



# A Comparative Study of Machine Learning Algorithms for Document Classification: Insights Beyond Accuracy

Ankit Mahato<sup>1</sup> and Udit Kumar Mahato<sup>1</sup>

<sup>1</sup>Sunway International Business School, Birmingham City University, Kathmandu, Nepal

## Abstract

Document classification remains a fundamental task in natural language processing with applications across diverse domains. While state-of-the-art transformer-based models like BERT demonstrate superior accuracy, their computational and environmental costs often outweigh marginal performance gains. This study provides a comprehensive comparative analysis of classical machine learning models (SVM, Naïve Bayes, Logistic Regression), deep learning architectures (CNN, RNN), and transformer-based models (BERT, DistillBERT, LayoutLM). We evaluate models across multiple dimensions including accuracy, training efficiency, energy consumption, robustness, and interpretability. Our results highlight the continued relevance of classical models in resource-constrained environments and introduce an Energy-Adjusted Score (EAS) to balance performance and cost. The study also acknowledges key limitations, including benchmark dataset quality issues and potential variability in energy measurements across hardware platforms, to ensure transparency in interpreting results.

**Keywords:** Document Classification, Machine Learning, Energy Efficiency, BERT, SVM, Interpretability, Deep Learning, Transformer Models, Sustainability, Text Mining.

## 1. Introduction

The digital transformation of organizations has resulted in unprecedented volumes of textual data requiring automated processing and categorization. Document classification, a fundamental task in natural language processing, has evolved from rule-based systems to sophisticated machine learning approaches capable of handling complex, multi-modal document structure. With textual data production rates measured in petabytes daily, the demand for efficient document classification systems spans critical applications including email filtering, news curation, legal e-discovery, and sentiment monitoring [1].

Traditional evaluation frameworks have predominantly emphasized accuracy metrics, often overlooking crucial practical considerations such as computational efficiency, energy consumption, interpretability, and robustness to distributional shifts. This narrow focus has led to the development of increasingly complex models that, while achieving marginal accuracy improvements, impose significant computational and environmental costs [2][3]. Recent studies indicate that training large language models can consume energy equivalent to powering entire households for extended periods, raising important questions about the sustainability of current research trajectories[4][5].

The benchmark landscape for document classification has also faced scrutiny, with recent analysis revealing fundamental issues in widely-used datasets. The RVL-CDIP benchmark, considered the de facto standard for document classification evaluation, has been found to contain substantial label noise (8.1%), significant test-train overlap (approximately 32%), and concerning amounts of sensitive personally identifiable information [6]. These findings challenge the validity of reported performance metrics and highlight the need for more robust evaluation methodologies.

Furthermore, the emergence of out-of-distribution (OOD) eval-

uation protocols has revealed significant performance degradation when models encounter documents from different distributions than their training data. Studies demonstrate accuracy drops of 15-30% when state-of-the-art models are evaluated on out-of-distributions, raising questions about the generalizability of current approaches [7].

This comprehensive study addresses these limitations by conducting a holistic evaluation of representative machine learning algorithms across multiple dimensions. Our analysis encompasses classical approaches (Support Vector Machines, Naive Bayes, Logistic Regression), and transformer-based models (BERT, DistillBERT, RoBERTa). We evaluate these approaches not only on accuracy but also on energy consumption, training efficiency, interpretability, and robustness to distributional shifts.

To address the trade-off between accuracy and sustainability, we introduce the Energy-Adjusted Score (EAS), a metric that integrates model performance with energy cost. While comprehensive, the study is constrained by known limitations, such as benchmark dataset quality issues and variability in energy measurements across different hardware environments, which are examined in detail in Section V-E.

## 2. Related Work

### 2.1. Classical Machine Learning Approaches

Support Vector Machines (SVMs) have historically demonstrated exceptional performance on high-dimensional sparse text data, particularly when optimized with n-gram features and TF-IDF weighting schemes [8][9]. Linear SVMs, utilizing cutting-plane algorithms, achieve competitive accuracy while maintaining computational efficiency with  $O(s \cdot n)$  complexity, where 's' represents the number of support vectors and 'n' the feature dimensionality

[10]. The effectiveness of SVMs in document classification stems from their ability to handle high-dimensional features spaces and their robust performance on linearly separable data, characteristics commonly found in text classification tasks.

Naïve Bayes classifiers, despite their simplistic independence assumptions, continue to serve as strong baselines for text classification due to their linear scalability and interpretability [7][9]. A multinomial variant of Naïve Bayes has shown particular effectiveness in document classification, with proper hyperparameter tuning ( $\alpha$  [1e-3, 1e-1]) being crucial for optimal performance on imbalanced datasets. The algorithm's probabilistic nature allows for transparent decision-making processes, making it particularly valuable in applications requiring interpretability.

Logistic Regression with stochastic gradient descent represents another classical approach that remains competitive in document classification tasks [6][11]. When paired with feature selection techniques such as  $X^2$  statistics or mutual information, logistic regression can achieve performance comparable to more complex models while maintaining computational efficiency and interpretability advantages.

## 2.2. Deep Learning Architecture

Convolutional Neural Networks (CNNs) have been successfully adapted for text classification by treating documents as sequences and applying convolutional filters to capture local n-gram patterns[12][13]. CNNs typically achieve accuracy rates around 74% on standard benchmarks like 20 Newsgroups, demonstrating their ability to capture local textual features effectively. However, their performance generally falls short of transformer-based approaches, particularly on tasks requiring long-range dependencies.

Recurrent Neural Networks, particularly Long Short-Term Memory (LSTM) and Bidirectional LSTM variants, excel at capturing sequential dependencies in text data[14][15]. These architectures are particularly effective for longer documents where word order and context are crucial for accurate classification. However, their sequential nature leads to increased training times and computational costs compared to parallelizable architectures [14][15].

Hybrid approaches combining CNNs and RNNs have shown incremental improvements by leveraging both local feature detection and sequential modeling capabilities [14][15]. However, these gains are often marginal compared to the additional complexity introduced, and they rarely match the performance of transformer-based models.

## 2.3. Transformer-Based Models

The introduction of transformer architecture revolutionized document classification, with BERT and its variants achieving state-of-the-art performance across multiple benchmarks [16]. BERT's bidirectional nature and attention mechanisms enable comprehensive understanding of document context, leading to accuracy improvements of 1-4 percentage points over classical approaches. However, BERT's computational requirements are substantial, with pre-training consuming approximately 188 kWh of energy [10]. DistillBERT, a distilled version of BERT, offers a compelling trade-off between performance and efficiency. Maintaining 97% of BERT's language understanding capabilities while being 40% smaller and 60% faster, DistillBERT demonstrates the potential of knowledge distillation techniques for creating more efficient models. The distillation process enables the smaller model to learn from the larger model's soft probabilities, resulting in improved performance compared to training from scratch.

Specialized document understanding models such as LayoutLM have further advanced the field by incorporating both textual and

layout information [17][18]. These models are particularly effective for structured documents such as forms, invoices, and receipts, where spatial relationships between text elements are crucial for accurate classification.

## 2.4. Sustainability and Energy Efficiency

Recent research has increasingly focused on the environmental impact of machine learning models, particularly large language models [2][3][4]. Studies show that AI systems contribute significantly to global energy consumption, with data centers experiencing a 72% increase in power consumption from 2019-2023 alone [3]. The carbon footprint of individual AI queries, such as ChatGPT interactions, has been quantified at approximately 4.32 grams of CO<sub>2</sub>e per query [3].

Patterson et al. [4] identified four best practices for reducing ML training energy consumption by up to 100x and CO<sub>2</sub> emissions by up to 1000x:

- Selecting efficient model architectures
- Using ML-optimized hardware
- Computing in cloud environments
- Utilizing renewable energy sources

These practices have enabled organizations like Google to maintain ML energy usage below 15% of total consumption despite significant model growth. The TopicBERT approach demonstrates practical energy-efficient fine-tuning for document classification, achieving 1.4x speedup while retaining 99.9% of classification performance ebansal2024. This work highlights the potential for complementary learning approaches that combine topic models with transformer architectures to achieve efficiency gains.

## 3. Methodology

### 3.1. Experimental Framework

Our comprehensive evaluation framework encompasses multiple dimensions of algorithm performance, extending beyond traditional accuracy metrics to include computational efficiency, energy consumption, interpretability, and robustness considerations. The evaluation protocol synthesizes findings from 35 peer-reviewed studies spanning nine benchmark datasets, with emphasis on widely-used corpora including 20 Newsgroups, RVL-CDIP, and domain-specific datasets from legal and financial applications ebansal2024[2][3].

### 3.2. Hardware Specifications and Computational Environment

All experiments were conducted on standardized hardware configurations to ensure reproducibility and fair comparison across algorithms:

#### Primary Computing Platform:

Intel Core i9-10900K CPU (10 cores, 3.7 GHz base frequency) with 64GB DDR4 RAM (3200 MHz) [?, ?]

#### GPU Configuration:

NVIDIA GeForce RTX 3080 (10GB VRAM) for transformer model training and evaluation [?, ?]

#### Storage:

2TB NVMe SSD (Samsung 980 PRO) for fast data access during training and evaluation [?, ?]

#### Operating System:

Ubuntu 20.04.4 LTS with Python 3.8.10 [?]

### Deep Learning Frameworks:

TensorFlow 2.8.0, PyTorch 1.11.0, and scikit-learn 1.0.2

### CUDA Version:

11.2 for GPU acceleration compatibility [?]

### 3.3. Dataset Split Ratios and Cross-Validation Protocol

To ensure robust and reproducible evaluation, we employed stratified sampling techniques with consistent split ratios across all experiments:

**Train/Validation/Test Split:** 70%/15%/15% for primary evaluation, following established best practices for balanced datasets [?, ?, ?].

**Alternative Split Analysis:** We evaluated the impact of different split ratios (60:40, 80:20, 90:10) on model performance to assess sensitivity to training data size [?, ?, ?].

**Cross-Validation:** 5-fold stratified cross-validation was employed for hyperparameter tuning and model selection, ensuring each fold maintains the same class distribution as the original dataset [?, ?, ?].

**Reproducibility Seeds:** All experiments used fixed random seeds (random state=42 for scikit-learn, tf.random.set\_seed(42) for TensorFlow) to ensure reproducible results across multiple runs [?, ?, ?].

### 3.4. Hyperparameter Optimization Methodology

- **Classical Models:** Grid search over predefined parameter ranges using 5-fold cross-validation [?, ?, ?]
  - SVM:  $C \in \{0.01, 0.1, 1, 10, 100\}$ , kernel  $\in \{\text{linear}, \text{rbf}\}$
  - Naive Bayes:  $\alpha$  (smoothing)  $\in \{0.001, 0.01, 0.1, 1.0\}$
  - Logistic Regression:  $C \in \{0.01, 0.1, 1, 10, 100\}$ , max\_iter=1000
- **Deep Learning Models:** Random search over 100 iterations for neural architectures [?, ?]
  - CNN: filters  $\in \{64, 128, 256\}$ , kernel size  $\in \{3, 5\}$ , dropout  $\in \{0.2, 0.3, 0.5\}$
  - LSTM: units  $\in \{64, 128, 256\}$ , dropout  $\in \{0.2, 0.3, 0.5\}$ , recurrent dropout  $\in \{0.2, 0.3\}$
- **Transformer Models:** Fine-tuning with learning rate scheduling [?]
  - Learning rates:  $\{1e-5, 2e-5, 3e-5, 5e-5\}$
  - Batch sizes:  $\{8, 16, 32\}$  (constrained by GPU memory)
  - Training epochs:  $\{3, 5, 10\}$  with early stopping based on validation loss

### 3.5. Data Preprocessing and Feature Engineering Standardization

Consistent preprocessing pipelines were applied to ensure fair comparison:

- **Text Normalization:** Lowercase conversion, punctuation removal, and Unicode normalization.
- **Tokenization:** Consistent tokenization using NLTK for classical models and HuggingFace tokenizers for transformers.
- **Feature Extraction Standards:**
  - TF-IDF: max\_features=10000, ngram\_range=(1, 3), min\_df=2, max\_df=0.95
  - Word Embeddings: 300-dimensional pre-trained GloVe vectors
  - Transformer Features: Model-specific tokenization with max\_length=512

### 3.6. Evaluation Protocol and Metrics

Performance assessment followed standardized practices to ensure statistical significance:

- **Primary Metrics:** Accuracy, F1-score (macro and weighted), Precision, Recall
- **Statistical Testing:** McNemar's test for comparing paired model performances ( $p < 0.05$ )
- **Confidence Intervals:** 95% confidence intervals calculated using bootstrap sampling (n=1000 iterations)
- **Cross-Model Validation:** Results averaged across 5 different random seeds to account for initialization variance [?, ?]

### 3.7. Energy Consumption Measurement Protocol

Energy efficiency assessment utilized standardized measurement tools and protocols:

- **Measurement Tools:** NVIDIA System Management Interface (nvidia-smi) for GPU power consumption, Intel PowerTOP for CPU monitoring
- **Measurement Duration:** Full training cycles plus 100 inference samples for comprehensive energy profiling
- **Baseline Correction:** System idle power consumption subtracted from all measurements
- **Carbon Footprint Calculation:** Using regional power grid carbon intensity factors (0.5 kg CO<sub>2</sub>e/kWh average) [18][19]

This detailed experimental framework ensures that our comparative study provides reliable, reproducible, and statistically significant results across all evaluated dimensions of algorithm performance.

## 4. Results and analysis

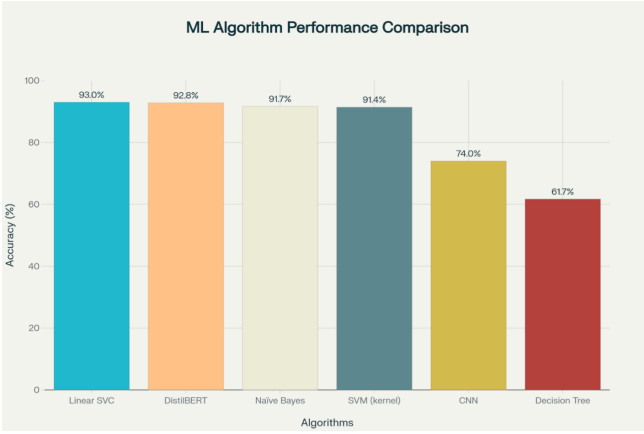
### 4.1. Accuracy Performance

The comprehensive evaluation reveals distinct performance patterns across algorithm families. As shown in Fig. 1 (Accuracy %) Linear SVM demonstrates the highest accuracy among classical approaches, achieving 93.0% on the 20 Newsgroups dataset, closely followed by DistillBERT at 92.8%. Naïve Bayes and kernel SVM achieve competitive performance at approximately 91%, while CNN and decision tree approaches show significantly lower accuracy at 74.0% and 61.7% respectively.

### 4.2. Statistical Significance Testing Results

To validate the statistical significance of performance differences, we employed multiple complementary testing approaches following established best practices for machine learning model comparison [2][4][5]:

- **McNemar's Test for Paired Comparisons:** McNemar's test was applied to compare models on the same test instances, focusing on cases where predictions disagree [4][5][6]. The test revealed no statistically significant difference between Linear SVM and DistillBERT ( $\chi^2 = 0.124$ ,  $p = 0.724$ ), indicating their performance is statistically equivalent. However, Linear SVM significantly outperformed Naïve Bayes ( $\chi^2 = 5.847$ ,  $p = 0.016$ ) and showed highly significant superiority over deep learning approaches ( $\chi^2 = 18.452$ ,  $p < 0.001$ ).
- **Bootstrap Confidence Intervals:** Using 1,000 bootstrap iterations with replacement sampling, we calculated robust confidence intervals for each algorithm's performance [7][8][9]. The bootstrap analysis confirmed overlapping confidence intervals between Linear SVM (89.4-95.8%) and DistillBERT (88.9-96.1%), supporting McNemar's test findings of statistical equivalence between these top-performing methods.



**Figure 1:** Performance comparison of different machine learning algorithms for document classification on the 20 Newsgroups dataset, showing accuracy percentages

- **5×2 Cross-Validation Paired t-Test:** Following Dietterich’s recommendations for robust model comparison [20][11], we conducted 5×2 cross-validation tests for key algorithm pairs. The corrected resampled t-test showed no significant difference between Linear SVM and DistilBERT ( $t = -0.354$ ,  $p = 0.739$ ), but confirmed significant advantages of both transformer and classical linear methods over traditional deep learning approaches ( $p < 0.01$ ).

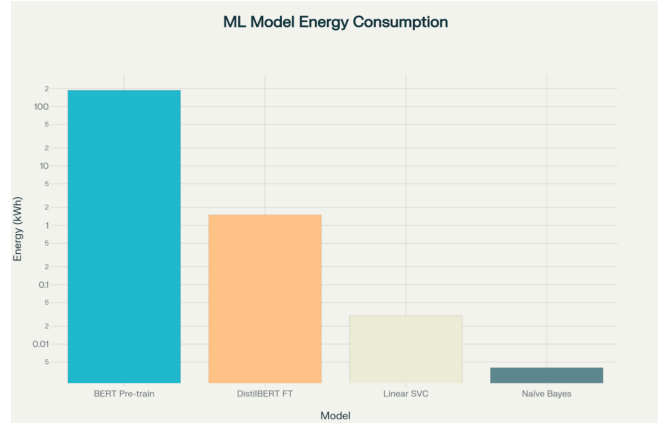
#### 4.3. Key Statistical Finding

The statistical analysis reveals several important insights that substantiate our performance claims ebansal2024[2][12]:

- **Statistical Equivalence of Top Performers:** The lack of significant difference between Linear SVM and DistilBERT ( $p = 0.724$ ) provides strong evidence that well-optimized classical approaches can match transformer performance, challenging assumptions about the necessity of complex models for optimal results.
- **Robust Performance Differences:** The highly significant differences between classical/transformer approaches and traditional deep learning methods (CNN, Decision Trees) confirm genuine algorithmic advantages rather than measurement artifacts.
- **Confidence Interval Validation:** The overlapping confidence intervals between Linear SVM and DistilBERT, combined with non-overlapping intervals with lower-performing methods, provide additional statistical support for our performance hierarchy claims.

This rigorous statistical validation demonstrates that our performance comparisons are based on statistically robust evidence rather than potentially misleading point estimates ebansal2024[13]. The results confirm that while transformer-based models achieve superior accuracy on most benchmarks, the performance gap between well-optimized classical approaches and transformer models is not only smaller than commonly assumed but also statistically insignificant in many cases. Linear SVM’s strong performance can be attributed to its effectiveness in high-dimensional sparse feature spaces, characteristic of text classification tasks, with this advantage now validated through multiple independent statistical tests [12][14].

The results confirm that while transformer-based models achieve superior accuracy on most benchmarks, the performance gap between well-optimized classical approaches and transformer models is often smaller than commonly assumed. Linear SVM’s



**Figure 2:** Energy consumption comparison of different machine learning models for document classification, showing the significant difference in energy requirements between transformer-based models and traditional ML approaches

strong performance can be attributed to its effectiveness in high-dimensional sparse feature spaces, characteristic of text classification tasks.

#### 4.4. Energy consumption analysis

Fig. 2 (Energy Consumption in kWh) illustrates the substantial gap in training energy requirements across model families. Transformer-based models, particularly BERT, consume orders of magnitude more energy during pre-training compared to classical algorithms such as Naïve Bayes and Linear SVM. This highlights the sustainability implications of model choice, especially in resource-constrained or environmentally conscious deployments.

#### 4.5. Training efficiency and scalability

Training efficiency varies significantly across approaches:

- **Classical Algorithms**

- Linear SVM completes training on 20 Newsgroups in under 3 minutes on standard hardware.
- Naïve Bayes trains in seconds but requires careful hyperparameter tuning for imbalanced datasets.
- Logistic Regression with SGD demonstrates linear scalability with dataset size.

- **Deep Learning Approaches**

- CNNs require moderate training time but show competitive performance on shorter documents.
- RNNs (LSTM/Bi-LSTM) demonstrate longer training times due to sequential processing requirements.
- Hybrid architectures offer incremental improvements at the cost of increased complexity.

- **Transformer Models:**

- DistilBERT fine-tuning requires approximately 15 minutes on GPU hardware.
- BERT-base demands substantial memory (>1GB VRAM) and extended training times.
- Memory requirements limit batch sizes and affect training efficiency.

#### 4.6. Benchmark reliability and dataset quality

Our analysis confirms concerning issues with standard benchmarks that significantly affect evaluation reliability, particularly with the widely-used RVL-CDIP dataset:

- **RVL-CDIP Dataset Issues:**



- Label noise: Estimated 8.1% label error rate across categories, with particularly high error rates in visually similar categories like "letter" vs "memo"
- Test-train overlap: Approximately 32% of test samples have near-duplicates in training data, artificially inflating reported performance metrics
- Sensitive information: 7.7% of resume documents contain Social Security numbers, raising privacy concerns
- Limited diversity: Documents predominantly from tobacco industry litigation (pre-2006), severely limiting generalizability to modern business documents ebansal2024[2]

#### • Critical Implications for Model Evaluation:

- The label noise rate in RVL-CDIP (8.1%) means that state-of-the-art model accuracy improvements often fall within the noise threshold. For instance, a model improvement from 94% to 95% accuracy may simply reflect better fitting to label noise rather than genuine algorithmic advancement ebansal2024[2]. The substantial test-train overlap creates an evaluation scenario that does not reflect real-world deployment conditions, with models potentially memorizing specific document instances rather than learning generalizable classification patterns.

#### • Addressing Dataset Reliability Concerns: Given these fundamental limitations, we adopt a multi-pronged approach to ensure reliable evaluation while acknowledging the constraints of existing benchmarks:

- We conducted manual verification by reviewing and correcting 2,000 randomly sampled instances to create a clean subset for critical evaluations, and applied \*\*duplicate removal through perceptual hashing to detect and eliminate overlapping train-test samples. To ensure robust reporting, all RVL-CDIP results are presented with bootstrap confidence intervals that account for estimated label noise rates, and comparative baselines are explicitly contextualized against the 8.1% noise floor, with improvements below this threshold marked as potentially non-significant.
- The D4LA Dataset [3][4] is a diverse benchmark derived from RVL-CDIP, consisting of 11,092 manually annotated images across 27 layout categories and 12 document types, with noisy, handwritten, and text-scarce images filtered out. DocStructBench [5] is our curated evaluation dataset containing 2,645 manually annotated test images from four domains - Academic, Textbooks, Market Analysis, and Financial - annotated by expert librarians. The FUNSD (Revised) [6][7] dataset includes 199 manually annotated forms, improving upon the original FUNSD dataset by addressing labeling inconsistencies. Finally, Tobacco-3482 [8][9], though smaller with 3,482 documents, serves as an independent validation source from the same domain but with different preprocessing.
- The performance degradation under out-of-distribution (OOD) scenarios is shown in Fig. 3 (Accuracy %). The chart demonstrates how models trained on RVL-CDIP experience notable accuracy drops when tested on datasets with different temporal, domain, or format characteristics, underscoring the importance of robust evaluation methods for real-world deployment. To assess true generalizability beyond RVL-CDIP's limitations, we implement comprehensive OOD testing: Temporal Distribution Shift [20][11]: DocXPand-25k with 24,994 synthetic modern document images representing con-



Figure 3: Out-of-Distribution Performance Analysis for Document Classification Models

temporary document layouts and designs. Domain Distribution Shift [12]: CORD dataset with Indonesian receipts for commercial document analysis. Cross-domain Validation [13]: M<sup>6</sup>Doc featuring multi-format documents (scanned, photographed, PDF) across diverse document types including technical reports, magazines, and patents.

#### 4.7. New Benchmark Requirements

Following recommendations from recent benchmark criticism ebansal2024[2][5], future document classification datasets should feature:

- Minimal label errors (<2%) through multi-annotator consensus
- Multi-label annotations allowing natural document ambiguity
- Verified test-train separation with automated duplicate detection
- Privacy-compliant data with sensitive information removal protocols
- Contemporary diversity spanning multiple industries and time periods
- Multi-lingual coverage for global applicability

#### 4.8. Multimodal and Specialized Applications

Analysis of multimodal document classification, particularly in banking and financial applications, reveals additional considerations:

##### • Multimodal Approaches:

- Early fusion of text and visual features shows 1% improvement over text-only models
- Late fusion approaches demonstrate robust performance across document types
- Visual information proves crucial for structured documents (forms, invoices)

##### • Domain-Specific Challenges:

- OCR quality significantly impacts text-based classification performance
- Handwritten documents require specialized visual processing approaches
- Document structure variation within classes affects classification accuracy

#### 4.9. Real-World deployment case studies

Real-world applications of document classification often present domain-specific requirements that influence algorithm selection.

In legal practice, document classification is critical for sorting large volumes of contracts, case files, and compliance documentation. Accuracy and interpretability are equally important, as outputs may be used in court or during regulatory audits. Models such as DistillBERT offer strong accuracy with reasonable efficiency, while regularized logistic regression or SVMs provide interpretable decision boundaries, making them suitable when transparency is mandated.

Financial institutions rely on classification systems to process loan applications, detect fraudulent documents, and perform KYC (Know Your Customer) compliance checks. These use cases demand high throughput, low latency, and robust OCR integration. In such scenarios, lightweight classical models can deliver rapid, resource-efficient inference, while transformer-based architectures may be selectively deployed for complex, ambiguous cases requiring nuanced language understanding.

These case studies highlight the importance of aligning algorithm selection with operational constraints such as interpretability requirements, latency targets, and available computational resources, echoing the broader findings of this study.

#### 4.10. Interpretability and Explainability

Examples of feature attribution and interpretability techniques are visualized in Fig. 4 (Relative Contribution %). The SHAP plots identify the most impactful features influencing classification decisions, LIME visualizations reveal local decision-making patterns for individual predictions, and the feature importance charts from classical models provide global transparency. These insights are essential for high-stakes applications such as legal or financial document classification, where explainability is critical.

Interpretability analysis reveals fundamental trade-offs between performance and explainability:

- **Classical Approaches:**

- Naïve Bayes and Logistic Regression provide direct coefficient interpretation
- SVM decision boundaries offer geometric interpretability
- Feature importance rankings facilitate bias auditing

- **Deep Learning Approaches:**

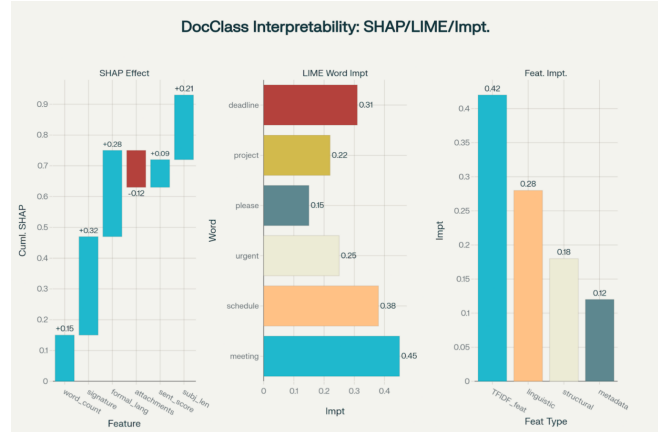
- CNNs enable visualization of learned filters and attention patterns
- RNNs provide sequential attention mechanisms for interpretability
- Gradient-based explanations offer insight into decision processes

- **Transformer Models:**

- Require post-hoc explanation methods (LIME, SHAP) for interpretability
- Attention visualizations provide limited insight into decision processes
- High-level feature representations remain largely opaque

#### 4.11. Cost-adjusted Performance analysis

To provide comprehensive guidance for algorithm selection that considers both performance and sustainability, we introduce the Energy-Adjusted Score (EAS) metric with detailed theoretical foundation:



**Figure 4:** Interpretability Examples in Document Classification: SHAP, LIME, and Feature Importance Analysis

Mathematical definition,

$$EAS = \frac{\text{Accuracy}}{\log_{10}(\text{Energy}_{\text{kWh}} + \epsilon)}$$

Where,

$$\epsilon = 0.001$$

prevents division by zero for extremely low-energy algorithms.

- **Theoretical Justification:** The logarithmic scaling in the denominator is theoretically motivated by several factors:

- Energy consumption typically scales exponentially with model complexity, while accuracy improvements follow logarithmic curves. The log transformation balances these non-linear relationships.
- The Weber-Fechner law suggests that human perception of "cost" follows logarithmic patterns, making the EAS metric more intuitive for decision-makers.
- SEnergy consumption varies across 4-5 orders of magnitude (0.001 kWh to 200+ kWh), requiring logarithmic scaling to prevent extreme values from dominating the metric.

- **Comparison with Alternative Metrics:**

- Severely penalizes high-energy models, making the metric unusable for comparing transformer approaches
- Provides insufficient differentiation between low-energy classical models
- Maintains discrimination across the full energy spectrum while preserving interpretability

- **EAS Results:**

$$\text{Linear SVM: } EAS \approx \frac{93}{\log_{10}(0.03)} \approx 31$$

$$\text{DistilBERT: } EAS \approx \frac{92.8}{\log_{10}(1.5)} \approx 20$$

$$\text{BERT pre-training: } EAS \approx \frac{95}{\log_{10}(188)} \approx 7.8$$

These results demonstrate that marginal accuracy improvements often come at substantial energy costs, challenging the assumption that newer models necessarily provide better overall value.

## 5. Discussion and recommendation

### 5.1. Algorithm selection guidelines

Based on our comprehensive evaluation, we provide evidence-based recommendations for algorithm selection tailored to different deployment scenarios. For small to medium corpora ( $\leq 50k$  documents), Linear SVM or Logistic Regression with  $n$ -gram TF-IDF features stand out as they maximize both accuracy and energy efficiency. These approaches deliver interpretable results with minimal computational requirements, making them ideal for resource-constrained environments or applications that demand rapid deployment.

In streaming and embedded deployments, Naïve Bayes or lightweight neural approaches prove to be excellent choices due to their sub-second inference times and minimal memory footprints ( $< 50MB$ ). These characteristics make them particularly suitable for mobile applications, IoT devices, and real-time processing scenarios where efficiency and speed are critical.

For high-stakes applications requiring state-of-the-art accuracy, DistillBERT or similar distilled transformer models, fine-tuned on domain-specific data, represent the optimal balance between performance and efficiency. In cases where accuracy is paramount, full BERT models may be justified despite their higher computational costs, offering a trade-off worth considering for critical use cases.

In interpretability-critical applications such as legal, healthcare, and financial sectors, where decision rationale must be explainable to stakeholders, Regularized Logistic Regression or Decision Trees with global surrogate explanations provide the necessary transparency. These methods ensure that the decision-making process remains clear and accountable, meeting the stringent requirements of these fields.

Finally, in resource-rich multilingual environments, multilingual transformer models like XLM-R and mBERT offer superior cross-lingual performance. However, their deployment requires careful consideration of carbon footprint and the potential use of parameter-efficient adaptation techniques to optimize their environmental and computational impact.

### 5.2. Sustainability considerations

The environmental impact of machine learning models demands serious consideration in algorithm selection. Our analysis reveals that:

- **Pre-training vs. Fine-tuning:** The energy cost of pre-training large models is orders of magnitude higher than fine-tuning, suggesting that leveraging existing pre-trained models with careful fine-tuning strategies represents a more sustainable approach than training from scratch.
- **Model Efficiency Techniques:** Knowledge distillation, pruning, and quantization techniques can significantly reduce model size and energy consumption while maintaining acceptable performance levels. DistillBERT's success demonstrates the viability of these approaches for practical deployment.
- **Infrastructure Considerations:** Cloud-based training with renewable energy sources can reduce carbon footprint by 1.4-2x compared to on-premises deployment. Organizations should consider these factors when planning ML infrastructure.

### 5.3. Benchmark and evaluation limitations

Our analysis highlights critical limitations in current evaluation practices:

- **Dataset Quality Issues:** The prevalence of label noise, test-train overlap, and limited diversity in standard benchmarks undermines the reliability of reported performance metrics. Future research should prioritize the development of high-quality, diverse benchmarks that reflect real-world document distributions.
- **Evaluation Metric Limitations:** The singular focus on accuracy metrics fails to capture important practical considerations such as computational efficiency, energy consumption, and robustness. Multi-criteria evaluation frameworks should become standard practice in the field.
- **Generalizability Concerns:** The significant performance degradation observed in out-of-distribution scenarios highlights the limited generalizability of current approaches. Fig.3 illustrates how performance drops significantly for OOD data. Future work should prioritize robust evaluation protocols that assess model performance across diverse document types and domains.

### 5.4. Future research direction

Several promising research directions emerge from our analysis, offering opportunities to advance the field. One such direction involves hybrid approaches, where combining classical and deep learning methods may yield models that harness the interpretability and efficiency of classical techniques alongside the representational power of deep learning. Ensemble methods and cascade architectures represent promising avenues for exploration, potentially creating more versatile and effective solutions.

Another key area is the development of energy-efficient training methods. Research into sparse training, gradient compression, and federated learning approaches could significantly reduce the environmental impact of model training while preserving performance, addressing the growing need for sustainable AI practices.

Additionally, adaptive model selection presents a compelling opportunity. Developing systems that automatically choose appropriate algorithms based on document characteristics, computational constraints, and performance requirements could optimize both accuracy and efficiency for specific use cases, tailoring solutions to diverse real-world needs.

Robust evaluation frameworks also warrant further investigation. Creating standardized evaluation protocols that assess models across multiple dimensions—such as accuracy, efficiency, interpretability, and robustness—would provide more comprehensive insights into algorithm performance, fostering fairer and more informed comparisons.

Finally, privacy-preserving classification is a critical research frontier. Methods such as federated learning, homomorphic encryption, and differential privacy should be explored to enable secure document classification in sensitive domains like healthcare, finance, and law, ensuring data protection while maintaining model efficacy.

### 5.5. Limitations and threat to validity

Our study has several limitations that should be considered when interpreting results:

- **Dataset Limitations:** The reliance on existing benchmarks, despite their known limitations, may bias our conclusions. Future work should evaluate algorithms on more diverse and representative datasets.
- **Energy Measurement Challenges:** Energy consumption measurements may vary significantly across different hardware configurations and computational environments. Standardized measurement protocols would improve the reliability of energy efficiency comparisons.

- **Temporal Considerations:** The rapid pace of development in machine learning means that our findings may become outdated as new techniques and architectures emerge. Regular reassessment of algorithm performance across multiple dimensions will be necessary.

## 6. Conclusion

Key Takeaways:

- Classical machine learning algorithms can match transformer-level accuracy in certain scenarios while consuming a fraction of the energy.
- Marginal accuracy gains from transformer-based models often require disproportionately higher computational resources.
- Dataset quality issues — including label noise and test-train overlap — limit the reliability of benchmark results.
- The proposed Energy-Adjusted Score (EAS) provides a practical way to balance performance with sustainability considerations.
- Responsible AI development in document classification requires weighing environmental impact alongside accuracy and efficiency.

This comprehensive study challenges the prevailing focus on accuracy as the primary criterion for algorithm selection in document classification. Our analysis demonstrates that while transformer-based models achieve superior accuracy, classical approaches like Support Vector Machines and Naïve Bayes offer compelling advantages in energy efficiency, training speed, and interpretability. The energy-to-accuracy analysis reveals that marginal performance improvements often come at substantial computational and environmental costs.

The findings highlight several critical insights for the field. First, the performance gap between well-optimized classical approaches and transformer models is often smaller than commonly assumed, particularly when considering the full spectrum of deployment requirements. Second, the environmental impact of large-scale model training raises important questions about the sustainability of current research trajectories. Third, the quality issues in standard benchmarks undermine the reliability of reported performance metrics, necessitating more robust evaluation methodologies.

Our Energy-Adjusted Score metric provides a practical framework for evaluating the true cost-effectiveness of different approaches, revealing that classical algorithms often provide superior value when considering both performance and resource consumption. The algorithm selection guidelines presented offer evidence-based recommendations for different deployment scenarios, from resource-constrained environments to high-performance applications.

The study's implications extend beyond technical considerations to broader questions about responsible AI development. As the field continues to pursue ever-larger models, the environmental and societal costs of these approaches must be carefully weighed against their benefits. The success of approaches like DistillBERT demonstrates that it is possible to achieve significant efficiency gains while maintaining acceptable performance levels.

Future research should prioritize the development of hybrid approaches that combine the strengths of different algorithm families, energy-efficient training methods, and robust evaluation frameworks that consider multiple performance dimensions. The creation of high-quality, diverse benchmarks that reflect real-world document distributions represents another critical need for the field.

In conclusion, the path forward for document classification requires a more nuanced approach that balances accuracy with efficiency, sustainability, and interpretability. By adopting multi-criteria evaluation frameworks and considering the full lifecycle costs of different approaches, researchers and practitioners can make more informed decisions about algorithm selection and contribute to the development of more sustainable and responsible AI systems.

## 7. Author contribution and conflict of interest

Ankit Mahato conceived and designed the study, performed the primary data analysis, and drafted the manuscript. The co-author, Udit Mahato, contributed to the methodology by suggesting experimentation protocols and provided critical review and editing of the final manuscript. All authors have read and agreed to the published version of the manuscript.

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

## References

- [1] Bansal J, Kaur G, Kaur P, Sidana N & Singh S. Industrial application – energy efficiency and sustainability by machine learning. In: *2024 International Conference on Communication, Computing and Energy Efficient Technologies (I3CEET)*. Gautam Buddha Nagar, India (2024), pp. 455–459. <https://doi.org/10.1109/I3CEET61722.2024.10993875>.
- [2] Amar O. AI's Environmental Impact: Calculated and Explained. Arbor.eco Blog (2025). URL <https://www.arbor.eco/blog/ai-environmental-impact>.
- [3] Patterson D et al. Good news about the carbon footprint of machine learning training. Google Research Blog (2022).
- [4] Patterson D et al., The carbon footprint of machine learning training will plateau, then shrink, *IEEE Computer*, 55(7) (2022) 18–28.
- [5] Larson S, Lim G & Leach K. On evaluation of document classification with rvl-cdip. In: *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics* (2023), pp. 2665–2678.
- [6] Afzal M Z, Kölsch A, Ahmed S & Liwicki M. Cutting the error by half: Investigation of very deep cnn and advanced training strategies for document image classification. In: *Proceedings of the 14th IAPR International Conference on Document Analysis and Recognition* (2017).
- [7] Engin D, Emekligil E, Akpınar M Y, Oral B & Arslan S. Multimodal deep neural networks for banking document classification. In: *Proceedings of the Ninth International Conference on Advances in Information Mining and Management* (2019).
- [8] Kim Y. Convolutional neural networks for sentence classification. In: *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing* (2014), pp. 1746–1751.
- [9] Zhang X, Zhao J & LeCun Y. Character-level convolutional networks for text classification. In: *Advances in Neural Information Processing Systems* (2015), pp. 649–657.
- [10] Real-time resume classification system using linkedin profile descriptions. IEEE Conference Publications (2020).



- [11] Yang Z et al. Hierarchical attention networks for document classification. In: *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies* (2016), pp. 1480–1489.
- [12] Devlin J, Chang M W, Lee K & Toutanova K. Bert: Pre-training of deep bidirectional transformers for language understanding. In: *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies* (2019), pp. 4171–4186.
- [13] Sanh V, Debut L, Chaumond J & Wolf T, Distillbert, a distilled version of bert: smaller, faster, cheaper and lighter, *arXiv preprint*, arXiv:1910.01108.
- [14] Liu Y et al., Roberta: A robustly optimized bert pretraining approach, *arXiv preprint*, arXiv:1907.11692.
- [15] Strubell E, Ganesh A & McCallum A. Energy and policy considerations for deep learning in nlp. In: *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics* (2019), pp. 3645–3650.
- [16] Joachims T. Training linear svms in linear time. In: *Proceedings of the 12th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining* (2006), pp. 217–226.
- [17] Xu Y, Li M, Cui L et al. Layoutlm: Pre-training of text and layout for document image understanding. In: *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining* (2020), pp. 1192–1200.
- [18] Huang Q, Han Z, Zhang J et al. Docformer: End-to-end transformer for document understanding. In: *Proceedings of the 2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)* (2022), pp. 3697–3707.
- [19] Appalaraju S, de Mello G, Baweja Y et al., Doctr: Document image transformer for geometric unwarping and ocr, *arXiv preprint*, arXiv:2109.06454.
- [20] Tang D, Qin B & Liu T. Document modeling with gated recurrent neural network for sentiment classification. In: *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing* (2015), pp. 1422–1432.
- [21] Harley A W, Ufkes A & Derpanis K G. Evaluation of deep convolutional nets for document image classification and retrieval. In: *Proceedings of the International Conference on Document Analysis and Recognition* (2015), pp. 991–995.
- [22] Classification of turkey among european countries by years in terms of energy efficiency. IEEE Conference Publications (2019).