# A Methodological Framework for High-Fidelity NFL Game Simulation via Monte Carlo Methods

## Executive Summary: A Framework for Realistic NFL Simulation

This report provides a comprehensive methodology for building a high-fidelity, play-by-play National Football League (NFL) game simulation using the nflfastR dataset. Achieving the core requirement of "realistic results" necessitates a multi-stage process that moves far beyond simple random sampling of historical outcomes. The framework begins by establishing a robust data foundation, identifying the critical variables within nflfastR that define the dynamic state of a football game. Second, it details the construction of a sophisticated, opponent-adjusted statistical model of team strength, which serves as the core parameter set for the simulation engine. This modeling phase is a crucial prerequisite, transforming noisy historical data into stable estimates of underlying team ability. Third, the report outlines the architecture of the Monte Carlo simulation engine itself, focusing on state management, probabilistic play-calling logic, stochastic outcome generation, and the often-overlooked but critical component of temporal realism through dynamic pace-of-play modeling. Finally, and most critically, it establishes a rigorous validation framework. This framework benchmarks the simulation's outputs against both historical game results and the highly efficient sports betting market, providing a quantitative and objective measure of the model's "realism." This structured, multi-layered approach ensures that the final product is not merely a descriptive tool but a powerful analytical engine capable of generating credible, probabilistic forecasts of NFL game outcomes.

## Section 1: The Anatomy of a Play - Foundational Data and Metrics with nflfastR

The empirical bedrock of any credible simulation is the data upon which it is built. This section deconstructs the nflfastR play-by-play dataset to identify the variables that constitute a complete "game state vector." It also introduces Expected Points Added (EPA) as the fundamental unit for measuring performance, which will serve as the simulation's primary currency for generating outcomes. A simulation cannot track every one of the 372+ columns available in the nflfastR dataset [1]; therefore, the initial analytical task is one of dimensionality reduction—selecting the subset of variables that form a sufficient and manageable game state vector.

## 1.1 Accessing and Structuring nflfastR Data

The first step in the simulation pipeline is the acquisition and structuring of a comprehensive historical dataset. The nflverse ecosystem provides powerful tools for this purpose. In the R programming language, the nflreadr package is the standard for efficiently downloading nflfastR's curated play-by-play data repositories.[1] For analysts working in Python, the

nfl-data-py library provides equivalent functionality, serving as an interface to the same underlying data sources.[3] To build robust statistical models, it is essential to load data from multiple seasons, creating a large sample of plays that captures a wide range of team strengths and strategic eras.[1]

The nflfastR data is structured with a "one row per play" philosophy, which is conceptually ideal for a play-by-play simulation.[6] However, this structure presents challenges with "abnormal" plays, such as those involving multiple turnovers or laterals, where a single row may not fully capture the sequence of events.[6] For the purpose of building the core simulation engine, an initial filtering step is necessary. Most of the underlying models for team strength and play outcomes should be trained on a dataset filtered for standard scrimmage plays, identified in

nflfastR with the binary indicators rush == 1 or pass == 1.[1] This excludes special teams plays, penalties that result in no play, and other non-scrimmage events. These excluded play types (e.g., field goals, punts) must be handled by dedicated, separate modules within the simulation, which will be discussed in Section 3.

## 1.2 Defining the Game State Vector

At any moment in an NFL game, its status can be described by a vector of key variables. This "game state vector" is the complete set of inputs required by the simulation engine to determine the context of the next play, model the likely strategic decisions, and generate a probable outcome. The granularity of nflfastR provides all the necessary components for this vector.

The core situational variables that define the tactical context of a play are: down, ydstogo (yards to gain for a first down), yardline_100 (distance in yards from the opponent's end zone), game_seconds_remaining, half_seconds_remaining, qtr (quarter), posteam (the team with possession), defteam (the defensive team), and score_differential.[5] These variables dictate everything from play-calling tendencies to the value of a given outcome. For player-level analysis or simulations involving individual statistics, it is critical to use the unique player identification codes provided (e.g.,

passer_id, rusher_id, receiver_id) rather than player names, which can be inconsistent or non-unique across seasons.[5] Table 1 provides a definitive data dictionary of the essential variables for constructing the simulation's state vector.

| Variable Name | Data Type | Description | Role in Simulation | Source |
|---|---|---|---|---|
| game_id | character | Unique identifier for each game. | Tracks individual simulation trials. | [5] |
| play_id | numeric | Unique identifier for each play within a game. | Sequential play tracking. | [5] |
| posteam | character | Abbreviation for the team on offense. | Defines possession. | [1] |
| defteam | character | Abbreviation for the team on defense. | Defines the defending team. | [1] |

| down | numeric | The down of the play (1, 2, 3, 4). | Core input for play-calling and outcome models. | [5] |
|---|---|---|---|---|
| ydstogo | numeric | Yards needed to achieve a first down. | Core input for play-calling and outcome models. | [11] |
| yardline_100 | numeric | Numeric yards from the opponent's end zone. | Core input for play-calling and outcome models. | [5] |
| game_seconds _remaining | numeric | Numeric seconds remaining in the game. | Input for pace-of-play and strategic models (e.g., 4th down). | [9] |
| score_different ial | numeric | Score of the possession team minus the score of the defensive team. | Key input for play-calling models (e.g., pass/run ratio). | [12] |
| pass | numeric | Binary indicator (1/0) for a pass play (includes sacks, scrambles). | Used for filtering and training play-type specific models. | [1] |
| rush | numeric | Binary indicator (1/0) for a rush play. | Used for filtering and training play-type specific | [1] |

| | | | | |
|---|---|---|---|---|
| | | | models. | |
| epa | numeric | Expected Points Added on the play. | The primary target variable for outcome generation. | [5] |
| home_team | character | Abbreviation for the home team. | Used to determine home-field advantage. | [1] |
| away_team | character | Abbreviation for the away team. | Used to determine home-field advantage. | [1] |

## 1.3 Expected Points Added (EPA): The Simulation's Currency

To generate realistic outcomes, the simulation must be grounded in a metric that accurately captures the value of each play. Traditional metrics like yards gained are context-dependent and fail to distinguish between a 5-yard gain on 3rd and 2 and a 5-yard gain on 3rd and 10. The superior metric for this purpose is Expected Points Added (EPA). The foundation of EPA is the concept of Expected Points (EP), which assigns a point value to any game state based on historical scoring outcomes from that state.[13] For example, a 1st and 10 from a team's own 25-yard line might have an EP of approximately 0.9 points, meaning that, on average, the possessing team is expected to score 0.9 more points than their opponent by the end of that drive.

EPA is simply the change in EP from the start of a play to the end of it.[15] A 20-yard pass on that 1st and 10 would move the ball to the 45-yard line, a state with a much higher EP value. The difference between the new EP and the starting EP is the EPA of that play. The

nflfastR dataset provides pre-calculated EPA values for every play, derived from its own robust and publicly documented models.[16]

For the purpose of simulation, EPA serves as a universal currency. It quantifies the value of every play in the common unit of points, allowing for direct comparison between a run, a pass,

a penalty, or a turnover.[18] This reframes the core task of the simulation engine. Instead of attempting to simulate "yards gained" and then inferring the point impact, a more robust approach is to directly simulate an "EPA outcome." The simulation's goal is to generate a sequence of plays whose resulting EPA distributions are consistent with the matchup of the two competing teams. The resulting yards gained, first downs, and scores become secondary outputs derived from this sequence of EPA values. This represents a fundamental shift in perspective from more traditional simulation approaches, grounding the engine in the metric that is most directly correlated with winning.

# Section 2: A Principled Approach to Modeling Team Strength

A realistic simulation is not a single model but a system of models, and the most critical precursor to a realistic simulation is the accurate quantification of the "true" underlying ability of each team. The simulation engine is merely a consumer of these strength parameters. Therefore, a separate, rigorous statistical modeling process must be undertaken *before* the simulation is constructed. This section details a principled approach to move beyond raw season averages to develop a robust, opponent-adjusted model of team strength that will parameterize the simulation engine.

## 2.1 The Fallacy of Raw EPA and the Need for Adjustment

Using a team's raw, season-long EPA per play as a measure of its strength is a common but deeply flawed approach. Raw metrics are highly susceptible to noise and, most importantly, are heavily influenced by the quality of a team's opponents.[14] A team that faces a schedule of porous defenses will naturally accumulate a higher offensive EPA per play, but this value is an inflated measure of their true offensive capability. Conversely, a strong defense that plays against a slate of elite offenses may appear mediocre by raw EPA allowed.

To obtain a more accurate estimate of true team ability, performance must be adjusted for the strength of the opponent. The core principle of opponent adjustment is that a team's performance on any given play should be evaluated relative to the quality of the unit it was facing.[19] This is the same foundational concept that powers more complex metrics like Football Outsiders' DVOA (Defense-adjusted Value Over Average), which contextualizes every play against a league-average baseline adjusted for situation and opponent.[14] For our

simulation to be realistic, its input parameters must be derived from such an adjusted framework.

## 2.2 A Hierarchical Model for Team Strength

Team strength is not a monolithic concept. To model it effectively, it must be deconstructed into distinct, interacting components. A standard and effective approach is to model four separate units for each team: Offensive Passing, Offensive Rushing, Defensive Passing, and Defensive Rushing.[21] This allows the simulation to capture specific matchup advantages, such as a strong rushing offense against a weak rushing defense.

A powerful and statistically sound method for estimating these component strengths is a hierarchical (or multilevel) regression model.[20] In this framework, the EPA of each play in a historical dataset is modeled as a function of the offensive ability of the possessing team and the defensive ability of the opposing team. A simplified representation of the model is:

$Play\_EPA_{ij} \sim N(\alpha_{offense,i} + \alpha_{defense,j}, \sigma_{play})$
Here, $Play\_EPA_{ij}$ is the EPA of a play where team i is on offense and team j is on defense. The terms $\alpha_{offense,i}$ and $\alpha_{defense,j}$ represent the latent "true" abilities of the respective units, which are estimated by the model.

A key advantage of the hierarchical structure is regularization, also known as shrinkage. The model assumes that each team's ability parameter is drawn from a common league-wide distribution (e.g., $\alpha_{offense,i} \sim N(0, \sigma_{offense})$). This has the practical effect of "pulling" or "shrinking" the estimates for teams with extreme or noisy results towards the league average. This process prevents overfitting, especially on the smaller sample sizes encountered early in a season, and produces more stable and predictive estimates of true talent.[20] The degree of shrinkage is learned from the data itself; if teams truly have widely varying abilities, the model will shrink them less. This data-driven regularization is crucial for generating reliable inputs for the simulation.

Furthermore, the stability of performance metrics over time should inform the model's structure. Historical analysis shows that offensive performance, particularly in the passing game, is more consistent and predictive from year to year ("stickier") than defensive performance.[18] This real-world observation implies that the hierarchical model should be structured with different prior variances for offense and defense. The model should place more confidence in a team's offensive strength, allowing for less variance around that estimate, while acknowledging the higher week-to-week and year-to-year volatility of defensive performance. This directly impacts the simulation: outcomes generated from the offensive passing matchup should have a lower intrinsic variance than those generated from

defensive matchups, mirroring the real-world observation that elite offenses tend to be more consistently elite than elite defenses.

## 2.3 Incorporating Key Contextual Factors

Beyond the four primary units, several overarching contextual factors must be modeled to ensure realism.

- **Home-Field Advantage (HFA):** HFA is a persistent and statistically significant factor in the NFL. It can be incorporated into the model as a fixed parameter that provides a small, consistent EPA boost to the home team on every play.[24] More advanced implementations could model HFA on a team-by-team basis, recognizing that some franchises have a more pronounced advantage than others.[24]
- **Quarterback (QB) Value and Injury Risk:** The quarterback is the single most influential player on the field, and his impact often transcends the general "Offensive Passing" rating of his team. A robust simulation must account for this in two ways. First, a separate rating for the starting quarterback, perhaps based on a composite metric like EPA + CPOE (Completion Percentage Over Expected), can be added as a powerful predictor.[14] Second, the model must include a parameter for the probability of a starting QB getting injured in any given game, and the significant drop-off in team performance that occurs when a backup enters the game.[27] Simulating this high-impact, stochastic event is essential for realistically capturing the variance of a team's season-long outcomes.

The output of this entire modeling phase is a set of clean, opponent-adjusted strength ratings for every team. These ratings, as shown in the conceptual Table 2, become the definitive input parameters for the Monte Carlo engine.

| Team | Adj. Pass Offense EPA | Adj. Rush Offense EPA | Adj. Pass Defense EPA | Adj. Rush Defense EPA | HFA (Points) | QB Value (EPA/play) |
|------|------|------|------|------|------|------|
| KC | +0.18 | +0.02 | +0.05 | +0.01 | 2.5 | +0.25 |
| SF | +0.15 | +0.08 | -0.10 | -0.12 | 2.2 | +0.18 |
| BAL | +0.12 | +0.15 | -0.15 | -0.18 | 2.7 | +0.20 |

| | | | | | | |
|---|---|---|---|---|---|---|
| PIT | -0.05 | -0.01 | -0.08 | -0.05 | 2.4 | -0.02 |
| CAR | -0.15 | -0.07 | +0.02 | +0.04 | 1.9 | -0.18 |

*(Note: Values are illustrative examples of model outputs.)*

# Section 3: The Engine Room - Constructing a Play-by-Play Monte Carlo Simulation

With a robust set of team strength parameters established, the next phase is to construct the simulation engine itself. This engine is responsible for transforming the static ratings into dynamic, play-by-play game narratives. The core of the engine is a loop that iteratively updates the game state by modeling play calls, generating stochastic outcomes, and managing the game clock. The entire process is then repeated thousands of times—the Monte Carlo method—to build a distribution of possible game outcomes. The decision to build a play-by-play, rather than a drive-by-drive, simulation is a deliberate one. While a drive-level simulation is computationally simpler [29], it cannot capture the intricate, path-dependent strategic decisions—most notably on fourth downs—that are critical to modern NFL game theory and essential for achieving a high degree of realism. [31] A play-by-play architecture is necessary to model these crucial inflection points explicitly. [33]

## 3.1 The Core Simulation Loop

The simulation of a single game is an iterative process that begins with an initial state and proceeds play by play until a terminal condition is met.

1. **State Initialization:** Each simulated game begins by setting the initial game state vector: Quarter 1, 15:00 on the clock (game_seconds_remaining = 3600), score 0-0, and the ball positioned for the opening kickoff.
2. **The Loop:** For each play, the engine executes a sequence of logical steps:
   - **Determine Game State:** The engine reads the current state vector (down, distance, field position, score, time, etc.).
   - **Model Play Call:** Based on the game state and pre-modeled team tendencies, a play type is probabilistically selected (e.g., pass, run, punt, field goal).

- - **Generate Play Outcome:** An outcome for the chosen play type is sampled from a probability distribution. The parameters of this distribution are determined by the specific matchup of the offensive and defensive unit strengths (from Section 2).
  - **Update Game State:** The generated outcome is applied to the state vector. This involves updating the down, distance, field position, score, and, critically, the game clock.
  - **Check for End Condition:** The engine checks if a terminal state has been reached (e.g., game_seconds_remaining is 0 at the end of a half or the end of regulation/overtime). If the game is not over, the loop repeats for the next play.
3. **Monte Carlo Repetition:** This entire simulated game constitutes one "trial." The Monte Carlo method leverages the law of large numbers by running thousands of these trials. The aggregated results from these trials form a probability distribution of final scores, win probabilities, and other game-level statistics.[25]

## 3.2 Modeling Play Calling and Pace of Play

To achieve realism, the simulation cannot assume that teams make decisions randomly or that time elapses uniformly. These strategic elements must be modeled explicitly.

- **Play Call Tendencies:** A team's decision to pass or run is heavily influenced by the game state. A simple but effective approach is to build a sub-model, such as a logistic regression or a gradient-boosted tree, trained on historical nflfastR data. This model predicts the probability of a pass (P(pass)) given inputs like down, ydstogo, score_differential, and game_seconds_remaining.[37] During the simulation, for each scrimmage play, the engine uses this model to calculate
P(pass) and then samples the play call (pass or run) based on that probability.
- **Pace of Play and the Game Clock:** This is a vital and often overlooked component of a realistic simulation. A naive model that subtracts a fixed amount of time for each play will fail to capture the strategic manipulation of the clock and will likely generate an unrealistic number of plays per game. An analysis of nflfastR data reveals that the time elapsed between plays is highly variable and conditional on the previous play's outcome and the game situation.[38] For instance, an incomplete pass stops the clock, resulting in minimal time runoff. A run in-bounds followed by a huddle consumes a significant amount of time (e.g., 35-40 seconds). A no-huddle offense can run plays with as little as 15-20 seconds of runoff. A realistic simulation must incorporate a "time off clock" sub-model that samples from a distribution of time elapsed, conditioned on factors like the previous play's result, the time remaining in the half (to account for two-minute drills), and the score differential (to account for four-minute offenses).[37] Without this dynamic clock management, the simulation will produce systematically biased game scripts and final

scores.

## 3.3 Generating Stochastic Play Outcomes

The core of the simulation's generative capability lies in sampling play outcomes from carefully parameterized probability distributions.

- **Sampling from Distributions:** For a given play call (e.g., a pass), the engine will not generate a deterministic outcome. Instead, it will sample an EPA value from a probability distribution, such as a Normal or a Student's t-distribution. The Student's t-distribution can be particularly useful as it has heavier tails, better capturing the higher frequency of extreme-outcome plays (like long touchdowns or turnovers) than a standard Normal distribution.[20]
- Conditioning the Distribution: The parameters of this sampling distribution (mean and variance) are not static; they are determined by the specific matchup based on the team strength ratings from Section 2. For a pass play between the home and away team, the mean of the EPA distribution would be calculated as:
  Mean_EPA = Home_Team_Pass_Offense_Rating + Away_Team_Pass_Defense_Rating + HFA_Effect + QB_Effect
  (Note: Defensive ratings are negative, so a good defense lowers the mean EPA). The variance of the distribution should also be modeled from historical data, allowing the simulation to capture the difference between consistently average teams and high-variance, boom-or-bust teams.
- **Translating EPA to Yards:** After an EPA value is sampled, it must be translated back into a yardage gain to update the game state's field position variables (yardline_100, ydstogo). This requires a "reverse lookup" model. A regression model can be trained on historical nflfastR data to predict yards_gained as a function of the play's epa and the situational context (down, ydstogo, yardline_100). This ensures that the generated yardage is consistent with the value of the play. For example, a high-positive EPA sample on 3rd and 15 will be translated into a large yardage gain, while a slightly negative EPA sample on 1st and 10 will translate into a short gain or loss.

## 3.4 Modeling Special Situations

The simulation must also include modules to handle non-scrimmage plays and critical strategic decisions.

- **Special Teams:**

- **Field Goals:** The outcome of a field goal attempt is primarily a function of its distance. A logistic regression model can be trained to predict the probability of a successful kick based on the kick_distance variable from nflfastR.[39] More advanced models can incorporate kicker-specific ratings or weather factors.[40]
  - **Punts:** The outcome of a punt is the opponent's resulting field position. This can be modeled by sampling from a historical distribution of net_punt_yards.
- **Turnovers:** Rather than being modeled as a distinct categorical outcome, turnovers can be elegantly handled as events in the tail of the EPA distribution. A very large negative EPA outcome sampled on a pass play corresponds to an interception. A large negative EPA on a run or completed pass can represent a lost fumble. The probability of these events occurring thus emerges organically from the overall EPA distribution for that matchup, correctly making them more likely when a poor offense faces an elite defense.
- **Fourth-Down Decisions:** This is a pivotal element of modern NFL strategy. The simulation should not rely on historical coaching tendencies, which are often suboptimal. Instead, it should incorporate a normative fourth-down decision model. This can be achieved by integrating a pre-built tool like the nfl4th package [2] or by constructing a dedicated sub-model based on the principles of win probability maximization.[31] When a fourth-down situation arises, the simulation consults this model, which recommends one of three actions: Go for it, attempt a Field Goal, or Punt. The engine then executes the recommended action. This ensures the simulated teams behave strategically optimally, reflecting the cutting edge of football analytics.

# Section 4: From Theory to Reality - Validation and Calibration

A simulation is only as valuable as its ability to replicate the real world. This section directly addresses the user's primary requirement for "realistic results" by outlining a multi-pronged strategy for validation and calibration. Realism is not a subjective assessment; it is a quantifiable and falsifiable claim that must be rigorously tested against objective benchmarks. The validation process is not a final step but an iterative loop, where discrepancies between simulated and real-world data are used to refine and improve the underlying models.

## 4.1 Internal and Face Validity

The first layer of validation involves ensuring the simulation is logically sound and produces

outputs that are plausible on the surface.

- **Sanity Checks:** This is the most basic form of validation, confirming that the simulation adheres to the fundamental rules of football. Do simulated statistics aggregate correctly (e.g., passing yards + rushing yards = total yards)? Are game states always logical (e.g., the simulation never produces a 5th down)? Does the clock always move forward? These checks ensure the internal consistency of the engine's code.
- **Distributional Checks:** The simulation should not just produce plausible individual games, but its aggregate output over many trials should mirror the statistical distributions of the real NFL. The distributions of key metrics from thousands of simulated games—such as total points per game, passing yards per game, plays per game, and turnover rates—should be compared to the historical distributions observed in the nflfastR dataset. If, for example, the simulated league average for points per game is 48, while the historical average is 44, it indicates a systemic issue, likely within the pace-of-play model, that must be addressed.

## 4.2 Backtesting Against Historical Results

The next layer of validation tests the simulation's predictive power against a held-out set of historical games.

- **Predictive Accuracy:** To perform this test, the team strength models from Section 2 are trained on data up to a certain season (e.g., through 2022). Then, the simulation is used to forecast the outcome of every game in a held-out season (e.g., 2023). For each game, thousands of trials are run to generate a win probability for the home team. The team with a win probability greater than 50% is the simulation's predicted winner. The model's straight-up accuracy—the percentage of actual winners correctly predicted—serves as a primary performance benchmark.[26] A well-constructed model should achieve an accuracy rate significantly higher than the 57-58% baseline of simply picking the home team.[43]
- **Probabilistic Calibration (Brier Score):** For a probabilistic model, accuracy alone is an incomplete measure of performance. It is crucial to assess whether the assigned probabilities are well-calibrated. For example, when the model predicts a set of games with a 70% win probability, do the favored teams actually win approximately 70% of those games? The Brier score is a proper scoring rule that measures the accuracy of probabilistic predictions. A lower Brier score indicates better calibration. This test ensures that the model's confidence levels are meaningful and reliable.

## 4.3 The Ultimate Benchmark: Calibration Against Betting Markets

The most rigorous and objective test of a simulation's realism is to compare its outputs to the sports betting market. The closing line of the betting market is widely considered the most accurate publicly available prediction of game outcomes, as it represents a consensus forecast refined by the "wisdom of the crowd," sophisticated modeling, and the allocation of vast amounts of capital.[36] A simulation that produces forecasts systematically different from the market cannot be considered truly realistic.

- **Comparing Spreads and Probabilities:**
  1. For each game in the historical backtest, a median projected score differential is generated from the thousands of simulation trials. This median differential is the model's intrinsic "point spread."
  2. This simulated spread is then compared to the actual closing line spread for that game. The Mean Absolute Error (MAE) between the simulated spreads and the market spreads is a key validation metric. A low MAE (e.g., in the range of the market's own error, typically around 10-11 points [26]) indicates that the simulation is well-aligned with the highly efficient market consensus.
  3. Similarly, the simulation's generated win probability can be converted into an equivalent "moneyline" price and compared directly to the closing odds offered by sportsbooks.
- **Against-the-Spread (ATS) Performance:** A powerful diagnostic test is to evaluate how the simulation would perform as a betting tool against the closing line. If the simulation's spread for a game is -7 and the market closes at -4, the simulation suggests a bet on the favorite. Over the entire held-out season, the simulation's ATS win rate is calculated. A result near 50% is a strong signal of a well-calibrated and realistic model, indicating that it has no systematic bias relative to the market. A result significantly below 50% points to a fundamental flaw in the model. A result consistently above 52.4% (the typical break-even point after accounting for vigorish) would be extraordinary and suggest a market inefficiency has been discovered, though this is a very high bar to clear.

This multi-faceted validation process can be summarized in a scorecard, providing a concise and powerful summary of the model's quantitative "realism."

| Validation Metric | Benchmark | Model Performance | Interpretation |
| --- | --- | --- | --- |
| **Straight-Up Accuracy** | >65% | 66.5% | The model has predictive power beyond simple |

| | | | heuristics. |
| --- | --- | --- | --- |
| **Brier Score** | Lower is better | 0.21 | The model's probabilities are well-calibrated. |
| **MAE vs. Closing Spread** | ~10.5 points | 10.2 points | The model's median outcome is aligned with the market consensus. |
| **ATS Win Rate vs. Close** | ~50% | 50.8% | The model exhibits no significant systematic bias relative to the market. |

(Note: Benchmark and performance values are illustrative examples based on high-performing public models.[26])

# Conclusion

This report has detailed a comprehensive, multi-stage framework for constructing a high-fidelity NFL game simulation capable of producing realistic results. The methodology emphasizes a principled approach that begins not with simulation code, but with rigorous statistical modeling. By leveraging the granular detail of nflfastR play-by-play data, it is possible to build a sophisticated, opponent-adjusted model of team strength that serves as the foundation for the entire system. This model deconstructs team ability into core offensive and defensive components, incorporating critical factors such as home-field advantage and the outsized impact of the starting quarterback.

The simulation engine itself is designed at the play-by-play level, a necessary complexity to capture the path-dependent strategic decisions, particularly on fourth downs, that define modern NFL football. Crucially, the engine moves beyond simplistic assumptions by incorporating dedicated sub-models for probabilistic play-calling and, most importantly, for the dynamic pace of play, ensuring that the temporal flow of the simulated game mirrors reality. The generation of play outcomes is grounded in the language of Expected Points Added (EPA), with stochastic results sampled from distributions whose parameters are

dictated by the specific matchup of team strengths.

Ultimately, the credibility of any such simulation rests on its validation. The framework presented establishes a stringent, multi-layered validation process. This process moves from internal sanity checks and distributional comparisons to rigorous backtesting against historical outcomes and, most critically, calibration against the efficient sports betting market. By measuring performance through metrics like predictive accuracy, Brier score, and Mean Absolute Error relative to the closing line, "realism" is transformed from a subjective goal into a quantifiable and objective benchmark. Building a truly realistic NFL simulation is an iterative and demanding process that lies at the intersection of statistical modeling, software engineering, and domain expertise. However, by following this structured and validation-driven approach, it is possible to create a powerful analytical tool capable of generating robust, credible, and probabilistic forecasts of NFL game dynamics.

## Works cited

1. A beginner's guide to nflfastR, accessed August 24, 2025, https://www.nflfastr.com/articles/beginners_guide.html
2. nflverse • Data and Tools for NFL Analytics • nflverse, accessed August 24, 2025, https://nflverse.nflverse.com/
3. nfl-data-py - PyPI, accessed August 24, 2025, https://pypi.org/project/nfl-data-py/
4. nflfastR to predict win probability : r/NFLstatheads - Reddit, accessed August 24, 2025, https://www.reddit.com/r/NFLstatheads/comments/18bhdl6/nflfastr_to_predict_win_probability/
5. nflfastR_python.md - GitHub Gist, accessed August 24, 2025, https://gist.github.com/Deryck97/dff8d33e9f841568201a2a0d5519ac5e
6. Get started with nflfastR, accessed August 24, 2025, https://www.nflfastr.com/articles/nflfastR.html
7. Get NFL Play by Play Data — fast_scraper - nflfastR, accessed August 24, 2025, https://www.nflfastr.com/reference/fast_scraper.html
8. nflfastR/vignettes/beginners_guide.Rmd at master - GitHub, accessed August 24, 2025, https://github.com/mrcaseb/nflfastR/blob/master/vignettes/beginners_guide.Rmd
9. Field Descriptions - nflfastR, accessed August 24, 2025, https://www.nflfastr.com/articles/field_descriptions.html
10. NFL Stats Variables - nflfastR, accessed August 24, 2025, https://www.nflfastr.com/articles/stats_variables.html
11. Examples • nflfastR, accessed August 24, 2025, https://mrcaseb.github.io/pages_dummy/articles/examples.html
12. Fourth Down Decision Calculator - RBSDM.com, accessed August 24, 2025, https://rbsdm.com/stats/fourth_calculator/
13. Expected Points and Relative Team Strength - inpredictable, accessed August 24, 2025, https://www.inpredictable.com/2012/12/expected-points-and-relative-team.html

14. NFL Advanced Metrics and Stats: DVOA, EPA, CPOE, aDOT & More, accessed August 24, 2025, https://www.covers.com/nfl/key-advanced-metrics-betting-tips
15. What are Expected Points Added (EPA) in the NFL - nfelo, accessed August 24, 2025, https://www.nfeloapp.com/analysis/expected-points-added-epa-nfl/
16. An R package to quickly obtain clean and tidy NFL play by play data • nflfastR - GitHub Pages, accessed August 24, 2025, https://mrcaseb.github.io/pages_dummy/
17. An R package to quickly obtain clean and tidy NFL play by play data • nflfastR, accessed August 24, 2025, https://www.nflfastr.com/
18. NFL EPA Tiers | Team Rankings by Expected Points Added | nfelo, accessed August 24, 2025, https://www.nfeloapp.com/nfl-power-ratings/nfl-epa-tiers/
19. Adjusting EPA for Strength of Opponent - Open Source Football, accessed August 24, 2025, https://opensourcefootball.com/posts/2020-08-20-adjusting-epa-for-strenght-of-opponent/
20. Estimating Team Ability From EPA - Open Source Football, accessed August 24, 2025, https://opensourcefootball.com/posts/2021-06-27-estimating-team-ability-from-epa/
21. Projection Model Picks - NFL - Shaving Points, accessed August 24, 2025, https://shaving-points.com/nfl/nfl-projection-model/
22. Modeling NFL game outcomes using Python … - Open Source Football, accessed August 24, 2025, https://opensourcefootball.com/posts/2021-01-21-nfl-game-prediction-using-logistic-regression/
23. NFL passing success rate - THE POWER RANK, accessed August 24, 2025, https://thepowerrank.com/nfl-passing-success-rate/
24. state-space model for National Football League scores - Mathematics & Statistics - Boston University, accessed August 24, 2025, http://math.bu.edu/people/mg/papers/nfl.pdf
25. The NFL Simulation: Monte Carlo Methods | Towards Data Science, accessed August 24, 2025, https://towardsdatascience.com/the-nfl-simulation-monte-carlo-methods-2f110424b1dd/
26. nfelo Model Performance, accessed August 24, 2025, https://www.nfeloapp.com/games/nfl-model-performance/
27. NFL Season Simulator - Unabated, accessed August 24, 2025, https://unabated.com/simulator-feature
28. How The Unabated NFL Simulator Works, accessed August 24, 2025, https://unabated.com/articles/how-the-unabated-nfl-simulator-works
29. Fast Drive Football, accessed August 24, 2025, https://www.fastdrivefootball.com/
30. NFL Game Simulator, accessed August 24, 2025, https://homes.cs.washington.edu/~kevinz/football-sim/
31. Analytics, have some humility: a statistical view of fourth-down decision making - arXiv, accessed August 24, 2025, https://arxiv.org/html/2311.03490v4

32. Go For It?: A Guide to Optimal Fourth Down Decision Making - Best Ball Stats, accessed August 24, 2025, https://bestballstats.com/2021/01/26/go-for-it-a-guide-to-optimal-fourth-down-decision-making/
33. Simulation-Based Decision Making in the NFL using NFLSimulatoR, accessed August 24, 2025, https://arxiv.org/pdf/2102.01846
34. Modeling Football: Combining ML models and Monte Carlo simulations | by Ian Dorward | DraftKings Engineering | Medium, accessed August 24, 2025, https://medium.com/draftkings-engineering/modeling-football-combining-ml-models-and-monte-carlo-simulations-170fd58cbe11
35. An Introduction and Step-by-Step Guide to Monte Carlo Simulations - Medium, accessed August 24, 2025, https://medium.com/@benjihuser/an-introduction-and-step-by-step-guide-to-monte-carlo-simulations-4706f675a02f
36. Using NFL Betting Simulations - WIN DAILY®, accessed August 24, 2025, https://windailysports.com/using-nfl-betting-simulations/
37. NFL Tendency Analysis and Basic Play Type Prediction | MLDS ..., accessed August 24, 2025, https://sites.northwestern.edu/msia/2020/01/31/nfl-tendency-analysis-and-basic-play-type-prediction/
38. How Madden Fails to Simulate Football: Quarter Length | Mega Bears Fan, accessed August 24, 2025, http://www.megabearsfan.net/post/2020/06/05/How-Madden-Fails-to-Simulate-Football-Quarter-Length.aspx
39. Releases · nflverse/nflfastR - GitHub, accessed August 24, 2025, https://github.com/nflverse/nflfastR/releases
40. Next Gen Stats: New advanced metrics you NEED to know for the 2024 NFL season, accessed August 24, 2025, https://www.nfl.com/news/next-gen-stats-new-advanced-metrics-you-need-to-know-for-the-2024-nfl-season
41. Fourth Down Decision Making: Challenging the Conservative Nature of NFL Coaches - DU Undergraduate Research Journal, accessed August 24, 2025, https://duurjportal.com/index.php/duurj/article/download/93/60/1023
42. The Anatomy of American Football: Evidence from 7 Years of NFL ..., accessed August 24, 2025, https://journals.plos.org/plosone/article?id=10.1371/journal.pone.0168716
43. Predicting NFL football games based on simulation and modeling, accessed August 24, 2025, https://www.mathsjournal.com/pdf/2018/vol3issue2/PartB/3-1-45-589.pdf
44. Sports Prediction Markets And Event Trading Explained - BettingUSA.com, accessed August 24, 2025, https://www.bettingusa.com/sports/prediction-markets/
45. How do betting analysts actually build their prediction models for NBA games? - Quora, accessed August 24, 2025, https://www.quora.com/How-do-betting-analysts-actually-build-their-prediction