# Vector-Based Opponent Modeling in Poker

Leonard Kinsman[1]

Rensselear Polytechnic Institute

**Abstract.** This research investigates a novel method for representing a poker player's complex playing style as a single, quantitative vector embedding. While traditional player analysis relies on a discrete set of statistics (e.g., VPIP, PFR), such metrics fail to capture the holistic and contextual nuances of a player's strategy. We propose a hierarchical deep learning architecture to encode this style. First, individual hand histories are encoded into feature vectors. A Transformer-based model then aggregates a player's entire set of hand vectors, trained using a deep metric learning (Triplet Loss) objective on a large-scale dataset. The central hypothesis is that proximity within this learned embedding space, measured by cosine similarity, will strongly correlate with true stylistic similarity. A successful model would provide a powerful, high-dimensional tool for quantitative player profiling, clustering, and comparative analysis, moving beyond the limitations of traditional statistical summaries.

**Keywords:** Representation Learning · Player Modeling · Game Theory · Imperfect Information

## 1 Introduction

### 1.1 Poker and Exploitability

For decades, game theory and artificial intelligence research have utilized complex games as benchmarks for developing decision-making agents. While perfect information games like chess and Go have seen superhuman solutions, imperfect information games present a far more complex and realistic challenge. Among these, Poker, and specifically No-Limit Texas Hold'em, has emerged as the quintessential problem. It is a multi-agent, stochastic environment where critical information—opponents' cards—is permanently hidden, forcing agents to reason under profound uncertainty. In poker, players combine private "hole cards" with public "community cards" to build the best five-card hand, but the game's strategic depth comes from its structured rounds of betting. Players use actions like betting, calling, and raising to maximize their expected value by either having the best hand at showdown or by forcing all opponents to fold.

This strategic landscape is governed by a fundamental duality: equilibrium play versus exploitative play. Game Theory Optimal (GTO) strategy refers to a balanced, unexploitable equilibrium. An agent playing a perfect GTO strategy is indifferent to its opponents' actions and cannot be beaten in the long

run, though it may not maximize its winnings against a flawed opponent. In contrast, exploitative strategy seeks to identify and attack the specific, systematic errors in a non-equilibrium opponent's play. The utility of exploitation is therefore immense, as the vast majority of human (and even many non-human) players deviate from GTO. A player who folds too often to aggression should be bluffed relentlessly; a player who calls too wide should be bet against only with strong value.

This powerful approach, however, is predicated on a critical prerequisite: accurate player modeling. To exploit an opponent, one must first possess a reliable model of their tendencies. Historically, this has been accomplished using a small set of discrete, hand-compiled statistics like VPIP (Voluntarily Put In Pot) and PFR (Pre-Flop Raise). While useful, these metrics provide a low-resolution, fragmented snapshot of a player's style. They fail to capture the contextual, sequential, and holistic nature of a player's strategy. This research addresses this gap. We posit that a player's complete strategic "fingerprint" can be learned and compressed into a high-dimensional vector embedding. This paper presents a novel deep learning framework to generate such a representation from a large corpus of hand histories, providing a quantitative tool for player profiling that moves far beyond traditional analytics.

### 1.2   Poker Crash Course

No-Limit Texas Hold'em is a game of skill and strategy where the primary objective is to win the "pot," which contains the sum of all chips bet during a hand. A player can win the pot in one of two ways: either by presenting the best five-card poker hand at the conclusion of all betting rounds (a "showdown"), or by making a bet that causes all other opponents to concede by folding their hands. The game progresses through a structured sequence, beginning with each player receiving two private "hole cards." These are followed by four rounds of betting, interspersed with the dealing of five "community cards" that all players can use. The first three community cards, dealt simultaneously, are known as "the flop," followed by a single "turn" card and a final "river" card. In each betting round, players can check, bet, call, raise, or fold. The "No-Limit" format imposes no maximum on the size of a bet, allowing a player to wager all of their chips at any point, which creates a highly complex and dynamic environment.

## 2   Motivation

The primary objective for any serious poker player is to maximize expected value. While Game Theory Optimal (GTO) play provides an unexploitable baseline, it does not maximize winnings against the vast majority of opponents who play a non-equilibrium, exploitable style. The ability to identify and attack the specific

strategic flaws of an opponent is the true hallmark of an expert player and the most direct path to increasing profit. However, effective exploitation is critically dependent on the accuracy of the underlying player model.

Current methods for player modeling almost exclusively rely on a dashboard of discrete, hand-compiled statistics, such as VPIP (Voluntarily Put In Pot), PFR (Pre-Flop Raise), and Aggression Factor. While these metrics are informative, they are fundamentally limited. They are low-resolution, non-contextual, and fail to capture the holistic or sequential nature of a player's strategy. For example, two players with an identical VPIP of 25 may have drastically different styles, one playing passively and another hyper-aggressively, a nuance lost in the single statistic.

This research is motivated by the limitations of these traditional analytics. We lack a comprehensive, high-dimensional, and quantitative representation of a player's complete "strategic fingerprint." Such a tool would unlock more powerful forms of analysis, allowing for nuanced player clustering, accurate similarity-based comparisons, and the discovery of complex stylistic patterns that discrete stats entirely miss. This paper seeks to create such a representation by moving from a small set of predefined statistics to a learned, data-driven vector embedding that captures the total essence of a player's style.

## 3   Related Work

The development of exploitative agents in poker is fundamentally a problem of opponent modeling. While Game Theory Optimal (GTO) play provides an unexploitable baseline, superior performance can be achieved by identifying and attacking the systematic weaknesses of a non-equilibrium opponent. Prior research has largely focused on creating explicit, statistical models of an opponent's strategy.

One prominent approach is Bayesian inference, as detailed by Southey et al. in "Bayes' Bluff" [1]. Their method constructs a probabilistic model to manage the uncertainty of an opponent's strategy. It begins with a prior distribution over the space of possible strategies and updates this to a posterior distribution as observations of the opponent's play are collected. This posterior is then used to compute an exploitative response, such as a Bayesian Best Response (BBR) or a response to a strategy sampled from the posterior (Thompson's Response). This method relies on a well-defined prior, which can range from a simple Dirichlet distribution in small games to a complex, expert-defined "informed prior" for larger games like Texas Hold'em.

A different, but related, method is the "Deviation-Based Best Response" (DBBR) proposed by Ganzfried and Sandholm [2]. This hybrid approach combines game-

theoretic reasoning with pure opponent modeling. It first computes an approximate equilibrium strategy offline, which serves as a baseline. During live play, the agent observes the opponent's action frequencies at various game states and models the *deviations* between these observed frequencies and the equilibrium's frequencies. A new opponent model is then constructed based on these observed deviations, and a best response is computed against this model in real-time.

Both of these foundational approaches aim to build an explicit model of the opponent's action probabilities. They are effective but rely on either a pre-computed equilibrium baseline [2] or a carefully constructed prior distribution [1]. Our work diverges from this paradigm. Instead of attempting to model an opponent's strategy at the level of individual game states, we propose a method to learn a single, holistic vector representation (an embedding) of a player's complete style directly from their raw hand histories. This deep metric learning approach is designed to capture a player's "strategic fingerprint" in a high-dimensional space, allowing for robust style-based clustering and similarity analysis without requiring a GTO solver or an expert-defined prior.

## 4   Methods and Data

This section details the dataset used for our research, the hierarchical model architecture designed to learn player representations, and the statistical method we will use to evaluate the quality of the resulting embedding space.

### 4.1   Data: Poker Hand History (PHH) Dataset

Our research is built upon the large-scale Poker Hand History (PHH) dataset from the University of Toronto. This dataset comprises several million anonymized hands played by over 100,000 unique players. Each hand history provides a complete, sequential log of all actions taken by all players at the table. This log includes the action (fold, check, call, bet, raise), the street (preflop, flop, turn, river), player positions, bet sizes, and the community cards. This rich, sequential data serves as the raw input for our representation learning models.

### 4.2   Methodology: A Hierarchical Learning Framework

Our methodology is a two-stage hierarchical process. First, we learn a fixed-size vector representation for *individual hands* using an LSTM autoencoder. Second, we aggregate these hand vectors into a single, comprehensive *player vector* using a Transformer-based deep metric learning model.

**Stage 1: LSTM Autoencoder for Hand Representation**   To create a numerical representation of a single hand, we first convert the sequence of events into a series of vectors. The **input vector** at each time step is a concatenation of the following features:

- 53 binary variables for Private Card 1 (with the 53rd variable representing "Unknown").
- 53 binary variables for Private Card 2 (with the 53rd variable representing "Unknown").
- 4 binary variables for the game state (Pre-Flop, Flop, Turn, River).
- A float representing normalized position ($i^{th}$ position / $N$ players).
- A float for the current pot size in Big Blinds (BBs).
- A binary vector for the action taken (Check, Bet, Call, Fold, Raise).
- A float for the bet size in Big Blinds (BBs), which is 0 for non-bet actions.
- A binary variable indicating if it is the player's turn to act.
- A binary variable indicating if the player has already folded.

We employ a Long Short-Term Memory (LSTM) based autoencoder, a type of recurrent neural network specifically designed to handle long-range dependencies in sequential data [3]. The **encoder** is an LSTM network that reads the input sequence of event vectors and compresses it into a single, fixed-size latent vector, $\boldsymbol{v}_h$. The **decoder** is a separate LSTM that attempts to reconstruct the original action sequence *only* from this latent vector. By training the model to minimize the reconstruction loss, the latent vector $\boldsymbol{v}_h$ is forced to capture the most salient, non-random strategic information about how the hand was played.

**Stage 2: Transformer-based Metric Learning for Player Representation** After processing all hands, each player $i$ is represented by a set of hand vectors $P_i = \{\boldsymbol{v}_{h1}, \boldsymbol{v}_{h2}, \ldots, \boldsymbol{v}_{hN}\}$. To aggregate this set into a single player vector, we use a Transformer encoder. We feed a sample of $N$ hand vectors into the Transformer and use the output embedding of the special `[CLS]` token as the final, aggregated player vector, $\boldsymbol{P}_i$.

Given the 100,000+ unique players, a standard classification objective is computationally infeasible. We therefore frame this as a **deep metric learning** problem [4], trained with a **Triplet Loss** function. For each training step, we sample a triplet:

1. **Anchor ($A$):** A set of $N$ random hands from a specific player, Player $i$.
2. **Positive ($P$):** A *different* set of $N$ random hands from the *same player*, Player $i$.
3. **Negative ($N$):** A set of $N$ random hands from a *different, randomly selected player*, Player $j$.

The Transformer model $f(\cdot)$ processes all three sets to produce three player vectors: $\boldsymbol{P}_A = f(A)$, $\boldsymbol{P}_P = f(P)$, and $\boldsymbol{P}_N = f(N)$. The Triplet Loss objective then trains the network to pull the Anchor and Positive vectors closer together while pushing the Anchor and Negative vectors further apart.

### 4.3   Evaluation Protocol

To validate the semantic meaning of our player embedding space, we must test our primary hypothesis: that vector proximity correlates with strategic similarity. We will use a "ground truth" measure of player similarity derived from

traditional poker statistics and compare it against the cosine similarity of our learned vectors.

First, we define a set of common, strategically-critical game states, $S$. This set includes, but is not limited to:

- $s_1$: Pre-flop, 2-Bet "Raise First In" (RFI).
- $s_2$: Pre-flop, facing a 2-Bet.
- $s_3$: Pre-flop, facing a 3-Bet.
- $s_4$: Flop, In-Position as Pre-flop Raiser on an Ace-High Monotone board.

For each player $i$ and each state $s \in S$, we compute their observable strategic distribution, $P_{s,i}$, by analyzing all their hands in the dataset. This distribution is the frequency of their actions in that state:

$$P_{s,i} = (\text{Fold \%}, \text{Check \%}, \text{Call \%}, \text{Bet \%}, \text{Raise \%})$$

To quantify the "ground truth" strategic distance $D(i,j)$ between any two players, $i$ and $j$, we compute the average **Kullback-Leibler (KL) Divergence** between their strategy distributions across all defined states:

$$D(i,j) = \frac{1}{|S|} \sum_{s \in S} D_{KL}(P_{s,i} \parallel P_{s,j})$$

This value $D(i,j)$ represents a robust, statistically-grounded measure of how differently two players play.

Our final validation will be a **Spearman rank correlation test** [5] between two ranked lists for a sample of players:

1. **List 1:** All other players ranked by their strategic distance $D(i,j)$ from a target player $i$.
2. **List 2:** All other players ranked by the cosine distance of their vector $\boldsymbol{P}_j$ from the target player's vector $\boldsymbol{P}_i$.

A high, positive rank correlation will validate our hypothesis, confirming that the learned embedding space successfully organizes players by their true, observable strategies.

## 5   Timeline

- Complete Training of LSTM (October 24th)
- Write Metric for Player Similarity (October 24th)
- Finish Training Transformer (November 7th)
- Calculate Correlations (November 14th)

# References

1. Southey, F., Bowling, M., Larson, B., Piccione, C., Burch, N., Billings, D., Rayner, C.: Bayes' Bluff: Opponent Modelling in Poker. In: Proceedings of the 21st Conference on Uncertainty in Artificial Intelligence (UAI 2005), pp. 550–557. AUAI Press, Arlington, Virginia (2005)
2. Ganzfried, S., Sandholm, T.: Game Theory-Based Opponent Modeling in Large Imperfect-Information Games. In: Tumer, K., Yolum, P., Sonenberg, L., Stone, P. (eds.) Proceedings of the 10th International Conference on Autonomous Agents and Multiagent Systems (AAMAS 2011), pp. 829–836. IFAAMAS, Taipei, Taiwan (2011)
3. Hochreiter, S., Schmidhuber, J.: Long Short-Term Memory. Neural Computation **9**(8), 1735–1780 (1997). https://doi.org/\BeginDoi
4. Schroff, F., Kalenichenko, D., Philbin, J.: FaceNet: A Unified Embedding for Face Recognition and Clustering. In: 2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp. 815–823. IEEE (2015). https://doi.org/\FaceNetDoi
5. Spearman, C.: The Proof and Measurement of Association between Two Things. The American Journal of Psychology **15**(1), 72–101 (1904). https://doi.org/\SpearmanDoi