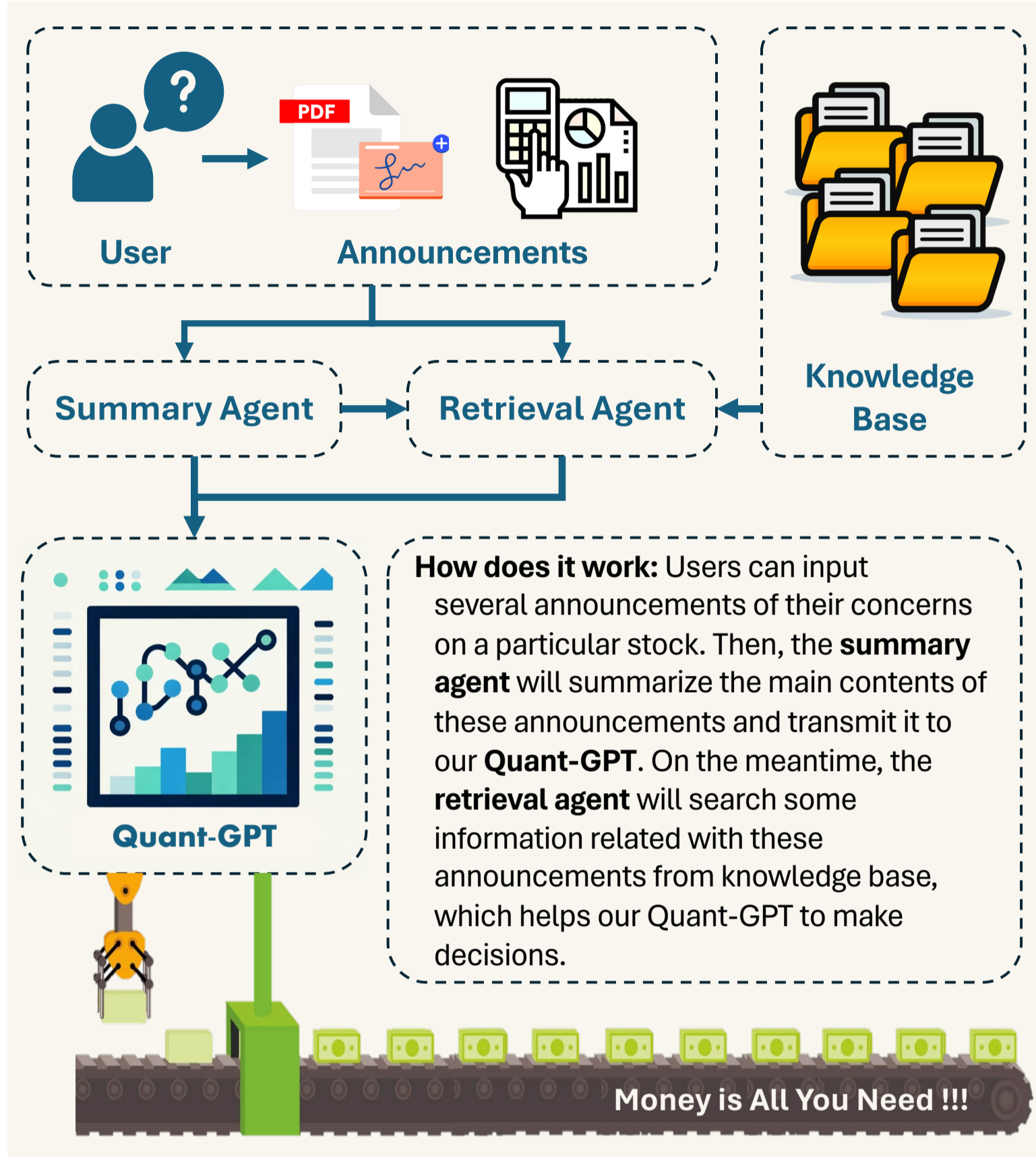
## Motivation



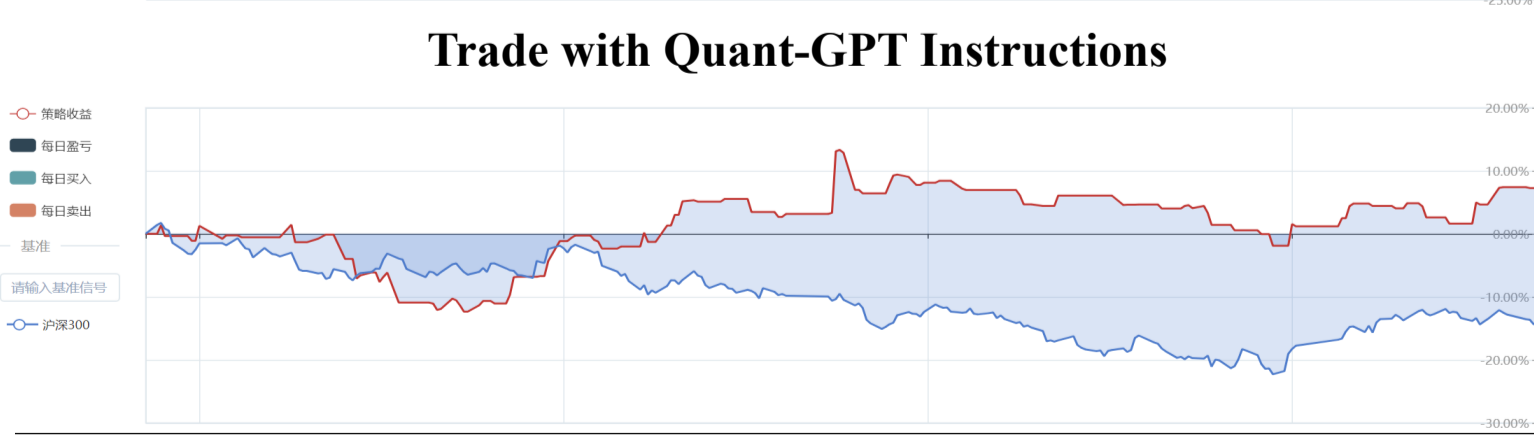
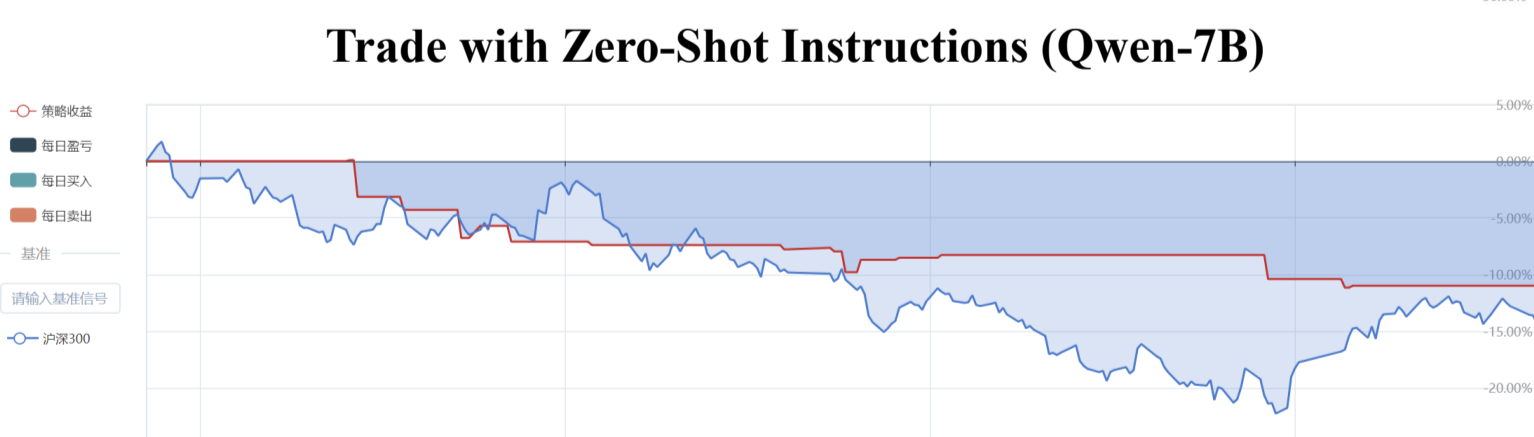
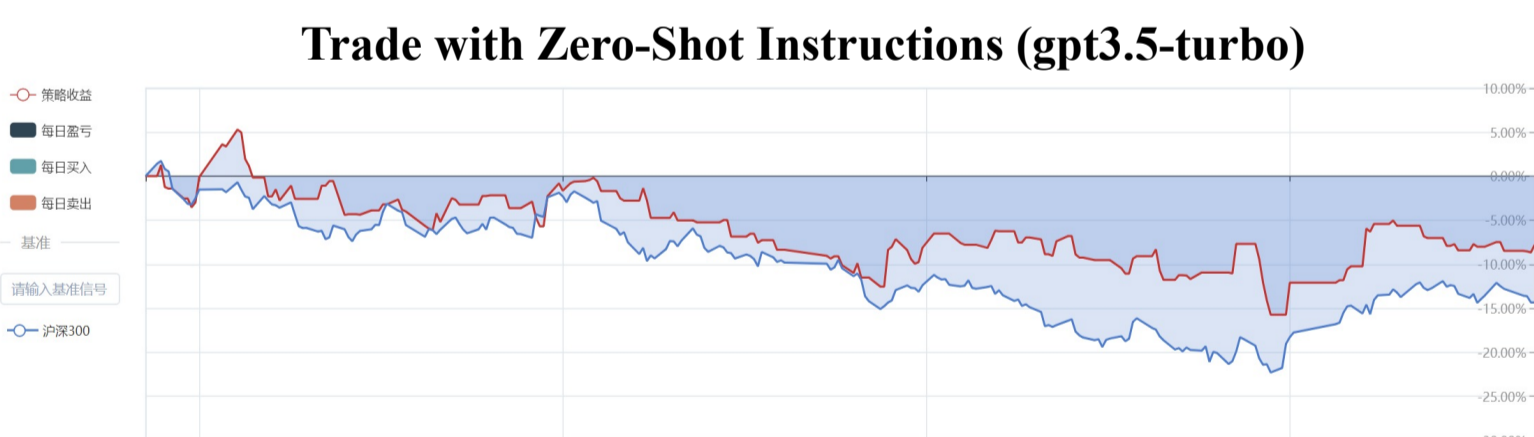
The RAG system is an agent that summarizes company announcements using a summary LLM trained on augment data from GPT. It is activated when a user inputs an announcement, focusing on key aspects such as the announcement time, company introduction, and relevant quantitative values. The system has embedded over 90,000 news articles using dense and lexical embeddings. **Dense embeddings** capture the overall semantics of the text by converting it into vectors, while **lexical embeddings** enhance the extraction of key information by assessing the importance of keywords in the text.

The top five news articles within a user-defined time window are retrieved through a hybrid similarity measure. Users can tailor the time window duration and news types in the UI, enabling the model to perform precise sentiment analysis of announcements.

## Methodology



## Results



Metrics	gpt-3.5-turbo	Qwen-7B	Quant-GPT
Annualized Return	-7.81%	-11.02%	7.26%
Max Drawdown	19.96%	11.26%	13.61%
Sharpe Ratio	-0.49	-2.02	0.40