# FinSight AI

## *Directional Sentiment Pressure (DSP) Engine for Equity Markets*

## Abstract

FinSight AI is a lightweight, production-oriented sentiment intelligence system that converts raw global news into actionable institutional signals. The engine ingests financial news, regulatory filings, and market metadata, applies deep relevance filtering, processes sentiment through LLM-based models, and transforms extracted signals into a normalized Directional Sentiment Pressure (DSP) score in the range [−3,+3].
This document outlines the architecture, mathematical formulation, data pipeline, weighting logic, relevance filtering, storage schema, and score computation methodology behind the DSP engine.

## Table of Contents

## 1. Introduction

FinSight AI provides an end-to-end pipeline that transforms noisy, unstructured financial news into stable sentiment-based alpha signals. The DSP metric aggregates per-article directional impact into a normalized score that is invariant to article count, source distribution, and short-term sentiment imbalance.

The system was designed for:

- equity screening,
- intraday and swing-trading overlays,
- risk-aware decision support, and
- factor-model augmentation.

## 2. System Overview

FinSight AI consists of six stages:

1. **Data Ingestion**
2. **Relevance Filtering**
3. **LLM Sentiment Analysis**
4. **Dynamic Weighting**
5. **DSP Engine**
6. **Aggregation & Serving Layer**

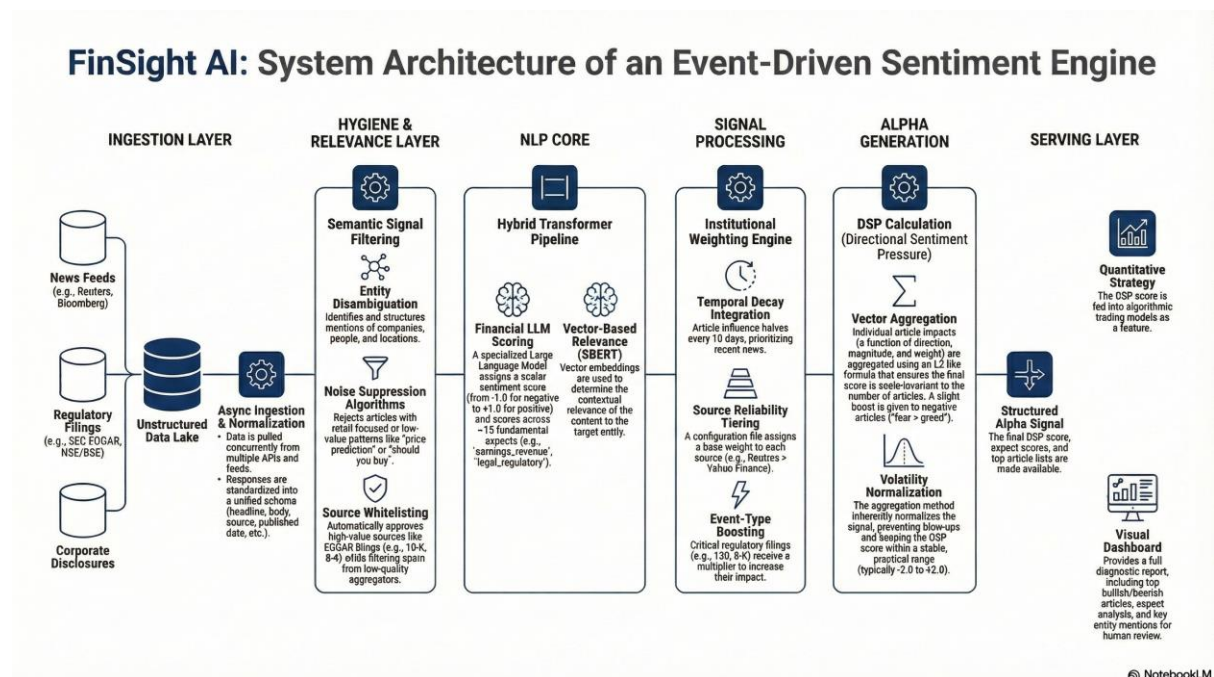Each stage is modular and independently replaceable.



**Figure 1:** End-to-End System Overview

## 3. Data Ingestion Layer

### 3.1 News Sources

Tiered ingestion from:

- Tier-1 media (Reuters, Bloomberg, WSJ)
- Tier-2 portals (Yahoo Finance, MarketWatch)
- Indian media (Moneycontrol, Economic Times, LiveMint)
- Alternative sources (Biztoc, Benzinga, aggregator feeds)

### 3.2 Regulatory Filings

Supported filings:

- SEC: 8-K, 10-K, 10-Q, 13D, 13G, Form-4
- NSE/BSE filings for Indian equities

### 3.3 Normalization

Every document is normalized into:

{ headline, body, source, published_at, symbol, entities }

### 3.4 LLM Summaries

Long-form articles are compressed using internal LLM summarization for consistent downstream text length.

## 4. Relevance Engine

The relevance engine removes noise aggressively before any scoring occurs.

### 4.1 Normalization

Lowercasing, punctuation removal, entity cleanup.

### 4.2 Company Matching

- Symbol match ("AAPL", "TSLA")
- Name match ("Tesla", "Apple Inc.")
- Alias match (e.g., "Alphabet" → GOOG)

### 4.3 Noise Filtering

Articles are discarded if:

- they exhibit SEO-spam patterns ("boosts stake", "should you buy"),
- they stem from low-quality sources (certain aggregator spam),
- they describe minor portfolio churn irrelevant to fundamentals.

### 4.4 Filing Override

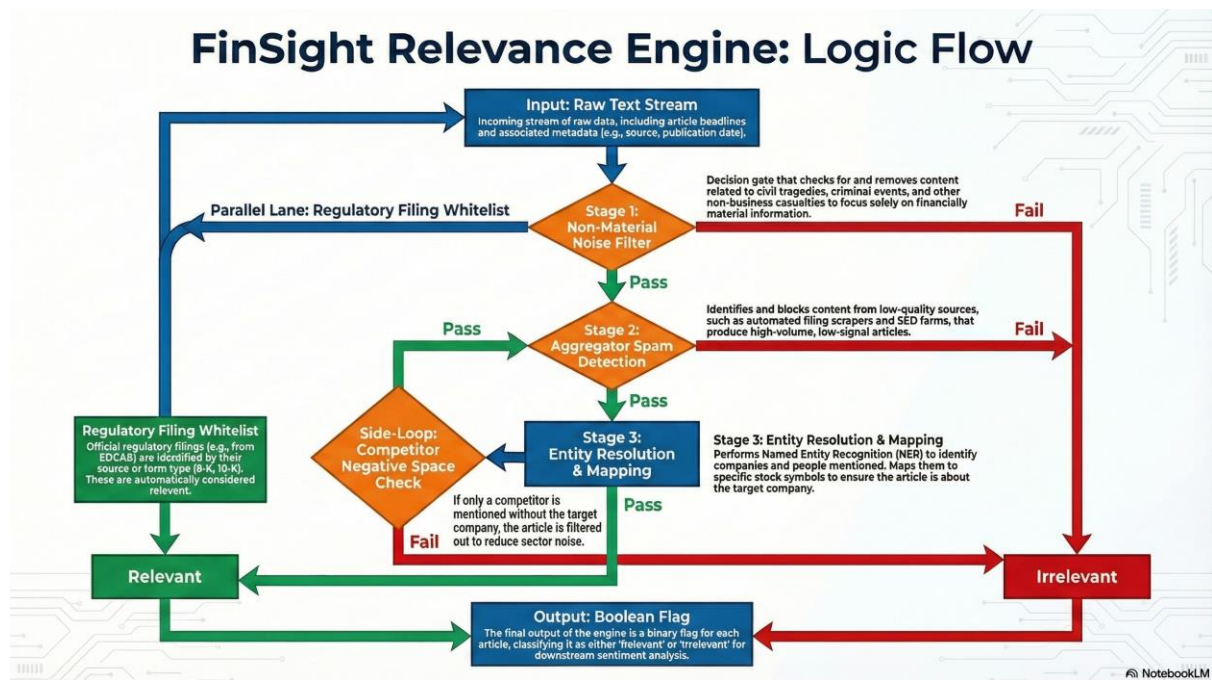Regulatory filings bypass relevance filters and are always included.

**Figure 2:** Relevance Engine

## 5. NLP & Sentiment Layer

### 5.1 LLM Sentiment

A fine-tuned sentiment model outputs:

sentiment_score $\in$ [-1.0, +1.0]

### 5.2 Aspect Sentiment

Fifteen business dimensions are computed:

- earnings & revenue
- costs & margins
- liquidity
- dividends & buybacks
- leadership
- strategy
- M&A
- legal & regulatory
- competition
- product
- technology
- supply chain
- customer experience
- brand
- labor force

### 5.3 Entity Extraction

NER (Named Entity Recognition) models extract organizations, persons, and locations.

# 6. Dynamic Weighting Engine

Each article receives a numerical weight:

final_weight = ( base_source_weight
        × recency_factor
        × engagement_factor
        × filing_adjustment)

### 6.1 Base Source Weight

Loaded from internal trust model. (Redacted in public repo.)

### 6.2 Recency Decay

Non-linear temporal decay calibrated to the swing trading horizon

### 6.3 Engagement Factor

Engagement impact is normalized using a logarithmic scaling function with a proprietary volatility dampening coefficient to marginalize viral noise while preserving signal magnitude.

# 7. DSP Computation

A normalized, direction-weighted, magnitude-adjusted aggregation of validated news sentiment over a fixed time window.

DSP is clipped to [−3.0, +3.0].

# 8. Data Storage Schema (SQLite)

The article table includes:

| Field Name | Data Type |
|------------|-----------|
| id | INTEGER PRIMARY KEY |
| symbol | TEXT |
| headline | TEXT |
| summary | TEXT |
| source | TEXT |
| published | TEXT (ISO-8601) |

| Field Name | Data Type |
|---|---|
| form_type | TEXT |
| entities | JSON |
| sentiment_score | REAL |
| sentiment_label | TEXT |
| aspect_sentiment | JSON |
| engagement | INTEGER |
| weight | REAL |
| weighted_score | REAL |
| fetch_date | TEXT |



**Figure 3:** SQLite Schema

## 9. Implementation Notes

- System is fully asynchronous (asyncio)
- Caching avoids API throttling
- LLM calls are rate-limited
- Normalization ensures stability across volatile news cycles
- L2 normalization protects against volume imbalances

## 10. Limitations

- Non-English sources are minimally supported
- Model depends on LLM quality
- Corporate synonym dictionary must be maintained manually

## 11. Future Work

- Multi-asset joint sentiment modeling
- Cross-symbol entity-flow mapping
- Predictive regressions for return forecasting
- Regime-aware DSP scaling

## 12. References

- SEC EDGAR documentation
- Financial NLP research papers
- Entity extraction models