# **INVESTORBENCH: A Benchmark for Financial Decision-Making Tasks** with LLM-based Agent

**Anonymous ACL submission** 

# Abstract

Recent advancements have underscored the po-001 tential of large language model (LLM)-based agents in financial decision-making. Despite this progress, the field currently encounters 005 two main challenges: (1) the lack of a comprehensive LLM agent framework adaptable to a variety of financial tasks, and (2) the ab-007 sence of standardized benchmarks and consistent datasets for assessing agent performance. To tackle these issues, we introduce INVESTOR-BENCH, the first benchmark specifically designed for evaluating LLM-based agents in diverse financial decision-making contexts. IN-VESTORBENCH enhances the versatility of LLM-enabled agents by providing a comprehensive suite of tasks applicable to different financial products, including single equities like 017 stocks, cryptocurrencies and exchange-traded funds (ETFs). Additionally, we assess the reasoning and decision-making capabilities of our agent framework using thirteen different LLMs as backbone models, across various market environments and tasks. Furthermore, we have curated a diverse collection of open-source, multimodal datasets and developed a comprehensive suite of environments for financial decisionmaking. This establishes a highly accessible platform for evaluating financial agents' performance across various scenarios.

#### 1 Introduction

037

041

The recent studies on large language model (LLM)based agents have demonstrated impressive performance across a range of decision-making tasks in complex and open-ended environments spanning various domains (Zhang et al., 2024b; Guo et al., 2024; Eigner and Händler, 2024; Wang et al., 2024). However, developing agentic frameworks tailored 038 specifically for financial decision-making remains a significant challenge. This complexity arises from the need for agents to acutely discern and prioritize decisive signals, and then make sequentially



Figure 1: General architecture of INVESTORBENCH.

high-quality decisions within the volatile and multifaceted financial markets, where information varies in time sensitivity and modality.

Furthermore, the design of financial agents becomes increasingly complex when applied across multiple decision-making tasks, due to the significant variation in key factors influencing financial decisions across different objectives and task types. For instance, single-equity tasks like stock trading require analyzing company-specific and industrywide data, including market metrics, sector trends, performance reports, and relevant news (Yi et al., 2022). In contrast, cryptocurrency trading is highly sensitive to crypto-specific news and sentiment due to its dynamic nature (Bhatnagar et al., 2023). ETFs, on the other hand, typically follow passive investment strategies, emphasizing long-term growth and cost efficiency (Madhavan, 2016).

The recent emergence of financial LLM-based agent frameworks such as FINMEM (Yu et al., 2024a), FINAGENT (Zhang et al., 2024a), CRYPTO-

TRADE (Li et al., 2024), FINROBOT (Yang et al., 063 2024), and FINCON (Yu et al., 2024b) has pre-064 sented a variety of architectural approaches tai-065 lored to specific financial tasks. This diversification has sparked substantial interest across both academic and industrial landscapes. FINROBOT is engineered specifically for market analysis, while FINMEM and FINAGENT are oriented towards trading individual equities like stocks and ETFs. CRYPTOTRADE focuses solely on cryptocurrency trading. FINCON pioneers in addressing portfolio management, although it currently handles only compact portfolios consisting of three stock assets. While these frameworks are effective within their respective niches, they generally focus on address-077 ing only limited types of financial decision-making tasks. This restricts them from further demonstrating the broader applicability of these frameworks and limits the comprehensive, comparative insights that could be drawn from their overall decisionmaking performance. Furthermore, the frequent reliance on proprietary financial data complicates the evaluation of these tools, obscuring their effectiveness and adaptability in broader contexts. 086 Therefore, there is a pressing need to develop innovative benchmarks specifically designed to evaluate LLM-based agents across a wider spectrum of financial decision-making scenarios. Such benchmarks would enable a more robust assessment of these technologies, facilitating advancements that could cater to various financial applications.

We introduce INVESTORBENCH, an opensource, LLM-based agent benchmark that generalizes across a broad range of financial decisionmaking tasks. Its detailed structure is illustrated in Figure 1. Further developed upon the foundational framework of FINMEM (Yu et al., 2024a), which focuses on single-stock investment decisions, our benchmark extends the scope to encompass an ensemble of diverse financial market environments for various financial tasks. INVESTOR-BENCH's cognitive architecture, similar to FIN-MEM, employs a layered memory processing mechanism with distinct decay rates, enabling the agent to store, retrieve, and consolidate insights and reflections more effectively than the pure similaritybased memory retrieval used in FINAGENT. This approach ensures that decisions are informed by timely and impactful data, a capability previously shown effective for single-asset trading. These features reflect how human traders draw sequential decisions upon investment signals from multiple

097

101

103

104

106

107

108

109

110

111

112

113

114

sources and varying time sensitivities, allowing the agent to naturally adapt to complex financial tasks. INVESTORBENCH expands its evaluation beyond the original stock trading tasks to encompass three decision tasks significant in the realm of financial investment: **stock trading, cryptocurrency trading**, and **ETF investing**. 115

116

117

118

119

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

In summary, we make three key contributions: 1) We establish INVESTORBENCH, an innovative and comprehensive financial agentic benchmark designed to evaluate the reasoning and sequential decision-making capabilities of LLM-based agents in complex, open-ended financial scenarios. This benchmark provides a realistic perspective for assessing the design and performance of such agents. 2) We provide a set of open-source, multi-source market environments that closely mirror real-world conditions. Furthermore, these environments also serve as a standardized platform for evaluating the decision-making performance of other LLM-based financial agents. 3) We present a unified, flexible language-agent framework that allows finance professionals to conveniently customize assess any LLMs serving as the agents reasoning core. In this paper, we conduct a holistic evaluation of 13 LLMs including recent, competitive, and domain-specific fine-tuned models (see Table 1) to provide a broad overview of their reasoning capabilities in sequential decision-making tasks within financial contexts.

# 2 LLM Trading Agents

In this section, we define a framework of the LLMbased agents in the INVESTORBENCH and formalize the financial decision-making tasks within the context of partially observable Markov decision process (POMDP) (Bertsekas and Shreve, 1996; Liu et al., 2020; Kabbani and Duman, 2022).

# 2.1 Definition

The LLM-based agent in INVESTORBENCH is structured as a large language model-modulo framework, designed to match or surpass the capabilities of professional human investors. This framework consists of several interconnected modules, each tailored to handle the distinct challenges presented by the financial markets volatility and complexity: **Brain/Backbone (LLM):** This module, which is the LLM itself, serves as the core of the LLMbased agent. It enhances the agent's capabilities by enabling it to understand, process, and generate

natural language. This module plays a crucial role 164 in supporting complex decision-making processes, 165 offering interpretations of market-related informa-166 tion, generating predictive analytics, and reflecting 167 on past investment decisions. 168

Perception: This module serves a critical function 169 by converting raw market data into a structured 170 format that is compatible with the LLM, specifying 171 what the agent perceives and observes, which includes numerical, textual, and visual information. 173

174

175

176

177

179

181

182

187

188

190

192

193

194

195

196

198

199

201

210

211

212

213

215

**Profile:** This module serves two functions articulated in natural language. Firstly, it describes the agent's role, highlighting its character as an experienced investor with expert-level knowledge 178 and a self-adaptive risk preference. This risk preference dynamically adjusts based on historical market momentum, allowing the agent to optimize its strategies in real time. Secondly, the module provides a detailed background of the decisionmaking task, specifying the key characteristics and pertinent information about the target assets involved in the trading decisions, such as equity historical performance, price fluctuations, and sector information. This dual-function module supports the agent's decisions with both the current market context and its historical performance.

> Memory: This module processes and retains essential market data and historical insights, allowing the agent to draw on a rich repository of knowledge for decision-making. Building upon the pioneering work of Yu et al. (2024a) in FINMEM, the memory architecture comprises two primary components: Working Memory and Layered Long-**Term Memory**, as depicted in Figure 2.

Working memory: This component maintains FINMEM's original functionalities: observation, summarization, and reflection. It incorporates two reflection mechanisms: immediate and extended. Immediate reflection produces the agent's reasoning outcomes by integrating current market indicators with the top-K ranked events from each long-term memory layer, which are significant during both warm-up and evaluation stages. In the warm-up stage, the emphasis shifts as the trading direction is predetermined, focusing on understanding market trends and enhancing predictive accuracy. In the evaluation stage, it outputs the trading direction (Buy, Sell, or Hold), the rationale for this decision, identifying the most influential memory events and their respective IDs from each layer.

Layered Long-Term Memory: Inspired by the human cognitive system's varying information decay speeds, Layered Long-Term Memory component structures financial insights across multiple layers. Each layer is represented by a vector database in the Long-Term Memory data warehouse, where information is prioritized and purged based on a specific decay rate. Deeper layers retain information longer with smaller decay rates, while shallower layers, dealing with more transient data, have larger decay rates. This tiered approach is critical as it allows the adaptation of the memory architecture to a broader range of financial tasks beyond single-asset decisions, accommodating an expanded variety of data sources and increasing overall system flexibility. Detailed mechanisms for ranking and decay in each layer are further elaborated in the Appendix A.

216

217

218

219

220

221

222

223

224

225

226

227

228

229

230

231

232

233

234

235

236

237

238

239

240

241

242

243

244

245

247

248

249

250

251

252

253

254

255

256

257

258

259

261

262

263

265

266

Action: This module executes trading and investment decisions based on the analysis provided by other modules. It directly outputs {"Buy", "Sell", "Hold"} for traded asset (stock, crypto, or ETF), as instructed by the backbone LLM. Action module synthesizes the operational outcomes from the Profile and Memory modules to facilitate precise and well-informed investment decisions. For its daily trading operations, the agent can choose from three specific actions for the traded asset: "Buy", "Sell", or "Hold". The functionality and input requirements of this module differ significantly between the warm-up and evaluation stages: during the warm-up stage, the agent observes daily adjusted price differences between consecutive days, which are critical for identifying potential "Buy" or "Sell" signals. This period allows the agent to calibrate and adjust its decisionmaking strategies based on near-term market movements; during the evaluation stage, access to future price data is restricted, compelling the agent to rely solely on available historical data and its cognitive processing capabilities. In response to trading inquiries, the module integrates historical Profit & Loss (PnL), outcomes from extended reflections, and the top-K retrieved memories. This integration ensures that each trading decision is grounded in a comprehensive analysis of past performance and current market conditions.

#### Modeling financial decision-making 2.2

Formally, we model a financial decision-making process as infinite horizon POMDP with time index  $\mathbb{T} = \{0, 1, 2, \cdots\}$  and discount factor  $\alpha \in (0, 1]$ .



Figure 2: (1) The language agent's memory module is crafted to interact with the market environment to conduct various financial decision-making tasks. It contains two core components – Working Memory and Layered Long-term Memory. (2) The outline of the agent's decision-making workflow for retrieving critical memory events and market observations to inform specific investment decisions.

This POMDP contains: (1) a state space  $\mathcal{X} \times \mathcal{Y}$ where  $\mathcal{X}$  is the observable component and  $\mathcal{Y}$  is unobservable component of the financial market; (2) the action space of the agent is A, which is modeled as {"Buy", "Sell", "Hold"}; (3) the reward function  $\mathcal{R}(o, b, a) : \mathcal{X} \times \mathcal{Y} \times \mathcal{A} \to \mathbb{R}$  uses daily profit & loss (PnL) as the output; (4) the observation process  $\{O_t\}_{t\in\mathbb{T}} \subseteq \mathcal{X}$  is a multi-dimensional process (5) the reflection process  $\{B_t\}_{t\in\mathbb{T}} \subseteq \mathcal{Y}$ represents the agent's self-reflection, which is updated from  $B_t$  to  $B_{t+1}$  on daily basis (Griffiths et al., 2023); (6) the action  $A_t \sim \pi(\cdot | \text{prompt})$  represents the way to make investment decision driven by the language conditioned policy  $\pi$ . By denoting daily profit & loss (PnLs) by  $R_t^{\pi} = \mathcal{R}(O_t, B_t, A_t)$ and the set of all admissible language conditioned policies as  $\Pi = \{\pi(\cdot | \text{prompt})\}\$ , the optimization objective of financial decision-making task is then:

$$\max_{\pi \in \Pi} \mathbb{E} \Big[ \sum_{t \in \mathbb{T}} \alpha^t R_t^{\pi} \Big] \tag{1}$$

# 3 InvestorBench

270

271

274

275

277

278

279

283

287

290

293

He we introduce the detailed architecture of InvestorBench, as illustrated in Figure 1.

## 3.1 Benchmark Composition

INVESTORBENCH is organized into four main components: (1) **Data Sources and Market Environments**: INVESTORBENCH utilizes a wide range of open-source data and incorporates third-party APIs, such as Yahoo Finance and SEC EDGAR, to create a comprehensive, multi-modal market environment data warehouse. (2) LLM Agent: INVESTOR-BENCH includes an advanced LLM-based agent equipped with modules for Brain, Perception, Profile, Memory, and Action. This agent is enhanced with external tools (such as tabular data readers and API callers) and data operations (including vector database management, information reinforcement, and retrieval). (3) Financial Decision-Making Tasks: INVESTORBENCH offers three distinct financial decision-making tasks, differentiated by their asset types. (4) Evaluation Metrics: The efficacy of all tasks within INVESTORBENCH is evaluated using a set of standard metrics in the quantitative finance field, providing a thorough evaluation of the decision-making capabilities of the LLMbased agent.

294

295

296

297

298

299

301

302

303

304

305

306

307

308

309

310

311

312

313

314

315

316

317

318

319

320

321

322

323

## 3.2 Trading Environments

We release three datasets, each curated from diverse sources, to construct tailored financial market environments for specific tasks. Our objective is to address the current gap in evaluation environments for financial decision-making agent frameworks and to offer a fully open platform for the comprehensive assessment of agents across various tasks. Below, we introduce each environment, categorized by task type, detailing its scope and the data sources it incorporates.

Stock market environment integrates information 324 from multiple sources, including: 1) Daily stock open, high, low, close, and volume (OHLCV) data acquired from Yahoo Finance. 2) Summarized insights from company quarterly and annual reports (Form 10-Q and 10-K) downloaded from the SEC EDGAR database. 3) News articles for 330 seven stocks collected daily between 2020-07-01, 331 and 2021-05-06. The news data for four of these companiesMicrosoft Corporation (MSFT), John-333 son & Johnson (JNJ), UVV Corporation (UVV), and Honeywell International Inc. 335 (HON)-are randomly selected from the pool with the most new 336 records (over five hundred) from the open-access 337 dataset provided by Zhou et al. (Zhou et al., 2021), while the news data for the remaining three companiesTesla, Inc. (TSLA), Apple Inc. (AAPL), and NIO Inc. (NIO)-are obtained from 341 Refinitiv Real-Time News, which primarily contains high-quality news information from Reuters. 4) The sentiment categories ('positive', 'negative', 'neutral') assigned to each news record are generated by gpt-3.5-turbo-0125.

**Cryptocurrency market environment** encompasses 1) the daily stock open-high-low-close-volume (OHLCV) acquired from CoinMarketCap; 2) the multisource cryptocurrency news data collected from cryptonews, cryptopotato, and cointelegraph(Vanhoucke, 2023); 3) news spanning from 2023-02-13 to 2023-11-05 collected by (Zhou et al., 2021) in daily frequency. 4) The sentiment categories generated by the same means.

**ETF market environment** is constructed using News-Informed Financial Trend Yield (NIFTY) dataset (Saqur et al., 2024). It contains the processed and curated daily news headlines from 2019-07-29 to 2020-09-21 and generated sentiment categories for each news headline.

In experimental use, we divide the dataset according to the date, with the train set used for the warmup phase to establish the memory database, and the test set used for the test phase to evaluate the model performance.

# **3.3 Evaluation metrics**

347

354

370

371

374

We employ four widely recognized financial metrics to evaluate and compare the investment performance of various LLMs serving as backbones across different tasks: : Cumulative Return (CR) (Hull, 2007), Sharpe Ratio (SR) (Sharpe, 1994), Table 1: INVESTORBENCH evaluates 13 proprietary or open-source LLMs on financial decision-making tasks.

Model	#Size	Form	Ver.	Model	#Size	Form	Ver.
gpt-4(Achiam et al., 2023)	N/A	api	0613	Qwen2.5-7b(Qwen team, 2024)	7B	open	Instruct
gpt-4o(OpenAI, 2022)	N/A	api	0806	Qwen2.5-32b(Qwen team, 2024)	32B	open	Instruct
gpt-o1-preview	N/A	api	0912	Qwen2.5-72b(Qwen team, 2024)	72B	open	Instruct
DeepSeek-v2(Xin et al., 2024)	15B	open	Lite	11ama3.1-8b(Llama team, 2024)	8B	open	Instruct
DeepSeek-11m(Xin et al., 2024)	67B	open	Chat	11ama3.1-70b(Llama team, 2024)	70B	open	Instruct
Yi-1.5-9b(Young et al., 2024)	9B	open	Chat	Palmyra-Fin(team, 2024)	70B	open	32K
Yi-1.5-34b(Young et al., 2024)	34B	open	Chat				

Annualized Volatility (AV) (Cochrane, 1988), and Maximum Drawdown(MDD) (Ang and Chen, 2003). Note that CR and the SR are often considered more essential than AV and MDD in evaluating asset trading performance due to their focus on long-term gains and risk-adjusted returns by their definition. Here, we regard these two metrics as primary metrics when evaluating the experiment outcomes. The detailed explanation is in Appendix A.

# 4 Experiment and Discussion





(b) By model parameter sizes across open-source LLMs. Note: Small-size models refer to models with no more than 10B parameters. Medium-size models refer to models with parameter numbers in the range of (10B, 65B]. Large-size models are those with more than 65B parameters.

Figure 3: Agent Performance Comparisons from two key perspectives. *The CR, SR, AV, and MDD represent the average values for each model type, expressed as a percentage relative to the Buy & Hold strategy.* 

To establish a baseline and assess the performance of LLM agents, we standardize experimental settings and evaluation metrics across various financial decision-making tasks. Results are pre375

376

377

378

379

382

383

384

385

Model	odel MSFT				JNJ				UVV				HON			
	CR↑	SR↑	AV↓	MDD↓	CR↑	SR↑	AV↓	MDD↓	CR↑	SR↑	AV↓	MDD↓	CR↑	SR↑	AV↓	MDD↓
Buy & Hold	15.340	1.039	24.980	9.428	13.895	1.343	17.500	9.847	36.583	2.112	29.299	15.406	33.256	2.347	23.967	9.195
						Financia	l Domain	Models								
Palmyra-Fin-70B	14.697	0.897	27.518	9.428	5.748	0.450	19.317	9.367	37.875	2.039	31.200	15.967	20.016	1.464	22.974	6.824
						Prop	ietary Me	odels								
GPT-o1-preview	17.184	0.962	30.000	9.428	13.561	1.086	20.864	9.847	41.508	2.147	32.479	9.633	13.162	0.776	28.511	11.558
GPT-4	16.654	0.932	30.022	9.428	13.712	1.103	20.894	9.860	31.791	1.640	32.567	10.434	34.342	2.005	28.779	9.195
GPT-40	12.461	0.924	22.653	6.647	9.099	0.875	17.471	7.169	8.043	0.496	27.241	14.889	38.540	2.418	26.782	8.979
						Open-	Source M	odels								
Qwen2.5-72B-Instruct	7.421	0.588	21.238	6.973	14.353	1.140	20.995	9.812	37.178	1.822	34.223	13.365	34.309	2.000	28.779	9.292
Llama-3.1-70B-Instruct	17.396	1.335	21.892	7.045	13.868	1.121	20.779	9.825	35.981	1.728	34.986	15.406	43.944	2.646	27.903	8.993
DeepSeek-67B-Chat	13.941	0.834	28.081	7.850	14.426	1.185	20.450	9.825	29.940	1.481	33.964	15.407	32.536	1.909	28.628	10.782
Yi-1.5-34B-Chat	22.093	1.253	29.613	9.428	14.004	1.180	19.938	9.847	20.889	1.020	34.417	14.936	30.743	1.823	28.335	9.195
Qwen2.5-32B-Instruct	-0.557	-0.041	22.893	8.946	2.905	0.292	16.725	7.169	-1.623	-0.097	27.973	17.986	26.332	1.980	22.348	5.261
DeepSeek-V2-Lite (15.7B)	11.904	0.694	28.796	16.094	-7.482	-0.670	18.773	17.806	33.560	1.703	33.099	12.984	16.686	0.974	28.771	16.806
Yi-1.5-9B-Chat	19.333	1.094	29.690	9.428	18.606	1.611	19.409	10.986	49.415	2.410	34.446	11.430	29.028	1.700	28.682	12.588
Llama-3.1-8B-Instruct	22.703	1.322	28.855	7.385	13.988	1.486	20.460	9.969	41.108	1.981	34.866	16.429	39.079	2.320	28.299	10.341
Qwen-2.5-Instruct-7B	-10.305	-0.724	23.937	23.371	21.852	0.980	37.425	9.573	11.752	0.853	22.988	15.451	4.291	0.285	24.933	14.156
									NIO							
Model		TS	LA			AA	PL			N	10			Ave	rage	
Model	CR↑	TS SR↑	AV↓	MDD↓	CR↑	AA SR↑	AV↓	MDD↓	CR↑	N SR↑	IO AV↓	MDD↓	CR↑	Ave SR↑	rage AV↓	MDD↓
Model Buy & Hold	CR↑ 39.244	TS SR↑ 0.869	AV↓ 75.854	<b>MDD</b> ↓ 37.975	<b>CR</b> ↑ 10.837	AA SR↑ 0.470	AV↓ 38.984	<b>MDD</b> ↓   19.119	<b>CR</b> ↑ 52.216	N SR↑ 0.858	IO AV↓ 107.197	<b>MDD</b> ↓ 47.766	<b>CR</b> ↑   34.099	Ave SR↑ 0.732	rage AV↓ 74.012	<b>MDD</b> ↓ 34.953
Model Buy & Hold	<b>CR</b> ↑ 39.244	TS SR↑ 0.869	AV↓ 75.854	<b>MDD</b> ↓ 37.975	<b>CR</b> ↑ 10.837	AA SR↑ 0.470 Financia	PL AV↓ 38.984 <i>l Domain</i>	MDD↓   19.119   Models	<b>CR</b> ↑ 52.216	N SR↑ 0.858	IO AV↓ 107.197	<b>MDD</b> ↓ 47.766	<b>CR</b> ↑   34.099	Ave SR↑ 0.732	rage AV↓ 74.012	<b>MDD</b> ↓ 34.953
Model Buy & Hold Palmyra-Fin-70B	CR↑ 39.244 -6.661	TS SR↑ 0.869 -0.222	AV↓ 75.854 50.379	MDD↓ 37.975 25.820	CR↑ 10.837 8.562	AA SR↑ 0.470 Financia 0.372	PL AV↓ 38.984 <i>l Domain</i> 38.891	MDD↓   19.119   <i>Models</i> 25.466	CR↑ 52.216 -3.261	N SR↑ 0.858 -0.057	IO AV↓ 107.197 101.711	MDD↓ 47.766 58.406	<b>CR</b> ↑   34.099   -0.453	Ave SR↑ 0.732 0.031	rage AV↓ 74.012 63.660	MDD↓ 34.953 36.564
Model Buy & Hold Palmyra-Fin-70B	CR↑ 39.244 -6.661	TS SR↑ 0.869 -0.222	ELA AV↓ 75.854 50.379	MDD↓   37.975   25.820	<b>CR</b> ↑ 10.837 8.562	AA SR↑ 0.470 Financia 0.372 Propr	PL AV↓ 38.984 <i>l Domain</i> 38.891 <i>ietary Me</i>	MDD↓   19.119   Models 25.466   odels	CR↑ 52.216 -3.261	N SR↑ 0.858 -0.057	IO AV↓ 107.197 101.711	MDD↓ 47.766 58.406	<b>CR</b> ↑   34.099   -0.453	Ave SR↑ 0.732 0.031	rage AV↓ 74.012 63.660	MDD↓ 34.953 36.564
Model Buy & Hold Palmyra-Fin-70B GPT-o1-preview	CR↑ 39.244 -6.661 34.499	TS SR↑ 0.869 -0.222 0.796	AV↓           75.854           50.379           72.822	MDD↓         37.975         25.820         35.490	CR↑           10.837           8.562           8.238	AA SR↑ 0.470 Financia 0.372 Propr 0.422	PL AV↓ 38.984 <i>l Domain</i> 38.891 <i>ietary Me</i> 33.045	MDD↓                 19.119                 25.466                 odels                 14.412	CR↑           52.216           -3.261           32.433	N SR↑ 0.858 -0.057 0.558	IO           AV↓           107.197           101.711           102.470	MDD↓ 47.766 58.406 54.016	CR↑   34.099   -0.453   25.057	Ave SR↑ 0.732 0.031 0.592	rage AV↓ 74.012 63.660 69.446	MDD↓ 34.953 36.564 34.639
Model Buy & Hold Palmyra-Fin-70B GPT-01-preview GPT-4	CR↑ 39.244 -6.661 34.499 45.246	TS SR↑ 0.869 -0.222 0.796 1.190	AV↓           75.854           50.379           72.822           63.896	MDD↓   37.975   25.820   35.490   25.031	CR↑           10.837           8.562           8.238           9.889	AA SR↑ 0.470 Financia 0.372 Propr 0.422 0.440	PL AV↓ 38.984 <i>l Domain</i> 38.891 <i>rietary Me</i> 33.045 38.008	MDD↓                 19.119                 25.466                 odels                 14.412                 19.119	CR↑ 52.216 -3.261 32.433 75.952	SR↑           0.858           -0.057           0.558           1.286	IO           AV↓           107.197           101.711           102.470           104.059	MDD↓ 47.766 58.406 54.016 37.867	CR↑           34.099           -0.453           25.057           43.696	Ave SR↑ 0.732 0.031 0.592 0.972	rage AV↓ 74.012 63.660 69.446 68.654	MDD↓ 34.953 36.564 34.639 27.339
Model Buy & Hold Palmyra-Fin-70B GPT-01-preview GPT-4 GPT-40	CR↑           39.244           -6.661           34.499           45.246           45.946	TS SR↑ 0.869 -0.222 0.796 1.190 1.348	AV↓           75.854           50.379           72.822           63.896           57.281	MDD↓                 37.975                 25.820                 35.490                 25.031                 21.631	CR↑           10.837           8.562           8.238           9.889           7.405	AA SR↑ 0.470 Financia 0.372 Prop 0.422 0.440 0.457	PL AV↓ 38.984 <i>l Domain</i> 38.891 <i>ietary Ma</i> 33.045 38.008 27.434	MDD↓                 19.119                 19.119                 Models                 25.466                 0dels                 14.412                 19.119                 12.824	CR↑           52.216           -3.261           32.433           75.952           63.743	SR↑           0.858           -0.057           0.558           1.286           1.318	AV↓           107.197           101.711           102.470           104.059           85.210	MDD↓ 47.766 58.406 54.016 37.867 29.220	CR↑           34.099           -0.453           25.057           43.696           39.031	Ave SR↑ 0.732 0.031 0.592 0.972 1.041	AV↓           74.012           63.660           69.446           68.654           56.642	MDD↓ 34.953 36.564 34.639 27.339 21.225
Model Buy & Hold Palmyra-Fin-70B GPT-01-preview GPT-4 GPT-4o	CR↑           39.244           -6.661           34.499           45.246           45.946	TS SR↑ 0.869 -0.222 0.796 1.190 1.348	AV↓           75.854           50.379           72.822           63.896           57.281	MDD↓         37.975         25.820         35.490         25.031         21.631	CR↑           10.837           8.562           8.238           9.889           7.405	AA SR↑ 0.470 Financia 0.372 Propu 0.422 0.440 0.457 Open-	PL AV↓ 38.984 <i>l Domain</i> 38.891 <i>rietary Ma</i> 33.045 38.008 27.434 Source M	MDD↓                 19.119                 Models       25.466         25.466                 odels                 14.412                 19.119                 12.824	CR↑           52.216           -3.261           32.433           75.952           63.743	N           SR↑           0.858           -0.057           0.558           1.286           1.318	AV↓           107.197           101.711           102.470           104.059           85.210	MDD↓ 47.766 58.406 54.016 37.867 29.220	CR↑           34.099           -0.453           25.057           43.696           39.031	Ave SR↑ 0.732 0.031 0.592 0.972 1.041	rage           AV↓           74.012           63.660           69.446           68.654           56.642	MDD↓ 34.953 36.564 34.639 27.339 21.225
Model Buy & Hold Palmyra-Fin-70B GPT-01-preview GPT-4 GPT-40 Qwen2.5-72B-Instruct	CR↑           39.244           -6.661           34.499           45.246           45.946           39.112	TS SR↑ 0.869 -0.222 0.796 1.190 1.348 1.075	ILA           AV↓           75.854           50.379           72.822           63.896           57.281           61.136	MDD↓                     37.975                     25.820                     35.490                     25.031                     21.631	CR↑           10.837           8.562           8.238           9.889           7.405           11.935	AA SR↑ 0.470 Financia 0.372 Propu 0.422 0.440 0.457 Open- 0.572	PL AV↓ 38.984 <i>l Domain</i> 38.891 <i>ietary Ma</i> 33.045 38.008 27.434 <i>Source M</i> 35.293	MDD↓                     19.119                     Models         25.466           25.466                     odels                     14.412                     19.119                     12.824                //odels                     19.119	CR↑           52.216           -3.261           32.433           75.952           63.743           87.412	N           SR↑           0.858           -0.057           0.558           1.286           1.318           2.181	AV↓           107.197           101.711           102.470           104.059           85.210           70.628	MDD↓ 47.766 58.406 54.016 37.867 29.220 12.464	CR↑           34.099           -0.453           25.057           43.696           39.031           46.153	Ave SR↑ 0.732 0.031 0.592 0.972 1.041 1.276	rage AV↓ 74.012 63.660 69.446 68.654 56.642 55.686	MDD↓ 34.953 36.564 34.639 27.339 21.225 19.523
Model Buy & Hold Palmyra-Fin-70B GPT-0-preview GPT-4 GPT-40 Qwen2.5-72B-Instruct Llama-3.1-70B-Instruct	CR↑           39.244           -6.661           34.499           45.246           45.946           39.112           37.545	TS SR↑ 0.869 -0.222 0.796 1.190 1.348 1.075 0.891	ILA           AV↓           75.854           50.379           72.822           63.896           57.281           61.136           70.815	MDD↓                     37.975                     25.820                     35.490                     25.031                     21.631                     26.985                     29.813	CR↑           10.837           8.562           8.238           9.889           7.405           11.935           12.772	AA SR↑ 0.470 Financia 0.372 Propu 0.422 0.440 0.457 Open- 0.572 0.583	PL AV↓ 38.984 <i>I Domain</i> 38.891 <i>ietary Ma</i> 33.045 38.008 27.434 Source M 35.293 37.057	MDD↓         I           19.119         I           Models         25.466           25.466         I           odels         14.412           19.119         12.824           Idels         19.119           12.824         I           Idels         19.119           16.021         I	CR↑           52.216           -3.261           32.433           75.952           63.743           87.412           66.522	N           SR↑           0.858           -0.057           0.558           1.286           1.318           2.181           1.118	IO           AV↓           107.197           101.711           102.470           104.059           85.210           70.628           104.848	MDD↓ 47.766 58.406 54.016 37.867 29.220 12.464 46.379	CR↑           34.099           -0.453           25.057           43.696           39.031           46.153           38.946	Ave SR↑ 0.732 0.031 0.592 0.972 1.041 1.276 0.864	rage AV↓ 74.012 63.660 69.446 68.654 56.642 55.686 70.907	MDD↓ 34.953 36.564 34.639 27.339 21.225 19.523 30.738
Model Buy & Hold Palmyra-Fin-70B GPT-01-preview GPT-4 GPT-40 Qven2.5-72B-Instruct Llama-3.1-70B-Instruct DeepSeek-67B-Chat	CR↑           39.244           -6.661           34.499           45.246           45.946           39.112           37.545           35.647	TS SR↑ 0.869 -0.222 0.796 1.190 1.348 1.075 0.891 0.885	ILA           AV↓           75.854           50.379           72.822           63.896           57.281           61.136           70.815           67.660	MDD↓         37.975           35.820         35.490           25.031         21.631           26.985         29.813           33.359         33.359	CR↑           10.837           8.562           8.238           9.889           7.405           11.935           12.772           14.213	AA SR↑ 0.470 Financia 0.372 Propu 0.422 0.440 0.457 Open- 0.572 0.583 0.666	PL AV↓ 38.984 <i>I Domain</i> 38.891 <i>ietary Ma</i> 33.045 38.008 27.434 Source M 35.293 37.057 36.118	MDD↓         I           19.119         I           Models         25.466           25.466         I           odels         14.412           19.119         12.824           Iodels         19.119           12.824         I           Iodels         19.119           16.021         10.876	CR↑           52.216           -3.261           32.433           75.952           63.743           87.412           66.522           30.963	N           SR↑           0.858           -0.057           0.558           1.286           1.318           2.181           1.118           0.599	IO           AV↓           107.197           101.711           102.470           104.059           85.210           70.628           104.848           91.146	MDD↓           47.766           58.406           54.016           37.867           29.220           12.464           46.379           45.855	CR↑           34.099           -0.453           25.057           43.696           39.031           46.153           38.946           26.941	Ave SR↑ 0.732 0.031 0.592 0.972 1.041 1.276 0.864 0.717	rage AV↓ 74.012 63.660 69.446 68.654 56.642 55.686 70.907 64.975	MDD↓ 34.953 36.564 34.639 27.339 21.225 19.523 30.738 30.030
Model Buy & Hold Palmyra-Fin-70B GPT-0-preview GPT-4 GPT-40 Qwen2.5-72B-Instruct Llama-3.1-70B-Instruct DeepSeek-67B-Chat Yi-1.5-34B-Chat	CR↑           39.244           -6.661           34.499           45.246           45.946           39.112           37.545           35.647           35.364	TS SR↑ 0.869 -0.222 0.796 1.190 1.348 1.075 0.891 0.885 0.808	ILA           AV↓           75.854           50.379           72.822           63.896           57.281           61.136           70.815           67.660           73.561	MDD↓         37.975           35.490         25.820           35.490         1.631           26.985         29.813           33.359         35.490	CR↑           10.837           8.562           8.238           9.889           7.405           11.935           12.772           14.213           14.227	AA SR↑ 0.470 Financia 0.372 Propu 0.422 0.440 0.457 Open- 0.572 0.583 0.666 0.623	PL AV↓ 38.984 <i>l Domain</i> 38.891 <i>ietary Ma</i> 33.045 38.008 27.434 Source M 35.293 37.057 36.118 38.596	MDD↓         I           19.119         I           Models         25.466           25.466         I           ndels         14.412           19.119         12.824           Iodels         I           10.2824         I           Iodels         I           10.876         19.432	CR↑           52.216           -3.261           32.433           75.952           63.743           87.412           66.522           30.963           64.307	N           SR↑           0.858           -0.057           0.558           1.286           1.318           2.181           1.118           0.599           1.063	IO           AV↓           107.197           101.711           102.470           104.059           85.210           70.628           104.848           91.146           106.597	MDD↓           47.766           58.406           54.016           37.867           29.220           12.464           46.379           45.855           48.042	CR↑           34.099           -0.453           25.057           43.696           39.031           46.153           38.946           26.941           37.966	Ave SR↑ 0.732 0.031 0.592 0.972 1.041 1.276 0.864 0.717 0.831	rage AV↓ 74.012 63.660 69.446 68.654 56.642 55.686 70.907 64.975 72.918	MDD↓ 34.953 36.564 34.639 27.339 21.225 19.523 30.738 30.030 34.321
Model Buy & Hold Palmyra-Fin-70B GPT-01-preview GPT-4 GPT-40 Qwen2.5-72B-Instruct Llama-3.1-70B-Instruct DeepSeek-67B-Chat Yi-1.5-34B-Chat Qwen2.5-32B-Instruct	CR↑           39.244           -6.661           34.499           45.246           45.946           39.112           37.545           35.647           35.364           21.336	TS SR↑ 0.869 -0.222 0.796 1.190 1.348 1.075 0.891 0.885 0.808 0.729	ILA           AV↓           75.854           50.379           72.822           63.896           57.281           61.136           60.815           67.660           73.561           49.157	MDD↓         37.975         37.975         35.490         35.490         35.490         35.490         35.490         30.311         30.312 </td <td>CR↑           10.837           8.562           8.238           9.889           7.405           11.935           12.772           14.213           14.227           13.220</td> <td>AA SR↑ 0.470 Financia 0.372 Prop 0.422 0.440 0.457 Open- 0.572 0.583 0.666 0.623 1.089</td> <td>PL AV↓ 38.984 <i>I Domain</i> 38.891 <i>ietary Mc</i> 33.045 38.008 27.434 <i>Source M</i> 35.293 37.057 36.118 38.596 20.522</td> <td>MDD↓         I           19.119         I           Models         25.466           25.466         I           odels         I           14.412         I           19.119         I           2.824         I           odels         I           19.119         I           16.021         I           19.432         8.943</td> <td>CR↑           52.216           -3.261           32.433           75.952           63.743           87.412           66.522           30.963           64.307           28.096</td> <td>N           0.858           -0.057           0.558           1.286           1.318           2.181           1.118           0.599           1.063           0.652</td> <td>IO           AV↓           107.197           101.711           102.470           104.059           85.210           70.628           91.146           106.597           72.344</td> <td>MDD↓ 47.766 58.406 54.016 37.867 29.220 12.464 46.379 45.855 48.042 37.975</td> <td>CR↑           34.099           -0.453           25.057           43.696           39.031           46.153           38.946           26.941           37.966           20.884</td> <td>Ave SR↑ 0.732 0.031 0.592 0.972 1.041 1.276 0.864 0.717 0.831 0.823</td> <td>rage AV↓ 74.012 63.660 69.446 68.654 56.642 55.686 70.907 64.975 72.918 47.341</td> <td>MDD↓ 34.953 36.564 34.639 27.339 21.225 19.523 30.738 30.030 34.321 22.541</td>	CR↑           10.837           8.562           8.238           9.889           7.405           11.935           12.772           14.213           14.227           13.220	AA SR↑ 0.470 Financia 0.372 Prop 0.422 0.440 0.457 Open- 0.572 0.583 0.666 0.623 1.089	PL AV↓ 38.984 <i>I Domain</i> 38.891 <i>ietary Mc</i> 33.045 38.008 27.434 <i>Source M</i> 35.293 37.057 36.118 38.596 20.522	MDD↓         I           19.119         I           Models         25.466           25.466         I           odels         I           14.412         I           19.119         I           2.824         I           odels         I           19.119         I           16.021         I           19.432         8.943	CR↑           52.216           -3.261           32.433           75.952           63.743           87.412           66.522           30.963           64.307           28.096	N           0.858           -0.057           0.558           1.286           1.318           2.181           1.118           0.599           1.063           0.652	IO           AV↓           107.197           101.711           102.470           104.059           85.210           70.628           91.146           106.597           72.344	MDD↓ 47.766 58.406 54.016 37.867 29.220 12.464 46.379 45.855 48.042 37.975	CR↑           34.099           -0.453           25.057           43.696           39.031           46.153           38.946           26.941           37.966           20.884	Ave SR↑ 0.732 0.031 0.592 0.972 1.041 1.276 0.864 0.717 0.831 0.823	rage AV↓ 74.012 63.660 69.446 68.654 56.642 55.686 70.907 64.975 72.918 47.341	MDD↓ 34.953 36.564 34.639 27.339 21.225 19.523 30.738 30.030 34.321 22.541
Model Buy & Hold Palmyra-Fin-70B GPT-01-preview GPT-4 GPT-40 Qwen2.5-72B-Instruct Llama-3.1-70B-Instruct DeepSeek-67B-Chat Yi-1.5-34B-Chat Qwen2.5-32B-Instruct DeepSeek-V2-Lite (15.7B)	CR↑ 39.244 -6.661 34.499 45.246 45.946 39.112 37.545 35.647 25.364 21.336 31.458	TS SR↑ 0.869 -0.222 0.796 1.190 1.348 1.075 0.891 0.885 0.808 0.729 0.724	ILA           AV↓           75.854           50.379           72.822           63.896           57.281           61.136           67.660           73.561           49.157           68.524	MDD↓         37.975         37.975         35.490         25.031         21.631         31.631         31.359         33.359         35.490         20.704         35.404 </td <td>CR↑           10.837           8.562           8.238           9.889           7.405           11.935           12.772           14.213           14.220           27.016</td> <td>AA SR↑ 0.470 Financia 0.372 Propu 0.422 0.440 0.457 Open- 0.572 0.583 0.666 0.623 1.089 1.221</td> <td>PL AV↓ 38.984 <i>I Domain</i> 38.891 <i>ietary Me</i> 33.045 38.008 27.434 <i>Source M</i> 35.293 37.057 36.118 38.596 20.522 11.860</td> <td>MDD↓         I           19.119         I           Models         25.466           25.466         I           models         14.412           19.119         12.824           Models         I           fodels         I           10.876         I           19.432         8.943           37.435         I</td> <td>CR↑           52.216           -3.261           32.433           75.952           63.743           87.412           66.522           30.963           64.307           28.096           27.762</td> <td>SR↑           0.858           -0.057           0.558           1.286           1.318           2.181           1.118           0.599           1.063           0.652           0.474</td> <td>IO           AV↓           107.197           101.711           102.470           104.059           85.210           70.628           104.848           91.146           106.597           72.344           103.193</td> <td>MDD↓ 47.766 58.406 54.016 37.867 29.220 12.464 46.379 45.855 48.042 37.975 48.478</td> <td>CR↑           34.099           -0.453           25.057           43.696           39.031           46.153           38.946           26.941           37.966           20.884           28.745</td> <td>Ave SR↑ 0.732 0.031 0.592 0.972 1.041 1.276 0.864 0.717 0.831 0.823 0.813</td> <td>rage AV↓ 74.012 63.660 69.446 68.654 55.686 70.907 64.975 72.918 47.341 61.192</td> <td>MDD↓ 34.953 36.564 34.639 27.339 21.225 19.523 30.738 30.030 34.321 22.541 40.439</td>	CR↑           10.837           8.562           8.238           9.889           7.405           11.935           12.772           14.213           14.220           27.016	AA SR↑ 0.470 Financia 0.372 Propu 0.422 0.440 0.457 Open- 0.572 0.583 0.666 0.623 1.089 1.221	PL AV↓ 38.984 <i>I Domain</i> 38.891 <i>ietary Me</i> 33.045 38.008 27.434 <i>Source M</i> 35.293 37.057 36.118 38.596 20.522 11.860	MDD↓         I           19.119         I           Models         25.466           25.466         I           models         14.412           19.119         12.824           Models         I           fodels         I           10.876         I           19.432         8.943           37.435         I	CR↑           52.216           -3.261           32.433           75.952           63.743           87.412           66.522           30.963           64.307           28.096           27.762	SR↑           0.858           -0.057           0.558           1.286           1.318           2.181           1.118           0.599           1.063           0.652           0.474	IO           AV↓           107.197           101.711           102.470           104.059           85.210           70.628           104.848           91.146           106.597           72.344           103.193	MDD↓ 47.766 58.406 54.016 37.867 29.220 12.464 46.379 45.855 48.042 37.975 48.478	CR↑           34.099           -0.453           25.057           43.696           39.031           46.153           38.946           26.941           37.966           20.884           28.745	Ave SR↑ 0.732 0.031 0.592 0.972 1.041 1.276 0.864 0.717 0.831 0.823 0.813	rage AV↓ 74.012 63.660 69.446 68.654 55.686 70.907 64.975 72.918 47.341 61.192	MDD↓ 34.953 36.564 34.639 27.339 21.225 19.523 30.738 30.030 34.321 22.541 40.439
Model Buy & Hold Palmyra-Fin-70B GPT-01-preview GPT-4 GPT-40 Qwen2.5-72B-Instruct Llama-3.1-70B-Instruct DeepSeet-67B-Chat Yi-1.5-34B-Chat Qwen2.5-32B-Instruct DeepSeet-V2-Lite (15.7B) Yi-1.5-9B-Chat	CR↑           39.244           -6.661           34.499           45.246           45.946           39.112           37.545           35.647           21.336           31.458           31.350	TS SR↑ 0.869 -0.222 0.796 1.190 1.348 1.075 0.891 0.885 0.808 0.729 0.744 0.703	ILA           AV↓           75.854           50.379           72.822           63.896           57.281           61.136           67.660           73.561           49.157           68.524           74.895	MDD↓         Image: Non-State	CR↑           10.837           8.562           8.238           9.889           7.405           11.935           12.772           14.213           14.227           13.220           27.016           3.640	AA SR↑ 0.470 Financia 0.372 0.422 0.440 0.457 Open- 0.572 0.583 0.666 0.623 1.089 1.221 0.162	PL AV↓ 38.984 <i>l Domain</i> 38.891 <i>ietary Mo</i> 33.045 38.008 27.434 <i>Source M</i> 35.293 37.057 36.118 38.596 20.522 11.860 37.947	MDD↓         I           19.119         I           19.119         I           25.466         I           odels         I           14.412         I           19.119         I           12.824         I           odels         I           19.119         I           16.021         I           10.876         I           19.432         8.943           37.435         I           17.578         I	CR↑           52.216           -3.261           32.433           75.952           63.743           87.412           66.522           30.963           64.307           28.096           27.762           33.748	SR↑           0.858           -0.057           0.558           1.286           1.318           2.181           1.118           0.652           0.474           0.569	IO           AV↓           107.197           101.711           102.470           104.059           85.210           70.628           104.848           91.146           106.597           72.344           103.193           103.193           104.502	MDD↓           47.766           58.406           54.016           37.867           29.220           12.464           46.379           45.855           48.042           37.975           48.478           55.284	CR↑           34.099           -0.453           25.057           43.696           39.031           46.153           38.946           26.941           37.966           20.884           28.745           22.913	Ave SR↑ 0.732 0.031 0.592 0.972 1.041 1.276 0.864 0.717 0.831 0.823 0.813 0.478	rage AV↓ 74.012 63.660 69.446 68.654 55.686 70.907 64.975 72.918 47.341 61.192 72.448	MDD↓ 34.953 36.564 34.639 27.339 21.225 19.523 30.738 30.030 34.321 22.541 40.439 36.946
Model Buy & Hold Palmyra-Fin-70B GPT-01-preview GPT-4 GPT-40 Qwen2.5-72B-Instruct Llama-3.1-70B-Instruct DeepSeek-67B-Chat Yi-1.5-34B-Chat Qwen2.5-32B-Instruct DeepSeek-V2-Lite (15.7B) Yi-1.5-9B-Chat Llama-3.1-8B-Instruct	CR↑ 39.244 -6.661 34.499 45.246 45.946 39.112 37.545 35.647 35.364 21.336 31.458 31.350 35.622	TS SR↑ 0.869 -0.222 0.796 1.190 1.348 1.075 0.891 0.885 0.808 0.729 0.744 0.703 0.832	ILA           AV↓           75.854           50.379           72.822           63.896           57.281           61.136           67.660           73.561           49.157           68.524           74.895           71.936	MDD↓         37.975           35.490         25.031           25.631         21.631           26.985         29.813           33.359         35.490           20.704         35.494           37.975         36.383	CR↑           10.837           8.562           8.238           9.889           7.405           11.935           12.772           14.213           14.227           13.220           27.016           3.640           7.079	AA SR↑ 0.470 Financia 0.372 Propy 0.422 0.440 0.457 Open- 0.572 0.583 0.666 0.623 1.089 1.221 0.162 0.309	PL AV↓ 38.984 <i>I Domain</i> 38.891 <i>ietary Ma</i> 33.045 38.008 27.434 <i>Source M</i> 35.293 37.057 36.118 38.596 20.522 11.860 37.947 38.757	MDD↓         I           19.119         I           Models         25.466           25.466         I           models         14.412           19.119         12.824           Iodels         I           19.119         12.824           Iodels         19.119           16.021         10.876           19.432         8.943           37.435         17.578           18.747         18.747	CR↑           52.216           -3.261           32.433           75.952           63.743           87.412           66.522           30.963           64.307           28.096           27.762           33.748           33.689	SR↑           0.858           -0.057           0.558           1.286           1.318           2.181           1.118           0.599           1.063           0.652           0.474           0.569           0.560	IO           AV↓           107.197           101.711           102.470           104.059           85.210           70.628           91.146           106.597           72.344           104.029           104.020           104.094	MDD↓           47.766           58.406           54.016           37.867           29.220           12.464           46.379           45.855           48.042           37.975           48.478           55.284           56.527	CR↑           34.099           -0.453           25.057           43.696           39.031           46.153           38.946           26.941           37.966           20.884           28.745           22.913           25.463	Ave SR↑ 0.732 0.031 0.592 0.972 1.041 1.276 0.864 0.717 0.823 0.813 0.478 0.567	rage           AV↓           74.012           63.660           69.446           68.654           55.686           70.907           64.975           72.918           47.341           61.192           72.448           72.262	MDD↓ 34.953 36.564 34.639 27.339 21.225 19.523 30.738 30.030 34.321 22.541 40.439 36.946 37.219

Table 2: Performance of stock trading with different LLMs as backbone model across seven stocks.

<sup>1</sup> The Buy & Hold strategy is a passive investment approach commonly used as a baseline strategy, where an investor purchases stocks and holds onto them for an extended period regardless of market fluctuations.

<sup>2</sup> An upward arrow ( $\uparrow$ ) next to a metric indicates that higher values signify better performance, while a downward arrow ( $\downarrow$ ) indicates that lower values are preferable.

<sup>3</sup> The numbers highlighted in red indicate the best-performing outcomes for the corresponding metrics.

sented on a task-by-task basis. We report the performance of INVESTORBENCH on three singleasset trading tasks: *stocks*, *cryptocurrencies*, and *ETFs* trading, using closed-source, open-source, and domain-specific LLMs.

# 4.1 Experiment Setup

391

392

400

401

402

403

404

405

406

407

408

409

410

411

412

Table 1 summarizes the performance of a comprehensive list of trading agents. For single equity tasks, the baseline is set up by Buy and Hold strategy, while for portfolio management task, it is set up by an equal-weight portfolio with the detailed rational explained in Appendix. In our experiments, the temperature parameter of all LLM-based agent systems is set at 0.6 to balance response consistency and reasoning creativity. The performance metrics are reported for the test trajectory with the median CR, SR, AV, and MDD from five repeated epochs. (If the median of these metrics does not belong to the same epoch, the performance is based on the trajectory with the median SR.)

Furthermore, the selection of warm-up and test periods differs across various tasks due to the varying time spans of data collected to construct the agent environment. For the single-asset trading tasks, the warm-up period of stock trading is from 2020-07-01 to 2020-09-30 and the test period is from 2020-10-01 to 2021-05-06. The warm-up period of cryptocurrency trading is from 2023-02-11 to 2023-04-04 and the test period is from 2023-04-05 to 2023-11-05. The warm-up period of ETF trading is from 2019-07-29 to 2019-12-30 and the test period is from 2020-01-02 to 2020-09-21.

413

414

415

416

417

418

419

420

421

422

423

424

425

426

427

428

429

430

431

432

433

434

435

For LLM deployment, we utilize vllm to deploy LLMs. For small-scale LLMs (under 10B parameters), we deploy models on two RTX A6000 GPUs, each with 48GB DRAM. For mid-scale LLMs (10B to 65B parameters), we use four RTX A6000 GPUs. For large-scale LLMs (over 65B parameters), models are deployed on eight A100 GPUs, each equipped with 80GB DRAM.

# 4.2 Result 1: Stock Trading

Table 2 presents the performance of thirteen backbone models across seven stocks, accompanied by the average of each metric for all stocks to offer a more comprehensive view of their overall performance. We outline three key insights as follows:

Superior stock trading performance is achieved 436 with proprietary LLMs as agent backbones 437 Compared to agents employing open-source or 438 financial-domain-specific fine-tuned LLMs, those 439 using the three proprietary LLMs demonstrated sig-440 nificantly higher and more consistent average CR 441 and SR, as shown in Figure 3a. Despite being fine-442 tuned with extensive financial contexts, domain-443 specific LLMs did not provide a decisive advantage 444 in sequential stock trading decision-making tasks. 445 This may be attributed to their primary training for 446 other functions, such as long financial report anal-447 ysis exemplified by Palmyra-Fin-70B, rather than 448 decision-making. 449

Model parameter size increment enhances agent 450 financial decision-making quality and robust-451 ness. In the category of open-source LLMs, those 452 exceeding 67B parameters displayed superior CRs 453 and SRs, along with markedly less variance within 454 their category, as illustrated in Figure 3b and Ta-455 ble 2. This underscores the prevailing belief that the 456 reasoning capabilities of LLMs are proportionate to 457 their parameter size, which holds also true in stock 458 trading, which is a sequential decision-making task 459 460 in an open-ended, volatile environment by nature. Proprietary models exhibit significantly 461 stronger decision-making capabilities compared 462 to even the largest open-source LLMs under 463 complex, mixed market conditions, though this 464 advantage is less evident in relatively monotone 465 market environments. During the test phase, 466 primarily influenced by the range of open-source 467 data collected, TSLA and NIO exhibited volatility 468 with mixed upward and downward stock price 469 trends, whereas the other five stocks generally 470 showed bullish trends. The investment signals 471 derived from such complex markets tend to be 472 noisy or delayed, as illustrated in Appendix C. We 473 observed that proprietary models possess a superior 474 ability to manage these challenging conditions and 475 consistently deliver better performance outcomes 476 than large-sized open-source LLMs. Their 477 reasoning capability enables them to effectively 478 utilize other decision-relevant information, such as 479 historical momentum, current holdings, and, most 480 critically, self-reflection outcomes from the agents, 481 thereby facilitating more accurate decisions. 482

# 4.3 Result 2 & 3: Cryptocurrency Trading and ETF Trading

483

484

485

486

In the test phases of both cryptocurrency and ETF trading tasks, market trends are mixed. Notably,

Table 3: Performance of cryptocurrency trading with different LLMs as backbone models across Bitcoin (BTC) and Ethereum (ETH).

Model		B	ГC		ETH				
	CR↑	SR↑	AV↓	MDD↓	CR↑	SR↑	AV↓	MDD↓	
Buy & Hold	21.821	0.683	37.426	20.796	4.528	0.146	41.817	29.889	
		Financia	al Domair	1 Models					
Palmyra-Fin-70B	-20.812	-1.212	20.036	27.782	4.795	0.240	26.924	16.405	
		Prop	rietary M	odels					
GPT-01-preview	34.060	1.114	35.846	17.075	2.496	0.085	39.641	27.692	
GPT-4	22.396	0.828	31.699	17.206	1.516	0.051	39.812	32.541	
GPT-40	14.330	0.532	31.304	17.278	4.666	0.190	33.051	22.539	
Average	23.595	0.825	32.950	17.186	2.893	0.109	37.501	27.591	
		Open-	Source M	lodels					
Qwen2.5-72B-Instruct	0.549	0.325	1.979	0.897	11.984	0.584	18.554	27.642	
Llama-3.1-70B-Instruct	20.440	0.758	31.604	17.813	-11.888	-0.410	39.047	36.416	
DeepSeek-67B-Chat	28.307	0.891	37.219	17.944	9.480	0.309	41.369	26.261	
Yi-1.5-34B-Chat	13.620	0.434	36.778	22.790	6.325	0.227	37.503	25.707	
Qwen2.5-32B-Instruc	11.566	0.869	15.608	7.984	2.823	0.194	19.571	7.883	
DeepSeek-V2-Lite (15.7B)	4.804	0.153	36.846	20.562	-9.504	-0.311	41.199	21.270	
Yi-1.5-9B-Chat	7.953	0.253	36.799	26.545	-3.684	-0.119	41.818	35.417	
Llama-3.1-8B-Instruct	20.521	0.646	37.240	21.104	4.939	0.163	40.928	29.466	
Qwen-2.5-Instruct-7B	19.477	0.612	37.289	20.796	-1.339	-0.075	24.124	-16.053	
Average	14.137	0.549	30.151	17.382	1.015	0.062	33.790	21.557	

the cryptocurrency task shows significantly smaller price fluctuations compared to the ETF task. We outline the key features of using an LLM-agent to make financial decisions across these two distinct markets as follows: 487

488

489

490

491

492

493

494

495

496

497

498

499

500

501

502

503

504

505

506

507

508

509

510

511

512

513

514

515

516

517

518

519

520

Large-sized open-source models and proprietary models are needed to effectively capture trading signals of cryptocurrency markets, which are highly sensitive to news and financial sentiment. As shown in Table 3, using midsized and small-sized open-source models as the decision-making agent backbone generally results in weaker performance than the market baseline with respect to CR and SR.

ETF investment requires proprietary models enriched with extensive pre-trained knowledge to serve as the agents brain and provide robust reasoning support. As shown in Table 4, proprietary models significantly outperform open-source and financial domain-specific models in this task. This advantage arises from the complexity of ETF trading, which necessitates interpreting actionable signals across diverse sectors, demanding more strategic, long-term decisions grounded in deep comprehension and reflection anchored by rich precontexts.

# 4.4 Discussion

Combining all the experimental results, we find that the performance of different LLM varies significantly in stock, cryptocurrency, and ETF trading. This variation not only reflects the inherent complexity of financial markets, but also highlights the importance of model selection or fine-tuning. For instance, proprietary LLM generally exhibit be

521

522

- 530 531 532
- 533
- 534

535 536

538

539 540

541

543

545

547

551

guage models (LMs) has stimulated the exploration of financial LMs, such as pre-trained LMs: FINBERT (Liu et al., 2021; Yang et al., 2020; Araci, 2019; Huang et al., 2023), FINBERT-MRC (Zhang and Zhang, 2023), FLANG (Shah et al., 2022), and several financial LLMs: FINGPT(Liu et al., 2023), FINMA (Xie et al., 2023), IN-VESTLM (Yang et al., 2023), BloombergGPT (Wu et al., 2023), which leverage extensive training on

making tasks.

5

5.1

**Related Work** 

Table 4: Performance of ETF trading with different LLMs as backbone models.

cial sector, such as FINMEM (Yu et al., 2024a), FINAGENT (Zhang et al., 2024a) and FINROBOT (Yang et al., 2024), characterized by their adaptability and openness. However, variations in framework design, task scope, and data types present challenges in uniformly evaluating the efficacy of LLM agents in financial scenarios.

552

553

554

555

556

557

558

559

561

563

564

565

566

567

568

569

570

571

572

573

574

575

576

577

578

579

580

581

582

583

584

585

586

587

588

590

591

592

593

594

595

596

597

598

600

#### 5.2 **Financial LLM Benchmarks**

In the realm of financial LLMs, several benchmarks have been developed: FLUE (Shah et al., 2022) introduces the first comprehensive benchmark with five financial NLP tasks, including sentiment analysis, headline classification, named entity recognition, structure boundary detection, and question answering. Pixiu (Xie et al., 2023) expands this benchmark to include financial document understanding and classification tasks, incorporating multimodal datasets. FinBen (Xie et al., 2024) encompasses 36 datasets covering 24 financial tasks. Despite these advancements, there remains a notable gap in benchmarks specifically designed for LLM-based agent applications within the financial sector.

#### Conclusion 6

INVESTORBENCH offers the community two distinct modes of engagement. The first mode allows participants to integrate their fine-tuned LLMs into the INVESTORBENCH's agent framework to undertake financial decision-making tasks. This setup enables them to benchmark the performance of their models against those previously experimented with by our work. The second mode permits users to directly incorporate the environment and evaluation metrics of INVESTORBENCH into their own designed agents, facilitating a comparative analysis of their agent design's effectiveness. This dual approach provides a flexible framework for testing and enhancing financial decision-making strategies within the INVESTORBENCH ecosystem.

Future research efforts will expand the benchmark by incorporating additional information modalities, such as audio (e.g., earnings call recordings) and graphs (e.g., K-lines, trade charts), to explore whether these data types can enhance decision-making quality. The foundational agent framework of INVESTORBENCH is designed to seamlessly accommodate these modalities, ensuring that the extended benchmark remains easy to use and scalable.

8

ETF	CR↑	SR↑	AV↓	MDD↓						
Buy & Hold	2.069	0.06	46.645	35.746						
Financial Domain Models										
Palmyra-Fin-70B	24.759	1.152	30.419	8.203						
Proprietary Models										
GPT-01-preview	21.224	0.849	43.766	20.054						
GPT-4	2.807	0.110	44.679	37.785						
GPT-40	12.292	0.377	46.150	32.678						
Average	12.108	0.445	44.865	30.172						
Open-Sour	ce Model	ls								
Qwen2.5-72B-Instruct	4.507	0.227	28.090	8.580						
Llama-3.1-70B-Instruct	9.895	0.464	30.184	12.759						
Yi-1.5-34B-Chat	4.996	0.322	21.986	12.858						
Qwen2.5-32B-Instruct	19.617	0.955	29.070	7.496						
DeepSeek-V2-Lite (15.7B)	1.389	0.063	31.371	31.831						
Yi-1.5-9B-Chat	-4.657	-0.228	28.907	15.545						
Llama-3.1-8B-Instruct	11.239	0.475	33.480	15.587						
Qwen-2.5-Instruct-7B	-0.384	-0.020	27.596	14.059						
Average	5.825	0.282	28.835	14.839						

performance in stock trading due to their strong

training on various financial datasets, while open-

source models struggle to achieve these results,

especially in more volatile environments such as

cryptocurrency trading. In addition, the effective-

ness of LLM-based agents depends heavily on their ability to adapt to market fluctuations. Agents

that incorporate advanced memory systems and dy-

namic risk assessment capabilities are better able to

cope with complex market situations, highlighting

the value of the complex architectural features of

LLM-based agent framework in financial decision-

The rapid development of general-domain lan-

diverse financial datasets (e.g. stock price data,

financial news and analyst reports) and adapt the

capabilities of LMs to the unique needs of financial

applications. Concurrently, the advancement of

LLMs has significantly enhanced the development

of language-based agent frameworks in the finan-

LLM for Financial Domain

# 601 Limitation

First, INVESTORBENCH is currently focusing on
single-asset financial decision-making task, without addressing multi-asset tasks such as portfolio
management. Second, copyright restrictions on financial domain data may compromise the quality
of the datasets we create, potentially limiting the
assessment of model performance.

# 9 Ethical Statement

The authors take full responsibility for the devel-610 opment of INVESTORBENCH, ensuring that the 611 publicly available part in dataset does not contain personal information, and conform to established 613 ethical guidelines. The data are shared under the 614 MIT license, requiring users to adhere to its terms. 615 INVESTORBENCH is intended for academic and educational purposes only and is not a substitute for 617 professional advice. While efforts have been made 618 to ensure its accuracy, the authors and their institutions disclaim liability for any outcomes arising from its use. Users agree to take responsibility for 621 ethical and lawful use and to indemnify the authors and their affiliates against any claims or damages 623 resulting from reliance on this Material.

# References

625

631

633

634

635

636

637

638

641

643

647

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.
- Andrew Ang and Joseph Chen. 2003. Downside risk. Journal of Portfolio Management, 29(4):103–112.
- Dogu Araci. 2019. Finbert: Financial sentiment analysis with pre-trained language models. *arXiv preprint arXiv:1908.10063*.
- Dimitri Bertsekas and Steven E Shreve. 1996. *Stochastic optimal control: the discrete-time case*, volume 5. Athena Scientific.
- Mukul Bhatnagar, Sanjay Taneja, and Ramona Rupeika-Apoga. 2023. Demystifying the effect of the news (shocks) on crypto market volatility. *Journal of Risk and Financial Management*, 16(2):136.
- John H. Cochrane. 1988. Volatility tests and efficient markets: A review essay. *Journal of Monetary Economics*, 22(3):463–485.
- Eva Eigner and Thorsten Händler. 2024. Determinants of llm-assisted decision-making. *arXiv preprint arXiv:2402.17385*.

Thomas L Griffiths, Jian-Qiao Zhu, Erin Grant, and R Thomas McCoy. 2023. Bayes in the age of intelligent machines. *arXiv preprint arXiv:2311.10206*. 649

650

651

652

653

654

655

656

657

658

659

660

661

662

663

664

665

666

667

668

669

670

671

672

673

674

675

676

677

678

679

680

681

682

683

684

685

686

687

688

689

690

691

692

693

694

695

696

697

698

699

700

701

- Taicheng Guo, Xiuying Chen, Yaqi Wang, Ruidi Chang, Shichao Pei, Nitesh V Chawla, Olaf Wiest, and Xiangliang Zhang. 2024. Large language model based multi-agents: A survey of progress and challenges. *arXiv preprint arXiv:2402.01680*.
- Allen H Huang, Hui Wang, and Yi Yang. 2023. Finbert: A large language model for extracting information from financial text. *Contemporary Accounting Research*, 40(2):806–841.
- John Hull. 2007. Risk Management and Financial Institutions. John Wiley & Sons.
- Taylan Kabbani and Ekrem Duman. 2022. Deep reinforcement learning approach for trading automation in the stock market. *IEEE Access*, 10:93564–93574.
- Yuan Li, Bingqiao Luo, Qian Wang, Nuo Chen, Xu Liu, and Bingsheng He. 2024. Cryptotrade: A reflective llm-based agent to guide zero-shot cryptocurrency trading. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 1094–1106.
- Xiao-Yang Liu, Guoxuan Wang, and Daochen Zha. 2023. Fingpt: Democratizing internet-scale data for financial large language models. *arXiv preprint arXiv:2307.10485*.
- Yang Liu, Qi Liu, Hongke Zhao, Zhen Pan, and Chuanren Liu. 2020. Adaptive quantitative trading: An imitative deep reinforcement learning approach. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pages 2128–2135.
- Zhuang Liu, Degen Huang, Kaiyu Huang, Zhuang Li, and Jun Zhao. 2021. Finbert: A pre-trained financial language representation model for financial text mining. In *Proceedings of the twenty-ninth international conference on international joint conferences on artificial intelligence*, pages 4513–4519.
- Meta Llama team. 2024. The llama 3 herd of models.
- Ananth N Madhavan. 2016. *Exchange-traded funds and the new dynamics of investing*. Oxford University Press.
- Jaap MJ Murre and Joeri Dros. 2015. Replication and analysis of ebbinghaus forgetting curve. *PloS one*, 10(7):e0120644.

OpenAI. 2022. Introducing chatgpt.

Joon Sung Park, Joseph O'Brien, Carrie Jun Cai, Meredith Ringel Morris, Percy Liang, and Michael S. Bernstein. 2023. Generative agents: Interactive simulacra of human behavior. In *Proceedings of the 36th Annual ACM Symposium on User Interface Software and Technology*, UIST '23, New York, NY, USA. Association for Computing Machinery.

704

- 705 706 710 711
- 712 713 714 715 716

717

- 718 719 720 721 722
- 727 730 732
- 733 736
- 739 740 741
- 742 743
- 744
- 745 746
- 747 748 749

751

752

755

- Alibaba Group Qwen team. 2024. Qwen2.5.
- Raeid Sagur, Ken Kato, Nicholas Vinden, and Frank Rudzicz. 2024. Nifty financial news headlines dataset. arXiv preprint arXiv:2405.09747.
- Raj Sanjay Shah, Kunal Chawla, Dheeraj Eidnani, Agam Shah, Wendi Du, Sudheer Chava, Natraj Raman, Charese Smiley, Jiaao Chen, and Divi Yang. 2022. When flue meets flang: Benchmarks and large pre-trained language model for financial domain. arXiv preprint arXiv:2211.00083.
- William F. Sharpe. 1994. The sharpe ratio. The Journal of Portfolio Management, 21(1):49-58.
- Writer Engineering team. 2024. Palmyra-Fin-70B-32k: a powerful LLM designed for Finance. https:// dev.writer.com.
- Olivier Vanhoucke. 2023. Crypto news dataset. Accessed: 2024-08-20.
- Lei Wang, Chen Ma, Xueyang Feng, Zeyu Zhang, Hao Yang, Jingsen Zhang, Zhiyuan Chen, Jiakai Tang, Xu Chen, Yankai Lin, et al. 2024. A survey on large language model based autonomous agents. Frontiers of Computer Science, 18(6):186345.
- Shijie Wu, Ozan Irsoy, Steven Lu, Vadim Dabravolski, Mark Dredze, Sebastian Gehrmann, Prabhanjan Kambadur, David Rosenberg, and Gideon Mann. 2023. Bloomberggpt: A large language model for finance. arXiv preprint arXiv:2303.17564.
- Qianqian Xie, Weiguang Han, Zhengyu Chen, Ruoyu Xiang, Xiao Zhang, Yueru He, Mengxi Xiao, Dong Li, Yongfu Dai, Duanyu Feng, et al. 2024. The finben: An holistic financial benchmark for large language models. arXiv preprint arXiv:2402.12659.
- Qianqian Xie, Weiguang Han, Xiao Zhang, Yanzhao Lai, Min Peng, Alejandro Lopez-Lira, and Jimin Huang. 2023. Pixiu: A large language model, instruction data and evaluation benchmark for finance. arXiv preprint arXiv:2306.05443.
- Huajian Xin, Z. Z. Ren, Junxiao Song, Zhihong Shao, Wanjia Zhao, Haocheng Wang, Bo Liu, Liyue Zhang, Xuan Lu, Qiushi Du, Wenjun Gao, Qihao Zhu, Dejian Yang, Zhibin Gou, Z. F. Wu, Fuli Luo, and Chong Ruan. 2024. Deepseek-prover-v1.5: Harnessing proof assistant feedback for reinforcement learning and monte-carlo tree search.
- Hongyang Yang, Boyu Zhang, Neng Wang, Cheng Guo, Xiaoli Zhang, Likun Lin, Junlin Wang, Tianyu Zhou, Mao Guan, Runjia Zhang, et al. 2024. Finrobot: An open-source ai agent platform for financial applications using large language models. arXiv preprint arXiv:2405.14767.
- Yi Yang, Yixuan Tang, and Kar Yan Tam. 2023. Investlm: A large language model for investment using financial domain instruction tuning. arXiv preprint arXiv:2309.13064.

Yi Yang, Mark Christopher Siy Uy, and Allen Huang. 2020. Finbert: A pretrained language model for financial communications. arXiv preprint arXiv:2006.08097.

756

757

758

759

760

761

762

763

765

768

769

770

771

772

773

774

778

779

780

781

782

783

784

785

786

787

788

789

790

791

792

793

794

795

796

797

798

799

800

801

802

803

804

805

- Ziruo Yi, Ting Xiao, Kaz-Onyeakazi Ijeoma, Ratnam Cheran, Yuvraj Baweja, and Phillip Nelson. Stock2vec: An embedding to improve 2022. predictive models for companies. arXiv preprint arXiv:2201.11290.
- Alex Young, Bei Chen, Chao Li, Chengen Huang, Ge Zhang, Guanwei Zhang, Heng Li, Jiangcheng Zhu, Jianqun Chen, Jing Chang, et al. 2024. Yi: Open foundation models by 01. ai. arXiv preprint arXiv:2403.04652.
- Yangyang Yu, Haohang Li, Zhi Chen, Yuechen Jiang, Yang Li, Denghui Zhang, Rong Liu, Jordan W Suchow, and Khaldoun Khashanah. 2024a. Finmem: A performance-enhanced llm trading agent with layered memory and character design. In Proceedings of the AAAI Symposium Series, volume 3, pages 595–597.
- Yangyang Yu, Zhiyuan Yao, Haohang Li, Zhiyang Deng, Yupeng Cao, Zhi Chen, Jordan W Suchow, Rong Liu, Zhenyu Cui, Zhaozhuo Xu, et al. 2024b. Fincon: A synthesized llm multi-agent system with conceptual verbal reinforcement for enhanced financial decision making. arXiv preprint arXiv:2407.06567.
- Wentao Zhang, Lingxuan Zhao, Haochong Xia, Shuo Sun, Jiaze Sun, Molei Qin, Xinyi Li, Yuqing Zhao, Yilei Zhao, Xinyu Cai, et al. 2024a. Finagent: A multimodal foundation agent for financial trading: Tool-augmented, diversified, and generalist. arXiv preprint arXiv:2402.18485.
- Yadong Zhang, Shaoguang Mao, Tao Ge, Xun Wang, Adrian de Wynter, Yan Xia, Wenshan Wu, Ting Song, Man Lan, and Furu Wei. 2024b. Llm as a mastermind: A survey of strategic reasoning with large language models. arXiv preprint arXiv:2404.01230.
- Yuzhe Zhang and Hong Zhang. 2023. Finbert-mrc: financial named entity recognition using bert under the machine reading comprehension paradigm. Neural Processing Letters, 55(6):7393–7413.
- Zhihan Zhou, Liqian Ma, and Han Liu. 2021. Trade the event: Corporate events detection for news-based event-driven trading. In Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021, pages 2114-2124, Online. Association for Computational Linguistics.

# Appendices

#### Memory Ranking Mechanism of Α **FINMEM**

Upon receiving an investment inquiry, FINMEM re-806 trieves the top-K critical memory events from each 807 layer and channels them to the immediate reflection 808

856

860

861

862

863

864

865

866

867

868

869

870

871

872

873

874

875

876

877

878

879

880

881

882

883

884

885

886

887

888

889

890

891

892

893

 layers compared to daily financial news, reflecting their extended timeliness, relevance, and impact on financial decision-making.

ports (e.g., Form 10-Ks) are assigned higher stabil-

ity values and categorized within deeper processing

$$S_{\text{Relevancy}_l}^E = \frac{\mathbf{m}_{\mathbf{E}} \cdot \mathbf{m}_{\mathbf{P}}}{\|\mathbf{m}_{\mathbf{E}}\|_2 \times \|\mathbf{m}_{\mathbf{P}}\|_2}$$
(4)

The relevancy score  $S_{\text{relevancy}l}^E$  quantifies the semantic similarity between a memory event E and the current query P using cosine similarity of their respective embedding vectors, **mE** and **m**<sub>P</sub>, as shown in Equation 4. These embeddings are generated from the event's textual content and the LLM prompt query (which includes trading inquiries and the agent's character setting) using OpenAI's "textembedding-ada-003" model.

The importance score  $S^E_{\text{Importancel}}$  for a memory event E in layer l is calculated as the product of a value  $v_l^E$  (derived from a uniform piecewise scoring function, Equation 5) and a degrading ratio  $\theta_l$ (Equation 6), as shown in Equation 7. This approach, adapted from (Park et al., 2023), is tailored to our stratified long-term memory structure. The likelihood of higher  $v_l^E$  values increases from shallow to deep layers, while  $\theta_l$  measures the diminishing importance of an event over time using layerspecific exponential functions. The base  $\alpha_l$  for each layer follows  $\alpha shallow < \alpha_{intermediate} < \alpha_{deep}$ (set to 0.9, 0.967, and 0.988 respectively), ensuring  $\theta_l$  decreases to a threshold of 5 after 30, 90, and 365 days for shallow, intermediate, and deep layers. This layered approach, implemented through three-piece-wise functions for both  $S^E_{\text{Importancel}}$  and  $SRecencyl^E$ , enables FinMem to process longterm memory in a stratified manner. Memory events are purged when  $SRecencyl^E$  falls below 0.05 or SImportance<sub>1</sub><sup>E</sup> is under 5 (pre-scaling), maintaining the relevance and efficiency of the memory store.

$$v_l^E = \begin{cases} 40 & \text{with probability } p_1 \\ 60 & \text{with probability } p_2 \\ 80 & \text{with probability } p_3 \end{cases}$$
(5)

$$\theta_l = (\alpha_l)^{\delta^E}, \quad l = \text{shallow, intermediate, deep,}$$
(6)

where  $p_1 + p_2 + p_3 = 1$ , but their values vary by shallow, intermediate, and 897

component of the working memory. These events are selected based on their information retrieval score,  $\gamma_l^E$ , where *l* represents the layer (shallow, intermediate, or deep), as defined in Equation 2.

809

810

811

812

813

814

815

816

817

818

821

823

833

834

837

838

839

841

847

851

855

$$\gamma_l^E = S_{\text{Recency}_l}^E + S_{\text{Relevancy}_l}^E + S_{\text{Importance}_l}^E, \quad (2)$$

where each memory event is only associated with one score and can only belong to a single layer.

Let E denote a given memory event. The scoring mechanism for E, adapted from Park et al. (Park et al., 2023) but with modified recency and importance computations, is tailored to handle data with various timelines and to achieve layered processing that represents the diverse periodicities of the financial environment. This score encapsulates three metrics: recency (how recently the event occurred), relevancy (the event's pertinence to the current context), and importance (the event's significance). Individual metric scores exceeding 1.0 are scaled to the [0,1] range before being summed, ensuring a balanced contribution from each component and preventing any single metric from dominating the overall score. The resulting composite score provides a comprehensive evaluation of the memory event's significance within the multi-layered, periodically varying financial landscape.

$$S_{\text{Recency}_l}^E = e^{-\frac{\delta^E}{Q_l}}, \quad \delta^E = t_{\text{P}} - t_E, \quad (3)$$

where  $\delta^E$  represents the time elapsed between a memory event's occurrence and the trading inquiry's arrival. The model utilizes three processing layers, each corresponding to a specific timeframe: shallow ( $Q_{\text{shallow}} = 14$ days), intermediate ( $Q_{\text{intermediate}} = 90$  days), and deep ( $Q_{\text{deep}} = 365$  days). These intervals represent two weeks, a quarter, and a year respectively.

When a trade inquiry P arrives in processing layer l via an LLM prompt, the agent calculates the recency score  $S_{\text{Recency}_l}^E$  for a memory event E using Equation 3. This score inversely correlates with the time elapsed between the inquiry and the event's memory timestamp, mapping to Ebbinghaus's forgetting curve (Murre and Dros, 2015). The stability term  $Q_l$  in Equation 3 modulates memory decay rates across layers, with higher values in deeper layers indicating longer memory persistence. For instance, in the trading context, company annual re-

11

894

895

901

902

903

904

905

907

908

909

910

911

912

913

914

915

916

917

919

920

921

923

925

926

927

931

932

933

937

938

940

942

deep processing. when shallow processing  $p_1, p_2, p_3 = \{0.8, 0.15, 0.05\}$ , intermediate processing,  $p_1, p_2, p_3 = \{0.05, 0.8, 0.15\}$  and deep processing,  $p_1, p_2, p_3 = \{0.05, 0.15, 0.8\}$ .

$$S_{\text{Importance}_{l}}^{E} = v_{l}^{E} * \theta_{l}, \tag{7}$$

Furthermore, FINMEM employs an access counter function to dynamically manage memory events across layers, ensuring that crucial events influencing trading decisions are elevated to deeper layers for extended retention and recurring access. This process, monitored by the LLM validation tool Guardrails AI, tracks critical memory IDs across layers. Events deemed pivotal for investment success receive a 5-point boost to their importance score  $(S_{\text{Importancel}}^E)$ . Upon meeting upgrade criteria for a deeper layer, an event's recency score  $(SRecency_l^E)$  is reset to 1.0, underscoring its significance and preventing rapid decay. Conversely, less relevant events gradually fade. This mechanism allows FINMEM to efficiently identify, prioritize, and retain key events based on their nature and retrieval frequency, while gradually phasing out less impactful information, thereby maintaining a dynamic and relevant memory structure for financial decision-making.

# **B** Details on Evaluation Metrics

Below is a brief overview of these metrics:
Cumulative Return (CR) % measures the total value change of an investment over time by summing daily logarithmic returns, shown in Equation 8. Higher values indicate better strategy effectiveness.

$$\mathbf{CR} = \sum_{t=1}^{n} r_i = \sum_{t=1}^{n} \left[ \ln\left(\frac{p_{t+1}}{p_t}\right) \cdot \operatorname{action}_t \right] \quad (8)$$

, where  $r_i$  is the logarithmic return from day t to t + 1,  $p_t$  and  $p_{t+1}$  are the closing prices on days t and t + 1, respectively, and action<sub>t</sub> is the model's trading decision for day t.

**Sharpe Ratio (SR)** assesses risk-adjusted returns by dividing the average excess return  $(R_p)$  over the risk-free rate  $(R_f)$  by its volatility  $(\sigma_p)$ , detailed in Equation 9. Higher ratios signify better performance.

$$\mathbf{SR} = \frac{R_p - R_f}{\sigma_p} \tag{9}$$

**Annualized Volatility (AV) % and Daily Volatility (DV) %** quantify return fluctuations; AV is derived by scaling DV (*standard deviation of daily logarithmic returns*) by the square root of the annual trading days (252), as in Equation 10. This metric highlights potential return deviations across the year.

$$\mathbf{AV} = \mathbf{DV} \times \sqrt{252} \tag{10}$$

943

944

945

946

947

948

949

950

951

952

953

954

955

956

957

958

959

960

961

962

963

964

965

966

967

968

969

970

971

972

973

974

975

976

977

978

979

980

981

982

983

984

985

986

987

988

989

**Max Drawdown (MDD)** % calculates the largest portfolio value drop from peak to trough, as given in Equation 11. Lower values indicate lesser risk and higher strategy robustness.

$$\mathbf{MDD} = \max(\frac{P_{\text{peak}} - P_{\text{trough}}}{P_{\text{peak}}})$$
(11)

Note that CR and the SR are often considered more essential than AV and MDD in evaluating asset trading performance due to their focus on long-term gains and risk-adjusted returns by their definition. Here, we regard these two metrics as primary metrics when evaluating the experiment outcomes.

# C An example of mixed and lagged market signals: Partial investment insights of TSLA on 2021-03-05

Here are some insights from the agent's memory module for TSLA as of 2021-03-05. A few memory records are omitted; these are either neutral or positive. Despite this, the stock price trend for TSLA is sharply downward, conflicting with the overall positive financial sentiments and market signals. Utilizing proprietary models such as GPT4 and GPT-01 as backbones, the financial decision-making agent can leverage other investment insights like historical momentum and self-reflection to consistently support a 'Sell' decision. In contrast, the largesized open-source models like Qwen2.5-72B and DeepSeek-67B-Chat exhibit instability in producing consistent actions across repeated experimental trials.

## Short-term Memory

1 Sentiment: Negative: The key insights from the news regarding Tesla Inc (NASDAQ: TSLA) losing market share to Ford Motor Company's (NYSE: F) Mustang Mach-E in the United States are as follows:1. \*\*Market Competition\*\*: Tesla is facing increased competition in the electric vehicle (EV) market, particularly from established automakers like Ford. The Mustang Mach-E's success indicates that other companies are effectively entering the EV space, which could impact Tesla's dominance.2. \*\*Market Share Impact\*\*: The loss of market share to Ford suggests that Tesla's growth trajectory in the U.S. may be challenged. Investors should consider the potential implications for Tesla's revenue and profitability if this trend continues...

991

992

995

996

997

1001

1002

1003

1004

1005

1006

1008

1009

1010

1011

1013

1014

1015

1016

1017

1018

1019

1020

1021

1022

1023

1024

1025

1026

1028

1029

1031

1032

1033

1034

1035

1036

1037

1039

- 2 Sentiment: Positive: The news about a Banksy artwork being burned and sold as a non-fungible token (NFT) highlights a few key insights relevant to investment decisions, particularly concerning innovative and disruptive technologies: 1. \*\*Emergence of NFTs\*\*: The transformation of physical art into digital assets through NFTs signifies a growing trend in the digital economy. This trend could influence sectors beyond art, including technology and finance, as more industries explore blockchain applications ...
  - 3 Sentiment: Neutral: To provide a summary of key insights relevant to making investment decisions about Tesla (TSLA) from the preopen movers news, I would focus on the following aspects: 1. \*\*Stock Performance\*\*: Look for any significant pre-market price movements for TSLA. If the stock is showing a notable increase or decrease, it could indicate investor sentiment or reaction to recent news.2. \*\*News Catalysts\*\*: Identify any specific news items or announcements related to Tesla that might be influencing its stock price. This could include earnings reports, product launches, regulatory news, or changes in leadership...
- 4 Sentiment: positive: The key insights from the news about Tesla's vehicle registrations in Germany are as follows: 1. \*\*Significant Growth in Registrations\*\*: Tesla experienced a 78% year-over-year increase in vehicle registrations in Germany as of January. This substantial growth indicates a strong demand for Tesla vehicles in one of Europe's largest automotive markets.2. \*\*Market Penetration\*\*: The surge in registrations suggests that Tesla is successfully penetrating the German market, which is traditionally dominated by local automakers. This could imply a growing acceptance and preference for electric vehicles (EVs) in Germany, benefiting Tesla as a leading EV manufacturer...

# **Mid-term Memory**

...

1 Sentiment: Negative: The news about U.S. 1042 Senator Elizabeth Warren's proposal for a 1043 Democratic wealth tax could have several im-1044 plications for investment decisions regarding 1045 Tesla (TSLA):1. \*\*Impact on Wealthy In-1046 vestors\*\*: The proposed wealth tax targets the 1047 wealthiest Americans, which could include 1048 major shareholders and investors in Tesla. If 1049 implemented, this tax could lead to changes in 1050 investment strategies among these individuals, 1051 potentially affecting Tesla's stock price due 1052 to shifts in large-scale holdings....5. \*\*Longterm Implications\*\*: While the immediate 1054 impact might be uncertain, the long-term im-1055 plications of a wealth tax could lead to a more 1056 equitable distribution of wealth, potentially 1057 increasing consumer spending power. This 1058 could benefit Tesla if it results in a larger mar-1059 ket for its products. Investors in Tesla should 1060 closely monitor the progress of this proposal 1061 and consider its potential effects on market 1062 dynamics and investor behavior. 1063

1040

1041

1083

1084

2 Sentiment: positive: The key insights from 1064 the news information regarding consumer 1065 preferences in China for electric cars, partic-1066 ularly focusing on price and driving range, 1067 can be relevant for making investment deci-1068 sions about Tesla (TSLA) in the following 1069 ways: 1. \*\*Market Demand Alignment\*\*: 1070 Tesla's ability to align its product offerings 1071 with the primary concerns of Chinese con-1072 sumers 2014price and driving range2014could 1073 significantly impact its market share and sales growth in China. Investors should monitor 1075 how Tesla's models compare to competitors 1076 in these aspects... Overall, Tesla's success in 1077 addressing the key consumer preferences in 1078 China 2014price and driving range 2014will 1079 be critical for its growth prospects in the re-1080 gion, making these factors important consid-1081 erations for investors. 1082

# Long-term Memory

1 Sentiment: Negative: The key insights from<br/>the news about Tesla Inc. (NASDAQ: TSLA)<br/>facing competition in China are as follows: 1.<br/>\*\*Increased Competition\*\*: Tesla's Model 31085<br/>1086<br/>1086

is experiencing significant competition from 1089 a budget electric vehicle produced through 1090 a joint venture involving General Motors 1091 (NYSE: GM). This indicates a more compet-1092 itive landscape in the Chinese electric vehi-1093 cle (EV) market, which could impact Tesla's 1094 market share and sales growth in the region... 1095 For investors, these insights suggest a need 1096 to closely watch Tesla's strategic responses 1097 to competition in China, its pricing strate-1098 gies, and any potential impacts on its finan-1099 cial performance. Additionally, understanding 1100 the broader competitive landscape and Tesla's 1101 ability to sustain its growth trajectory will be 1102 crucial for making informed investment deci-1103 sions. 1104

2 Sentiment: positive: The news about Bill 1105 1106 Gates' concerns regarding Bitcoin's impact on climate change highlights a broader issue 1107 of environmental sustainability in the tech 1108 and financial sectors. Here are the key in-1109 sights relevant to making investment decisions 1110 about Tesla (TSLA): 1. \*\*Environmental Im-1111 pact Awareness\*\*: Bill Gates' concerns un-1112 derscore the growing awareness and scrutiny 1113 of the environmental impact of technology 1114 and financial products. This is relevant for 1115 Tesla, as the company positions itself as a 1116 leader in sustainable energy and electric ve-1117 hicles (EVs)... Overall, the emphasis on en-1118 vironmental impact and sustainability in the 1119 tech sector could reinforce Tesla's strategic ad-1120 vantages and appeal to investors prioritizing 1121 green investments. 1122

> D Case Study on Single Stock Trading: Forecast for TSLA on 2022-10-25 to Predict Trading Decision on 2022-10-26

1123

1124

1125

1126

...

## Initialize Profile

### **1. Operations:**

- Provide a performance overview of the trading stock based on available data.

- Set up the risk inclination as the key character of the trading agent.

2. Range: Financial information such as the financial sectors, historical performance, and previous stock trends of the trading stock.

**3. Prompts:** You are an experienced trading manager and investment firm. Your task is to make informed decisions on the given stock based on the provided information.

Under Self-Adaptive Risk Character Setting: When historical momentum is positive, you are a risk-seeking investor. But when historical momentum is negative, you are a risk-averse investor.

4. General background setting:

You have accumulated a lot of information about the following sectors, so you are especially good at trading them: 1)Electric Vehicles (Automotive Sector). 2) Energy Generation and Storage...From year 2021 to 2022 September, Tesla's continued growth and solid financial performance over the defined period ...

### Summarize

### **1. Operations:**

- Summarize different types of input information.
- Distribute them to corresponding layers of the long-term memory database.
- 2. Range: Daily market news, Long Documents such as company 10-K and 10-Q reports
- 3. Prompts:
- (1). Summarize the contents: Summarize the following documents into 1000 words.

- (2). Comprehend the investment sentiment of news insights: The positive, neutral and negative scores are for understanding the investment sentiments, opinions, or emotions. For example, positive news about a company can lift investor sentiment, encouraging more buying activity, which in turn can push stock prices higher...

# 4. Outputs:

## (1). To Shallow Memory Layer:

- [News (ID: 261)] Here's How Much You Would Have Made Owning Tesla Stock In The Last 10 Years Tesla (NASDAQ:TSLA) has outperformed the market over the past 10 years by 50.69% on an annualized basis producing an average annual return of 60.76%. Currently, Tesla has a market capitalization of \$683.54 billion.... The sentiment is {positive}.

- [News (ID: 278)] Tesla Q3 Earnings Are Imminent. Can Nio Foreshadow What's To Come? What To Know Before The Print Tesla Inc (NASDAQ: TSLA) shares were trading down slightly Wednesday afternoon ahead of the automakers' third-quarter report, but the stock is up 6% over the last five sessions... The sentiment is {positive}.

### (2). To Intermediate Memory Layer:

- [Form 10-Q (ID: 222)] Tesla Q3 2022 revenues were \$21.5 billion, up 56% year-over-year. Automotive sales revenue grew 56% to \$17.8 billion driven by higher Model 3/Y and Model S/X deliveries. Gross automotive margin declined to 27.9% due to cost inflation and factory ramps. Net income was \$3.3 billion, up 102% year-over-year. Positive free cash flow was \$6.1 billion...

- [News (ID: 275)] Tesla Q3 Earnings Highlights: Record Revenue, Operating Margin And Free Cash Flow, Tesla Semi Deliveries Coming In December Electric vehicle leader Tesla Inc (NASDAQ: TSLA) reported third-quarter financial results after market close Wednesday...The sentiment is {neutral}.

- [News (ID: 274)] Tesla Preps For 2023 Cybertruck Launch, Will Make Battery Packs In California The Cybertruck is one of Tesla Inc. (NASDAQ: TSLA) most hotly anticipated, but also most delayed, products. - ...The sentiment is {negative}.

## (3). To Deep Memory Layer:

- [News (ID: 161)] Tesla Whale Trades Spotted A whale with a lot of money to spend has taken a noticeably bearish stance on Tesla. Looking at the options history for Tesla (NASDAQ:TSLA) we detected 477 strange trades. The sentiment is {positive}.

- [Self-reflection (ID: 226)] Given the short-term positive news score in the market for TSLA and a positive cumulative return, there is a high probability of continued growth in the short term. However, investor should be aware of potential threats in the mid-term market with competitors like General Motors, and Nio...

### Observe

1. Operations: Access and interpret market indicators such as current stock prices and historical momentum data.

**2. Range:** Stock's daily adjusted closing price, historical momentum in the past k days (k = 3 in this case), etc. **3. Prompts:** 

- The information below provides a summary of stock price fluctuations over the previous few days, which is the "momentum" of a stock. It reflects the trend of a stock. Momentum is based on the idea that securities that have performed well in the past will continue to perform well, and conversely, securities that have performed poorly will continue to perform poorly.

4. Outputs:

- (1). The daily adjusted closing price of TSLA on {2022-10-25} is {\$222.42}.

- (2). **Train:** On {2022-10-25}, the momentum of TSLA, indicated by the price difference between the current and the next trading

day, is {\$2.22}.

**Test:** On {2022-10-25}, the historical momentum of TSLA, as measured by its cumulative logarithmic returns in the past

{3} days, was {7.05%}.

15 Figure 4: First section of FINMEM's workflow for perceiving and processing multi-sourced information from market environment.

### Reflect

### 1. Operations:

**Train:** - Infer the reasoning from the retrieved memories insights for the recommended trading actions the increase or decrease of the trading stock price.

**Test:** - Generate recommended trading actions. - Infer the reasoning for making such trading action.

Train & Test: - Track the memory insight IDs that contributed to the current trading decision.

**2.** Range: Interpretation of observed market indicators and top k insights from each layer of the long-term memory database.

### 3. Prompts:

*Train:* Given the following information, can you explain why the financial market fluctuation from the current day to the next day behaves like this? Summarize the reason for the decision. You should provide summary information and the IDs of the information to support your summary.

*Test:* Given the information, can you make an investment decision? Just summarize the reason for the decision. Consider only the available short-term information, the mid-term information, the long-term information, and the reflection-term information. Consider the momentum of the historical stock price. Also, consider how much share of the stock the investor holds now. You should provide exactly one of the following investment decisions: buy or sell. You could go with the 'hold' option if necessary. You also need to provide the IDs of the information to support your decision. **4. Outputs:** 

### Train:

- Recommended Action Reason: The recommended trading action is 'Buy'. This can be attributed to positive market sentiment driven by a range of favorable factors. Notably, Tesla's strong battery supply chain, effective market expansion plans in Germany, and its ability to maintain market strength despite potential competitive threats positively influenced investor confidence, contributing to the stock price rise.

- Reflected Contributed Insights IDs: {261, 278, 275, 274, 161,...}

Test:

- Trading Action: {Buy}

- Recommended Action Reason: Considering the short-term positive revenue growth, higher Model 3/Y and Model S/X deliveries, positive free cash flow, and increased cash and cash equivalents, Tesla appears to have strong financial health. The momentum being positive, it also encourages a risk-seeking strategy. The significant positive signs from news such as the imminent Q3 earnings add to the favorable view of the investment. Despite the mid-term neutral sentiment, long-term information shows only marginal drops and the stock is still looking strong overall. These, along with positive reflection-term information, form the reasoning behind the 'buy' decision...

- Reflected Contributed Insights IDs: {261, 278, 275, 274, 161,...}

Figure 5: Second section of FINMEM's workflow for generating trading action, reasoning and reflection.