

# AlphaManager: A Data-Driven-Robust-Control Approach to Corporate Finance\*

Murillo Campello<sup>†</sup>      Lin William Cong<sup>‡</sup>      Luofeng Zhou<sup>§</sup>

First draft: Feb 2023; this draft: March 2025.

## Abstract

Corporate decision-making entails complex, high-dimensional, non-linear stochastic control during which managers learn and adapt via dynamic interactions with the market environment. We propose a data-driven-robust-control (DDRC) framework to complement traditional theory, reduced-form models, and structural estimations in corporate finance research, emphasizing both empirical explanation and prediction of firm outcomes while delivering policy recommendations for a variety of business objectives. Specifically, we develop a predictive environment module using supervised deep learning and integrate a decision-making module based on generative deep reinforcement learning. By incorporating model ambiguity and robust control techniques, our framework not only better explains and predicts corporate outcomes in- and out-of-sample but also prescribes key managerial actions that significantly outperform historical ones. We document rich heterogeneity in model prediction performance, ambiguity, and policy efficacy in the cross section of U.S. public firms and across time regimes. Importantly, DDRC helps delineate where theory and causal analysis should concentrate, integrate fragmented knowledge (e.g., via transfer learning), and reveal managerial preferences (through an extension combining inverse reinforcement learning).

**Keywords:** AI, Ambiguity, Big Data, Corporate Finance, Deep Learning, Generative Modeling, Inverse RL, Offline RL, Stochastic Control, Transfer Learning.

---

\*We thank Wei Jiang and Tom Sargent for detailed feedback. We also thank Serdar Dinc, Jianqing Fan, Zhiguo He, Markus Pelger, Lea Stern, conference and seminar participants at AI and Big Data in Finance Research (ABFR) Webinar, 8th Annual Global Quantitative and Macro Investment Conference (Wolfe Research), the 10th Annual Workshop of Business Financing and Banking Research Group (USyd), 5th Big Data Econometrics Theory and Applications Conference (CJoE), China Meeting of Econometric Society (CMES 2021, Shanghai), Chinese University of Hong Kong, City University of Hong Kong, Columbia Business School, Cornell Tech, 2024 EFMA Annual Meeting, 2023 GSU-RFS FinTech Conference, Harvest Fund 25th Anniversary Ceremony, INFORMS Annual Conference, 4th International FinTech Research Forum (RUC), Johns Hopkins University, National University of Singapore, NYU Stern, Peking University Guanghua, Rutgers Business School, 6th Shanghai Financial Forefront Symposium, Tsinghua University SEM, UBC Sauder, 4th Workshop on Big Data Econometric Theory and Application, and 2023 XJTLU AI and Big Data in Accounting and Finance Conference for helpful comments and discussions. We are grateful to Matthew Wong, Logan Kraver, and Atharwa Pandey for exceptional research assistance. Send correspondence to Cong at will.cong@cornell.edu.

<sup>†</sup>University of Florida Warrington College of Business and NBER.

<sup>‡</sup>Cornell SC Johnson College of Business (Johnson), ABFER, IC3, and NBER.

<sup>§</sup>New York University Stern School of Business

# 1 Introduction

Corporate finance research studies firm decision-making and outcomes using theoretical models and archival data. The literature is built on simplified and tractable representations of the corporate environment where rational agents optimize their utilities and interact. Notably, scholars have raised concerns about the low predictability and limited explainability of corporate outcomes, underscoring a systematic gap in our understanding of corporate decision-making (see, e.g., Graham, 2022, Presidential Address of the American Finance Association). Many theoretical models in corporate finance are static, analyze partial-equilibrium, or overlook the interactions with the environment, while empirical studies typically report evidence that is merely “consistent with” *ex-ante* a priori theories without providing a realistic alternative (Spiegel, 2023). In contrast, real-world managerial actions and firm outcomes are now recorded with exceptional granularity and timeliness. Our paper demonstrates that a data-driven perspective — leveraging abundant data, algorithm advances, and powerful computation — can reveal novel empirical patterns and deliver valuable economic insights to guide both academic research and managerial practice.

Implementing a data-driven approach to uncover empirical patterns and advise corporate decision-making presents several challenges. First, managers face a combination of complex decisions in a high-dimensional action space, often contingent on numerous state variables. Their interdependent actions produce highly nonlinear effects that go beyond the reach of conventional econometric models and low-dimensional causal inference. Second, managerial decisions interact dynamically with the economic environment, with market feedback further influencing subsequent choices (e.g., Bond et al., 2010, 2012; Edmans et al., 2012, 2015). Such feedback loops not only impose significant costs but also complicate empirical analyses.<sup>1</sup> Moreover, because optimal corporate decisions are rarely “labeled” in historical data, standard supervised learning techniques are of limited use. Third, unlike physical laws, financial markets evolve rapidly; consequently, concerns about data distributional shifts become both relevant and pressing.<sup>2</sup>

---

<sup>1</sup>More generally in social science research, real-time online interaction to generate new data is impractical, either because unfiltered, continuous data collection is expensive (e.g., in high-frequency trading), unethical (e.g., hiring and firing employees), or possibly dangerous (e.g., law enforcement). Even in domains where online interaction is feasible, we might still want to utilize previously collected data instead — for example, if the domain is complex and effective generalization requires large datasets. Therefore, as explained later, we pursue an offline reinforcement learning (RL) approach using historical data similar to Cong et al. (2020).

<sup>2</sup>Data shifts, also known as distributional shifts, occur when the joint distribution of inputs and outcomes differs between training and test samples.

Given that managerial decision-making is essentially robust control problems characterized by high dimensionality, nonlinearity, dynamic learning, and evolving complexity, existing models often fall short of providing practical policy recommendations for corporate executives. For example, reduced-form empirical models can illuminate the economic mechanisms underlying a particular policy for a selective group of agents; however, their focus is typically local and low-dimensional. As a result, they lack the capacity to explain broader empirical outcomes, generate comprehensive counterfactuals, or yield generalizable recommendations beyond isolated causal effects. Structural estimations excel at managing model complexity to generate counterfactuals by modeling environments in a holistic manner that retains interpretability and a clear economic rationale. However, these models confine themselves to analytically tractable theories and conventional Markov decision processes (MDPs) with pre-specified transition probabilities, largely neglecting dynamic feedback and continuous learning about the environment. Doing so limits their ability to accurately explain observed data, generate reliable out-of-sample predictions, or ultimately provide practical guidance for managerial decision-making.

To advance predictions of firm outcomes under various counterfactual managerial decisions — and ultimately to guide corporate decision-making — we integrate deep learning, offline reinforcement learning, and robust control based on ambiguity, thereby framing corporate decision-making as a data-driven-robust-control (DDRC) problem. The framework we propose (“AlphaManager model”) serves as a data-driven counterpart to structural estimations by searching a broader modeling space that lends theoretical insights with data-driven patterns. In doing so, it provides a comprehensive depiction of the economic system, explain corporate outcomes, and generate high-dimensional, effective recommendations for enterprise decisions. In our DDRC framework, managerial decision-making is modeled as a robust control problem in which managers maximize their utilities based on contemporaneous states. AlphaManager comprises two modules: (1) the predictive environment module (PEM) and (2) the decision-making module (DMM). In PEM, we leverage supervised deep learning to capture the nonlinear, high-dimensional features inherent in financial big data. In DMM, assuming PEM as given, we apply offline deep reinforcement learning (RL) to reduce search costs and incorporate of flexible managerial objectives along with dynamic feedback.<sup>3</sup> Fi-

---

<sup>3</sup>Reinforcement learning (RL) is “learning how to map situations to actions so as to maximize a numerical reward signal.” It is one of the three paradigms of modern machine learning together with supervised learning and unsupervised learning. In RL, an agent learns about states of its environment and takes actions that potentially affect states going forward as well as its objective function to maximize. RL is particularly well suited for this task because it learns optimal actions through sequential decision-making and iterative experience ac-

nally, the robust control techniques such as ambiguity aversion are introduced as constraints, ensuring that AlphaManager remains conservative in the face of high model uncertainty (for instance, due to overfitting, data shifts, or endogeneity concerns), thus mitigating risks.

AlphaManager first constructs predictive environments to test counterfactual via PEM, which generates counterfactual forecasts that both explain and predict variations in corporate outcomes. This process circumvents the need for costly experimentation or explicit causal identification, while still capturing environmental feedback effects. Despite apparent complexity, the DDRC framework enhances traditional reduced-form and structural approaches in several important ways. First, by simulating outcomes from a range of counterfactual managerial decisions, PEM uncovers empirical patterns that can both challenge and refine existing theories — informing new theoretical research. Second, it signals scenarios in which historical data and theory alone are insufficient, underscoring the necessity of reduced-form or structural models for effective learning. To assess model uncertainty, PEM deploys an ensemble of deep neural networks with identical architectures but different initializations; disagreement among these networks indicates that the model is extrapolating beyond its training data, rendering counterfactual predictions less reliable. Finally, our framework integrates insights from *both* reduced-form and structural models through ambiguity-guided transfer learning.<sup>4</sup> In situations where ambiguity is high, empirical causal identification and theoretical modeling play a joint, critical role in enriching the base for counterfactual predictions.

AlphaManager then incorporates a decision-making module (DMM) for managerial policy optimization, marking our work the first to combine RL with robust control in corporate finance. At its core, corporate finance problems can be viewed as involving a “system” (the firm and its environment) and a “controller” (the manager).<sup>5</sup> The controller’s objective is to optimally manage the system, constrained by (1) the manager’s knowledge of the

---

cumulation (Sutton and Barto, 1998). Unlike dynamic programming in structural estimation – which relies on known transition probabilities and fixed rewards structures — RL addresses a more general MDP where transition probabilities and rewards are unknown. RL does it either via a model-based approach (e.g., AlphaManager which learns the a model of the environment from data) or a model-free approach (e.g., Q-learning).

<sup>4</sup>Transfer learning has proven valuable across many domains; for instance, the pre-training used in large language models (LLMs) is a form of transfer learning. This approach can be tailored to economics and finance (see Chen et al. (2023)).

<sup>5</sup>Traditionally control theory is generally designed to solve linear (or linear quadratic) systems with well-defined objectives, law of motions for states, and constraints. However, real-world systems are nonlinear, and techniques that linearize these systems are limited to specific cases. Moreover, accurately modeling such system is challenging, which is why traditional model-driven control approaches — such as structural models in corporate finance — often exclude systems whose underlying dynamics are not fully known. The complexity of these systems increases when addressing hyper-scale issues, such as financial contagion or climate change responses, where unknown nonlinearities and unobserved environmental states render standard dynamic programming and simulation methods inadequate.



system’s state and environment — information that is gleaned via “system sensors” like accounting, auditing, and reporting — and (2) the limited set of parameters that can be directly controlled. In AlphaManager, these two aspects are operationalized by first training a deep supervised learning based PEM from historical data that allows effective counterfactual analysis and data generation. We then train DMM using RL to identify the most effective combination of managerial actions for a given business objective. RL generates the optimal control trajectory based on an exogenously defined reward structure (e.g., market capitalization appreciation) using unlabeled data and by interacting with the environment. In doing so, it offers normative recommendations for managerial decision-making. Finally, by integrating robust control theory — particularly through ambiguity aversion — we guide the training of AlphaManager to extract optimal decisions even under high model uncertainty (Hansen and Sargent, 2023).

Our AlphaManager application is trained on Compustat, which contains a long panel of fundamental variables of US-listed firms. We study nearly 20 thousand distinct firms with over half-a-million unique firm-quarter observations, ranging from 1976 to 2023. We supplement our dataset with additional stock market data from CRSP, and incorporate macro-level data; e.g., the National Financial Conditions Index (NFCI), from the Chicago Fed. We define two sets of variables: state variables and managerial decision variables. State variables describe the state that a firm faces in a given period of time. Typical examples involve fundamentals of firms (internal states), and macroeconomic or market conditions (external states). Managerial decision variables mediate the future state dynamics and also influence the utility functions of managers.

In PEM, a deep neural network is trained to minimize the mean-squared error (MSE) between real future states and predicted ones, with the input of current states and current managerial decisions. Our PEM achieves high accuracy in predicting and explaining the evolution of firm outcomes (the state variables) with the help of information on managerial decisions, i.e., the managerial planning one period forward. For instance, AlphaManager produces a 64.7% cross-sectional  $R^2$  for book asset growth and 3.2% for market cap growth, both out of sample. To make it comparable to other empirical research in corporate finance, we also calculate the predictability without managerial planning information. The predictability and explainability remain high for many variables, but more interestingly, for outcomes such as book asset, market capitalization, and enterprise value, the out-of-sample  $R^2$ s become negative. This finding is consistent with empirical asset pricing research where the market environment is simply too noisy. Critically for our analysis, however, these results

inform us which managerial decision variables matter most for firm dynamics going forward.

Based on PEM, we are able to analyze the heterogeneous outcomes and model ambiguity under counterfactual managerial decisions by state variable, sector, book-to-market decile, and macroeconomic regimes. PEM performs particularly well in trade and transportation, education and healthcare, and manufacturing; predictability is also higher during expansions. Firms with higher firm-level states in the cross section have lower MSE and ambiguity. Notably, model ambiguity is highest for high book-to-market firms.

In DMM, with an exogenously specified utility function of managers, we obtain optimal policies that generate managerial decisions that maximize utilities based on current states. When optimizing the short-term market cap growth, the quarterly outperformance of optimal decisions compared to real managerial decisions is 10.1%. When optimizing the long-term market cap growth, the quarterly outperformance is 8.7%. We find similar patterns when the objective is set as enterprise value growth, with quarterly outperformance of 4.4% and 2.7%, for short and long-term growth, respectively.

We contrast DMM-suggested policy under short-term versus long-term growth in firm value, as well as the term structure of this growth under DMM suggested optimal managerial policy. The contrast between short- and long-term oriented managerial decisions and their implications for firm value are receiving renewed interest in the corporate finance literature (see, e.g., Almeida et al., 2024). Short-term AlphaManager RL outperforms long-term RL in the long run but with higher ambiguity. After ambiguity adjustment, the performance of short-term RL is decreasing in the evaluation time horizon compared to the long-term RL. We further discuss the long-term implication of firm valuation under managerial short-termism from the point of view of board members. Notably, this approach does not require any knowledge of real decisions, which is the key difference between RL and commonplace supervised machine learning algorithms.

Naturally, the action and information space of a manager is vast. One might wonder about the set of variables needed for a given application and about when should one stop collecting data and constructing new variables. Conveniently, our proposed approach is not hampered by these concerns. Among other things, most variables one uses in applied corporate finance research already convey information about many others. For example, the firm’s capital structure already conveys information about its access to credit, risk, asset mix (collateral), and more. The same applies to firm size, or as to whether it is publicly traded, hence the amount of information available to investors. Data “confoundedness” is a plus — not a challenge — to our data-driven approach. The approach we propose is less prone to criticism re-

garding inconsistencies in variable selection made by researchers in existing studies (see Mitton, 2022) — inconsistencies leading to practices like *P*-hacking and *ex-post* theory-fitting.

Given the large counterfactual tests we do using the predictive environmental module, AlphaManager necessarily makes extrapolation and interpolation from historical data. Given that historical data are endogenously generated, potential model misspecification poses a particular challenge. The literature traditionally examines parameter stability using Chow-type tests (Andersen et al., 2015). One can also detect static misspecification’s by testing moment conditions under a GMM framework (Pan, 2002). For dynamic misspecification tests, Jarrow and Kwok (2015) provide an intuitive exact calibration approach utilizing analytical theoretical models. Our use of deep neural networks goes beyond economic theory and is data-driven. While it explores a much larger modeling–functional space, model misspecification is generally dynamic and hard to detect with conventional econometric tools. Therefore, to guard against and mitigate misspecifications, we adapt the new approach of utilizing ambiguity aversion and entropy-based measures to a data-driven setting. The model ambiguity measure also informs us of scenarios where additional knowledge from theories and causal identifications is important and the ones where predictive models trained on historical observations suffice. It not only guides corporate finance researchers on the dimensions to focus on, but also allows the integration of fragmented knowledge through ambiguity-guided transfer learning.

Our study contributes to the emerging literature on AI in finance. Machine learning and natural language processing have been widely applied in investment and asset pricing (e.g. Cong et al., 2020, 2021; Gu et al., 2020; Feng et al., 2020). They have also seen applications in studies on corporate finance or financial market risk (e.g., Cong et al., 2018; Li et al., 2020; Bellstam et al., 2020; Hanley and Hoberg, 2019; Campello et al., 2024); but their applications beyond creating or improving the measure of some explanatory variables are rather limited. Exceptions to this line of work include Erel et al. (2018) on predicting board director performance, Cao et al. (2023) on how machines and managers interact in the context of conference calls and earnings announcements, Cao et al. (2021) on how an AI analyst can provide additional insights on stock market forecasts to human analysts. We note that these studies employ standard — often rudimentary — models designed for prediction or focus exclusively on supervised learning (learning through examples) without applying the core paradigm (i.e., deep RL) in AI innovations in the past two decades. In contrast, we present the first AI and robust control application in finance. In fact, we are among the first to apply model-based offline RL in economics to offer a data-driven alternative to reduced-form models and structural estimations.

Our paper also contributes to the literature on model uncertainty (“ambiguity”) and robust control. Ambiguity represents a source of uncertainty where an economic agent is not confident in which prior models to choose for prediction tasks Hansen and Sargent (2023).<sup>6</sup> Its applications in finance are rare and mostly focus on asset pricing and investment (e.g., Wang, 2005; Mamaysky et al., 2007; Dicks and Fulghieri, 2021), with Dicks and Fulghieri (2019) (financial intermediation), Garlappi et al. (2017) (corporate investment), Izhakian et al. (2022) (capital structure), and Malenko and Tsoy (2020) (security design) as exceptions. Meanwhile, robust control has been put forth in Hansen and Sargent (2001) where a rational agent solves a stochastic control problem under ambiguity aversion. Studies about robust control are mostly theoretical, except for Barnett et al. (2020) which applies robust control to the context of climate change risks. Our work contributes to the literature on model uncertainty as the first empirical/methodological application of the ambiguity concept in corporate finance. We also add to robust control studies by estimating ambiguity with the help of deep learning and then approximate solutions to robust control problems using estimated ambiguity and offline RL.

Finally, our work can be placed in the context of the emerging computer science literature on offline RL (a.k.a, batched RL, e.g., Fujimoto et al., 2019; Kidambi et al., 2020).<sup>7</sup> A number of papers have illustrated the power of such an approach in enabling data-driven learning of policies for dialogue (Jaques et al., 2019), robotic manipulation behaviors (Ebert et al., 2018; Kalashnikov et al., 2018), and robotic navigation skills (Kahn et al., 2021). While Cong et al. (2020) is the first finance paper that applies offline RL (with online updates) to portfolio management, the authors do not fully optimize the environment module to mimic the real environment. We instead follow Kidambi et al. (2020) to build and optimize an environment module to calculate the transition probability across states without requiring experimenting via a simulator or costly interactions with the actual corporate or market environment. Together with Chen et al. (2023), we are the first studies introducing transfer learning in finance. Transfer learning leverages knowledge from one task or domain and apply

---

<sup>6</sup>Hansen and Sargent (2023) defines three sources of uncertainty: risk (in-model innovation), ambiguity (model uncertainty), and misspecification (model class uncertainty). Campello and Kankanhalli (2024) provide a comprehensive review of the research of uncertainty in corporate finance.

<sup>7</sup>Unlike online RL where real-time interactions with the actual environment or environment simulators are possible and counterfactual statements can be evaluated directly, offline RL often works only on historical data without any online interactions with the environment to generate additional data for model training (e.g., Fu et al., 2020). In a pure data-driven scheme, offline RL enables researchers to explore fields that are considered infeasible by classical online RL algorithms, especially those closely related to human behaviors where environment interaction is costly, infeasible, or dangerous (Levine et al., 2020).

it to improve learning in a different but related task or domain. Our DDRC framework is compatible to incorporate inferences from existing empirical causalities and theoretical models to our PEM for better internal validity. Our innovation also lies in using ambiguity to guide the choice of transfer learning and the application in corporate finance. More broadly, our study adds to emerging studies utilizing AI for goal-oriented search, which involves both heuristic search using RL and greed search using panel trees (Cong et al., 2022, 2023).

The rest of the paper is organized as follows. Section 2 describes the DDRC framework. Section 3 details data and model training. Section 4 reports the functionality and empirical results from PEM, whereas Section 5 investigates the optimal managerial actions DMM recommends under various given objectives. Section 6 discusses novel research questions that DDRC is particularly suited for. Section 7 concludes.

## 2 The Data-Driven-Robust-Control Framework

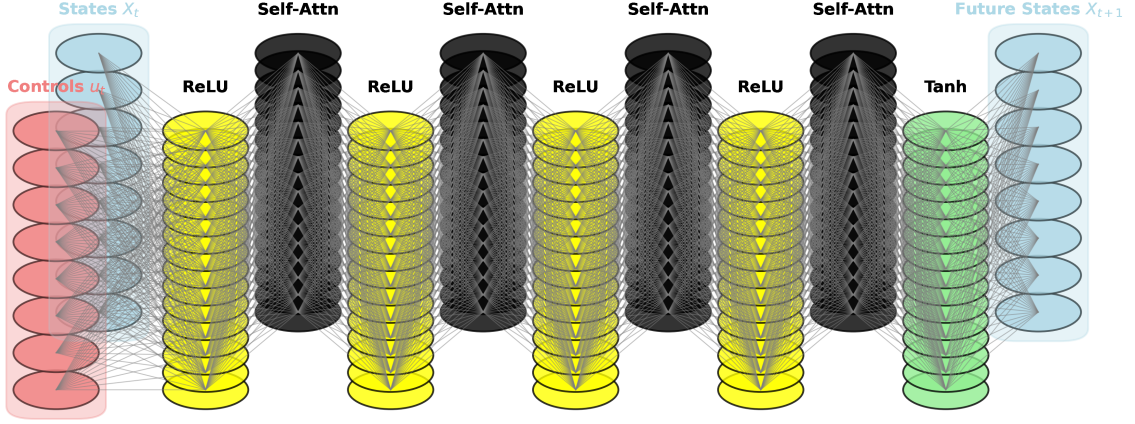
The DDRC approach is comprised of two modules, one utilizes a collection of deep learners to describe the market environment and how the outcomes of interest respond to managerial actions, and the other involves a reinforcement learner offering an optimal policy while dynamically interacting with the market environment and incorporating feedback. The resulting AlphaManager architecture is illustrated in Figure 1. The uses of the rectifier or *ReLU* (rectified linear unit) and *tanh* activation functions are standard in deep learning; the self-attention mechanism is widely seen in many deep learning models, such as transformer models; the state variable vector  $X_t$  and managerial control vector  $u_t$  are introduced shortly.

### 2.1 Predictive Environment Module (PEM)

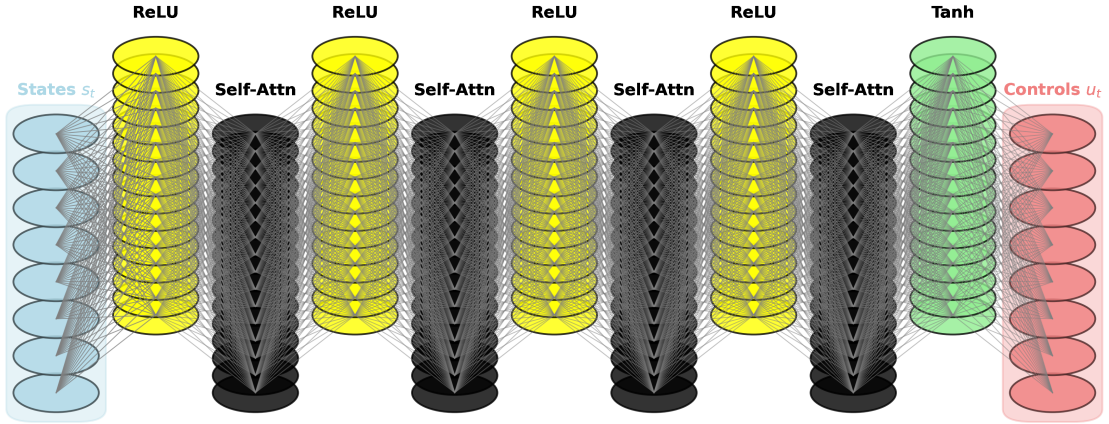
To conduct counterfactual analyses in our context, we need to model the market environment. Conventional causal analyses (e.g., instrumental variables and regression discontinuity designs) do not provide a comprehensive solution because these identification strategies are hard to come by or costly to establish (e.g., through experiments). Even when we have them, causal links are only identified locally and typically when varying a single treatment dimension. We are interested in learning about the outcomes of not only large corporate events, such as mergers or bankruptcies, but also any marginal decisions, such as combinations of increasing dividend payouts and reducing investment spending — decisions that are far more frequent and high-dimensional in the practice of financial management. By

Figure 1: AlphaManager (AM) Architecture

Panel A: Predictive Environment Module (PEM)



Panel B: Decision-making Module (DMM)



Note: This figure visualizes the AlphaManager architecture. Panel A displays the predictive environment module (PEM) with an input layer of a state vector  $\mathbf{X}_t$  concatenated with a decision (control) vector  $\mathbf{u}_t$ . The state vector  $\mathbf{X}_t$  (in blue color) includes firm-specific state variables such as firm fundamentals, macroeconomic variables, etc. The decision vector  $\mathbf{u}_t$  (in red color) contains leverage, cash holdings, equity financing, etc. Appendix A contains a full list of firm-specific state variables and control variables. The second, third, and fourth layers are fully-connected layers with rectified linear unit (*ReLU*) activation functions with a self-attention mechanism. The fifth layer of *tanh* function transforms the output spectrum to interval  $[-1, 1]$  and the last layer is the output layer with predictions of the system state evolution. Similarly, Panel B shows the structure of the decision-making module (DMM) mapping current state vector  $\mathbf{X}_t$  to the optimal control  $\mathbf{u}_t$  for a given managerial objective. When implemented, AlphaManager uses DMM to generate the optimal managerial actions and then predicts system state evolution using PEM.

building a market environment, we are taking a model-based RL approach that is similar to a structural estimation, except that the structure is data-driven rather than based on theory with more restrictive assumptions and closed-form solutions.

In AlphaManager, we use a deep neural network to model how the internal state variables and external market states going forward react to potentially high-dimensional managerial actions. Neural networks are designed to capture non-linear dependencies in high-dimensional spaces, and they have the potential to approximate any functional form (referred to as the *universal approximation theorem* in computer science, see, e.g., Hornik et al., 1989). The large number of corporate finance variables and decisions available and the potential nonlinear interactions among them necessitate such a module, which we refer to as the Predictive Environment Module (PEM).

Our approach recognizes that corporate decisions involve the manager taking a sequence of high-dimensional actions to optimize some given economic objective, which makes them stochastic control problems. Suppose the dynamics of the state variables in the system follow:

$$\Delta X_{t+1} = f(X_t, u_t) + \varepsilon_{t+1}, \quad \varepsilon_{t+1} \sim N(0, \Sigma), \quad (1)$$

where  $X_t$  denotes a firm’s internal state variables (e.g., accounting fundamental) and external state variables (e.g., inflation and unemployment) which are out of the manager’s control at time  $t$ ,  $u_t$  stands for the vector of managerial actions that are made by managers at time  $t$ ,  $f$  describes how the expected value of state change,  $\Delta X_{t+1}$ , corresponds to current state and managerial actions, and  $\varepsilon_{t+1}$  is the normally distributed risk term which represents the in-model uncertainty (i.e., the uncertainty that has been captured by the model). We can then use a fully connected neural network to approximate the function  $f$ .

PEM is expected to generate reliable out-of-sample predictions for counterfactual evaluations — tasks that require both extrapolation and interpolation. This process is inherently problematic because historical data are generated by endogenous managerial decisions, leaving gaps in some regions of the action-state space. Similar to reduced-form models, two primary challenges arise: overfitting and data shift. In structural estimations, these concerns are mitigated by assuming that the theoretical model perfectly represents the underlying system. Analogously, if the state variables we include fully capture the system’s evolution, the risks associated with extrapolation and interpolation would be minimized.

Various statistical procedures can mitigate these concerns without relying on such strong assumptions. For instance, to address overfitting, one common approach is to include L1-

norm or L2-norm penalties in the loss function during the training of a supervised deep neural network (e.g., Kaniel et al., 2023). Furthermore, by leveraging the virtue of complexity (e.g., Kelly et al., 2024), our over-parameterized PEM in the empirical implementation is designed to reduce overfitting and enhance out-of-sample performance. However, even when overfitting is controlled, PEM inherently introduces model ambiguity (Hansen and Sargent, 2023) because the data-generating process may shift, leading to uncertainty about whether a properly trained model will generalize to new datasets. In other words, neural network models with different parameters might perform similarly on the training data, yet yield divergent predictions due to a combination of residual overfitting and data shifts. As such, selecting the appropriate model *ex ante* adds an additional layer of uncertainty. In Section 2.3, we demonstrate how to empirically estimate ambiguity for any point in the variable space and address this issue by incorporating ambiguity aversion into the objective function.<sup>8</sup>

## 2.2 Decision-making Module and Offline Reinforcement Learning

Managers take a sequence of actions to solve the following stochastic control problem:

$$\max_{\{u_{t_0}, \dots, u_{t_0+T}\}} \mathbb{E}_{t_0} \sum_{t=t_0}^{t_0+T} r(X_t, u_t) \quad s.t. \Delta X_{t+1} = f(X_t, u_t) + \varepsilon_{t+1}, \varepsilon_{t+1} \sim N(0, \Sigma), \quad (2)$$

where  $r(X_t, u_t)$  is the instantaneous reward function given the state  $X_t$  and the managerial decision  $u_t$ , and  $T$  is the length of the optimizing time period associated with the managerial objective. An example of instantaneous reward function is the enterprise value growth for the next period given current state and managerial decisions.

The conventional approach to stochastic control problems involves deriving HJB equations to obtain closed-form solutions. When analytical solutions are unattainable, numerical algorithms can yield reliable, low-variance empirical approximations — provided the problem is low-dimensional. In our setting, the dynamics lack an analytical characterization, and the environment evolves in a high-dimensional, potentially nonlinear manner, which motivates our use of reinforcement learning.

In order to functionally represent managerial decisions of interest, we follow the canonical policy gradient algorithm in RL (Sutton and Barto, 1998) and assume that decisions are made

---

<sup>8</sup>This approach essentially adds a penalty based on relative entropy (also known as Kullback-Leibler divergence in computer science literature).



contingent on the state variable  $X$ . Specifically,

$$u_t = g(X_t), \quad (3)$$

where the functional representation  $g$  is called decision-making module (DMM), serving as a decision-making device contingent on the state vector. By substituting (3) into (2), we can reformulate the stochastic control problem as:

$$\begin{aligned} \max_{g(\cdot)} \quad & \mathbb{E}_{t_0} \sum_{t=t_0}^{t_0+T} r(X_t, u_t) \\ \text{s.t.} \quad & \Delta X_{t+1} = f(X_t, u_t) + \varepsilon_{t+1}, \quad \varepsilon_{t+1} \sim N(0, \Sigma) \\ & u_t = g(X_t). \end{aligned}$$

AlphaManager continuously learns about the environment through interactions and refines her understanding using the predictive environment module (PEM), which models the function  $f(X_t, u_t)$ . At the same time, she maximizes her expected cumulative payoff by optimizing over a broad space of dynamic policies, represented by  $g(X_t)$ , using the decision-making module (DMM).

## 2.3 Robust Control and Ambiguity Aversion

Neural networks usually have more free parameters than available training samples, leading to under-identified parameter estimates. Even when fewer parameters are used, endogeneity issues in the training data can undermine a model’s out-of-sample predictive power because data-driven training does not distinguish causal relationships or separate supply-side shocks from demand-side shocks. Consequently, variations in parameter initialization or model architecture may produce multiple models that perform similarly in-sample yet yield different out-of-sample predictions. Traditional model selection focuses on validation set performance, but when the first moment of predictions is comparable across models, additional metrics to guide model choices become necessary. Data shifts further compound the problem; the dispersion of predictions on the test set not only signals potential overfitting but also reflects model uncertainty arising from these shifts. This “across-model” uncertainty is referred to as *ambiguity* in robust control literature (Hansen and Sargent, 2023). The literature typically handles ambiguity problems in two steps. First, researchers estimate ambiguity through some metrics (e.g., relative entropy) to summarize the discrepancy among

multiple models defined as statistical distributions indexed by different parameters. If this discrepancy for a given test sample is large, the ambiguity for this specific sample is high. Second, when modeling a decision-maker’s problem, a punishment for ambiguity is added to the original objective function, making the agent ambiguity-averse.

To estimate ambiguity, we apply the idea of boosting to train a set of neural networks with similar structures for PEM, under the same loss function but with different initial coefficients or slight variations in the model specification.<sup>9</sup> In the supervised learning scheme, these networks perform the same functionality in-sample, as errors would be punished during the training process consistently. In contrast, out-of-sample — especially for unprecedented scenarios or unusual events that are not covered in the training sample —, the performances of these models might diverge. Accordingly, the disagreement level of these models is a natural proxy of whether a given input is likely within the training set span or not. When the disagreement level of these models is high, it is either because the current states are “unusual” or the current managerial decisions are “uncommon.” Formally, we trained a set of environment modules  $i = 1, 2, \dots, I$  mapping  $(X_t, u_t)$  to  $\hat{X}_{t+1}^i$ . The boosting error, motivated by the relative entropy which is a theoretical metric of ambiguity, for  $(X_t, u_t)$  is calculated as:

$$\text{BoostingError}(X_t, u_t) = \frac{1}{D} \sum_{d=1}^D \left( \max_{i=1,2,\dots,I} \hat{X}_{t+1,d}^i - \min_{i=1,2,\dots,I} \hat{X}_{t+1,d}^i \right)^2, \quad (4)$$

where  $D$  is the number of dimensions for state vector  $X$  and  $u_t$  is an arbitrary managerial control attempted. Note that there are alternative specifications of boosting errors one can use, which still capture the concept of ambiguity.

In addition to punishing ambiguity directly, a common practice in economics and computer science is to consider an agent solving a max-min problem.<sup>10</sup> Under model uncertainty, the agent tries to maximize the lowest reward generated by a set of models. This pessimism effectively serves as an additional device, together with ambiguity aversion, to control for model uncertainty in suggesting decisions.

To train DMM and solve the RL problem, we use fully connected neural networks to parameterize the decision-making device  $g$ , the mapping from state to control. We also

---

<sup>9</sup>The choice of ambiguity metrics could vary, as shown in Hansen and Sargent (2023). We choose to use boosting error in our context given its simplicity of calculation and its intuition of maximum dispersion of predictions for models around our main PEM in the sense of model structure and parameter initialization under the same training process.

<sup>10</sup>See Barnett et al. (2020) in economics and Jin et al. (2021) in computer science for modeling pessimistic agents solving max-min problems to mitigate ambiguity.

incorporate the idea of Kidambi et al. (2020) and Yu et al. (2020) on how to incorporate pessimistic reward and punishment on ambiguity in the training process. Specifically, in training DMM, for any given state  $X_{t_0}$ , the empirical objective function is calculated as:

$$J(X_{t_0}) := \sum_{t=t_0}^{t_0+T} \min_{i=1,2,\dots,I} r^i(X_t, g(X_t)) - \delta \cdot \max \left\{ 0, \max_{t=t_0, \dots, t_0+T} \text{BoostingError}(X_t, g(X_t)) - \theta(X_{t_0}) \right\} \quad (5)$$

In Eq (5),  $r^i(X_t, g(X_t))$  denotes the instantaneous reward function by environment module  $i$ , given state PEM-predicted time- $t$  state,  $X_t$ , and the decision-making device  $g$ , and we take the lowest instantaneous reward across  $I$  different environment modules. The ambiguity punishment parameter,  $\delta$ , is chosen to be very large so that we guarantee that the boosting error is always small enough. Empirically we pick  $\delta = 10^7$ .  $\text{BoostingError}(X_t, g(X_t))$  is the boosting error at time  $t$ , given PEM-predicted time- $t$  state,  $X_t$ , and the decision-making device  $g$ , calculated as Eq (4). The benchmark error,  $\theta(X_{t_0})$ , is calculated as a constant  $\alpha$  plus  $\beta$  times the boosting error of  $X_{t_0}$  and the lagged real managerial decisions  $\bar{u}_{t_0-1}$  which is known at time  $t$ :

$$\theta(X_{t_0}) = \alpha + \beta \cdot \text{BoostingError}(X_{t_0}, \bar{u}_{t_0-1}), \quad (6)$$

where  $\beta > 0$  and  $\bar{u}_{t_0-1}$  is the actual managerial decision in the previous period just a simple benchmark for how ambiguous the states correspond to. What we do here is essentially prioritizing the reduction in ambiguity over the baseline rewards, if the ambiguity exceeds some threshold. When considering different optimizing period  $T$ , given the monotonicity of  $\max_{t=t_0, \dots, t_0+T} \text{BoostingError}(X_t, g(X_t))$ , we choose higher values of  $\alpha$  and  $\beta$  when optimizing under a higher  $T$ . The empirical choices of  $\alpha$  and  $\beta$  under different optimizing period  $T$  are detailed in Appendix C.

## 3 Data and Training

### 3.1 Data and Variables

Our panel data is fairly standard in corporate finance research and is based on quarterly CRSP–Compustat merged database, where we are able to include 20,485 different firms ranging from 1976:Q1 to 2023:Q2, with 784,460 firm-quarter observations. The variables in this paper fall into two classes: state variables ( $X_t$ ) and managers’ decision variables ( $u_t$ ). In each period, managers make their decisions and their decisions will impact state variable

dynamics for the next period. Examples of state variables include fundamental variables that are not entirely determined by managers, such as the firm’s market capitalization, interest expenses, sales growth, earnings forecasts, and Tobin’s Q. They also include past managerial decisions. Decision variables  $u_t$  reflect managerial decisions directly (or at least largely so). Examples include corporate investment spending, cash savings, and capital structure.

In order to train our supervised learning-based PEM, we require the existence of observations for two consecutive fiscal quarters in standard Compustat dataset (with “datafmt” code to be “STD”). We only focus on domestic (with “popsrc” code to be “D”) industrial companies (with “indfmt” code to be “INDL”), and we exclude financial and utility firms (with “naics2” code to be 22 or 52) from our sample. We also require book assets to be positive. AlphaManager has a rolling training for RL. We divide our dataset by fiscal quarter to have roughly 30% initial training set, which implies that the initial burn-in is until 1991. We winsorize each variable at 1% and 99%, and subsequently normalize them to be bounded inside  $[-1, 1]$  for training convenience. We also consider macroeconomic states for the same time period, using the Chicago Fed National Financial Condition Index (NFCI) subindices (risk, credit, financial leverage, and non-financial leverage) as macro covariates. The detailed variable selection is listed in Appendix A. Summary statistics for firm-level variables are reported in Table 1.

### 3.2 Training and Computation

We describe the training of the AlphaManager model here and list the hyperparameter choices in Appendix B. PEM has 268 inputs, including 12 firm fundamental variables (e.g., log book assets), 12 variables capturing their growth (first difference in variable values). We further incorporate the corresponding variables lagged 1 through 4 quarters to capture patterns in the time series. Involving lagged variables provides additional information on time dependence, and incorporating the changes in values helps neural networks better extract more information than typically learned from the levels of state variables alone. We additionally include lagged decisions and their growth, which are publicly observable at the current time point, as well as their 1-4 quarter lagged versions. We also have 2 stock market states and their current growth, together with the 1-4 quarter lagged values. For the 4 macro states, we only consider their current levels and their levels lagged 1 through 4 quarters. For managerial decisions, we have 9 decision variables (e.g., leverage) and their correspondingly changes (growth) from their previous values. The output of PEM has 14 dimensions, namely

Table 1: Summary Statistics for Firm-specific States and Managerial Decisions

variable	count	mean	median	std
Leverage	784,460	0.2449	0.2020	0.2333
Acquisition	784,460	0.0047	0.0000	0.0227
Investment	784,460	0.0144	0.0068	0.0227
Cash	784,460	0.2041	0.0932	0.2651
Dividend	784,460	0.0019	0.0000	0.0051
DebtIssue	784,460	0.9122	0.0000	1.8047
EquityIssue	784,460	0.5231	0.0010	1.0551
RDExp	784,460	0.0130	0.0000	0.0309
Repurchase	784,460	0.0026	0.0000	0.0095
Total Assets	784,460	5.2788	5.0987	2.2292
Current Assets	784,460	4.2719	4.2478	2.1785
Sales	784,460	3.7052	3.6178	2.2490
Payables	784,460	2.7043	2.3038	2.0113
COGS	784,460	3.3015	3.1079	2.1523
InterestExp	784,460	0.5467	0.0000	1.1486
Inventory	784,460	2.4818	2.1510	2.2778
CurrentLiability	784,460	3.5870	3.3725	2.1935
Receivables	784,460	3.1032	2.9350	2.1419
Revenue	784,460	3.6999	3.6127	2.2525
MarketCap	784,460	5.2181	5.0620	2.2858
EnterpriseValue	784,460	5.7726	5.6021	2.2511
Volume	784,460	10.3049	10.4611	2.6685
Return	784,460	-0.0140	0.0000	0.2872

Note: This table documents the summary statistics of our firm-level state variables and decision variables. Our sample starts from 1976:Q1 to 2023:Q2, and covers 784,460 unique firm-quarter observations. In the sample selection, we keep domestic industrial firms (with “popsrc” code to be “D” and “indfmt” code to be “INDL”) from standard Compustat (with “datafmt” code to be “STD”). We exclude financial and utility firms (with “naics2” code to be 22 or 52) from the sample. We also require book assets to be positive. We winsorize our each variable at 1% and 99%. The detailed variable definitions and formulations are in Appendix A.

12 next-quarter changes in fundamental and 2 next-quarter changes in return and volume.

For training PEM, we first initialize a neural network and train the model in mini-batch mode (with batch size to be 2048 observations) using data until 1991:Q4, which is ready to predict system states in 1992:Q1 “out-of-sample.” In this stage, we train our neural network for 30 epochs. Next, we include 1992:Q1 data into our training set to update the environment module and use the updated model to predict the system states in 1992:Q2 for additional 5 epochs, and so on. We also train 10 auxiliary models with different realizations in the initialization phase and with slightly different model designs at the same time for ambiguity

estimation. We use four hidden layers in the baseline PEM, with each layer having 512 neurons. In total, we have 1,186,304 parameters.

Our DMM has 250 inputs, namely the inputs of PEM excluding current decisions and their growth versions ( $9 + 9 = 18$ ), and has 9 outputs, which are changes in current decisions compared to the last-quarter ones. For training DMM, we initialize a neural network and train the RL model using data until 1991 to recommend managerial decisions in 1991. We train the model for 64 epochs first, and then when the loss function (after ambiguity punishment) turns positive, we continue to train the model for 5 epochs and move on to the next stage. Then, we focus on 1992 data to recommend 1992 managerial decisions, and follow the same training scheme going forward in time. As baseline specifications, we consider two objectives, market capitalization growth and enterprise value growth. We do this for two horizons: one quarter ahead (“short horizon”) and eight quarters ahead (“long horizon”). We use a four-hidden layer network in training our DMM, with each layer having 256 neurons. In total, we have 265,993 parameters.

We train our neural networks on the Red Cloud platform, with 16 CPUs, 55G memory, and an A100 GPU.<sup>11</sup> The training process takes 30-35 hours for PEM (including the auxiliary models), 8 hours for the DMM with short-term objectives, and 75-80 hours for the DMM under long-term objectives. We use Adam optimizer (Kingma and Ba, 2014) in training neural networks.

## 4 The Economic Environment for U.S. Public Firms

### 4.1 Supervised PEM: A Corporate Finance “World Model”

When implemented on U.S. equities, PEM predicts next period system evolution, including reactions to current period managerial actions. In other words, environment module predicts  $X_{t+1}$  given  $X_t$  and  $u_t$ . As discussed earlier, a deep neural network representation could easily fit the training sample, but is not guaranteed to produce sensible out-of-sample counterfactuals, especially in the presence of distributional shifts. Let us discuss this in turn.

Suppose there is a rarely observed — or an unseen — event in the training sample, which is the general focus of event studies and causal identification strategies. One can use a dummy variable to denote whether this event occurs or not. In the training sample,

---

<sup>11</sup>Red Cloud is a subscription-based Infrastructure as a Service cloud that provides root access to virtual servers and storage on-demand under Cornell University Center for Advanced Computing ([link](#)).

this variable is almost always labeled 0 and only rarely 1, which makes it a phenomenon too “weak” to pick up. In this case, the ambiguity for our test sample is infinite, i.e., any models trained on our training set could provide unconstrained estimations on the treatment effect. To address this concern of ambiguity, we follow the idea of boosting and train a set of auxiliary neural networks to estimate the ambiguity by their dispersion in predictions.<sup>12</sup>

Acting as a metric of ambiguity, boosting error differentiates scenarios where all models coincide in predictions versus ones where model predictions are highly dispersed. By punishing the boosting error in the training of our decision-making module later on, the ambiguity averse agent avoids high-ambiguity actions which we have little reference from prior empirical data, and only explores the action domain with low ambiguity. Similar ideas are also seen in computer science literature related to offline reinforcement learning, such as Kidambi et al. (2020) and Yu et al. (2020).

## 4.2 Predicting Firm Fundamentals and Market Reactions

To train PEM, we use standard mean-squared error (MSE) as our loss function. We standardize input variables to be in the range  $[-1, 1]$  to fit the  $\tanh(\cdot)$  function in the neural networks, rather than normalizing them to be with zero mean and unit variance.

To better understand the economic implication of training loss, we convert the technical loss to pseudo  $R^2$  by dependent variables. The expression of pseudo  $R^2$  for the  $d^{th}$  state for a given sample  $S_d$  with corresponding predictions  $\hat{S}_d$  is:

$$\text{Pseudo } R^2(S_d, \hat{S}_d) = 1 - \frac{\sum_{X_d \in S_d} (X_d - \hat{X}_d)^2}{\sum_{X_d \in S_d} (X_d - \bar{X}_d)^2}, \quad (7)$$

where  $X_d \in S_d$  is a sample of the  $d^{th}$  dimension of state vector, and  $\hat{X}_d \in \hat{S}_d$  is its corresponding prediction from PEM;  $\bar{X}_d$  is the initial training set (i.e., sample before 1991:Q4) average of  $X_d$ . It is not surprising that our environment module gains very high pseudo  $R^2$  for most variables, as decision variables are used as network inputs. The pseudo  $R^2$  for stock return growth is not as high as other variables, since it is driven by investors’ demand — which incorporates factors like sentiment and liquidity — that is out of managers’ control, at least before the control decisions are revealed to the investors in subsequent quarters.

---

<sup>12</sup>The main difference between our boosting method and commonly used algorithms such as AdaBoost lies in how we use the auxiliary predictors. These auxiliary networks share a similar structure with our neural network but with different initialization. Traditional applications combine all the auxiliary predictors to be a stronger predictor, while we use auxiliary predictors to jointly form a natural gauge of model ambiguity.

Table 2: Pseudo  $R^2$  for the Growth of State Variables Without and With Control Information

State Variable	Ignoring Control Information		Including Control Information	
	Training $R^2$	Test $R^2$	Training $R^2$	Test $R^2$
Total Assets	-4.09%	-8.15%	55.44%	62.56%
Current Assets	-3.58%	-7.10%	44.49%	51.21%
Sales	29.54%	28.68%	31.33%	30.88%
Payables	21.46%	24.43%	24.40%	27.64%
COGS	25.68%	26.76%	27.00%	28.56%
InterestExp	73.26%	77.17%	73.36%	77.28%
Inventory	12.78%	13.71%	17.04%	18.92%
CurrentLiability	8.88%	7.72%	21.89%	22.69%
Receivables	17.52%	18.77%	21.59%	23.20%
Revenue	29.51%	28.59%	31.31%	30.80%
MarketCap	1.32%	-3.33%	9.32%	7.07%
EnterpriseValue	-0.97%	-5.73%	14.61%	13.14%
Volume	12.81%	16.53%	15.77%	20.75%
Return	47.90%	45.27%	50.04%	48.19%

Note: This table documents the pseudo  $R^2$  for the growth of state variables without and with the current control ( $u_t$ ). PEM intends to predict future firm-specific state variables given current state and control. Since the current control vector is private information (for instance the manager’s planning on leverage for the next period), it is not observable by outsiders such as investors or econometricians. In predicting the future, we first ignore all controls, by setting them to be 0, calculate PEM predictions, and then the pseudo  $R^2$  (ignoring control) for each firm-specific state variable. Then, we incorporate control variables as inputs, calculate PEM predictions, and then the pseudo  $R^2$  (with control) for each firm-specific state variable.

Table 2 shows the pseudo  $R^2$  of future fundamentals ignoring or with managerial decision information for training sets (in-sample) and test sets (out-of-sample). With managerial decision information, we exploit the full capacity of PEM to generate counterfactuals for future firm-level states by feeding in current state variables as well as managerial decisions to be made. To evaluate the importance of managerial planning in predictions, we ignore the managerial decisions by setting them to be their corresponding last-quart values — as if firms stick with their latest decisions — and generate counterfactuals.

By comparing pseudo  $R^2$ s of future fundamentals ignoring or with managerial planning, we are able to identify pivotal state variables that managerial decisions tend to influence more. For outcomes such as gross revenue or net income, whether we consider the managerial action does not matter much in predictions. However, managerial actions seem to have a significant impact on state variables such as total asset and market capitalization. This reveals that some corporate outcomes are influenced by managerial actions while others are not, at least not immediately. We consider both the case with and without controls because



from outside investors’ perspective, managerial controls are only disclosed with a delay and constitute “insider” information that outsiders may not use for trading. Another way to interpret this “controllability” result is to look at the outcome variables less affected by the manager. For example, growth in trading volume is likely driven by traders in the secondary market rather than directly by the manager of the firm.

Table 3 reports PEM’s heterogeneous performance (measured by OOS MSE) for each state variable by its cross section ranking. We report results for three subsets of cross sections: pre-dot com bubble, between the dot com bubble and the great financial crisis (GFC), and post-GFC, where cutoff points of these two events are based on the NBER recession end months. For each state variable, we rank firms in each cross section as “low” and “high” by comparing the state variable level to its cross-sectional median, and calculate average and standard deviation of PEM MSE within each subsample. For most firm-level state variables, the lower half generally has higher prediction error, with the exceptions of book value of current asset, interest paid and book value of current liability. For macro state variables, the dot com to GFC period is the hardest to predict; the same applies to our two objectives (market capitalization and enterprise value). Table 4 reports PEM’s heterogeneous performance by sector, defined as the first digit of North American Industry Classification System (NAICS) code. Trade and Transportation sector (“naics1” = 4), Education and Healthcare (“naics1” = 6), and Other Services sector (“naics1” = 8) always have the lowest prediction error across three episodes. Table 5 reports PEM’s performance by book-to-market decile, where 1 indicates the lowest decile and 10 indicates the highest decile. The best performers with regard to prediction error of PEM oscillate among middle deciles and the distribution across deciles exhibits a U-shape MSE, while the bottom decile always has the worst performance.

Table 6 results reveal heterogeneous ambiguity for each state variable (measured by boosting error defined in the previous section), following the same logic. For most firm-level state variables, the lower half generally has lower ambiguity. The dot com to GFC period usually has the highest ambiguity, and the pre dot com period has the lowest in general. In Table 7, we do the same exercise for sector. Manufacturing (“naics1” = 3), Education and Healthcare (“naics1” = 6) and Other Services (“naics1” = 8) achieve the lowest ambiguity across three episodes. Mining and Construction (“naics1” = 2) has low ambiguity in the pre dot com period, while the ambiguity for Education and Healthcare sector (“naics1” = 6) gradually goes up. Table 8 shows that the highest three deciles of book-to-market has the lowest ambiguity over the three episodes.

Even though PEM is high-dimensional, we can analyze low dimensional action combi-

nations, which are the focus of conventional reduced-form models. Appendix D uses firm recapitalization as an example to showcase how to use PEM to analyze low-dimensional policy counterfactuals.

### 4.3 Predicting the Evolution of Macroeconomic States

Macroeconomic state variables are distinct from firm-specific fundamentals or stock market reactions. These variables operate on a broader scale, influence firm fundamentals, and are largely insulated from the direct influence of individual firm decisions. For example, changes in an individual firm’s leverage — whether it issues more debt or equity — are unlikely to substantially affect macroeconomic indicators like aggregate non-financial leverage or credit risk. These macroeconomic variables tend to be co-integrated and exhibit minimal cross-sectional variability, making them difficult to predict with precision.

In this context, the National Financial Conditions Index (NFCI) subindices, including risk, credit, financial leverage, and non-financial leverage, provide an important summary of macroeconomic conditions. However, predicting these macroeconomic variables poses a significant challenge due to their overlapping information content and their strong interrelationships. To investigate whether neural networks can overcome this, we compare the predictive performance of neural networks and vector autoregressive (VAR) models, which have long been considered a gold standard in time series forecasting for macroeconomic variables.

Our neural network approach achieves much better performance (measured by out-of-sample pseudo  $R^2$ s) in all four dimensions, especially for risk and credit where VAR fails to attain positive pseudo  $R^2$ s while our neural network gets 11.4% and 27.3% out-of-sample pseudo  $R^2$ s respectively. We detail the time series results in Appendix E.

## 5 Optimal Policy Recommendations to Managers

### 5.1 Reinforcement Learning and the Decision-Making Module

The decision-making module, built on the predictive environment module, intends to seek optimal managerial decision-making processes given a certain utility function (i.e., objective) of the manager. Managers make decisions contingent on the current state of the firm including the current macroeconomic status, a neural network could projection the decision-making process as a function of current states:  $u_t = g(X_t)$ . The neural network is trained

Table 3: Heterogeneous PEM Performance for System States

Variable		Full Sample		Pre-Dotcom		Dotcom-GFC		Post-GFC	
		Mean	Std	Mean	Std	Mean	Std	Mean	Std
Total Assets	high	1.60%	6.28%	1.82%	6.58%	1.79%	6.70%	1.32%	5.75%
	low	2.95%	9.34%	3.33%	9.89%	2.96%	9.07%	2.66%	9.04%
Current Assets	high	2.76%	8.20%	2.68%	7.62%	2.88%	8.36%	2.76%	8.54%
	low	4.04%	11.33%	4.31%	11.33%	3.78%	10.83%	3.99%	11.62%
Sales	high	3.17%	10.00%	3.33%	10.01%	3.10%	9.94%	3.09%	10.03%
	low	6.92%	16.94%	6.56%	15.42%	6.30%	15.35%	7.56%	18.80%
Payables	high	5.43%	12.57%	6.25%	13.54%	5.74%	13.06%	4.61%	11.40%
	low	6.39%	13.21%	5.22%	11.08%	6.14%	12.65%	7.42%	14.82%
COGS	high	2.75%	10.55%	3.09%	11.64%	2.66%	10.25%	2.53%	9.82%
	low	4.12%	12.40%	4.18%	12.10%	4.01%	12.25%	4.14%	12.70%
InterestExp	high	1.72%	5.59%	1.25%	4.69%	1.79%	6.05%	2.03%	5.90%
	low	1.84%	7.63%	1.45%	6.62%	1.89%	8.18%	2.11%	8.00%
Inventory	high	5.06%	12.85%	5.48%	13.68%	5.39%	13.40%	4.55%	11.82%
	low	4.85%	14.89%	4.82%	14.28%	4.58%	14.40%	5.03%	15.61%
CurrentLiability	high	5.23%	13.24%	5.65%	13.76%	5.46%	13.75%	4.78%	12.50%
	low	6.55%	15.11%	6.52%	14.43%	6.40%	14.87%	6.66%	15.74%
Receivables	high	3.29%	9.37%	3.47%	9.61%	3.41%	9.43%	3.07%	9.15%
	low	6.87%	16.00%	5.92%	14.21%	6.60%	15.24%	7.76%	17.58%
Revenue	high	3.16%	9.99%	3.34%	10.13%	3.07%	9.89%	3.07%	9.94%
	low	6.90%	16.93%	6.56%	15.46%	6.30%	15.39%	7.52%	18.74%
MarketCap	high	7.68%	15.70%	9.92%	18.70%	8.30%	16.43%	5.55%	11.96%
	low	12.53%	21.42%	12.99%	21.63%	13.62%	22.77%	11.55%	20.38%
EnterpriseValue	high	6.06%	13.67%	8.59%	17.98%	6.14%	12.80%	4.05%	9.24%
	low	10.34%	18.96%	11.53%	20.69%	11.55%	20.10%	8.73%	16.65%
Volume	high	4.01%	8.49%	4.89%	8.76%	3.68%	8.21%	3.52%	8.38%
	low	7.74%	14.11%	8.83%	14.05%	8.21%	14.96%	6.65%	13.57%
Return	high	5.05%	10.92%	5.88%	12.14%	5.43%	11.49%	4.19%	9.41%
	low	6.77%	14.92%	7.48%	15.94%	7.51%	15.95%	5.79%	13.37%
MacroRisk	high	6.00%	6.71%	6.42%	6.80%	5.99%	6.54%	4.84%	5.88%
	low	4.91%	5.78%	4.62%	5.33%	4.63%	5.37%	4.97%	6.11%
MacroCredit	high	5.86%	6.64%	6.59%	6.90%	6.26%	6.70%	5.01%	6.10%
	low	4.70%	5.53%	4.74%	5.47%	4.44%	5.18%	4.81%	5.89%
MacroFinLev	high	5.43%	6.28%	5.43%	6.10%	5.53%	6.26%	5.26%	6.22%
	low	5.05%	5.91%	5.64%	6.28%	5.14%	5.81%	4.55%	5.74%
MacroNonfinLev	high	5.75%	6.36%	6.35%	6.74%	5.16%	5.96%	5.06%	6.26%
	low	4.82%	5.83%	4.58%	5.33%	5.47%	6.09%	4.77%	5.75%

Note: This table shows the heterogeneous PEM performance for each state variable. For each cross section, we divide firms in low and high groups for each state variable, and we calculate mean and standard deviation of PEM MSE for the focal state variable within three subsamples: pre dot com, between dot com and GFC, and post GFC.

Table 4: Heterogeneous PEM Performance for Sectors

MSE	full sample		pre-dotcom		dotcom–GFC		post-GFC	
sector	mean	std	mean	std	mean	std	mean	std
agriculture	7.00%	7.03%	7.40%	7.59%	7.04%	6.58%	6.62%	6.74%
mining & construction	6.84%	7.00%	6.20%	6.32%	7.23%	7.20%	7.10%	7.30%
manufacturing	4.95%	5.93%	4.98%	5.77%	5.17%	6.00%	4.78%	6.02%
trade & transportation	4.35%	5.17%	4.60%	5.36%	4.33%	5.09%	4.15%	5.05%
information & professional services	5.14%	5.92%	6.28%	6.76%	5.18%	5.71%	4.24%	5.15%
education & healthcare	4.37%	5.22%	4.99%	5.57%	4.05%	4.76%	4.02%	5.11%
recreation & accommodation	4.76%	5.44%	5.09%	5.22%	4.66%	5.03%	4.52%	5.85%
other services	3.89%	4.66%	4.27%	5.05%	3.52%	4.07%	3.60%	4.39%
public administration	7.00%	7.66%	7.52%	7.99%	7.09%	7.63%	5.05%	6.15%

Note: This table shows the heterogeneous PEM performance for each sector, defined by the first digit of North American Industry Classification System (NAICS) code. For each firm-quarter observation, we calculate average MSE of PEM for state variables to get the average MSE of PEM for each firm-quarter observation, and then for each cross section, we calculate mean and standard deviation of PEM MSE within three subsamples: pre dot com, between dot com and GFC, and post GFC, and we take the average across state variables. The top three sectors (with the lowest MSEs) within each time period is marked in red.

Table 5: Heterogeneous PEM Performance for “Value” Deciles

MSE	full sample		pre-dotcom		dotcom–GFC		post-GFC	
book-to-market decile	mean	std	mean	std	mean	std	mean	std
1	6.45%	6.98%	7.40%	7.34%	6.41%	6.77%	5.72%	6.73%
2	5.68%	6.45%	6.53%	6.94%	5.55%	6.07%	5.08%	6.19%
3	5.19%	6.10%	5.89%	6.51%	5.20%	5.95%	4.64%	5.78%
4	4.75%	5.72%	5.41%	6.13%	4.88%	5.63%	4.16%	5.38%
5	4.68%	5.65%	5.14%	5.91%	4.77%	5.56%	4.27%	5.48%
6	4.68%	5.60%	4.96%	5.72%	4.95%	5.69%	4.31%	5.43%
7	4.72%	5.66%	4.86%	5.52%	4.96%	5.78%	4.47%	5.69%
8	4.82%	5.70%	4.78%	5.37%	5.13%	5.93%	4.67%	5.80%
9	5.13%	5.85%	4.97%	5.63%	5.29%	5.99%	5.15%	5.92%
10	6.03%	6.55%	5.37%	5.98%	6.06%	6.59%	6.50%	6.88%

Note: This table shows the heterogeneous PEM performance for each book-to-market decile. For each firm-quarter observation, we calculate average MSE of PEM for state variables to get the average MSE of PEM for each firm-quarter observation, and then for each cross section, we calculate mean and standard deviation of PEM MSE within three subsamples: pre dot com, between dot com and GFC, and post GFC, and we take the average across state variables. The top three deciles (with the lowest MSEs) within each time period is marked in red.

Table 6: PEM: Heterogeneous Ambiguity for System States

variable		full sample		pre-dotcom		dotcom-GFC		post-GFC	
		mean	std	mean	std	mean	std	mean	std
Total Assets	high	9.35%	4.19%	8.45%	3.85%	9.56%	4.17%	9.93%	4.33%
	low	7.96%	3.84%	7.04%	3.19%	7.78%	3.57%	8.76%	4.24%
Current Assets	high	10.68%	4.74%	8.99%	3.95%	10.85%	4.68%	11.90%	4.93%
	low	9.56%	4.59%	8.07%	3.62%	9.57%	4.31%	10.68%	5.07%
Sales	high	9.15%	4.35%	7.66%	3.48%	9.44%	4.31%	10.13%	4.66%
	low	8.81%	4.85%	7.13%	3.25%	8.79%	4.46%	10.10%	5.63%
Payables	high	9.51%	4.34%	7.84%	3.52%	9.60%	4.19%	10.74%	4.58%
	low	8.65%	4.38%	7.01%	3.23%	8.33%	3.96%	10.07%	4.88%
COGS	high	9.52%	4.30%	8.50%	3.92%	9.90%	4.36%	10.09%	4.40%
	low	9.07%	4.25%	8.07%	3.61%	9.10%	4.02%	9.82%	4.64%
InterestExp	high	7.43%	3.64%	6.11%	2.82%	7.58%	3.65%	8.33%	3.88%
	low	6.96%	3.47%	6.18%	2.85%	7.26%	3.53%	7.40%	3.77%
Inventory	high	8.87%	4.02%	7.51%	3.35%	9.39%	4.11%	9.60%	4.18%
	low	9.10%	4.32%	7.82%	3.57%	9.21%	4.12%	10.00%	4.71%
CurrentLiability	high	11.21%	4.93%	9.49%	4.26%	11.51%	5.06%	12.37%	4.96%
	low	10.31%	5.00%	8.77%	4.03%	10.22%	4.99%	11.52%	5.34%
Receivables	high	9.24%	3.96%	7.99%	3.44%	9.64%	3.95%	9.96%	4.11%
	low	9.08%	4.34%	7.69%	3.39%	9.12%	3.91%	10.10%	4.89%
Revenue	high	9.00%	4.36%	7.52%	3.48%	9.24%	4.31%	10.00%	4.67%
	low	8.69%	4.87%	7.03%	3.30%	8.71%	4.44%	9.95%	5.68%
MarketCap	high	13.33%	6.97%	11.10%	5.06%	15.98%	8.97%	13.50%	6.31%
	low	12.96%	8.25%	9.56%	4.60%	15.10%	9.74%	14.24%	8.63%
EnterpriseValue	high	12.39%	6.32%	10.08%	4.63%	14.63%	7.83%	12.86%	5.87%
	low	11.81%	7.66%	8.57%	4.13%	13.69%	8.70%	13.15%	8.26%
Volume	high	11.07%	4.64%	9.77%	4.07%	11.98%	5.18%	11.55%	4.51%
	low	9.50%	4.46%	7.98%	3.48%	10.40%	5.21%	10.10%	4.33%
Return	high	11.13%	5.86%	8.93%	4.24%	12.94%	7.11%	11.74%	5.60%
	low	11.62%	6.57%	8.92%	4.23%	13.57%	8.11%	12.53%	6.36%
MacroRisk	high	10.82%	4.17%	8.85%	3.28%	11.45%	4.33%	10.68%	3.93%
	low	9.47%	3.72%	7.53%	2.60%	9.44%	3.51%	10.82%	4.04%
MacroCredit	high	10.55%	4.07%	8.95%	3.31%	11.18%	4.43%	10.81%	4.03%
	low	9.29%	3.66%	7.64%	2.68%	9.79%	3.59%	10.70%	3.93%
MacroFinLev	high	10.46%	4.12%	8.14%	2.99%	11.33%	4.50%	10.88%	4.01%
	low	9.37%	3.65%	8.26%	3.08%	9.69%	3.49%	10.62%	3.94%
MacroNonfinLev	high	9.78%	3.85%	8.81%	3.27%	11.49%	4.19%	10.59%	4.04%
	low	9.91%	3.94%	7.48%	2.56%	9.56%	3.75%	10.89%	3.93%

Note: This table shows the heterogeneous ambiguity for each state variable measured by the greatest different in predictions among PEM and its auxiliary models for that specific state variable. For each cross section, we divide firms in low and high groups for each state variable, and we calculate mean and standard deviation of ambiguity for the focal state variable within three subsamples: pre dot com, between dot com and GFC, and post GFC.

Table 7: PEM: Heterogeneous Ambiguity for Sectors

Ambiguity	full sample		pre-dotcom		dotcom-GFC		post-GFC	
sector	mean	std	mean	std	mean	std	mean	std
agriculture	10.77%	4.04%	9.59%	3.48%	11.55%	4.33%	11.40%	4.08%
mining & construction	10.55%	4.02%	8.45%	2.76%	11.72%	4.39%	11.44%	4.02%
manufacturing	9.42%	3.94%	7.70%	2.90%	10.10%	4.04%	10.49%	4.15%
trade & transportation	9.63%	3.76%	7.73%	2.65%	10.35%	4.04%	10.84%	3.77%
information & professional services	10.07%	3.70%	8.82%	3.28%	10.53%	3.89%	10.73%	3.65%
education & healthcare	9.21%	3.35%	7.88%	2.70%	9.81%	3.49%	9.99%	3.42%
recreation & accommodation	9.81%	3.65%	8.35%	2.78%	10.41%	3.81%	10.79%	3.82%
other services	8.62%	3.22%	7.22%	2.55%	9.42%	3.22%	10.19%	3.24%
public administration	9.81%	3.83%	8.76%	3.24%	10.66%	4.26%	11.84%	3.67%

Note: This table shows the heterogeneous ambiguity for each sector, defined by the first digit of North American Industry Classification System (NAICS) code. The ambiguity is measured by the greatest different in predictions among PEM and its auxiliary models for each state variable. For each firm-quarter observation, we calculate average ambiguity of PEM for state variables to get the average ambiguity of PEM for each firm-quarter observation, and then for each cross section, we calculate mean and standard deviation of PEM ambiguity within three subsamples: pre dot com, between dot com and GFC, and post GFC, and we take the average across state variables. The three sectors with the lowest ambiguity within each time period is marked in red.

Table 8: PEM: Heterogeneous Ambiguity for “Value” Deciles

Ambiguity	full sample		pre-dotcom		dotcom-GFC		post-GFC	
book-to-market decile	mean	std	mean	std	mean	std	mean	std
1	11.50%	4.30%	9.86%	3.40%	11.99%	4.38%	12.51%	4.50%
2	10.56%	3.98%	8.97%	3.09%	11.11%	4.10%	11.50%	4.14%
3	10.17%	3.92%	8.55%	3.08%	10.78%	4.13%	11.07%	3.98%
4	9.90%	3.78%	8.43%	3.05%	10.37%	3.94%	10.79%	3.86%
5	9.79%	3.74%	8.21%	2.99%	10.48%	3.96%	10.59%	3.74%
6	9.55%	3.69%	8.01%	2.92%	10.29%	3.98%	10.30%	3.67%
7	9.40%	3.66%	7.83%	2.86%	10.18%	3.93%	10.11%	3.67%
8	9.36%	3.74%	7.69%	2.81%	10.07%	3.95%	10.20%	3.81%
9	9.27%	3.73%	7.40%	2.59%	9.85%	3.95%	10.29%	3.80%
10	8.96%	3.75%	6.94%	2.43%	9.38%	3.81%	10.16%	3.89%

Note: This table shows the heterogeneous ambiguity for each book-to-market decile. The ambiguity is measured by the greatest different in predictions among PEM and its auxiliary models for each state variable. For each firm-quarter observation, we calculate average ambiguity of PEM for state variables to get the average ambiguity of PEM for each firm-quarter observation, and then for each cross section, we calculate mean and standard deviation of PEM ambiguity within three subsamples: pre dot com, between dot com and GFC, and post GFC, and we take the average across state variables. The three deciles with the lowest ambiguity within each time period is marked in red.

to maximize the expected objective value predicted by PEM. To mitigate the uncertainty of predictive environment module under unusual decisions, we punish directly the boosting error defined in Equation (4) and train a model based on the discrete version of utility function in Equation (5).

We start from a randomly initialized neural networks whose inputs are current states and outputs are current decisions to be made. We optimize an exogenously specified objective function using RL. The value of the objective function is given by the counterfactual statements from PEM. For each epoch, the decision-making module first gives out a set of decisions to be made, then together with state variables, these decisions are fed into PEM to generate counterfactual market reactions in the next quarter. Finally, the difference between counterfactual market capitalization (alternatively, enterprise value) and current market capitalization (enterprise value) is treated as the objective function value for this trait. After the feed-forward process is complete, backward propagation takes place and the parameters of the DMM are updated, given the loss function to be negative objective value.

We consider a set of reasonable objective functions and also a spectrum of horizons for managers’ objectives. Specifically, we consider both short-term (1 quarter) and long-term (2 years) market capitalization and enterprise value as our objectives. We ignore agency problems and take managers as maximizing shareholders’ value or overall investors’ value. The DDRC framework is flexible to allow other objectives and these two are useful benchmarks to consider.

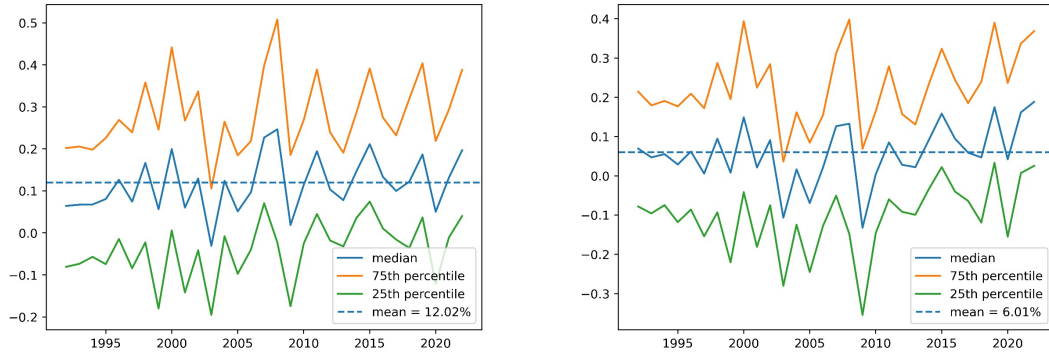
## 5.2 Performance Under Various Managerial Objectives

**Long-term objectives maximizing valuation growth in 2 years.** To outperform in the longer run, AlphaManager needs to have a longer horizon in its objective. We therefore examine how AlphaManager performs if the object is to maximize increases in the market capitalization or enterprise value over the subsequent 8 quarters instead of just over the next quarter. For long-term model with capitalization growth as the objective function, the average quarterly return is 8.73% while the first quarter performance is on average only 3.09%. When using enterprise value growth as the objective function, these two numbers become 4.43% and 1.28%, respectively.

**Short-term objectives centering on next quarter valuation growth.** We then consider a somewhat myopic objective of increasing market equity valuation over the next

quarter. We estimate the corresponding short-term ambiguity constraint parameters under optimal long-term AM agents to make the results comparable. The overall average out-performance is 12.02% quarterly in the short-term when using market capitalization growth as the objective. When the enterprise value growth serves as the objective, the out-performance becomes 6.01%. The cross-sectional dispersion is non-trivial, as seen in the 75th and 25th percentile firms in Figure 2.

Figure 2: Out-performance of AlphaManager with Short-term Objectives



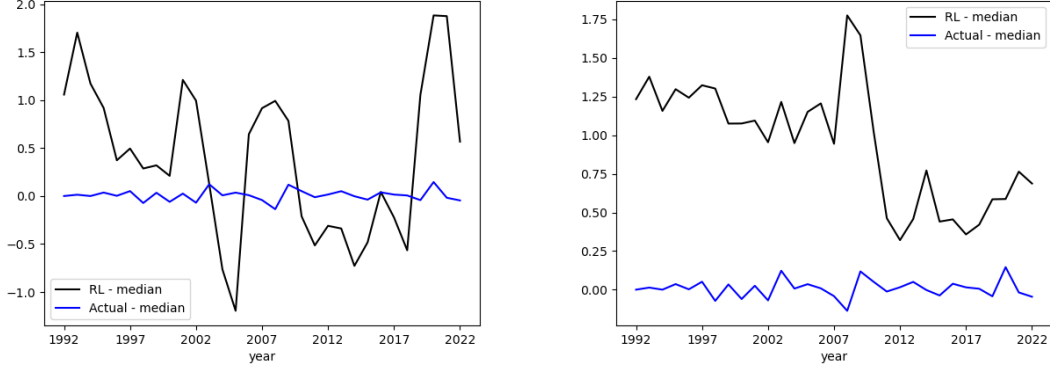
Note: This figure shows the out-performance of AlphaManager trained with short-term market cap growth (left) and enterprise value growth (right) as objective functions. The solid blue line shows the median performance, the orange line and the green line show the 75th and 25th percentiles respectively, and the dotted blue line shows the median for each year. The y-axis is in percentage of the change in objective for the next quarter predicted by the PEM.

Figure 3 depicts what myopic behavior (a focus on next quarter outcomes) implies for the long-term (the next two-year) performance of the firms. This relationship is tracked over 30 years of data. *Ex-ante*, several outcomes could appear. On the one hand, one could argue that managers historically may not be myopic, and as such, historical actions may generate better long-term performance if short-termism hurts firm fundamentals in the intermediate or long run. On the other hand, even with potentially longer-term goals, managers could have taken very suboptimal actions historically. The 1-quarter-market-cap-focused myopic AlphaManager can end up perform much better. In the data, we observe that though AlphaManager outperforms under the specified objective, over the intermediate term (i.e., for a different objective with longer horizon) it underperforms. Algorithmic predictions do not capture cross-sectional “risk” (in-model volatility), so the algorithmic predictions (black lines) have much smaller variations.

Figure 4 shows the term structure for short-term and long-term AM under the objective



Figure 3: Long-Term Performance Under Short-Termism



Note: This figure shows the long-term performance of AlphaManager trained with short-term market cap growth (left) and enterprise value growth (right) as objective functions. The solid blue line shows the median average performance of actual managerial decisions, while the solid black line shows the median performance of actions suggested by the RL module.

of market cap growth. In the left panel, the short-term AM is out performing the long-term AM in longer terms' rewards, but this out performance is at the cost of high ambiguity. In the right panel, we adjust reward (i.e., market cap growth) using the formula:

$$\text{Adj. Reward}_t = \frac{1 + \text{Reward}_t}{\sqrt{\max\{1, \frac{\text{BoostingError}_t}{\text{BoostingError}_{t_0}}\}}} - 1, \quad (8)$$

where  $\text{Reward}_t$  is the reward that AM achieves under the PEM at time  $t$ ,  $\text{BoostingError}_t$  is the boosting error for the AM-suggested managerial decisions at time  $t$ , and  $\text{BoostingError}_{t_0}$  is the benchmark boosting error at time  $t_0$ . When the ambiguity is as low as the benchmark ambiguity, there is no punishment on the reward; the adjusted reward is decreasing in ambiguity of AM. After the adjustment, even though the short-term AM outperforms in the short-term, the long-run performance is worse than the long-term AM, and further into the future, the gap between adjusted rewards from short-term AM and long-term AM is also wider.

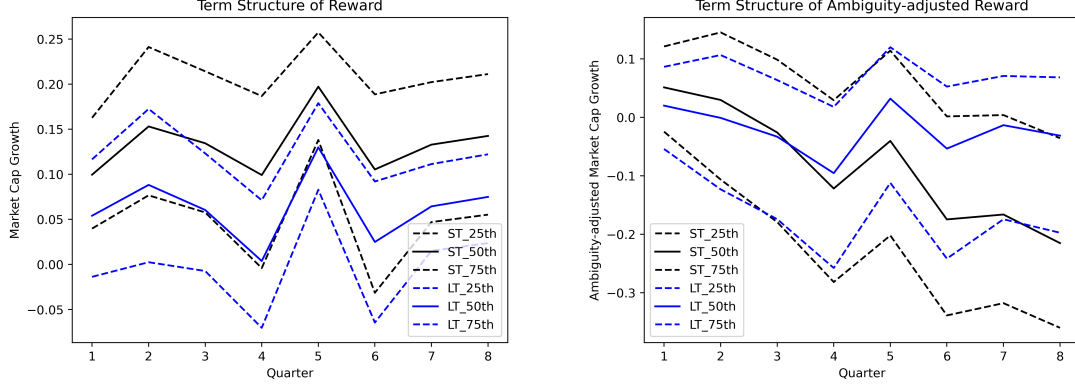
Table 9 reports heterogeneous predicted performance under the objective of 2-year enterprise value growth for the RL-recommended decisions by each state's cross-sectional ranking. It is worth noting that AM's impressive performance is driven by small and illiquid firms during high-risk episodes. The pre-dot com period enjoys the highest performance of AM. Table 10 documents the heterogeneous performance of RL. Manufacturing sector has high performance across three episodes. Table 8 shows the heterogeneity by book-to-market deciles. The best performers are gradually shifting down as time goes by.

Table 9: Heterogeneous Performance of AlphaManager Across Firms

variable		full sample		pre-dotcom		dotcom-GFC		post-GFC	
		mean	std	mean	std	mean	std	mean	std
Total Assets	high	7.13%	12.91%	6.67%	10.42%	3.61%	16.67%	10.27%	11.01%
	low	9.35%	4.15%	7.54%	13.18%	-1.19%	21.68%	11.17%	18.34%
Current Assets	high	6.54%	14.13%	6.56%	11.28%	2.15%	17.91%	10.15%	12.39%
	low	10.14%	4.56%	8.06%	13.21%	-0.98%	22.18%	11.59%	18.80%
Sales	high	6.78%	12.98%	6.67%	10.16%	2.84%	16.68%	10.03%	11.26%
	low	8.88%	4.20%	7.68%	13.75%	-1.25%	22.53%	11.56%	19.01%
Payables	high	6.95%	12.99%	6.45%	10.24%	3.31%	16.81%	10.26%	11.33%
	low	9.40%	4.31%	7.72%	13.32%	-0.90%	21.59%	11.16%	18.07%
COGS	high	6.60%	13.16%	6.53%	10.25%	2.59%	16.85%	9.84%	11.52%
	low	9.26%	4.20%	7.82%	13.69%	-1.03%	22.50%	11.81%	18.90%
InterestExp	high	6.74%	15.01%	7.59%	11.89%	1.71%	19.12%	9.88%	13.59%
	low	7.12%	3.51%	6.73%	12.40%	0.10%	20.43%	11.47%	16.79%
Inventory	high	6.37%	13.89%	6.25%	10.45%	1.93%	17.72%	10.14%	12.56%
	low	8.62%	3.96%	8.25%	13.72%	-0.39%	21.98%	11.47%	18.18%
CurrentLiability	high	6.35%	13.88%	6.43%	10.89%	1.96%	17.65%	9.86%	12.26%
	low	10.77%	4.86%	8.22%	13.60%	-0.69%	22.47%	11.97%	18.87%
Receivables	high	6.90%	12.90%	6.63%	10.09%	2.93%	16.82%	10.19%	11.09%
	low	9.16%	3.92%	7.57%	13.34%	-0.56%	21.55%	11.20%	18.09%
Revenue	high	6.64%	13.09%	6.67%	10.37%	2.60%	16.73%	9.94%	11.39%
	low	8.63%	4.19%	7.81%	13.97%	-1.30%	22.94%	11.73%	19.17%
MarketCap	high	6.43%	14.13%	6.08%	11.78%	2.32%	17.96%	10.06%	11.78%
	low	13.14%	6.94%	8.19%	12.37%	-0.62%	21.48%	11.50%	18.53%
EnterpriseValue	high	6.85%	13.51%	6.34%	11.24%	3.10%	17.27%	10.21%	11.32%
	low	12.29%	6.23%	7.81%	12.72%	-0.94%	21.52%	11.26%	18.39%
Volume	high	6.22%	16.26%	6.02%	12.92%	1.25%	19.87%	10.55%	15.04%
	low	10.74%	4.57%	8.52%	11.03%	0.38%	19.80%	11.00%	15.94%
Return	high	5.44%	15.88%	6.04%	12.02%	-0.92%	19.56%	9.98%	14.42%
	low	10.81%	5.76%	8.22%	12.16%	2.63%	19.96%	11.53%	16.40%
MacroRisk	high	6.13%	14.24%	6.36%	12.97%	-0.39%	22.55%	10.36%	14.79%
	low	4.53%	14.18%	8.02%	11.18%	2.11%	16.58%	11.28%	16.30%
MacroCredit	high	6.87%	14.10%	6.74%	12.94%	0.54%	21.97%	13.37%	15.83%
	low	4.02%	14.20%	7.51%	11.52%	1.15%	17.62%	8.92%	14.94%
MacroFinLev	high	5.05%	15.23%	6.37%	12.10%	-0.17%	23.13%	12.02%	15.79%
	low	5.38%	13.36%	7.87%	12.13%	1.77%	16.30%	9.60%	15.07%
MacroNonfinLev	high	5.50%	13.95%	6.57%	12.90%	0.43%	23.55%	12.17%	16.63%
	low	5.04%	14.42%	7.88%	11.19%	1.23%	15.89%	8.92%	13.59%

Note: This table shows the heterogeneous performance under RL-recommended decisions by each state variable. For each cross section, we divide firms in low and high groups for each state variable, and we calculate mean and standard deviation of the performance of RL-recommended decisions within three subsamples: pre dot com, between dot com and GFC, and post GFC.

Figure 4: Term Structure of Reward and Ambiguity-adjusted Reward



Note: This figure shows the term structure for AlphaManager trained with short-term market cap growth with reward (left) and ambiguity-adjusted reward (right). The ambiguity adjustment follows Eq.(8). The solid blue lines show the median performance of long-term AM, while the solid black lines show the median performance of short-term AM. Dashed lines show the 25th and 75th percentiles.

Table 10: Heterogeneous Performance of AlphaManager Across Industries

RL performance (quarterly avg)	full sample		pre-dotcom		dotcom-GFC		post-GFC	
	mean	std	mean	std	mean	std	mean	std
agriculture	8.08%	15.11%	7.83%	12.73%	2.83%	16.35%	11.30%	15.50%
mining & construction	7.87%	16.41%	6.97%	11.83%	0.40%	19.85%	11.91%	16.05%
manufacturing	7.94%	17.66%	6.76%	12.23%	0.71%	20.87%	13.58%	17.63%
trade & transportation	8.76%	14.04%	7.27%	10.45%	4.12%	16.60%	12.46%	14.23%
information & professional services	7.55%	17.03%	7.86%	12.94%	-0.23%	19.64%	12.39%	16.13%
education & healthcare	7.40%	16.24%	7.58%	11.84%	-0.04%	19.05%	11.76%	16.09%
recreation & accommodation	7.67%	13.95%	8.03%	11.02%	2.13%	16.42%	10.68%	13.76%
other services	6.45%	13.34%	6.35%	10.94%	1.12%	16.34%	11.03%	12.44%
public administration	6.17%	16.39%	7.49%	13.75%	-0.04%	19.18%	13.18%	15.21%

Note: This table shows the heterogeneous performance under RL-recommended decisions by each sector, defined by the first digit of North American Industry Classification System (NAICS) code. For each cross section, we calculate mean and standard deviation of the performance of RL-recommended decisions within three subsamples: pre dot com, between dot com and GFC, and post GFC. The top three sectors (with the highest average rewards) within each time period is marked in red.

Table 11: Heterogeneous Performance of AlphaManager Across Value Deciles

RL performance (quarterly avg)	full sample		pre-dotcom		dotcom-GFC		post-GFC	
book-to-market decile	mean	std	mean	std	mean	std	mean	std
1	3.98%	19.72%	4.22%	15.90%	-3.89%	23.80%	9.99%	17.76%
2	5.37%	17.42%	5.63%	13.74%	-1.00%	20.94%	10.27%	16.36%
3	5.87%	16.05%	6.28%	11.79%	-0.26%	19.94%	10.42%	15.07%
4	6.53%	15.47%	7.28%	11.31%	0.43%	18.94%	10.67%	14.83%
5	7.17%	15.04%	7.81%	10.90%	1.73%	18.68%	10.90%	14.42%
6	7.72%	14.69%	8.40%	10.81%	2.54%	18.10%	11.32%	14.05%
7	8.10%	14.91%	8.67%	10.72%	3.46%	18.16%	11.48%	14.84%
8	8.15%	15.10%	8.56%	10.63%	3.36%	18.55%	11.50%	15.24%
9	8.05%	15.54%	9.06%	10.18%	2.88%	19.60%	10.67%	15.73%
10	7.88%	15.38%	8.88%	10.72%	2.20%	17.65%	10.51%	16.63%

Note: This table shows the heterogeneous performance under RL-recommended decisions by each book-to-market decile. For each cross section, we calculate mean and standard deviation of the performance of RL-recommended decisions within three subsamples: pre dot com, between dot com and GFC, and post GFC. The top three deciles (with the highest average rewards) within each time period is marked in red.

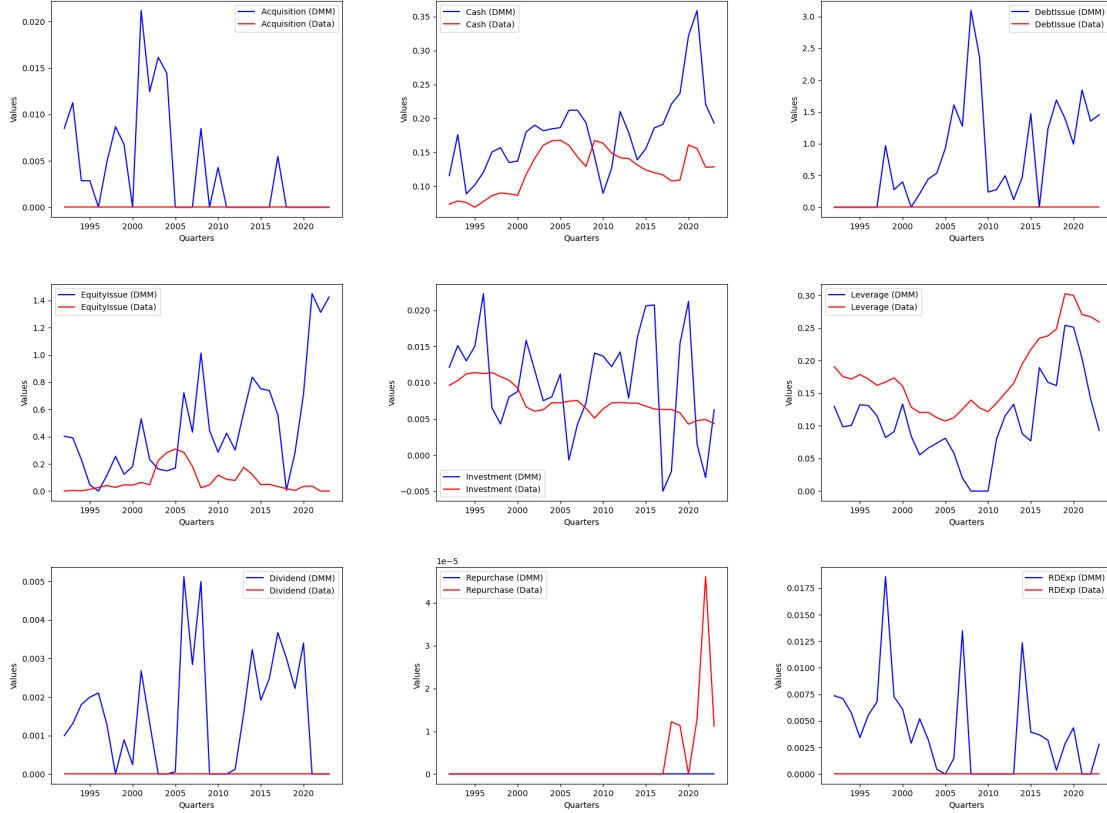
The analysis shows that AlphaManager again outperforms historical actions and there is heterogeneity in the cross-section. The average outperformance is much lower, reflecting that over the longer horizon, we just cannot predict much. Notably, these patterns also emerge in papers based on CEO surveys (e.g., Graham, 2022).

### 5.3 AI-Recommended Actions

We analyze the managerial actions recommended by AlphaManager under the four specified objective functions. Our analysis compares actual historical decisions (e.g., leverage levels) with model-generated recommendations for changes in these variables (e.g., leverage growth rates). When model-implied decisions exceed feasible bounds (such as negative leverage), we apply boundary constraints (setting them to zero, for instance).

Focusing first on the objective of next-quarter market capitalization maximization, as shown in Figure 5, AlphaManager’s recommendations diverge from historical practices in several key aspects: (1) pursuing more aggressive acquisition strategies, (2) maintaining higher cash reserves, (3) adopting a more equity-heavy capital structure while reducing debt (thereby decreasing leverage), (4) increasing dividend payouts, and (5) boosting investment expenditures, particularly in R&D. The model further suggests managers should tolerate greater investment variability and implement more share repurchases during economic downturns.

Figure 5: Empirical and Optimal Decisions under 1-QTR Market Cap Growth

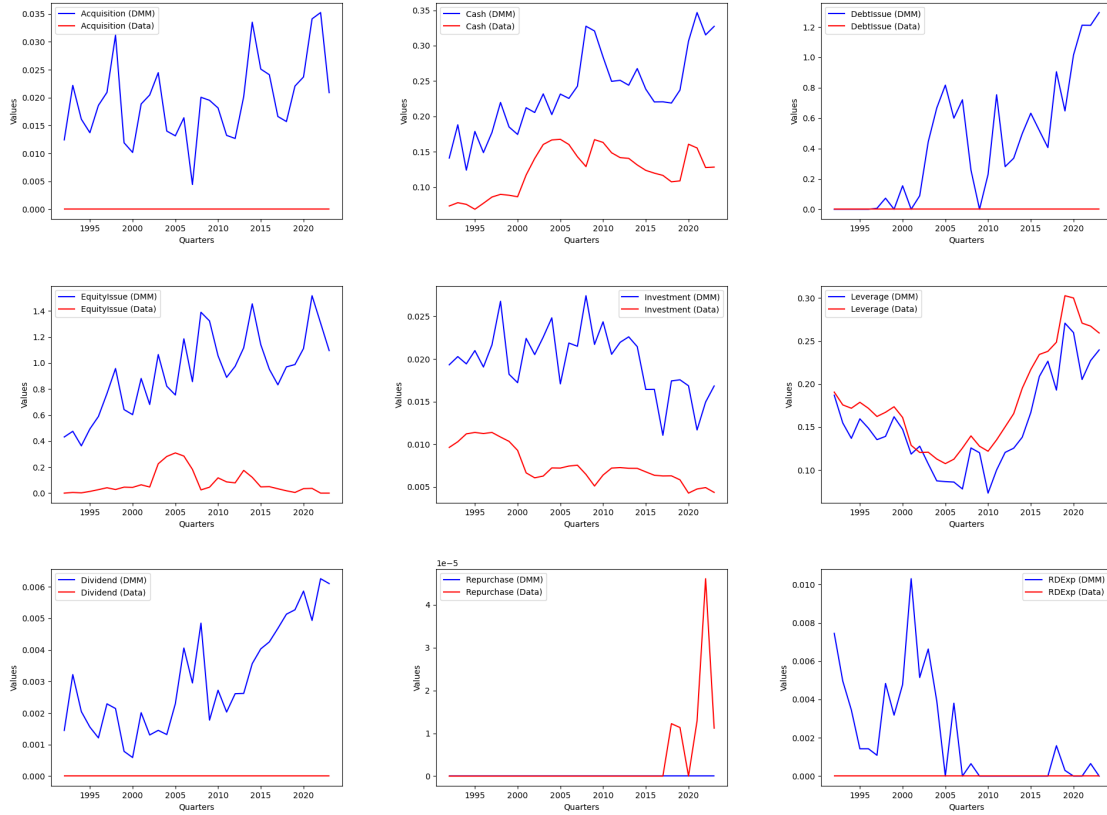


Note: This figure shows the empirical (red) and optimal decisions (blue) suggested by AM under 1-quarter market cap growth as the objective to maximize. We plot median to represent cross-sectional variation. The decisions are acquisitions (Faqc\_fund), cash holdings (Fcash\_hold), debt issuance (Fdebt\_issue\_log), equity issuance (Fequity\_issue\_log), investment ratio (Finv\_ratio2), leverage (Flev\_market), dividend (Fpay-out\_div), repurchase (Fpayout\_repurchase), and R&D expenses (Frd\_exp), with formal definitions detailed in Appendix A.

In Figure 6, when optimizing for short-term enterprise value maximization, AlphaManager generates recommendations largely consistent with the market capitalization objective, with one notable exception: our DMM suggests maintaining or even moderately increasing debt issuance rather than reducing it. This recommendation reflects the general benefits of leverage, particularly during the pre-financial crisis period, where higher leverage appears to have been value-enhancing.

In Figure 7, under the objective of maximizing long-term market capitalization growth (over eight quarters), AlphaManager recommends several strategic adjustments relative to historical practices: (1) pursuing more acquisitions, particularly for small firms; (2) maintaining higher cash reserves; (3) moderately increasing debt issuance while also issuing more

Figure 6: Empirical and Optimal Decisions under 1-QTR Enterprise Value Growth

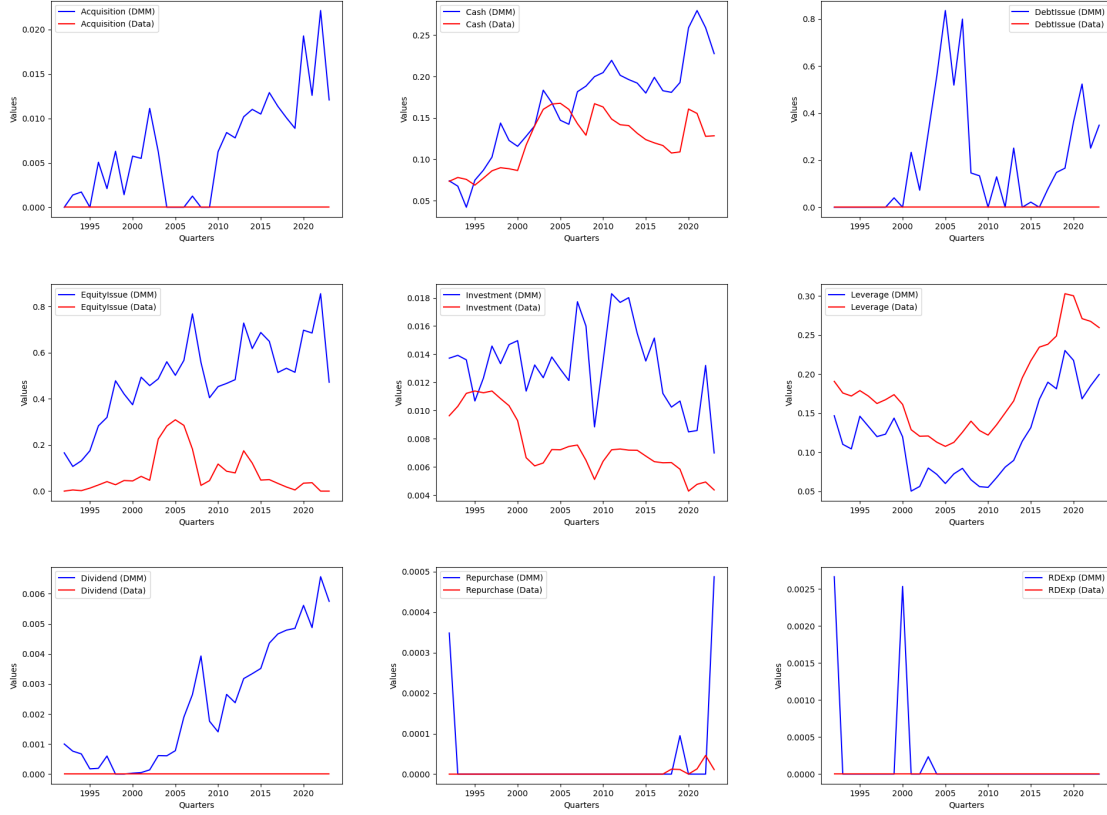


Note: This figure shows the empirical (red) and optimal decisions (blue) suggested by AM under 1-quarter enterprise value growth as the objective to maximize. We plot median to represent cross-sectional variation. The decisions are acquisitions (Faqc\_fund), cash holdings (Fcash\_hold), debt issuance (Fdebt\_issue\_log), equity issuance (Fequity\_issue\_log), investment ratio (Finv\_ratio2), leverage (Flev\_market), dividend (Fpayout\_div), repurchase (Fpayout\_repurchase), and R&D expenses (Frd\_exp), with formal definitions detailed in Appendix A.

equity (with the latter being especially beneficial for small firms); and (4) reallocating investment by reducing R&D expenditures while increasing capital expenditures, dividend payouts, and share repurchases. Additionally, the model suggests a net reduction in leverage despite the slight increase in debt, as the equity issuance more than offsets it.

Finally, in Figure 8, for the objective of maximizing long-term enterprise value (over eight quarters), AlphaManager recommends pursuing more aggressive acquisition strategies while maintaining moderately higher cash reserves. The model suggests increasing both debt and equity issuance, resulting in slightly higher leverage during normal economic conditions but reduced leverage following the onset of the GFC. Additional recommendations include moderately decreasing dividend payouts, expanding share repurchase programs, and reducing

Figure 7: Empirical and Optimal Decisions under 8-QTR Market Cap Growth



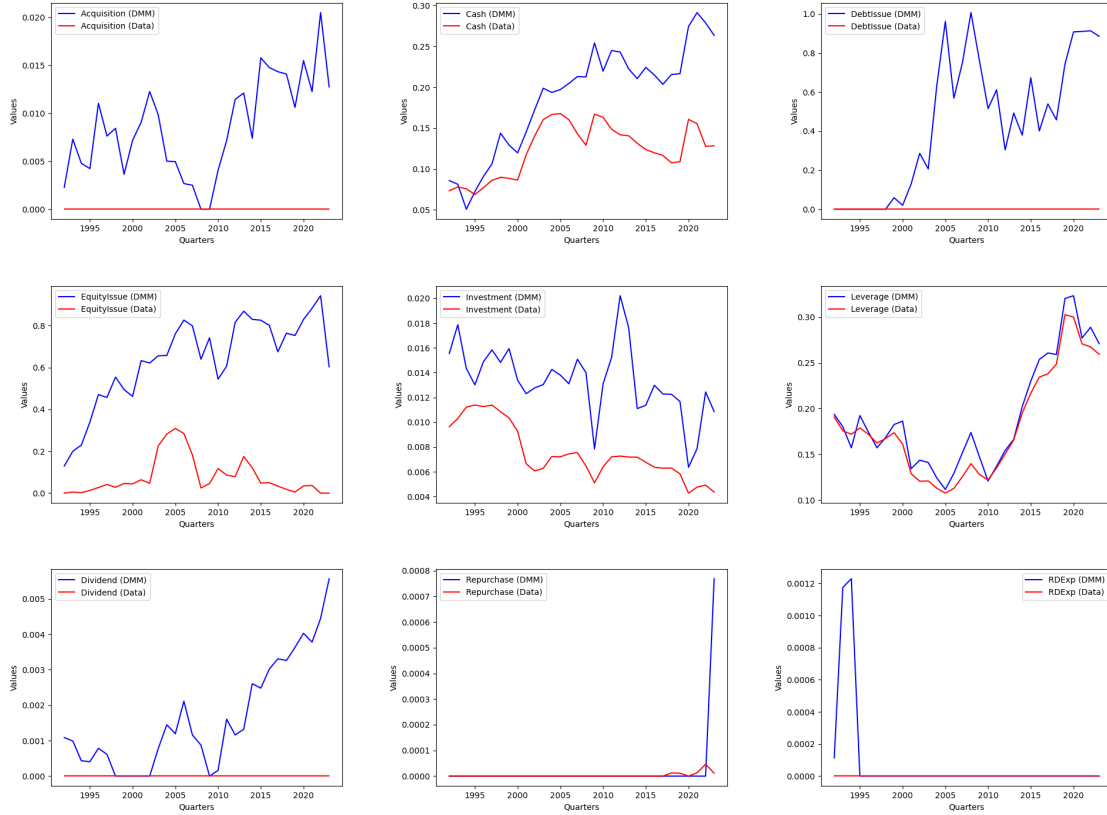
Note: This figure shows the empirical (red) and optimal decisions (blue) suggested by AM under 8-quarter market cap growth as the objective to maximize. We plot median to represent cross-sectional variation. The decisions are acquisitions (Faqc\_fund), cash holdings (Fcash\_hold), debt issuance (Fdebt\_issue\_log), equity issuance (Fequity\_issue\_log), investment ratio (Finv\_ratio2), leverage (Flev\_market), dividend (Fpay-out\_div), repurchase (Fpayout\_repurchase), and R&D expenses (Frd\_exp), with formal definitions detailed in Appendix A.

RD expenditures. These strategic adjustments collectively aim to optimize capital structure while reallocating resources toward growth-oriented activities.

## 6 New Research Questions in Corporate Finance

Our DDRC approach focuses more on producing reliable empirical predictions and policy recommendations. It does not serve as a substitute for theory, reduced-form models, or structural estimations, because it does not offer insights on the underlying economic mechanisms. The associated cost of the many connections AlphaManager establishes in the data is that they are sometimes difficult to rationalize — they lack the usual notion of *ex-ante*

Figure 8: Empirical and Optimal Decisions under 8-QTR Enterprise Value Growth



Note: This figure shows the empirical (red) and optimal decisions (blue) suggested by AM under 8-quarter enterprise value growth as the objective to maximize. We plot median to represent cross-sectional variation. The decisions are acquisitions (Faqc\_fund), cash holdings (Fcash\_hold), debt issuance (Fdebt\_issue\_log), equity issuance (Fequity\_issue\_log), investment ratio (Finv\_ratio2), leverage (Flev\_market), dividend (Fpayout\_div), repurchase (Fpayout\_repurchase), and R&D expenses (Frd\_exp), with formal definitions detailed in Appendix A.

theory predictions or causality statements — and they can be computationally expensive.

However, both problems can be mitigated by incorporating theoretical guidance (multiple concurrent models). DDRC is still a useful approach for understanding the underlying economic mechanisms by complementing other paradigms in corporate finance research. In a way, DDRC is a data-driven version of the structural approach that allows the incorporation of knowledge from theory and reduced-form models, as well as from financial big data. To wit, DDRC is not at odds — nor compete — with reduced causal exercises and theory-driven structural estimations. Instead, they give us guidance as to when we should use those other approaches to data analysis.

We next discuss how the measure of model ambiguity in AlphaManager reveals the types



of corporate decisions, firms, and market environment that would benefit more from theoretical guidance or causal identifications, as well as the types for which big data and powerful algorithms suffice. We then discuss how to incorporate the insights from other approaches into DDRC through ambiguity-guided transfer learning. Finally, we demonstrate how DDRC enables us to study managerial preferences.

## 6.1 Ambiguity and the Need for Theory/Reduced-Form Models

Our measure of ambiguity or relative entropy among the boosted PEMs informs how much we can rely on the data-driven approach. First, when the dispersion in predictions is high, one has to go back to theory or reduced-form models (with or without causal identification) in order to fully explain or predict firm outcomes, or to make counterfactual recommendations. Second, robust control makes DMM avoid exploring policies with more dispersed responses from PEMs. AlphaManager is thus conservative in the offline learning at the expense of missing out learning more profitable policies, the knowledge of which has to be derived from other conventional approaches. Figure 9 plots the estimated average ambiguity time series for each one of our state variables. For example, at the turns of macroeconomic regimes, ambiguity is estimated to be high.

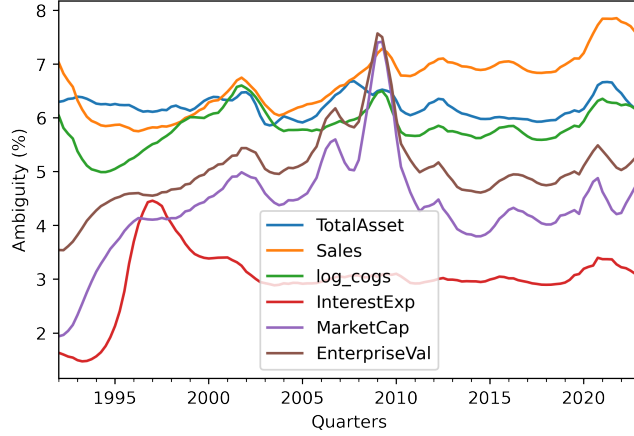
This apparent limitation of AlphaManager is its strength in disguise. The corporate finance literature has documented and studied a large collection of phenomena. Yet findings are typically siloed without a unified framework. DDRC, in addition to guiding managers in practice, provides insights about what areas in corporate finance research deserve more emphasis and attention going forward. This is so because along the dimensions that AlphaManager has greater ambiguity, theoretical insights and causal identification would be particularly informative and thus deserve more attention.

## 6.2 Ambiguity-Guided Transfer Learning from Other Models

Even after identifying the areas firms, and macroeconomic conditions where theory, reduced-form models, or structural estimations can be particularly informative, the question remains how one can incorporate the insights gained into a DDRC framework. A modified transfer learning technical comes to the rescue.

Transfer learning is a powerful machine learning technique that leverages the knowledge gained while solving one problem to solve a different but related problem. This approach is particularly useful when dealing with scenarios where the data for the new task is scarce or

Figure 9: Average Ambiguity Time Series by State Variables



Note: This figure shows average estimated ambiguity time series for key state variables—total assets, sales, COGS, interest expenses, market cap, and enterprise value. Ambiguity is calculated as follows: we first compute the square-rooted boosting error for each observation, aggregate these values at the state-quarter level, and then apply a 4-quarter moving average to the aggregated series.

when training a model from scratch would be computationally expensive. By transferring and adapting pre-trained models to new tasks, researchers and practitioners can achieve significant improvements in model performance with less data and in shorter training times. Transfer learning has found widespread applications across various domains, including natural language processing, computer vision, and speech recognition, demonstrating its versatility and efficiency in enhancing learning processes and outcomes in diverse settings.

In our setting, if one wants to incorporate any theoretical predictions (e.g., Hennessy and Whited, 2005; Bolton et al., 2011), or relationships identified in reduced-form models, or counterfactual predictions from structural estimations, one can use them to generate data to be added to the training sample for PEM. Note that using the historical data in reduced-form empirical models is unlikely to work. For example, a causal identification may concern only one particular firm and the effect would be too “weak” to be picked up in AlphaManager. Therefore, one has to necessarily generate more such cases by extending the reduced-form results to the cross section and broader ranges of the treatments, in order to generate sufficient observations.

Transfer learning is traditionally applied to one related dataset. Yet as we discuss earlier, we have many insights drawn from extant corporate finance research. How to incorporate them in a unified framework? Where to draw the line? That is where the ambiguity metric

can be used to guide the weights we put on the data for the transfer learning. If a pattern is rarer or more ambiguously described by PEM, one should increase the weights of data concerning that pattern in the transfer learning.

### 6.3 Revealed Managerial Preferences

In our DMM, we derive optimal decisions by optimizing exogenous objective functions. However, empirical managers often consider additional factors beyond our model’s specifications, such as ESG commitments or various forms of compensation. This divergence in objective functions naturally leads to differences between the theoretically optimal decisions suggested by our model and actual managerial decisions observed in practice.

Our key insight emerges from analyzing performance differentials across different objectives: the objective function yielding the smallest performance gap between model predictions and empirical decisions likely best approximates true managerial preferences. Motivated by this observation, we formulate the revealed preference problem as a min-max optimization challenge. To solve this, we employ advanced techniques like generative adversarial networks (GANs) to efficiently explore a broad space of potential managerial objective functions. The optimal function is identified as the one that minimizes the performance differential between our model’s recommendations and observed decisions.

This methodology effectively projects the true managerial preferences onto a space spanned by observable candidate objectives, providing a data-driven approach to understanding managerial decision-making processes.

## 7 Concluding Remarks

For any given managerial objective, corporate decision-making can be modeled as a high-dimensional, nonlinear stochastic control problem where managers learn about and interact to the evolving economic environment while formulating optimal dynamic policies. Conventional approaches in corporate finance and stochastic control often fall short of explaining or predicting empirical firm outcomes and have limited practical adoption. We introduce an AI-assisted, data-driven framework that prioritizes empirical performance and the generation of the useful counterfactuals to deliver effective policy recommendations for diverse business objectives.

Our approach first constructs a predictive environment module using supervised neural

networks and then integrates a decision-making module based on deep reinforcement learning. This dual-module structure goes beyond mere hypothesis testing on historical data or simulations by incorporating model ambiguity and robust control techniques. As a result, our framework not only improves in-sample and out-of-sample prediction of corporate outcomes but also identifies critical managerial decisions and offers dynamic policies that adapt to market evolution and feedback. Moreover, it distinguishes scenarios where theoretical insights and causal identification are essential from those where data-driven predictive models suffice. The framework’s flexibility is further enhanced by its ability to incorporate insights from theory, reduced-form, and structural models through ambiguity-guided transfer learning, making it a promising unified approach for corporate finance research.

We believe that the DDRC approach opens the door to investigating a range of important, yet previously unexplored, topics. For instance, while we assume a given managerial objective in this paper, techniques exist to learn historical managerial objectives from data. Additionally, our framework provides new insights into the heterogeneity of managerial actions, firm controllability, and their interactions with macroeconomic conditions – areas that warrant further study.

## References

- Almeida, Heitor, Nuri Ersahin, Vyacheslav Fos, Rustom M Irani, and Mathias Kronlund**, “How Do Short-Term Incentives Affect Long-Term Productivity?,” *The Review of Financial Studies*, 2024, p. hhae064.
- Andersen, Torben G, Nicola Fusari, and Viktor Todorov**, “Parametric inference and dynamic state recovery from option panels,” *Econometrica*, 2015, 83 (3), 1081–1145.
- Barnett, Michael, William Brock, and Lars Peter Hansen**, “Pricing uncertainty induced by climate change,” *The Review of Financial Studies*, 2020, 33 (3), 1024–1066.
- Bellstam, Gustaf, Sanjai Bhagat, and J Anthony Cookson**, “A text-based analysis of corporate innovation,” *Management Science*, 2020.
- Bolton, Patrick, Hui Chen, and Neng Wang**, “A unified theory of Tobin’s q, corporate investment, financing, and risk management,” *The journal of Finance*, 2011, 66 (5), 1545–1578.
- Bond, Philip, Alex Edmans, and Itay Goldstein**, “The real effects of financial markets,” *Annu. Rev. Financ. Econ.*, 2012, 4 (1), 339–360.
- , **Itay Goldstein, and Edward Simpson Prescott**, “Market-based corrective actions,” *The Review of Financial Studies*, 2010, 23 (2), 781–820.
- Campello, Murillo and Gaurav Kankanhalli**, “Corporate decision-making under uncertainty: review and future research directions,” *Handbook of Corporate Finance*, 2024, pp. 548–590.

- , – , and **Pradeep Muthukrishnan**, “Corporate hiring under Covid-19: Financial constraints and the nature of new jobs,” *Journal of Financial and Quantitative Analysis*, 2024, *Forthcoming*.
- Cao, Sean, Wei Jiang, Baozhong Yang, and Alan L Zhang**, “How to talk when a machine is listening: Corporate disclosure in the age of AI,” *The Review of Financial Studies*, 2023, *36* (9), 3603–3642.
- , – , **Junbo L Wang, and Baozhong Yang**, “From man vs. machine to man+ machine: The art and AI of stock analyses,” *Columbia Business School Research Paper*, 2021.
- Chen, Hui, Yuhang Cheng, Yanchu Liu, and Ke Tang**, “Teaching Economics to the Machines,” *Available at SSRN 4642167*, 2023.
- Cong, Lin William, Guanhao Feng, Jingyu He, and Junye Li**, “Sparse Modeling Under Grouped Heterogeneity with an Application to Asset Pricing,” 2023.
- , – , – , and **Xin He**, “Growing the Efficient Frontier on Panel Trees,” *NBER Working Paper*, 2022, (w30805).
- , **Ke Tang, Jingyuan Wang, and Yang Zhang**, “AlphaPortfolio: Direct Construction Through Reinforcement Learning and Interpretable AI,” *Working Paper*, 2020.
- , – , – , and – , “Deep Sequence Modeling: Development and Applications in Asset Pricing,” *The Journal of Financial Data Science*, 2021, *3* (1), 28–42.
- , **Tengyuan Liang, and Xiao Zhang**, “Textual Factors: A Scalable, Interpretable, and Data-driven Approach to Analyzing Unstructured Information,” Technical Report, resubmission requested. 2018.
- Dicks, David and Paolo Fulghieri**, “Uncertainty, investor sentiment, and innovation,” *The Review of Financial Studies*, 2021, *34* (3), 1236–1279.
- Dicks, David L and Paolo Fulghieri**, “Uncertainty aversion and systemic risk,” *Journal of Political Economy*, 2019, *127* (3), 1118–1155.
- Ebert, Frederik, Chelsea Finn, Sudeep Dasari, Annie Xie, Alex Lee, and Sergey Levine**, “Visual foresight: Model-based deep reinforcement learning for vision-based robotic control,” *arXiv preprint arXiv:1812.00568*, 2018.
- Edmans, Alex, Itay Goldstein, and Wei Jiang**, “The real effects of financial markets: The impact of prices on takeovers,” *The Journal of Finance*, 2012, *67* (3), 933–971.
- , – , and – , “Feedback effects, asymmetric trading, and the limits to arbitrage,” *American Economic Review*, 2015, *105* (12), 3766–3797.
- Erel, Isil, Léa H Stern, Chenhao Tan, and Michael S Weisbach**, “Selecting directors using machine learning,” Technical Report, National Bureau of Economic Research 2018.
- Feng, Guanhao, Stefano Giglio, and Dacheng Xiu**, “Taming the Factor Zoo: A Test of New Factors,” *The Journal of Finance*, 2020, *75* (3), 1327–1370.
- Fu, Justin, Aviral Kumar, Ofir Nachum, George Tucker, and Sergey Levine**, “D4rl: Datasets for deep data-driven reinforcement learning,” *arXiv preprint arXiv:2004.07219*, 2020.

- Fujimoto, Scott, David Meger, and Doina Precup**, “Off-policy deep reinforcement learning without exploration,” in “International Conference on Machine Learning” PMLR 2019, pp. 2052–2062.
- Garlappi, Lorenzo, Ron Giammarino, and Ali Lazrak**, “Ambiguity and the corporation: Group disagreement and underinvestment,” *Journal of Financial Economics*, 2017, 125 (3), 417–433.
- Graham, John R.**, “Presidential address: Corporate finance and reality,” *The Journal of Finance*, 2022, 77 (4), 1975–2049.
- Gu, Shihao, Bryan Kelly, and Dacheng Xiu**, “Empirical asset pricing via machine learning,” *The Review of Financial Studies*, 2020, 33 (5), 2223–2273.
- Hanley, Kathleen Weiss and Gerard Hoberg**, “Dynamic interpretation of emerging risks in the financial sector,” *The Review of Financial Studies*, 2019, 32 (12), 4543–4603.
- Hansen, Lars Peter and Thomas J Sargent**, “Robust control and model uncertainty,” *American Economic Review*, 2001, 91 (2), 60–66.
- **and Thomas J. Sargent**, “Risk, ambiguity, and misspecification: Decision theory, robust control, and statistics,” *Journal of Applied Econometrics*, 2023, n/a (n/a).
- Hennessy, Christopher A and Toni M Whited**, “Debt dynamics,” *The journal of finance*, 2005, 60 (3), 1129–1165.
- Hornik, Kurt, Maxwell Stinchcombe, and Halbert White**, “Multilayer feedforward networks are universal approximators,” *Neural networks*, 1989, 2 (5), 359–366.
- Izhakian, Yehuda, David Yermack, and Jaime F Zender**, “Ambiguity and the tradeoff theory of capital structure,” *Management Science*, 2022, 68 (6), 4090–4111.
- Jaques, Natasha, Asma Ghandeharioun, Judy Hanwen Shen, Craig Ferguson, Agata Lapedriza, Noah Jones, Shixiang Gu, and Rosalind Picard**, “Way off-policy batch deep reinforcement learning of implicit human preferences in dialog,” *arXiv preprint arXiv:1907.00456*, 2019.
- Jarrow, Robert and Simon Sai Man Kwok**, “Specification tests of calibrated option pricing models,” *Journal of Econometrics*, 2015, 189 (2), 397–414.
- Jin, Ying, Zhuoran Yang, and Zhaoran Wang**, “Is pessimism provably efficient for offline rl?,” in “International Conference on Machine Learning” PMLR 2021, pp. 5084–5096.
- Kahn, Gregory, Pieter Abbeel, and Sergey Levine**, “Badgr: An autonomous self-supervised learning-based navigation system,” *IEEE Robotics and Automation Letters*, 2021, 6 (2), 1312–1319.
- Kalashnikov, Dmitry, Alex Irpan, Peter Pastor, Julian Ibarz, Alexander Herzog, Eric Jang, Deirdre Quillen, Ethan Holly, Mrinal Kalakrishnan, Vincent Vanhoucke et al.**, “Scalable deep reinforcement learning for vision-based robotic manipulation,” in “Conference on Robot Learning” PMLR 2018, pp. 651–673.
- Kaniel, Ron, Zihan Lin, Markus Pelger, and Stijn Van Nieuwerburgh**, “Machine-learning the skill of mutual fund managers,” *Journal of Financial Economics*, 2023, 150 (1), 94–138.

- Kelly, Bryan, Semyon Malamud, and Kangying Zhou**, “The virtue of complexity in return prediction,” *The Journal of Finance*, 2024, 79 (1), 459–503.
- Kidambi, Rahul, Aravind Rajeswaran, Praneeth Netrapalli, and Thorsten Joachims**, “MOReL: Model-based offline reinforcement learning,” in “NeuIPS 2020” 2020.
- Kingma, Diederik P and Jimmy Ba**, “Adam: A method for stochastic optimization,” *arXiv preprint arXiv:1412.6980*, 2014.
- Levine, Sergey, Aviral Kumar, George Tucker, and Justin Fu**, “Offline reinforcement learning: Tutorial, review, and perspectives on open problems,” *arXiv preprint arXiv:2005.01643*, 2020.
- Li, Kai, Feng Mai, Rui Shen, and Xinyan Yan**, “Measuring corporate culture using machine learning,” *The Review of Financial Studies*, 2020.
- Malenko, Andrey and Anton Tsoy**, “Asymmetric information and security design under Knightian uncertainty,” *Available at SSRN 3100285*, 2020.
- Mamaysky, Harry, Matthew Spiegel, and Hong Zhang**, “Improved forecasting of mutual fund alphas and betas,” *Review of Finance*, 2007, 11 (3), 359–400.
- Mitton, Todd**, “Methodological variation in empirical corporate finance,” *The Review of Financial Studies*, 2022, 35 (2), 527–575.
- Pan, Jun**, “The jump-risk premia implicit in options: Evidence from an integrated time-series study,” *Journal of financial economics*, 2002, 63 (1), 3–50.
- Spiegel, Matthew**, “For corporate finance to truly advance we need more genuinely testable models,” *Financial Review*, 2023.
- Sutton, Richard S and Andrew G Barto**, *Introduction to reinforcement learning*, Vol. 135, MIT press Cambridge, 1998.
- Wang, Zhenyu**, “A shrinkage approach to model uncertainty and asset allocation,” *Review of Financial Studies*, 2005, pp. 673–705.
- Yu, Tianhe, Garrett Thomas, Lantao Yu, Stefano Ermon, James Zou, Sergey Levine, Chelsea Finn, and Tengyu Ma**, “Mopo: Model-based offline policy optimization,” *arXiv preprint arXiv:2005.13239*, 2020.

# Appendix A - Variable Selection and Definitions

Table 12 shows how we construct our firm-specific state and decision variables using Compustat and CRSP.

Table 12: Variable Selection and Definitions

variable	explanation	category	formula	source
lev	leverage	Decision	$(dtllq + dlcq) / L.atq$	Compustat
aqc	acquisition	Decision	$aqcq / L.atq$	Compustat
inv_ratio	investment	Decision	$capxq / L.atq$	Compustat
cash_hold	cash holding	Decision	$cheq / L.atq$	Compustat
div_payout	dividend	Decision	$devq / L.atq$	Compustat
log_debt_issue	debt issuance	Decision	$\log(1 + dltisq)$	Compustat
log_eq_issue	equity issuance	Decision	$\log(1 + sstkq)$	Compustat
rd_exp	R&D expenses	Decision	$xrdq / L.atq$	Compustat
repurchase	repurchases	Decision	$prstkq / L.atq$	Compustat
log_at	book asset	State	$\log(1 + atq)$	Compustat
log_act	current asset	State	$\log(1 + actq)$	Compustat
log_sale	gross revenue	State	$\log(1 + saleq)$	Compustat
log_ap	accounts Payables	State	$\log(1 + apq)$	Compustat
log_cogs	cost of good sold	State	$\log(1 + cogsq)$	Compustat
log_intpn	net interest paid	State	$\log(1 + intpnq)$	Compustat
log_inv	inventories	State	$\log(1 + invtq)$	Compustat
log_lct	current liabilities	State	$\log(1 + lctq)$	Compustat
log_rect	Receivables	State	$\log(1 + rectq)$	Compustat
log_revt	net revenue	State	$\log(1 + revtq)$	Compustat
log_market_cap	market cap	State	$\log(1 + csho * prccq)$	Compustat
log_enterprise_val	enterprise value	State	$\log(1 + atq + csho * prccq - ceqq - txdbq)$	Compustat
log_vol	trading volume	State	$\log(1 + VOL)$	CRSP
log_ret	equity return	State	$\log(1 + RET)$	CRSP
macro1	risk	Macro	risk	ChicagoFed
macro2	credit	Macro	credit	ChicagoFed
macro3	leverage	Macro	leverage	ChicagoFed
macro4	non-financial leverage	Macro	nonfinancial_leverage	ChicagoFed



## Appendix B - Technical Details

### B.1 Predictive Environment Module Hyper-Parameters

Table 13: Predictive Environment Module: Neural Network Hyperparameters

Hyperparameter	Value (main)	Value (macro)
optimizer	Adam	Adam
learning rate	0.000003	0.0005
batch size	2048	full
epoch (first)	40	50
epoch (rolling)	12	3
depth (main & aux)	4	2
neurons (main)	(512, 512, 512, 512)	(300, 300)
aux network num	10	-
dropout	0.3	-
l2 norm	0.00001	0.01

### B.2 Decision-making Module Hyperparameters

Table 14: Decision-making Module: Neural Network Hyperparameters

Hyperparameter	Value
optimizer	Adam
periods in obj. func.	{1, 8}
learning rate	0.00003
batch size	full
epoch	64
depth	4
neurons	(256, 256, 256, 256)

## Appendix C - Ambiguity Constraint Parameters

We first estimate parameters in the ambiguity constraints, namely  $\alpha$  and  $\beta$  in Eq.(6), for the long-term models. We follow an iterative process to search for proper constraints that could give out positive and reasonable long-term performance. For each cross section, we start from a tight lower-bound for these parameters, with 0 intercept and 1 slope. Essentially, this constraint requires AM managerial decisions as certain as the last set of observable decisions (the benchmark). Instead of search for optimal decisions which maximize long-term firm values, AM in this stage tries to minimize ambiguity of suggested managerial decisions. After 10 training epochs, we use a dynamic programming algorithm to estimate  $\alpha$  and  $\beta$  so that the constraint never binds and the area under the curve is minimized. In this way, the ambiguity punishment for the current AM equals to 0, and the ambiguity constraint is tight enough that AM performance is not likely to be irrationally high. After training for 5 more epochs, we examine whether the reward after ambiguity punishment is positive on average. If yes, we continue to train AM for 5 epochs and move on to the next cross section; if not, the current ambiguity constraint is likely too tight, so that AM is too ambiguity-averse and does not have sufficient incentive to explore the environment and reach for higher firm value. In this case, we re-estimate  $\alpha$  and  $\beta$  to loose the ambiguity constraint again and try to train the model for 5 more epochs, until it has a positive reward after ambiguity punishment.

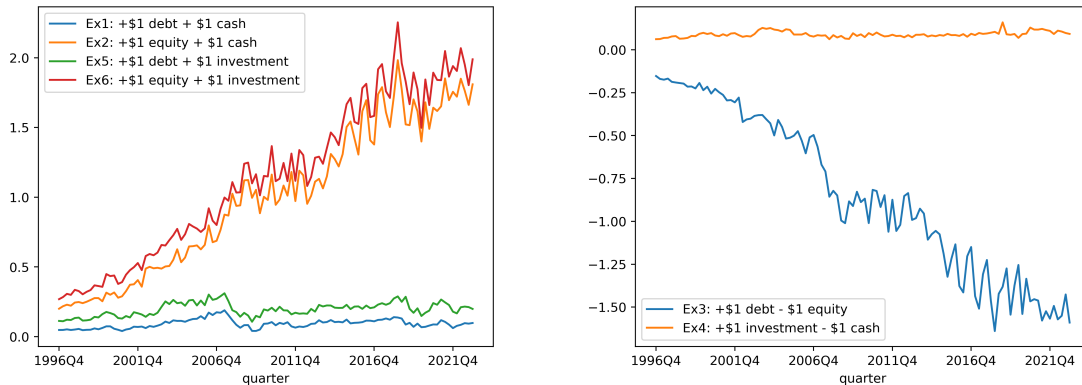
We then estimate  $\alpha$  and  $\beta$  for the short-term models. To make short-term models comparable to the corresponding long-term models in each cross section, we first calculate the long-term AM's short-term boosting error, and then we use the dynamic programming algorithm to estimate  $\alpha$  and  $\beta$ . Throughout the training of the short-term AM for that cross section, we stick with this ambiguity constraint and try to optimize short-term objective. The short-term AM is expected to behave better than the long-term AM, because under the long-term AM, the suggestion managerial decisions fall inside the constraint, while short-term objective values are suboptimal.

## Appendix D - An Illustration of Firm Recapitalization

Even though PEM is high-dimensional, we can analyze low dimensional action combinations, which are the focus of conventional reduced-form models. Firm recapitalization that shows up in every corporate finance textbook serves as an excellent illustration. Specifically, we ask how the enterprise value changes when a manager (firm):

- (1) raises \$1 in debt and put that \$1 into its cash savings
- (2) raises \$1 in equity and put that \$1 into its cash savings
- (3) raises \$1 in debt and \$1 less equity
- (4) puts \$1 cash into investment
- (5) raises \$1 in debt and put that \$1 into investment
- (6) raises \$1 in equity and put that \$1 into investment

Figure 10: Recapitalization Exercises: Capital Structure and Financing Investments



Note: This figure shows the recapitalization exercises where counterfactuals are generated by PEM. The x-axis is the calendar year, and the y-axis is dollar change in enterprise value. Solid lines plot the median value of the change in enterprise value for each cross section, and their legends explain the content of the recapitalization exercises.

Exercises 1–3 primarily address capital structure adjustments, whereas Exercises 4–6 focus on financing investments. Overall, Exercise 6 emerges as the optimal recapitalization strategy for increasing enterprise value. Notably, the dispersion in performance among the six strategies has grown over time.

We observe several patterns. First, Exercise 4 corresponds to the return on asset investment and is only marginally above zero. Second, in Exercise 3, altering the capital structure while keeping enterprise value constant results in a reduction in enterprise value — suggesting that firms generally carry excessive leverage and should decrease their debt proportion.

Third, Exercises 1, 2, 5, and 6 show that debt financing is less effective in boosting enterprise value compared to equity financing.

## Appendix E - Detailed Results on PEM for Macroeconomic Predictions

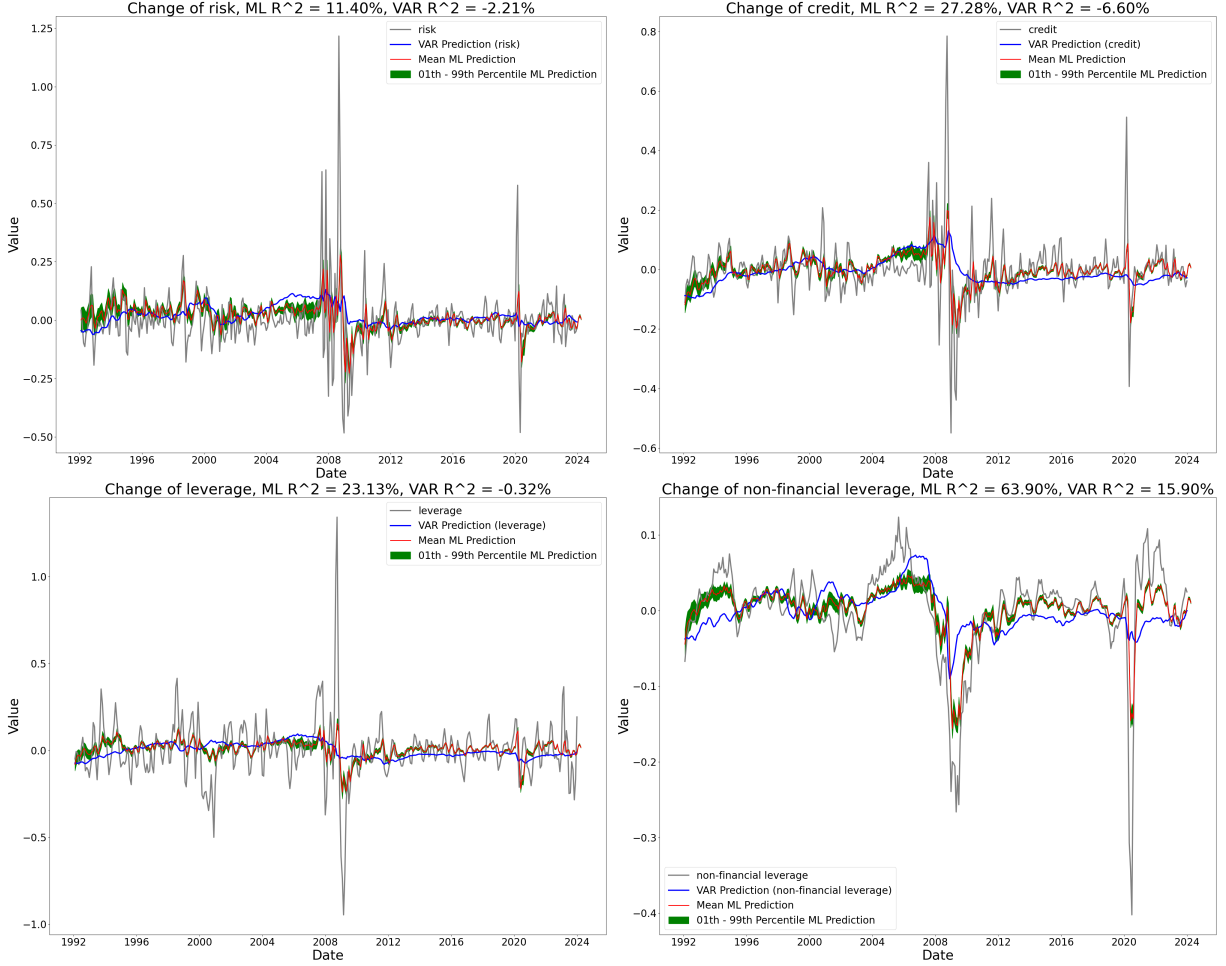


Figure 11: Predicting macroeconomic variables using neural network and VAR

Figure 11 shows the predicted change in next month for macroeconomic state variables using current and one-month lagged macroeconomic states generated by a neural network model and a VAR model. Neural networks, with their capability to handle collinear variables, are enhanced by adding current growth rates of macroeconomic state variables as additional inputs, along with their historical data. VAR models are generally used in time series modelling, which are widely used to model macroeconomic variables. In our setting, the neural network model is comparable to the VAR model because they are using the same information set as well as the same training paradigm when predicting the macroeconomic growth. Consistent with PEM, we update the VAR model every year using all available

observations till that year. We find that using simple VAR models to predict these four time series generates much lower (or even negative) out-of- sample pseudo  $R^2$  than PEM.

**Risk Prediction** The neural network achieves an  $R^2$  of 11.40% when predicting changes in risk, demonstrating moderate predictive power. The model captures significant risk fluctuations, particularly during periods of financial distress like the 2008 crisis, though it still leaves some noise in the actual values. On the other hand, the VAR model performs poorly, with a negative  $R^2$  of -2.21%, indicating that its predictions are worse than simply using the mean of the data. The VAR model struggles to capture key peaks and troughs, especially in volatile periods.

**Credit Prediction** The neural network fares better in predicting credit changes, achieving an  $R^2$  of 27.28%, with the model capturing general trends during financial disruptions. The VAR model, however, performs significantly worse with an  $R^2$  of -6.59%, highlighting its inability to accurately model credit risk in this context.

**Leverage Prediction** For leverage, the neural network attains an  $R^2$  of 23.13%, demonstrating reasonable predictive capability. The model captures key variations in leverage, especially during volatile periods like 2008, though it struggles during calmer periods. The VAR model continues to underperform, with an  $R^2$  of -0.32%, reflecting its failure to model significant leverage shifts accurately.

**Non-Financial Leverage Prediction** The neural network excels in predicting non-financial leverage, achieving an  $R^2$  of 63.90%, the highest across all variables. The model effectively tracks the changes in non-financial leverage across economic expansions and contractions, particularly during the financial crisis and the COVID-19 pandemic. The VAR model performs better than for other variables, achieving a positive  $R^2$  of 15.90%, though it still significantly underperforms relative to the neural network.

Overall, the neural network model consistently outperforms the VAR model in predicting macroeconomic conditions across all variables. The VAR model, while traditionally used for macroeconomic forecasting, struggles to capture the complex non-linear relationships inherent in these time series data, especially during periods of economic volatility. The neural network, with its ability to process both current and historical growth rates and

handle collinear variables, demonstrates a clear advantage in forecasting macroeconomic variables, particularly non-financial leverage.