



White Paper

# EntitySense Intelligence

The next-gen entity analysis platform

Dilip Singh, Sukanksha Totade  
Equifax D&A

July 8, 2025

# Contents

EntitySense Model .....	3
Background.....	4
How Adverse Media Screening Works.....	5
Introduction .....	6
Motivation .....	7
Solution.....	8
How language models work .....	8
EntitySense Model Framework .....	9
File reading and text extraction .....	9
Named entity recognition .....	9
Sentence splitting .....	10
Sentiment analysis .....	10
Crime classification .....	10
Entity verification .....	10
Data aggregation .....	10
EntitySense Model Module Performances .....	11
Name entity recognition solution.....	11
Sentiment analysis solution.....	13
Summary & Conclusion.....	14

## EntitySense Model

This Paper introduces EntitySense Intelligence System, a next-generation entity analysis platform designed to empower organizations with actionable insights from unstructured text data. This **comprehensive solution provides named entity recognition to identify and categorize entities like names and locations, dates, organizations, sentiment analysis to gauge emotional tone, criminal recognition for risk assessment, and identity verification for security purposes.** By integrating these diverse functionalities, EntitySense offers a holistic approach to data analytics, converting complex unstructured text into easily interpretable, structured information for better decision-making.

- Model boasts over **90% accuracy** in entity, sentiment, and crime detection, whereas humans have achieved only a 15% rate, according to our own research sample.
- The model can process **approximately 6,171 documents per hour**, each averaging 3,600 words, dwarfing average human reading speeds of 200-300 words per minute.<sup>1</sup>
- Our solution offers **unbiased** assessments across **17+ languages**, while the average person is fluent in merely 1-2 languages.
- Enabled with **concurrent processing**, the model offers **near-instantaneous insights**, empowering decision-making in real-time.
- The model ensures standardized output quality regardless of volume, **eliminating the inconsistencies** of human interpretation.
- Our framework is prepped for **evolution**. With deep learning codes set for NER and sentiment analysis, a simple import statement adaptation ensures integration with the latest models, positioning us **ahead in the AI curve**.

<sup>1</sup> <https://www.sciencedirect.com/science/article/abs/pii/S0749596X19300786>

## Background

Adverse media screening, or negative media detection, has become an important aspect of financial institutions to identify any potential financial crime risks. According to the Association of Certified Anti-Money Laundering Specialists (ACAMS), nearly half of all organizations use adverse media screening as a part of their due diligence. Out of those organizations, **37% have detected adverse media information that was not detected through any other due diligence methods.**

In this proliferated connected world, information travels with such a high speed that an adverse media report can turn into a global compliance concern within a few minutes. Adverse media can cause profound financial, operational, and reputational damage.

According to 2021 Boston Consulting Group (BCG)'s report, many financial institutions globally have paid more than \$321 billion in fines since the 2007-2008 financial crisis for non-compliance with anti-money laundering, know your customer, and sanctions regulations.

Many research studies say that public companies undergo a dip in stock prices in the event of adverse media reports concerning security and financial frauds. According to a 2012 Association of Certified Fraud Examiners survey that investigated cases between January 2010 and December 2011, organizations around the world **lose an estimated 5% of their annual revenues to fraud.**

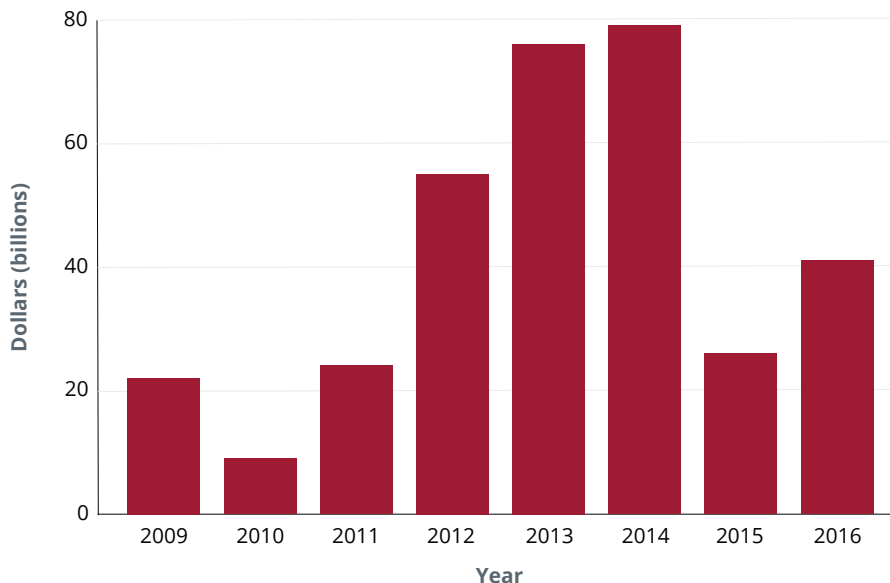
Consumer trust is another casualty. A survey by the Edelman Trust Barometer reveals that a high percentage of consumers lose trust in a brand after negative media exposure, directly impacting sales and customer retention.

Financial institutions globally have paid more than \$321 billion in fines since the 2007-2008 financial crisis for non-compliance with anti-money laundering.

- 2021 BCG Report

Figure 1: Bank penalties over the years

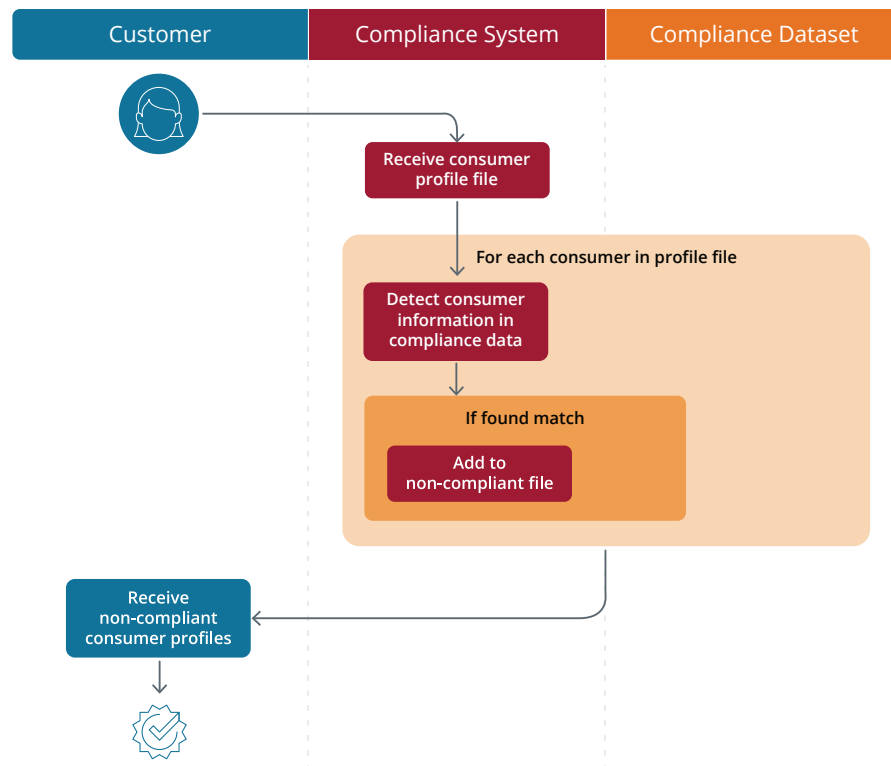
Financial institution paid \$321 billion in charges since financial crisis



### How Adverse Media Screening Works

- The customer sends a portfolio file/files with details of consumers.
- The system search is initiated to all trustable data sources .
- The system sends the search results.
- Sentiment analysis, risk detection and crime detection are identified manually by humans/analysts.
- The result is sent to the customer.

Figure 2: Adverse media screening process



## Introduction

In this rapidly evolving landscape of regulatory compliance, the ability to quickly and accurately identify risk using adverse media screening is necessary to support corporate integrity and adherence to legal standards. **Traditionally and even today adverse media screening is done manually, which is often inadequate.**

Currently, Equifax has a compliance process that screens criminal records. The process aims to find out if we should do business with an individual or organization based on adverse media. Without this process, our customers could spend billions of dollars per year to combat financial crime like money laundering and Ponzi schemes.

The emergence of Natural Language Processing (NLP) technologies has revolutionized adverse media detection, providing companies with a critical edge in compliance efforts. This whitepaper talks about our innovative approach to solve compliance adverse media screening problems. Our technology includes a Name Entity Recognition (NER) system, an advanced text analysis framework employing Textblob, and a cutting-edge sentiment analysis model, all underscored by a specialized crime detection algorithm. Together, these robust mechanisms are designed to anatomize, understand, and identify relevant information.

The customized NER is the anchor of our solution, which identifies and extracts not only main entities but relevant entities associated with compliance risks from unstructured text. It is fine-tuned to detect subtle references and connections to illegal activities, individuals, or entities that may indicate potential compliance infringements. With NER we have a textblob module that offers subtle understanding of text, adept at capturing complex nuances that may elude conventional screening methods. When amalgamated with sentiment analysis, this system provides a multifaceted examination of media tone, intent, and context, which is critical for distinguishing between neutral mentions and genuinely adverse implication.

Moreover, our crime detection algorithm is the final layer of inspection, scanning extracted entities and sentiment evaluations to identify potential criminal activities or connections. This algorithm is a vital tool for compliance departments, ensuring that organizations stay ahead of potential risks by preemptively addressing issues that could otherwise escalate into significant legal and reputational damage.

In a nutshell, our comprehensive solution for compliance adverse media screening leverages the combined strengths of bespoke NER, Textblob analysis, sentiment analysis, and crime detection of an entity, representing a significant leap forward in the domain of regulatory compliance and risk management.



Traditionally and even today adverse media screening is done manually, which is often inadequate.

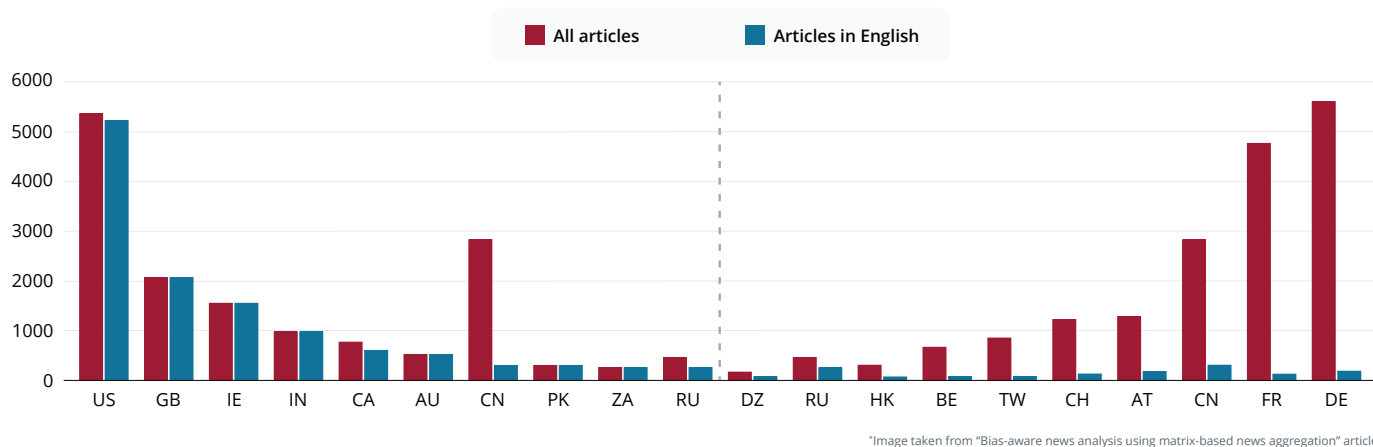
## Motivation

The current process for adverse media screening involves pulling credible media and news articles from wall street journal and Don Jones Factiva on a daily basis. The articles are then loaded into a system called Graycon that searches and highlights the keywords from a predefined keyword list including offense, action and penalty words (CDC Manual). Next, the highlighted articles are passed to analysts for screening.

For each article, the analyst records the PII and article abstract with respect to the keywords and marks the article as acceptable, duplicated, or irrelevant. The optional last step is to confirm the information extracted from those articles with official criminal records.

**This compliance process is highly manual, averaging 3k articles processed per day.** However, in this era of massive data, there are ~3 million news articles published per day globally, not to mention Equifax itself takes in ~30K adverse media per day.

Figure 3: Adverse media screening of articles



We wanted to save analysts time and resources by automating the process using the state of art language models. The problem seems simple, but in reality it's very complex.

First we need to identify the entity or entities from the news article. If there are two entities, we need to decipher if they are a person or a company/organization. Then we must extract other PIIs belonging to each entity and accurately identify the sentiment of that entity.

From there, we have to answer a few questions, such as: If the sentiment is negative does that mean we cannot do business with that entity? What if the article is subjective? We then need to take into account crimes related to the entity, both financial and non-financial.

## Solution

Equifax developed a product named EntitySense that enables our customers, like banks or broker dealer organizations, to screen and monitor their consumers using our compliance database that this automated process generated.

Our goal was to build an automated process with low false positive rates, good feedback, and high business performance. We made sure that decisions are mutable, traceable, and explainable. Unlike the current manual process, EntitySense pulls the credible media and news articles without human intervention and feeds the articles into our language models. Our language models then extract PII information from each full article.

## How language models work

To show how our language models work, we will look at an example. This article is from [Wikipedia](#).

*Bernard Lawrence Madoff (/ˈmeɪdɔːf/ MAY-dawf;<sup>[2]</sup> April 29, 1938 – April 14, 2021) born in NYC, NY, U.S. was an American financial criminal and financier who was the admitted mastermind of the largest known Ponzi scheme in history, worth an estimated \$65 billion.<sup>[3][4]</sup> He was at one time chairman of the Nasdaq stock exchange.<sup>[5]</sup> Madoff's firm had two basic units: a stock brokerage and an asset management business; the Ponzi scheme was centered in the asset management business.*

The model will extract:

Name	Bernard Lawrence Madoff
Person or organization	Person
Important dates	April 29, 1938; April 14, 2021
Important addresses	NYC, NY, U.S.
Other entities found	Nasdaq stock exchange
Sentiment	Positive
Crime committed	Yes
Financial crime	Ponzi scheme



## EntitySense Model Framework

In this section, we will describe how our model framework functions step-by-step.

Figure 4: Our model framework, step-by-step

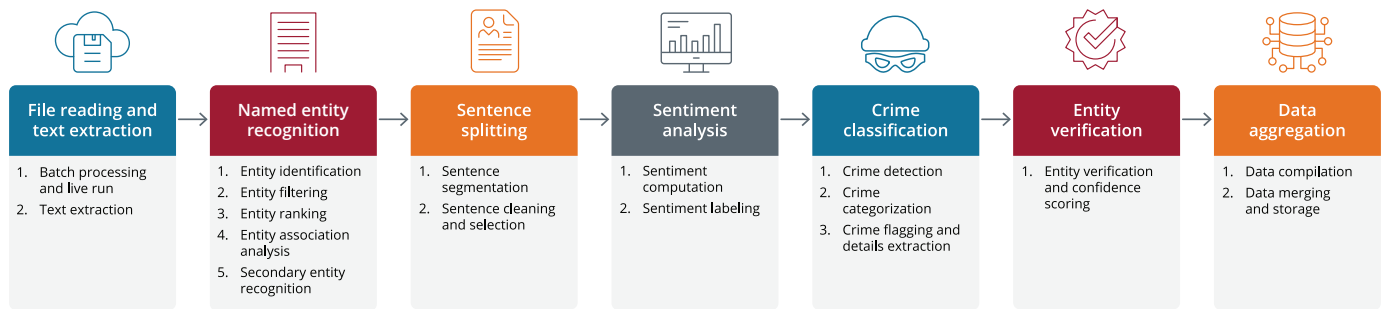


Figure 5: Our model uses the below equation to get the entities

$$Entity = NERScore_{CustomizedNER} ( ContentExtractor_{NLP} SentenceCleaning_{NLP} ( SentimentScore_{NLP} SocialMedia_{+} CrimeDetectorScore_{CustomizedNER} + EntityVerification_{ScoringModel} ) )$$

### File reading and text extraction

#### Step 1: Batch processing and live run

In this initial step, the system identifies the optimal way to process files, be it in a batch process for large quantities of data or a live run for real-time analysis. Moreover, the feature to directly extract data from document links enhances the flexibility of the data acquisition phase.

#### Step 2: Text extraction

Leveraging advanced text mining techniques, the system efficiently extracts text from the sourced files. This step ensures the raw data is prepared and primed for the subsequent stages of data analysis and processing.

### Named entity recognition

#### Step 3: Entity identification

Utilizing powerful NLP capabilities, the system meticulously scans each article to identify and extract various named entities, fostering a rich dataset for deeper analysis.

#### Step 4: Entity filtering

To streamline the data, any entity identified as a character of a movie is filtered out, thereby focusing the analysis on more relevant data segments.

#### Step 5: Entity ranking and separation

A sophisticated algorithm ranks and lists the top five entities, distinctly categorizing them into persons and organizations. Furthermore, the algorithm discerns the main person and organization entities, paving the way for detailed analysis.

#### Step 6: Entity association analysis

This step embarks on an intricate process where the main identified entities are analyzed in tandem with associated important addresses and dates, creating a comprehensive profile for each entity.

#### Step 7: Secondary entity recognition

The system also recognizes secondary entities, categorized under "other\_entities", thus providing a holistic view of the data spectrum.

## Sentence splitting

### *Step 8: Sentence segmentation*

Deploying powerful NLP context understanding methods, the system divides the article into individual sentences, creating a structured format that facilitates efficient processing in the subsequent steps.

### *Step 9: Sentence cleaning and selection*

Post-segmentation, sentences undergo a cleaning process to remove noise and are meticulously selected if they pertain to the main entity identified in the earlier steps.

## Sentiment analysis

### *Step 10: Sentiment computation*

The refined data is then channeled into the sentiment analysis engine, where it calculates a polarity sentiment score, providing a quantitative measure of the sentiment expressed in the textual data.

### *Step 11: Sentiment labeling*

This phase involves the generation of sentiment labels, offering analysts an intuitive understanding of the sentiment landscape pertaining to the entities in question.

## Crime classification

### *Step 12: Crime detection*

In this critical stage, the system utilizes a bespoke crime classifier to identify potential references to criminal activities within the data set.

### *Step 13: Crime categorization*

Following identification, crimes are categorized into financial and non-financial buckets, facilitating detailed analysis and reporting.

### *Step 14: Crime flagging and details extraction*

If the entity is associated with crimes that are of interest, the system flags it appropriately and extracts detailed information about accusations and criticisms, if any.

## Entity verification

### *Step 15: Entity verification and confidence scoring*

Here, the system cross-verifies the main entity with internal databases to ascertain the accuracy of the identification, supplementing this with a confidence score to gauge the reliability of the match.

## Data aggregation

### *Step 16: Data compilation*

This phase marks the convergence of all processed data streams, aggregating them to form a comprehensive data set.

### *Step 17: Data merging and storage*

Finally, the aggregated data is seamlessly merged into existing storage mediums such as CSV files, Excel files, or BigQuery tables, thus completing the data processing pipeline and readying the data for insightful analysis and reporting. This meticulous and multi-faceted approach ensures a high level of accuracy and depth in data analysis, providing a robust foundation for insightful decision-making and reporting.

## EntitySense Model Module Performances

Now, let's dive into the more technical details of ensemble models. We are dealing with two NLP problems: name entity recognition (NER) sentiment identification. There is a third problem, which is entity verification, but it is a simple classification problem.

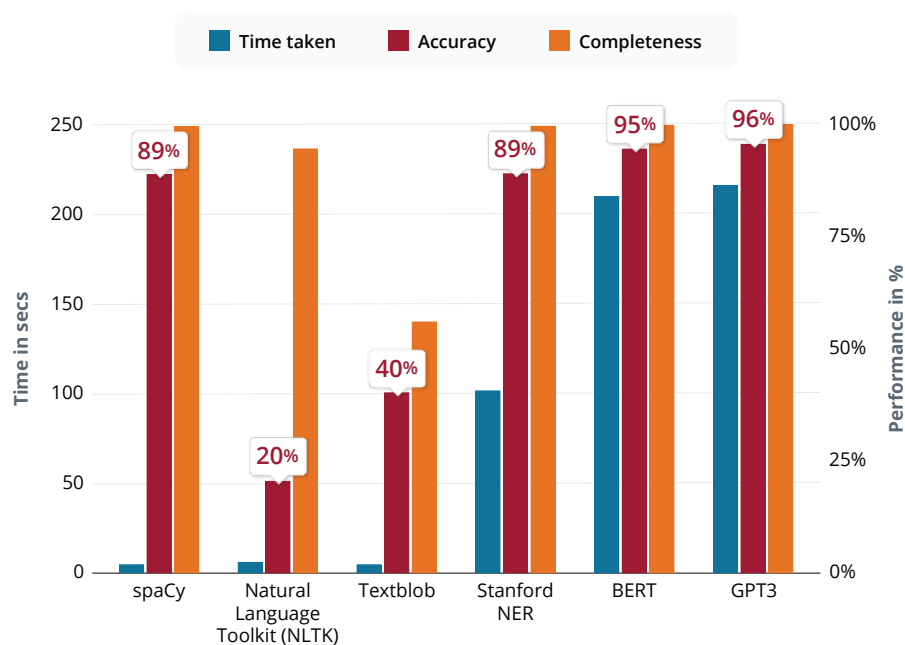
### Name entity recognition solution

KPI definitions:

- **Accuracy:** This is a score based on the correctness of the entities recognized. For this, we are comparing ground truth or gold standard to the outputs. The name is used as the ground truth.
- **Completeness:** Measures the depth of information captured by each model for a recognized entity. E.g., for Leonardo Wilhelm DiCaprio, if one model captures just the name, while another captures name, address, and DOB, the latter is more complete.
- **Multilingual:** Processing multiple language documents.

Figure 6: NER model performances

All pretrained models



During our experiments, we used various methods to solve NER and sentiment analysis problems. A quick overview on our data: we had 100 articles written in English with ~6500 average words. The following table provides a comparative overview of different NER models that are measured using accuracy, completeness and multilingual metrics. All these models are pre-trained.

- **spaCy:** This is a speedy performer, clocking in at 5 seconds. It provides an accuracy of 89% and impressively captures all personally identifiable information (PIIs). It also works across multiple languages.
- **Natural Language Toolkit (NLTK):** While taking a similar time as spaCy, its accuracy dips to 20%. It mainly identifies names, making it a more superficial extractor.
- **Textblob:** A bit faster at 4 seconds, it doubles the accuracy of NLTK to 40%. But, like NLTK, it focuses mainly on names.

Transitioning from these relatively simple models, we delve into more complex ones.

- **Stanford NER:** This model showcases the power of deeper models, with an accuracy comparable to spaCy. However, its time consumption of 101 seconds indicates its thorough processing. It's comprehensive and multilingual.
- **BERT and GPT3:** These represent the zenith in deep learning models for NER. While they take longer (around 3.5 minutes), their accuracies of 95% and 96% speak for their efficacy. Their depth is evident as they capture all PIs and operate in multiple languages.

In essence, as we progress from simpler to deep learning models, there's a marked improvement in performance, depth, and accuracy, albeit at the cost of processing time.

#### Why Spacy over BERT/Transformer/Complex models?

Models	Spacy	Deep Learning/BERT/Transformer/....
Efficiency and speed	<ul style="list-style-type: none"> <li>• Fast and efficient processing for many NLP tasks.</li> <li>• Allows for easy and efficient text preprocessing.</li> </ul>	<ul style="list-style-type: none"> <li>• Slower at inference, especially on CPU, because of the large model.</li> <li>• Require significant computational resource.</li> </ul>
Ease of use	<ul style="list-style-type: none"> <li>• Provides a simple, high-level API.</li> <li>• Comprehensive documentation and a large community provide ample resource.</li> </ul>	<ul style="list-style-type: none"> <li>• Working directly with transformer models like BERT can be more complex and may require a deeper understanding of the underlying model architecture.</li> </ul>
Memory footprint	<ul style="list-style-type: none"> <li>• Generally has a smaller memory footprint, making it more suitable for edge devices or environments with limited computational resources.</li> </ul>	<ul style="list-style-type: none"> <li>• Models generally have a larger memory footprint, requiring more memory and computational resources.</li> </ul>
Customizability and training	<ul style="list-style-type: none"> <li>• Allows for customization and the integration of other models, but may not provide as much flexibility as working directly with transformer models.</li> <li>• Provides pre-trained models for various languages that are optimized for a variety of NLP.</li> </ul>	<ul style="list-style-type: none"> <li>• Provides more flexibility for fine-tuning and customizing models for specific tasks or domains.</li> <li>• State-of-the-art performance on many NLP tasks, but requires more data and computational resources for training.</li> </ul>
Pretrained models and language support	<ul style="list-style-type: none"> <li>• Offers pre-trained models for various languages and tasks, but may not cover as many languages as some transformer models.</li> </ul>	<ul style="list-style-type: none"> <li>• Many pre-trained models are available for a wide variety of languages and tasks, offering state-of-the-art performance on many benchmarks.</li> </ul>
NLP task coverage	<ul style="list-style-type: none"> <li>• Provides out-of-the-box support for many common NLP tasks, such as tokenization, part-of-speech tagging, named entity recognition, and dependency parsing.</li> </ul>	<ul style="list-style-type: none"> <li>• Excels in tasks such as text classification, question answering, and language generation, but may require more setup for specific tasks.</li> </ul>

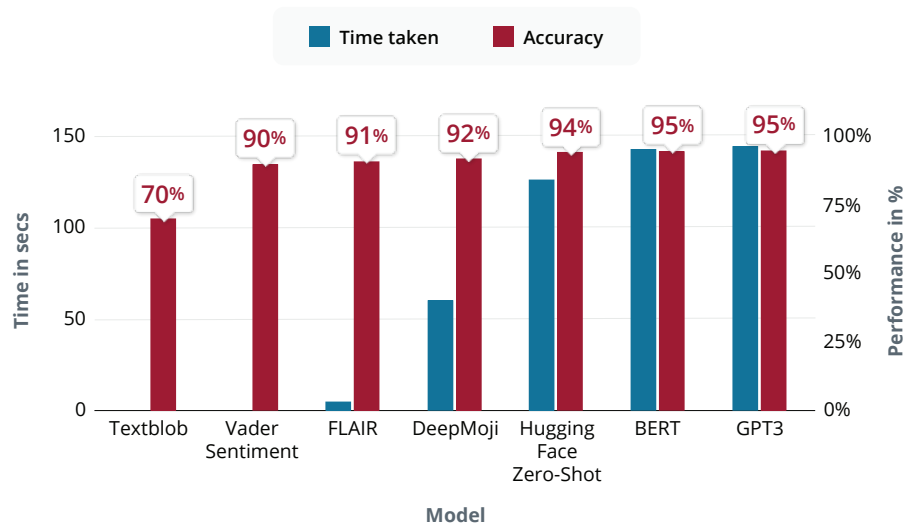
## Sentiment analysis solution

KPI definition:

- **Accuracy:** This is a score based on the correctness of the entities recognized. For this, we are comparing ground truth or gold standard to compare the outputs. The name is used as the ground truth.

Figure 7: Sentiment analysis model performances

All pretrained models



This section showcases a comparative analysis of various sentiment analysis pre-trained models.

- **Textblob:** It's one of the fastest models, completing its task in approximately 1 second. Although it has a 90% accuracy rate, it struggles with intricate nuances in sentiment, possibly overlooking subtleties.
- **Vader Sentiment:** Another swift model at around 1 second. While it shares the same accuracy as Textblob, it's particularly optimized for social media content, making it adept at discerning sentiments in tweets or posts.

Diving into more sophisticated models.

- **FLAIR:** With a slightly elongated processing time of 5 seconds, it boasts a commendable accuracy of 91%. Its underlying mechanism is a character-level LSTM network, allowing it to deeply understand the text structure and sentiment.
- **DeepMoji:** Taking a minute to process, its accuracy hovers at 92%. Unique in its approach, it's trained on emojis, decoding the emotional subtext behind them.

Moving into state-of-the-art deep learning.

- **Hugging Face Zero-Shot:** Clocking in at 2 minutes, it achieves an impressive 94% accuracy. Its strength lies in using transformers and contextual embeddings, which allow it to grasp the underlying sentiment in diverse contexts.
- **BERT and GPT3:** Both are high-end models that demand around 2 minutes but ensure a striking 95% accuracy. They utilize similar technologies, with GPT3 mirroring BERT's strategies.

In summation, as we progress from basic models to deep learning ones, there's a clear increment in understanding and accuracy. While the processing time might increase, the depth and precision these models offer are unparalleled.

## Why Vader over BERT/Transformer/Complex models?

Models	Vader Sentiment	BERT/Transformer/....
Simplicity and efficiency	<ul style="list-style-type: none"><li>• It is a lexicon and rule-based sentiment analysis tool that is specifically attuned to sentiments expressed in social media.</li><li>• It's fast and requires minimal resources since it does not require GPU and heavy computational power.</li></ul>	<ul style="list-style-type: none"><li>• More complex models that take longer time for training and inference.</li></ul>
No need for training	<ul style="list-style-type: none"><li>• VADER is pre-trained and does not require additional training data, making it ready to use.</li></ul>	<ul style="list-style-type: none"><li>• May require fine-tuning on a specific dataset to achieve optimal performance for sentiment analysis, necessitating additional data and training time.</li></ul>
Interpretability	<ul style="list-style-type: none"><li>• Being a rule-based model, VADER offers higher interpretability, allowing users to understand how the sentiment scores are derived.</li></ul>	<ul style="list-style-type: none"><li>• As deep learning models, they are generally considered black boxes, making it harder to interpret their decisions.</li></ul>
Handling of short text	<ul style="list-style-type: none"><li>• Due to its efficiency, VADER is suitable for real-time sentiment analysis.</li></ul>	<ul style="list-style-type: none"><li>• May not be as suitable for real-time analysis due to higher computational requirements.</li></ul>
Low resource requirement	<ul style="list-style-type: none"><li>• Requires minimal computational resources, making it more accessible for small to medium-scale projects.</li></ul>	<ul style="list-style-type: none"><li>• Require significant computational resources and memory.</li></ul>
Real-time analysis	<ul style="list-style-type: none"><li>• Due to its efficiency, VADER is suitable for real-time sentiment analysis.</li></ul>	<ul style="list-style-type: none"><li>• Require significant computational resources and memory.</li></ul>

## Summary & Conclusion

This whitepaper has presented a detailed analysis of the complexities around compliance adverse media detection and outlined a sophisticated solution leveraging customized NER, textblob, sentiment analysis, and crime detection algorithms. We looked at what is happening today with a microscope and demonstrated how our approach can address limitations.

The world of compliance and risk management is an evolving landscape, with adverse media posing significant threats to organizational integrity, operational, and financial stability. Our product analysis indicates that using NLP and AI technologies to enhance adverse media screening not only streamlines the compliance process but also significantly reduces oversight risk.

EntitySense's customized NER efficiently extracts primary entities, relevant entities to the primary, and other entities present in the article. Textblob gives us a high advantage in text interpretation and sentiment analysis, which helps to identify the tone and intent of an article for the entity in question. The crime detection algorithm adds a crucial layer, vetting the identified risks for potential criminal implications. This piece makes our model less biased as we provide fact rather than subjectivity of the article, leading to better business decisions.

Moving forward, we will be refining and training our models to keep pace with the evolving nature of language used in media. As entities and regulatory bodies increasingly acknowledge the significance of robust compliance mechanisms, the integration of advanced technologies like ours will become a staple in risk management strategies.