



WEEK 1 & 2 Detailed Implementation Plan



WEEK 1: Data Processing & Strong Baseline

Day 1: Environment Setup & Data Infrastructure

Objective: Establish robust development environment and data pipeline foundation

A. Development Environment Configuration

1. Hardware Setup Optimization

- **MacBook Pro Configuration:** Optimize for 16GB RAM usage and thermal management
- **Virtual Environment:** Python 3.9+ with isolated dependencies
- **GPU Considerations:** Prepare for future Colab integration needs
- **Storage Management:** Allocate 20GB+ for datasets and model artifacts

2. Software Stack Installation

- **Core ML Libraries:** pandas, numpy, scikit-learn, lightgbm, xgboost
- **Deep Learning:** tensorflow/pytorch, transformers (for future BERT integration)
- **Visualization:** matplotlib, seaborn, plotly for comprehensive plotting
- **Statistical Analysis:** scipy, statsmodels for hypothesis testing preparation
- **Development Tools:** jupyter, git, pre-commit hooks for code quality

3. Project Structure Initialization

```
meme_stock_prediction/  
├── data/  
│   ├── raw/                # Original datasets  
│   ├── processed/          # Cleaned and merged data  
│   ├── features/            # Engineered features  
│   └── external/            # Additional data sources  
├── src/  
│   ├── data/                # Data processing modules  
│   ├── features/            # Feature engineering  
│   ├── models/              # Model implementations  
│   └── evaluation/           # Evaluation frameworks
```

```
|   └─ utils/                # Utility functions
|   └─ notebooks/           # Jupyter notebooks for exploration
|   └─ models/              # Trained model artifacts
|   └─ results/             # Output files and reports
|   └─ tests/               # Unit tests
|   └─ docs/                # Documentation
```

B. Data Source Acquisition Strategy

1. Reddit WSB Dataset Processing

- **Data Validation:** Verify dataset integrity and completeness
- **Quality Assessment:** Check for spam, duplicate posts, and data anomalies
- **Privacy Compliance:** Ensure user data anonymization and ethical usage
- **Sample Data Generation:** Create realistic sample data for testing if original unavailable

2. Stock Price Data Integration

- **API Setup:** Configure Yahoo Finance or Alpha Vantage for reliable data access
- **Historical Data Validation:** Verify price accuracy against multiple sources
- **Missing Data Strategy:** Implement forward-fill and interpolation for holidays/weekends
- **Multiple Asset Support:** Ensure pipeline handles GME, AMC, BB simultaneously

3. Mention Count Data Preparation

- **Extraction Methodology:** Develop robust ticker symbol detection algorithms
- **False Positive Filtering:** Distinguish between stock mentions and common word usage
- **Temporal Alignment:** Ensure consistent daily aggregation across datasets
- **Validation Sampling:** Manual verification of mention detection accuracy

C. Deliverables

- Fully configured development environment with all dependencies
- Project structure with initial documentation and README
- Data loading pipeline with error handling and validation
- Sample data generation system for testing and development

Day 2: Data Quality Assessment & Integration

Objective: Ensure data quality and create unified dataset for modeling

A. Comprehensive Data Exploration

1. Reddit Data Analysis

- **Temporal Distribution:** Analyze posting patterns across time periods
- **Content Quality:** Assess text length, engagement metrics, and content diversity
- **User Behavior:** Understand posting frequency and engagement patterns
- **Language Analysis:** Identify common themes, keywords, and sentiment patterns

2. Stock Data Validation

- **Price Consistency:** Verify OHLC relationships and detect anomalous movements
- **Volume Patterns:** Analyze typical trading volumes and identify outliers
- **Corporate Actions:** Account for stock splits, dividends, and other adjustments
- **Market Hours:** Properly handle pre-market and after-hours data

3. Cross-Dataset Temporal Alignment

- **Date Range Analysis:** Identify optimal time period with maximum data overlap
- **Missing Data Patterns:** Understand systematic vs. random missing data
- **Timezone Handling:** Ensure consistent temporal reference across data sources
- **Weekend/Holiday Treatment:** Develop strategy for non-trading days

B. Data Cleaning and Preprocessing Pipeline

1. Reddit Text Preprocessing

- **Content Standardization:** Unicode normalization, case handling, special characters
- **Spam Detection:** Remove promotional content, bot posts, and irrelevant discussions
- **Language Filtering:** Focus on English content with appropriate language detection
- **Content Categorization:** Separate stock-relevant vs. general discussion posts

2. Financial Data Cleaning

- **Outlier Detection:** Identify and handle extreme price movements and volume spikes
- **Data Consistency:** Ensure price relationships (high \geq low, etc.) are maintained
- **Corporate Action Adjustment:** Apply appropriate adjustments for stock splits/dividends
- **Currency Normalization:** Ensure consistent currency representation

3. Unified Dataset Creation

- **Temporal Aggregation:** Convert all data to consistent daily frequency
- **Feature Alignment:** Ensure all datasets share common date index
- **Missing Value Strategy:** Implement forward-fill, interpolation, or imputation as appropriate
- **Data Validation:** Comprehensive checks for logical consistency and completeness

C. Initial Data Statistics and Insights

1. Descriptive Statistics Generation

- **Reddit Metrics:** Post volume, engagement rates, sentiment distribution over time
- **Market Metrics:** Price volatility, trading volume patterns, return distributions
- **Cross-Correlation Analysis:** Preliminary relationships between social and market data

2. Data Quality Report

- **Completeness Assessment:** Percentage of missing data by variable and time period
- **Consistency Validation:** Logical relationship verification across datasets
- **Outlier Documentation:** Catalog and justify treatment of extreme values
- **Temporal Coverage:** Document final date range and data availability

D. Deliverables

- Clean, integrated dataset ready for feature engineering
 - Comprehensive data quality report with statistics and visualizations
 - Data preprocessing pipeline with full documentation
 - Initial exploratory analysis identifying key patterns and relationships
-

Day 3-4: Comprehensive Feature Engineering

Objective: Create robust feature set combining social, financial, and temporal signals

A. Reddit-Based Feature Engineering (25 features)

1. Basic Engagement Metrics (8 features)

- **Volume Indicators:** Daily post count, comment count, unique user count
- **Engagement Quality:** Average score per post, score-to-comment ratios

- **Temporal Patterns:** Posting velocity, engagement acceleration
- **Weekend Effects:** Weekend vs. weekday posting pattern differences
- **Activity Concentration:** Gini coefficient for post distribution across users

2. Sentiment Analysis Features (10 features)

- **Polarity Metrics:** Basic positive/negative sentiment ratios using VADER or TextBlob
- **Sentiment Momentum:** Rate of change in sentiment over 1, 3, 7-day windows
- **Sentiment Volatility:** Standard deviation of sentiment within rolling windows
- **Extreme Sentiment:** Proportion of posts with highly positive/negative sentiment
- **Sentiment Consensus:** Measure of agreement vs. polarization in community sentiment

3. Content Analysis Features (7 features)

- **Keyword Density:** Frequency of stock-specific and trading-related terminology
- **Linguistic Complexity:** Average sentence length, vocabulary diversity
- **Urgency Indicators:** Presence of time-sensitive language and calls to action
- **Emotional Intensity:** Caps usage, exclamation marks, emotional language
- **Information vs. Opinion:** Ratio of fact-based vs. opinion-based content

B. Financial Market Features (35 features)

1. Price-Based Features (15 features per stock: GME, AMC, BB)

- **Returns:** 1-day, 3-day, 7-day, 14-day price returns
- **Volatility:** Rolling standard deviation of returns (5, 10, 20 day windows)
- **Price Momentum:** Rate of change and acceleration in price movements
- **Relative Performance:** Performance vs. market indices and sector peers
- **Technical Levels:** Distance from recent highs/lows, support/resistance levels

2. Volume-Based Features (10 features per stock)

- **Volume Patterns:** Raw volume, volume moving averages, volume ratios
- **Volume-Price Relationship:** Volume-weighted average price (VWAP) deviations
- **Unusual Activity:** Volume spikes relative to historical patterns
- **Liquidity Indicators:** Bid-ask spread proxies, market impact measures

3. Market Microstructure Features (10 features)

- **Volatility Clustering:** GARCH-based volatility modeling
- **Jump Detection:** Identification of unusual price movements

- **Market Regime:** Bull/bear market indicators, trend strength measures
- **Cross-Asset Correlations:** Relationships between different meme stocks

C. Temporal and Cross-Modal Features (19 features)

1. Time-Based Features (9 features)

- **Calendar Effects:** Day of week, month, holiday proximity effects
- **Market Session:** Pre-market, regular hours, after-hours indicators
- **Seasonal Patterns:** Quarterly earnings seasons, options expiration cycles
- **Event Windows:** Time relative to significant market or company events

2. Cross-Modal Interaction Features (10 features)

- **Sentiment-Price Correlations:** Rolling correlations between sentiment and returns
- **Volume-Mention Synchronization:** Alignment between social activity and trading volume
- **Prediction Lag Effects:** Sentiment predicting future price movements at various horizons
- **Feedback Effects:** Price movements influencing subsequent social sentiment

D. Feature Engineering Pipeline Implementation

1. Automated Feature Generation

- **Modular Design:** Separate feature generators for each category
- **Scalable Architecture:** Easy addition of new features without breaking existing pipeline
- **Error Handling:** Robust handling of missing data and edge cases
- **Performance Optimization:** Efficient computation for large datasets

2. Feature Validation and Quality Control

- **Statistical Properties:** Distribution analysis, outlier detection, correlation assessment
- **Temporal Stability:** Ensure features are stable across different time periods
- **Predictive Power:** Initial univariate analysis of feature-target relationships
- **Redundancy Assessment:** Identify and handle highly correlated features

E. Deliverables

- Complete feature engineering pipeline generating 79 features
- Feature documentation with mathematical definitions and business interpretations
- Feature quality report including distributions, correlations, and predictive power analysis
- Engineered dataset ready for model training and validation

Day 5-6: Baseline Model Development

Objective: Establish competitive baseline models for performance benchmarking

A. Model Architecture Selection and Implementation

1. LightGBM for Classification Tasks

- **Target Variables:** 1-day and 3-day price direction prediction for GME, AMC, BB
- **Model Configuration:** Gradient boosting with early stopping and cross-validation
- **Hyperparameter Space:** num_leaves (10-100), learning_rate (0.01-0.3), feature sampling rates
- **Regularization:** L1/L2 penalties, minimum child samples, bagging parameters

2. XGBoost for Regression Tasks

- **Target Variables:** 3-day and 7-day price magnitude prediction for all stocks
- **Architecture:** Extreme gradient boosting with regularization
- **Parameter Optimization:** max_depth, learning_rate, subsample, colsample_bytree
- **Loss Functions:** Squared error with custom evaluation metrics

3. LSTM for Sequential Pattern Recognition

- **Architecture:** Multi-layer LSTM with dropout and batch normalization
- **Sequence Length:** 30-60 day lookback windows for temporal pattern capture
- **Features:** Time series of engineered features with proper scaling
- **Training Strategy:** Early stopping, learning rate scheduling, gradient clipping

B. Training and Validation Framework

1. Time Series Cross-Validation

- **Methodology:** Walk-forward validation with expanding window
- **Split Strategy:** 70% training, 15% validation, 15% testing with temporal ordering
- **Data Leakage Prevention:** Strict temporal boundaries, no future information
- **Performance Metrics:** Accuracy, F1-score, AUC-ROC for classification; RMSE, MAE for regression

2. Model Training Procedures

- **Feature Scaling:** StandardScaler for continuous features, appropriate encoding for categorical
- **Class Imbalance:** Handling for direction prediction using class weights or sampling
- **Overfitting Prevention:** Early stopping, regularization, dropout, cross-validation monitoring
- **Computational Efficiency:** Parallel processing, memory optimization, progress tracking

3. Hyperparameter Optimization

- **Search Strategy:** Grid search for initial exploration, random search for refinement
- **Validation Approach:** Nested cross-validation to avoid overfitting to validation set
- **Computational Budget:** Balance between thorough search and practical time constraints
- **Documentation:** Track all hyperparameter experiments and results

C. Model Evaluation and Analysis

1. Performance Metrics Calculation

- **Classification Metrics:** Accuracy, precision, recall, F1-score, AUC-ROC, confusion matrices
- **Regression Metrics:** RMSE, MAE, MAPE, directional accuracy, correlation coefficients
- **Business Metrics:** Sharpe ratio estimation, maximum drawdown, profit factor
- **Statistical Significance:** Confidence intervals, statistical tests vs. random baseline

2. Feature Importance Analysis

- **Model-Specific Importance:** Native feature importance from tree-based models
- **Permutation Importance:** Model-agnostic importance through feature shuffling
- **SHAP Values:** Detailed feature contribution analysis for individual predictions
- **Partial Dependence:** Understanding feature effects across their ranges

3. Error Analysis and Model Diagnostics

- **Residual Analysis:** Pattern identification in prediction errors
- **Temporal Performance:** Model performance across different time periods
- **Market Condition Analysis:** Performance during high/low volatility periods
- **Failure Case Study:** Detailed analysis of worst predictions

D. Baseline Performance Benchmarking

1. Target Performance Levels

- **Classification Accuracy:** Target >70% for direction prediction (vs. 50% random)

- **Regression Performance:** Target RMSE <0.6 for magnitude prediction
- **Consistency:** Stable performance across different stocks and time periods
- **Business Relevance:** Positive risk-adjusted returns in trading simulation

2. Comparative Analysis

- **Simple Baselines:** Moving average, momentum, mean reversion strategies
- **Technical Analysis:** RSI, MACD, Bollinger Bands-based predictions
- **Sentiment-Only Models:** Using only Reddit features for comparison
- **Price-Only Models:** Using only financial features for comparison

E. Deliverables

- Trained baseline models for all prediction tasks with saved weights/parameters
 - Comprehensive performance evaluation report with statistical validation
 - Feature importance analysis with business interpretations
 - Model comparison framework ready for Week 2 enhancements
-

Day 7: Documentation and Week 1 Summary

Objective: Consolidate Week 1 achievements and prepare for Week 2 development

A. Comprehensive Documentation Creation

1. Technical Documentation

- **Code Documentation:** Docstrings, type hints, inline comments for all functions
- **API Reference:** Complete function and class documentation with examples
- **Pipeline Documentation:** Step-by-step data flow and processing procedures
- **Configuration Management:** Parameter files and environment setup instructions

2. Experimental Documentation

- **Model Architecture:** Detailed specifications and design rationale
- **Hyperparameter Logs:** Complete record of optimization experiments
- **Performance Tracking:** Systematic results logging with timestamps and versions
- **Error Analysis:** Documentation of challenges encountered and solutions implemented

B. Week 1 Performance Summary and Analysis

1. Achievement Summary

- **Data Pipeline:** Successfully processed and integrated 3 distinct data sources
- **Feature Engineering:** Created 79 meaningful features with validation
- **Model Performance:** Achieved baseline accuracy targets with statistical validation
- **Infrastructure:** Established robust development and evaluation framework

2. Key Performance Metrics

- **Best Classification Performance:** Report highest accuracy achieved and for which target
- **Best Regression Performance:** Report lowest RMSE achieved and for which target
- **Feature Importance Rankings:** Top 10 most important features across all models
- **Computational Efficiency:** Training times and resource utilization metrics

3. Statistical Validation

- **Performance Confidence:** Statistical significance of results vs. random baselines
- **Cross-Validation Stability:** Consistency across different validation folds
- **Temporal Robustness:** Performance stability across different time periods
- **Error Analysis:** Common failure patterns and prediction uncertainty quantification

C. Week 2 Preparation and Planning

1. Identified Enhancement Opportunities

- **Feature Engineering:** Gaps in social signal capture and advanced sentiment analysis
- **Model Architecture:** Opportunities for ensemble methods and deep learning integration
- **Data Sources:** Additional data that could improve prediction accuracy
- **Evaluation Framework:** More sophisticated metrics and validation procedures

2. Technical Debt and Optimization Opportunities

- **Code Refactoring:** Areas for improved modularity and maintainability
- **Performance Optimization:** Bottlenecks in data processing or model training
- **Memory Management:** Opportunities for more efficient resource utilization
- **Testing Coverage:** Areas needing additional unit tests and validation

3. Week 2 Success Criteria Definition

- **Performance Targets:** Specific improvement goals for accuracy and other metrics

- **Feature Development:** Planned advanced features and their expected contributions
- **Model Innovation:** Advanced architectures and ensemble methods to implement
- **Validation Requirements:** Enhanced statistical testing and robustness analysis

D. Deliverables and Knowledge Transfer

1. Complete Week 1 Package

- **Source Code:** Clean, documented, and tested codebase
- **Data Assets:** Processed datasets and feature engineering pipelines
- **Model Artifacts:** Trained models with performance metrics and documentation
- **Results Reports:** Comprehensive analysis and performance documentation

2. Week 1 Summary Report

- **Executive Summary:** High-level achievements and key metrics
- **Technical Details:** Methodology, implementation details, and results analysis
- **Lessons Learned:** Challenges overcome and insights gained
- **Week 2 Roadmap:** Planned enhancements and success criteria

E. Week 1 Deliverables Summary

- Robust data processing pipeline handling 3 data sources
- 79 engineered features with comprehensive documentation
- Baseline models achieving 75%+ accuracy for direction prediction
- Complete development framework ready for advanced enhancements
- Comprehensive documentation and performance analysis

WEEK 2: Meme-Specific Features & Advanced Models

Day 8-9: Advanced Meme Feature Engineering

Objective: Develop sophisticated features capturing meme stock-specific behaviors

A. Viral Pattern Detection System

1. Viral Growth Modeling (15 features)

- **Exponential Growth Detection:** Mathematical modeling of mention/engagement acceleration
- **Viral Velocity Indicators:** Rate of change in social media activity and engagement
- **Cascade Analysis:** User participation patterns and influence propagation
- **Saturation Detection:** Identification of peak viral moments and decline phases
- **Cross-Platform Amplification:** Correlation between Reddit activity and broader social media trends

2. Viral Lifecycle Classification

- **Growth Phase Identification:** Early-stage viral pattern recognition
- **Peak Detection:** Maximum attention capture and engagement identification
- **Decline Phase Analysis:** Post-peak engagement pattern characterization
- **Resurrection Patterns:** Secondary viral waves and revival detection

3. Implementation Strategy

- **Mathematical Foundations:** Epidemiological models adapted for social media viral spread
- **Feature Validation:** Statistical significance testing of viral indicators vs. price movements
- **Temporal Sensitivity:** Multi-timeframe viral pattern detection (hourly, daily, weekly)
- **Robustness Testing:** Validation across different viral events and market conditions

B. Advanced Sentiment Analysis Architecture

1. Multi-Model Sentiment Fusion (20 features)

- **Financial BERT Integration:** FinBERT for financial domain-specific sentiment analysis
- **Emotion Classification:** Multi-dimensional emotional state detection (joy, fear, anger, surprise)
- **Confidence Scoring:** Prediction confidence and uncertainty quantification
- **Contextual Understanding:** Situation-aware sentiment interpretation
- **Temporal Sentiment Dynamics:** Sentiment momentum, acceleration, and volatility measures

2. Meme-Specific Language Analysis

- **Diamond Hands Detection:** "Hold" sentiment strength and conviction measurement
- **Paper Hands Identification:** "Sell" pressure and weak conviction indicators
- **FOMO/FUD Analysis:** Fear of missing out vs. fear/uncertainty/doubt balance
- **Moon Expectation Modeling:** Price target optimism and expectation quantification
- **Tribal Language Intensity:** Community-specific terminology and identity markers

3. Advanced NLP Techniques

- **Semantic Similarity:** Word embeddings and contextual meaning analysis
- **Sarcasm Detection:** Irony and sarcasm identification in financial context
- **Influence Scoring:** Author credibility and post influence measurement
- **Topic Modeling:** Latent topic discovery and trend identification

C. Social Network Dynamics Quantification

1. Community Behavior Analysis (10 features)

- **Echo Chamber Measurement:** Opinion homogeneity and diversity quantification
- **Influential User Tracking:** High-karma user activity and influence patterns
- **New User Integration:** Fresh participant conversion and retention analysis
- **Community Fragmentation:** Sub-group formation and consensus breakdown detection
- **Information Cascade Strength:** Follow-the-leader behavior quantification

2. Network Effect Modeling

- **Coordinated Behavior Detection:** Synchronized posting and voting pattern identification
- **Brigading Analysis:** External influence and manipulation detection
- **Organic vs. Artificial Growth:** Distinguishing natural from manufactured viral patterns
- **Community Leadership Changes:** Shift in influential voices and opinion leaders

D. Cross-Modal Feature Innovation

1. Social-Financial Signal Integration (14 features)

- **Sentiment-Price Correlation Evolution:** Dynamic relationship tracking over time
- **Volume-Mention Synchronization:** Trading activity and social activity alignment
- **Prediction Lead-Lag Analysis:** Temporal precedence between social signals and price movements
- **Feedback Loop Detection:** Price movement influence on subsequent social sentiment
- **Cross-Asset Contagion:** Meme stock interconnection and influence spillover

2. Advanced Interaction Features

- **Regime-Dependent Correlations:** Relationship changes during different market conditions
- **Volatility-Sentiment Coupling:** Volatility impact on community behavior and vice versa
- **Options Flow Integration:** Social sentiment relationship with derivatives activity
- **Institutional vs. Retail Sentiment:** Different participant behavior pattern separation

E. Implementation and Validation Framework

1. Feature Engineering Pipeline Enhancement

- **Scalable Architecture:** Efficient processing of additional complex features
- **Real-Time Capability:** Streaming data processing for live feature computation
- **Quality Assurance:** Automated testing and validation of new feature calculations
- **Performance Monitoring:** Computational efficiency and memory usage optimization

2. Feature Validation Methodology

- **Statistical Significance:** Individual feature predictive power assessment
- **Information Content:** Mutual information and correlation analysis with targets
- **Temporal Stability:** Feature behavior consistency across different time periods
- **Business Logic Validation:** Economic and behavioral interpretation verification

F. Deliverables

- Advanced feature engineering pipeline generating 45+ new meme-specific features
 - Comprehensive validation report demonstrating feature quality and predictive power
 - Documentation of viral pattern detection algorithms with mathematical foundations
 - Integration framework ready for advanced model development
-

Day 10-11: Advanced Model Architecture Development

Objective: Implement sophisticated models leveraging new features and advanced architectures

A. Multi-Modal Transformer Architecture

1. BERT Integration for Text Processing

- **Model Selection:** Financial BERT (FinBERT) for domain-specific language understanding
- **Text Preprocessing:** Tokenization, encoding, and attention mask generation for Reddit posts
- **Fine-Tuning Strategy:** Domain adaptation for meme stock terminology and context
- **Computational Optimization:** Efficient batch processing and memory management for MacBook Pro

2. Transformer Encoder for Temporal Sequences

- **Architecture Design:** Multi-head attention for temporal feature sequences
- **Positional Encoding:** Time-aware position encoding for financial time series
- **Feature Fusion:** Integration of text embeddings with numerical features
- **Multi-Task Learning:** Simultaneous prediction of direction and magnitude

3. Advanced Attention Mechanisms

- **Cross-Modal Attention:** Attention between social sentiment and financial signals
- **Temporal Attention:** Dynamic weighting of different time periods
- **Feature Group Attention:** Selective focus on different feature categories
- **Ensemble Attention:** Model confidence and uncertainty-aware attention weighting

B. Enhanced LSTM Architecture

1. Bidirectional LSTM with Attention

- **Architecture:** Forward and backward temporal processing with attention pooling
- **Feature Integration:** Multi-scale temporal features with different lookback windows
- **Regularization:** Dropout, batch normalization, and gradient clipping
- **Memory Optimization:** Efficient implementation for limited computational resources

2. LSTM Variants Exploration

- **GRU Comparison:** Gated Recurrent Units for faster training and comparable performance
- **ConvLSTM:** Convolutional LSTM for spatial-temporal pattern recognition
- **Attention-LSTM:** Attention mechanism integration for improved long-term dependencies
- **Ensemble LSTM:** Multiple LSTM models with different configurations

C. Advanced Ensemble System Design

1. Multi-Level Ensemble Architecture

- **Base Model Diversity:** LightGBM, XGBoost, Transformer, LSTM with different strengths
- **Meta-Learning:** Second-level models learning optimal combination strategies
- **Dynamic Weighting:** Market condition-aware ensemble weight adjustment
- **Confidence Integration:** Prediction uncertainty incorporation in ensemble decisions

2. Adaptive Ensemble Strategies

- **Market Regime Detection:** Volatility, volume, and sentiment-based regime classification
- **Time-Varying Weights:** Temporal adaptation of model contributions
- **Performance-Based Weighting:** Historical performance-driven weight adjustment

- **Bayesian Model Averaging:** Uncertainty quantification in ensemble predictions

D. Model Training and Optimization Strategy

1. Advanced Training Techniques

- **Mixed Precision Training:** FP16 optimization for memory efficiency on available hardware
- **Gradient Accumulation:** Effective batch size increase through gradient accumulation
- **Learning Rate Scheduling:** Adaptive learning rate with warmup and decay
- **Early Stopping:** Overfitting prevention with patience and performance monitoring

2. Multi-Task Learning Framework

- **Shared Representations:** Common feature extraction for multiple prediction tasks
- **Task-Specific Heads:** Specialized output layers for classification and regression
- **Loss Function Balancing:** Optimal weighting of different task losses
- **Performance Evaluation:** Multi-task performance assessment and optimization

E. Model Validation and Testing Framework

1. Comprehensive Evaluation Methodology

- **Time Series Cross-Validation:** Rigorous temporal validation preventing data leakage
- **Out-of-Sample Testing:** Reserved test set for unbiased performance assessment
- **Robustness Testing:** Performance evaluation across different market conditions
- **Ensemble Validation:** Individual model and ensemble performance comparison

2. Advanced Metrics and Analysis

- **Prediction Confidence:** Uncertainty quantification and confidence interval estimation
- **Feature Attribution:** Model interpretability through attention weights and SHAP analysis
- **Error Analysis:** Systematic study of prediction failures and model limitations
- **Business Impact:** Trading simulation with transaction costs and slippage

F. GPU Training Requirements and Colab Integration

1. Colab Training Strategy (Days 10-11)

- **Transformer Training:** BERT fine-tuning and multi-modal transformer training requiring GPU
- **Hyperparameter Optimization:** Efficient search using GPU acceleration
- **Ensemble Training:** Parallel training of multiple models with GPU resources

- **Model Validation:** Comprehensive testing and performance evaluation

2. Local-Colab Workflow

- **Development:** Architecture design and small-scale testing on MacBook Pro
- **Training:** Heavy computational tasks on Colab with GPU acceleration
- **Integration:** Model weights and results integration back to local environment
- **Deployment:** Final model packaging for inference on local hardware

G. Deliverables

- Advanced multi-modal transformer architecture with BERT integration
 - Enhanced LSTM models with attention mechanisms and advanced regularization
 - Sophisticated ensemble system with adaptive weighting and meta-learning
 - Comprehensive training and validation framework with GPU optimization
-

Day 12-13: Model Training and Integration

Objective: Train advanced models and integrate into comprehensive prediction system

A. Systematic Model Training Execution

1. Individual Model Training Schedule

- **Day 12 Morning:** Enhanced LightGBM and XGBoost with new features
- **Day 12 Afternoon:** BERT sentiment analysis pipeline training (Colab GPU)
- **Day 12 Evening:** Multi-modal transformer architecture training (Colab GPU)
- **Day 13 Morning:** Advanced LSTM variants training and optimization
- **Day 13 Afternoon:** Ensemble system training and meta-model development

2. Training Monitoring and Quality Control

- **Performance Tracking:** Real-time monitoring of training progress and metrics
- **Overfitting Detection:** Validation loss monitoring and early stopping implementation
- **Resource Management:** Memory usage and computational efficiency optimization
- **Error Handling:** Robust training procedures with automatic restart and checkpointing

3. Hyperparameter Optimization

- **Automated Search:** Grid search and Bayesian optimization for model parameters
- **Cross-Validation:** Nested CV for unbiased hyperparameter selection
- **Computational Budget:** Efficient allocation of training time across models
- **Performance Documentation:** Systematic recording of hyperparameter experiments

B. Advanced Model Integration Framework

1. Ensemble Architecture Implementation

- **Model Combination Logic:** Weighted averaging, voting, and stacking strategies
- **Dynamic Weight Optimization:** Market condition-based ensemble weight adjustment
- **Confidence Integration:** Prediction uncertainty incorporation in final decisions
- **Performance Monitoring:** Real-time ensemble performance tracking and adjustment

2. Multi-Task Learning Integration

- **Shared Feature Extraction:** Common representation learning across prediction tasks
- **Task-Specific Optimization:** Individual loss functions and performance metrics
- **Joint Training Strategy:** Simultaneous optimization of all prediction objectives
- **Transfer Learning:** Knowledge transfer between related prediction tasks

C. Model Performance Validation and Comparison

1. Comprehensive Performance Assessment

- **Individual Model Evaluation:** Standalone performance of each model architecture
- **Ensemble Performance:** Combined system performance vs. individual components
- **Baseline Comparison:** Performance improvement over Week 1 baseline models
- **Statistical Significance:** Formal testing of performance improvements

2. Advanced Evaluation Metrics

- **Classification Performance:** Accuracy, F1-score, AUC-ROC, precision-recall curves
- **Regression Performance:** RMSE, MAE, directional accuracy, correlation analysis
- **Business Metrics:** Sharpe ratio, maximum drawdown, profit factor estimation
- **Robustness Metrics:** Performance stability across different market conditions

D. System Integration and End-to-End Testing

1. Complete Pipeline Integration

- **Data Flow Validation:** End-to-end testing from raw data to final predictions

- **Feature Engineering Integration:** Seamless integration of new features with models
- **Prediction Pipeline:** Real-time prediction capability with appropriate latency
- **Error Handling:** Robust error recovery and graceful degradation

2. Performance Optimization

- **Computational Efficiency:** Optimization of inference time and memory usage
- **Scalability Testing:** Performance with larger datasets and extended time periods
- **Memory Management:** Efficient resource utilization for production deployment
- **Code Quality:** Refactoring and optimization of critical performance bottlenecks

E. Model Interpretability and Analysis

1. Feature Importance and Attribution

- **Global Importance:** Overall feature rankings across all models and tasks
- **Local Explanations:** Individual prediction explanations using SHAP and LIME
- **Attention Analysis:** Transformer attention weight interpretation and visualization
- **Business Insight:** Translation of model insights into actionable business understanding

2. Model Behavior Analysis

- **Prediction Confidence:** Understanding when models are confident vs. uncertain
- **Error Pattern Analysis:** Systematic study of when and why models fail
- **Market Condition Sensitivity:** Model performance under different market regimes
- **Temporal Stability:** Model behavior consistency over time

F. Deliverables

- Fully trained advanced model ensemble with optimized hyperparameters
- Comprehensive performance evaluation demonstrating improvements over baseline
- Integrated prediction system with end-to-end testing and validation
- Model interpretability analysis with business insights and recommendations

Day 14: Week 2 Integration and Performance Analysis

Objective: Finalize Week 2 developments and prepare comprehensive performance assessment

A. Final System Integration and Testing

1. End-to-End System Validation

- **Complete Pipeline Testing:** Verification of entire system from data input to predictions
- **Performance Consistency:** Ensuring reproducible results across multiple runs
- **Error Handling Validation:** Testing system robustness under various failure scenarios
- **Documentation Completeness:** Ensuring all components are properly documented

2. Production Readiness Assessment

- **Inference Performance:** Measuring prediction latency and computational requirements
- **Memory Efficiency:** Optimizing system for available hardware resources
- **Scalability Validation:** Testing with larger datasets and extended time periods
- **Deployment Preparation:** Packaging for easy deployment and maintenance

B. Comprehensive Performance Evaluation

1. Week 1 vs Week 2 Comparison

- **Statistical Testing:** Formal hypothesis testing of performance improvements
- **Effect Size Analysis:** Quantifying practical significance of improvements
- **Confidence Intervals:** Uncertainty quantification for performance metrics
- **Multiple Comparison Adjustment:** Proper statistical handling of multiple models

2. Advanced Performance Metrics

- **Risk-Adjusted Returns:** Sharpe ratio, Sortino ratio, maximum drawdown analysis
- **Prediction Quality:** Calibration analysis and prediction confidence assessment
- **Market Condition Performance:** Performance breakdown by volatility, volume, sentiment regimes
- **Temporal Robustness:** Performance consistency across different time periods

C. Model Analysis and Insights

1. Feature Contribution Analysis

- **Ablation Studies:** Individual feature group contribution assessment
- **Feature Interaction:** Analysis of feature combinations and synergies
- **Marginal Improvement:** Quantifying improvement from each new feature category
- **Business Value:** Translation of technical improvements into business impact

2. Model Behavior Understanding

- **Prediction Patterns:** Analysis of when models perform best and worst
- **Market Regime Adaptation:** How models adapt to different market conditions
- **Social Signal Integration:** Effectiveness of social media signal incorporation
- **Ensemble Contributions:** Individual model contributions to ensemble performance

D. Documentation and Knowledge Transfer

1. Technical Documentation Update

- **Architecture Documentation:** Complete system design and implementation details
- **API Documentation:** Function and class documentation with usage examples
- **Configuration Guide:** Parameter settings and tuning recommendations
- **Troubleshooting Guide:** Common issues and resolution procedures

2. Research Documentation

- **Methodology Documentation:** Detailed explanation of novel approaches and techniques
- **Experimental Results:** Comprehensive results analysis with statistical validation
- **Lessons Learned:** Key insights and recommendations for future development
- **Reproducibility Package:** Complete instructions for result reproduction

E. Week 3 Preparation and Planning

1. Performance Gap Analysis

- **Target Achievement:** Assessment of Week 2 goals and remaining gaps
- **Optimization Opportunities:** Identification of areas for further improvement
- **Technical Debt:** Areas requiring refactoring or optimization
- **Statistical Validation Needs:** Requirements for formal statistical testing

2. Week 3 Strategy Development

- **Statistical Testing Plan:** Comprehensive hypothesis testing and validation strategy
- **Optimization Priorities:** Focus areas for hyperparameter and ensemble optimization
- **Ablation Study Design:** Systematic analysis of component contributions
- **Business Impact Assessment:** Framework for quantifying practical value

F. Week 2 Deliverables Summary

- Advanced prediction system with 45+ new meme-specific features

- Multi-modal ensemble achieving 78%+ accuracy (target: >Week 1 + 5%)
 - Comprehensive performance analysis with statistical validation
 - Complete documentation package ready for Week 3 optimization
-

Week 1 & 2 Success Metrics

Week 1 Completion Criteria

- ☐ Successfully integrate 3 data sources with quality validation
- ☐ Generate 79 engineered features with comprehensive documentation
- ☐ Achieve 75%+ accuracy on direction prediction tasks
- ☐ Establish robust evaluation framework with time series CV
- ☐ Complete baseline model training with performance benchmarking

Week 2 Completion Criteria

- ☐ Implement 45+ advanced meme-specific features with validation
- ☐ Deploy multi-modal transformer and ensemble architectures
- ☐ Achieve 78%+ accuracy representing 5%+ improvement over Week 1
- ☐ Complete system integration with end-to-end testing
- ☐ Demonstrate statistical significance of improvements

Overall Technical Achievements

- **Data Pipeline:** Robust processing of 50,000+ Reddit posts and financial data
 - **Feature Engineering:** 124+ total features across social, financial, and cross-modal categories
 - **Model Performance:** >75% baseline accuracy with statistically significant improvements
 - **Architecture Innovation:** Multi-modal transformer and adaptive ensemble systems
 - **Code Quality:** Production-ready codebase with comprehensive documentation
-

Implementation Guidelines

Daily Schedule Recommendations

- **Morning (4-5 hours):** Core development and implementation work

- **Afternoon (2-3 hours):** Testing, validation, and documentation
- **Evening (1-2 hours):** Planning, research, and next-day preparation

Resource Management Strategy

- **Local Development:** MacBook Pro for development, testing, and analysis
- **GPU Training:** Colab for BERT fine-tuning and transformer training
- **Data Storage:** Local storage with cloud backup for important artifacts
- **Version Control:** Git repository with regular commits and branching

Quality Assurance Framework

- **Code Quality:** Regular refactoring, documentation, and testing
- **Performance Monitoring:** Continuous tracking of metrics and computational efficiency
- **Reproducibility:** Fixed random seeds, documented procedures, version control
- **Error Handling:** Robust error recovery and graceful degradation

This comprehensive plan provides detailed daily guidance for implementing a competition-winning meme stock prediction system within the first two weeks, establishing the foundation for advanced statistical validation and academic paper preparation in Weeks 3 and 4.