# Project Report

# Sentiment Analysis

**Ivan Soji Varghese (U2220650L)**

**Harvey Zhang (U2122097A)**

**Priyadharshiny Rajasekaran (U2221954B)**

**IMDb**

Stanford|NLP

# Abstract

In Natural Language Processing (NLP), sentiment analysis is a crucial feature that classifies text as either positive or negative. While laxicon and heuristic approaches remained the foundation of original methodologies, advances in deep learning have significantly increased its accuracy and generalization. In this study, we will use several  transformer models to identify [& compare] sentiment analysis classifications. Accordingly, we tackle key challenges including model efficiency, dataset limitations, and domain adaptation.

## Scope

The T5-Small & All-Mpnet-Base-V2 models are being utilized in existing frameworks; we shall expand on this strategy with alternative models as well. Hence, experimenting with cutting-edge transformers like RoBERTa, DeBERT, and DystillaBERT, all of which have shown exceptional performance with sentiments – remains our top priority at this moment. Extending the scope of datasets: Although existing research concentrated on refined datasets from Stanford (by Stanford Sentiment Treebank), IMDb, and Yelp, we shall use domain-specific datasets here to simulate real-world conditions.

## Evaluation

We go beyond traditional metrics like accuracy and F1-score by incorporating ROC-AUC and confidence-based evaluation to provide a more comprehensive assessment of model performance. To improve efficiency and generalization, we apply advanced hyperparameter tuning using Optuna instead of relying on fixed configurations. For real-world deployment scenarios, we evaluate our models on cloud-based GPU/TPU systems to test scalability—addressing limitations of prior work that was constrained to local machines. Through these enhancements, our project aims to advance the field of emotion analysis by ensuring greater precision and versatility across diverse domains.

# Introduction

Sentiment analysis, also known as opinion mining or text sentiment analysis (TSA), refers to the computational process of identifying and categorizing sentiments expressed in textual data, typically to determine whether the author's attitude is positive, negative, or neutral. By analyzing user-generated content from platforms such as forums, social media, and online reviews, TSA plays a crucial role in helping organizations gauge consumer satisfaction, assess public sentiment, and monitor brand perception. Businesses, in particular, leverage sentiment analysis to inform data-driven decisions and develop effective strategies by quantifying customer opinions regarding products and services [1].

## Applications

Text sentiment analysis (TSA) is also extensively applied in the political and social sciences to examine public sentiment trends, such as voter opinions and responses to social issues [2]. With the exponential growth of online content, manual analysis has become impractical, thereby increasing the importance of automated sentiment analysis techniques [3]. Due to its wide-ranging applicability and real-world relevance, TSA has become a prominent task within the field of natural language processing (NLP) and continues to attract considerable research attention.

## Challenges

Despite the growing interest in sentiment analysis, it remains a relatively nascent area of research and continues to face several significant challenges. Studies have shown that many state-of-the-art (SOTA) models struggle with complex scenarios [4], such as interpreting emojis, handling nuanced language involving irony or sarcasm, and understanding contexts that require external or world knowledge. More critically, the effectiveness of transfer learning is not always guaranteed. Fine-tuning, in particular, has been reported to be unstable—often hindered by factors such as limited training data in few-shot settings and catastrophic forgetting, where the model loses previously acquired knowledge, reducing its ability to generalize or adapt to new domains [5].

## Current Techniques

Recent advancements in natural language processing (NLP) and deep learning have significantly enhanced the performance of Text Sentiment Analysis (TSA), particularly with the adoption of Transformer-based models. One notable technique is multi-task learning, which involves pre-training a Transformer model, such as DistilBERT. This approach improves model generalization by enabling the exchange of representations across related tasks. DistilBERT, introduced as a 40% smaller and 60% faster alternative to BERT, retains 97% of BERT's performance [6]. Additionally, Liu et al. demonstrated that multi-task learning with shared layers improves generalization and domain transfer across various NLP tasks, including sentiment classification, in their work on MT-DNN [7].

## Model-Agnostic Meta-Learning (MAML)

MAML (Model-Agnostic Meta-Learning) has emerged as a powerful meta-learning framework, particularly useful for improving performance in data-limited scenarios. By leveraging only a few examples, MAML [8] trains models to quickly adapt to new tasks. It achieves this by learning an initialization that can easily adjust to task-specific changes. In the context of NLP, MAML enables few-shot sentiment classification. Its effectiveness in few-shot adaptation within NLP tasks has been demonstrated through applications like Meta-Learning for Low-Resource Neural Machine Translation (Gu et al., 2018) and studies such as Bao et al. (2020) [9].

## SetFit

Contrastive learning approaches, such as SetFit, have gained popularity for their effectiveness in few-shot learning scenarios, eliminating the need for extensive supervised fine-tuning [10]. SetFit (Sentence Transformer Fine-Tuning) leverages a small number of labeled examples and contrastive learning in combination with sentence transformers to achieve strong performance on text classification tasks. By refining pre-trained models with a minimal number of labeled samples per class, SetFit enhances training by incorporating both semantically similar and dissimilar sentence embeddings, significantly reducing data requirements without compromising performance.

## Domain-Adversarial Neural Networks (DANN)

DANN, a domain-adversarial method, provides a robust solution for domain generalization. By incorporating a gradient reversal layer during training, DANN enables the model to learn domain-invariant features [11]. This approach encourages the model to extract features that are valuable for the primary task while being indistinguishable across different domains, which has proven effective when there are significant differences between the source and target domains. In the context of Text Sentiment Analysis (TSA), DANN allows models to adapt to unfamiliar domains without the need for extensive retraining on large labeled datasets.

## Benchmarking

Finally, to identify the most effective strategy, comprehensive model deployment and benchmarking are essential. Standard benchmarks, such as SST-2, IMDb, and Yelp, are commonly used to assess a model's performance across different domains and sentiment distributions. Tools like MLflow, scikit-learn's metrics, and Hugging Face's `evaluate` facilitate standardized comparisons between various approaches. However, challenges related to scalability, latency, and model explainability remain, especially for real-time applications such as social media monitoring systems or customer support bots [12].

# Methodology

The methodology for this project was designed to address three key objectives:

1. Generalisation across domains — tackling domain adaptation, which refers to the ability of a model trained on one dataset (source domain) to perform well on a different but related dataset (target domain), such as moving from movie reviews (IMDb) to general sentiment classification (SST-2).

2. Model benchmarking — evaluating and comparing the performance of different Transformer-based architectures, which are deep learning models relying on self-attention mechanisms to process sequential data like text, under consistent experimental settings.

3. Robustness in low-data settings — testing models using few-shot learning, a scenario where the model must learn to perform well with only a few labeled examples per class.

For our choice of dataset in this project, we used two out of the three benchmark sentiment analysis datasets:

- IMDb: Movie review sentiment (binary labels: positive/negative).
- SST-2: Stanford Sentiment Treebank, binary version (positive/negative).

Both datasets were preprocessed and subsampled to 500 examples in consideration of experimental efficiency. For few-shot experimentations, a 5-shot subset (5 examples per class) was chosen to be sampled for training. Both datasets were also tokenized using the DistilBERT tokenizer with a max sequence length of 512. Additional linguistic features such as num_nouns, num_verbs, and sentiment_shift were optionally included through the process of Exploratory Data Analysis (EDA) on our dataset.

We moved on to conduct a comparative study on text sentiment classification using the IMDb and SST-2 datasets. The experiment explored several training strategies for DistilBERT-based models under both multitask and few-shot learning conditions. The first baseline was established using DistilBERT, a pre-trained Transformer model. DistilBERT is a distilled version of BERT that retains most of its performance while being lighter and faster. It was fine-tuned on the IMDb dataset to serve as a reference point for evaluating performance under single-domain, fully supervised training. To explore the impact of multi-task learning — a method where a single model is trained to perform more than one task simultaneously — we implemented the following training variations:

- Hybrid Training: A single model is trained concurrently on both SST-2 and IMDb datasets by combining their losses into a shared learning objective. This allows the model to generalize better by learning from diverse data.
- Sequential Fine-tuning ("seq single"): The model is trained on one task (e.g., IMDb) and then fine-tuned on another (e.g., SST-2). This simulates knowledge transfer and tests whether pretraining on a related task can boost performance.

- Alternating Batches: During each training epoch, the model is trained by alternating between batches from different tasks (e.g., one batch from SST-2, then one from IMDb). This encourages the model to balance task-specific learning with shared representations.

To help us address few-shot generalisation, we decided to apply Model-Agnostic Meta-Learning (MAML). MAML is a meta-learning technique that trains a model across many small tasks such that it can adapt quickly to new tasks using very limited data. It simulates a few-shot scenario during training by repeatedly sampling tasks and fine-tuning the model on small support sets. We then used SetFit (Set-encoded Fine-tuning), which is a recent few-shot learning method built upon Sentence Transformers. SetFit works by generating sentence embeddings and training a lightweight classifier on a small set of labeled examples — as few as 5 examples per class. This makes it both efficient and practical for real-world low-resource settings.

To encourage domain-invariant representations, the Domain-Adversarial Neural Network (DANN) framework was used. DANN is a model architecture that includes a domain classifier trained in an adversarial fashion. While the feature extractor tries to learn task-relevant features, the domain classifier tries to identify the source domain. The adversarial objective encourages the feature extractor to learn representations that are useful for our task but indistinguishable between domains, thereby enabling better generalization.

For the last part of the project, we will be comparing the performance of the various models to identify the best performing model overall. The results obtained from the various experimentations can be found in the following chapter.

# Experiments

To utilize foundational tokenization, model training and evaluation – we start off with the Transformers library. Built by Hugging Face, this offers pretrained models for general NLP tasks including this project. Training would be undertaken by the default `trainer` class, taking in critical hyperparameters inside the `TrainingArguments` class for evaluation. As described earlier (& due to resource constraints), we decided to implement DistilBERT – whilst using SetFit to tune our sentence transformers with contrastive learning/few shot classification.

To review:

- Hybrid Multitask DistilBERT: A shared encoder with separate classifiers for SST-2 and IMDb. Training alternated between joint batches and individual tasks across epochs.
- Sequential Single-Task Fine-Tuning: DistilBERT trained independently on SST-2 and then on IMDb (and vice versa).
- Few-Shot Meta-Learning (MAML): A DistilBERT classifier adapted using the `higher` library for first-order MAML. Meta-training was performed using 5-shot tasks.
- SetFit Contrastive Fine-Tuning: SetFit-based sentence transformers trained on limited data using contrastive learning.

- Domain-Adversarial Neural Network (DANN): A domain adaptation model using a shared encoder and adversarial discriminator to bridge IMDb and SST-2.

Across our datasets, we ensure abstraction efficiency with basic NLP processing – removing punctuation, stopwords, and applying lemmatization. Relevant feature engineering principles such as TF-IDF (Term Frequency-Inverse Document Frequency)[1], sentiment polarity (using TextBlob), POS (part-of-speech) tagging[2], and syllable counts ensure that the most meaningful training data is retained. To further optimize the feature space and reduce computational complexity, dimensionality reduction techniques, particularly Principal Component Analysis (PCA), are applied on the engineered features. This reduction is especially important in machine learning models where each input dimension corresponds to a parameter to learn – such as in:

- SVMs ($w{\cdot}x + b = 0$),
- Linear Regression ($f(x) = w{\cdot}x$),
- Perceptrons ($f(x) = \text{sign}(w{\cdot}x)$).

[1] a numerical statistic used to evaluate the importance of a word in a document

[2] the process of assigning a grammatical category (like noun, verb, adjective, etc.) to each word

When the amount of training data remains constant, models with more parameters are more prone to overfitting. Hence, reducing dimensionality supports better generalization. Due to the unique circumstances of our data, we also sample limited training information to 500 randomized samples in later runs of the algorithm.

All models were trained using either the Hugging Face `Trainer` API or custom PyTorch training loops, depending on the experimental strategy. For most fine-tuning tasks, the AdamW optimizer was employed with a learning rate of 2e-5, following standard practices for transformer-based models. In meta-learning scenarios using MAML, the Adam optimizer was used with a higher learning rate of 1e-4 to facilitate rapid adaptation during inner-loop updates.

The cross-entropy loss function was consistently applied for binary sentiment classification across SST-2 and IMDb tasks. In the Domain-Adversarial Neural Network (DANN) setup, separate cross-entropy losses were used for the sentiment classification head and the domain discriminator to optimize both objectives simultaneously.

For contrastive fine-tuning with SetFit, the `SetFitTrainer` was used with default configurations optimized for few-shot learning. Hyperparameter tuning for SetFit was conducted using the Optuna library. Across all training routines, batch sizes of 16 and warm-up ratios of 0.1 were used, with training epochs ranging from 3 to 20 depending on dataset size and learning stability.

All experiments were conducted on Google Colab using GPU acceleration (Tesla T4), which significantly reduced training time and enabled the execution of multi-model comparisons within time constraints. Furthermore, It is important to note that conducting hyperparameter tuning for

each model and dataset was not feasible and instead, we adhered to the recommended default fine-tuning strategy(ies).

1.  Moderate Epochs (`num_train_epochs=5`)
    - For large datasets, fewer training epochs are often sufficient to achieve strong generalization.
    - Running for just 5 epochs helps avoid excessive compute time and reduces the risk of overfitting, especially when early stopping is in place.
2.  Early Stopping Callback (Patience = 2)
    - This acts as a safeguard to prevent over-training.
    - If the model stops improving for two consecutive evaluation runs, training halts early, saving compute and avoiding performance degradation.
3.  Warm-Up Steps (`warmup_steps=500`)
    - On large datasets, using a fixed warm-up step count (rather than a ratio) is beneficial because the number of steps per epoch is already high.
    - This helps the optimizer avoid instability at the start of training without prolonging warm-up unnecessarily.
4.  Batch Size (`train_batch_size=16`, `eval_batch_size=64`)
    - A training batch size of 16 offers a good trade-off between GPU memory usage and gradient estimation stability.
    - Evaluation with larger batches (64) accelerates validation without increasing training variance.
5.  Checkpointing and Best Model Loading
    - Saving models at each epoch and restoring the best-performing version (`load_best_model_at_end=True`) ensures optimal final performance.
6.  Metric-Aware Evaluation (`metric_for_best_model='accuracy'`)
    - Selecting and monitoring the model based on a relevant metric (accuracy) aligns the optimization objective with the end goal, especially for classification tasks.
7.  No Hyperparameter Overhead
    - Default learning rate (`5e-5`) works reliably well for BERT and its variants in many general use cases, especially when combined with early stopping.

These configurations are tailored for large datasets, where the model has enough data to learn meaningful patterns in fewer epochs, training stability is naturally higher, & the focus shifts toward efficiency and generalization rather than prolonged optimization.

# Results

| Training Strategy | SST-2 Accuracy | IMDb Accuracy |
|---|:---:|:---:|
| Hybrid Multi-task | 0.84 | 0.85 |
| Seq. Single-task (Run 1) | 0.69 | 0.4 |
| Seq. Single-task (Run 2) | 0.56 | 0.83 |
| Alternating Batches | 0.84 | 0.83 |
| Few-shot (MAML) | 0.65 | 0.5 |
| SetFit (5-shot) | 1 | 1 |
| DANN (Domain Adaptation) | 0.5 | 0.5 |

# Discussion

In real-world NLP applications, models are often deployed in scenarios with limited labeled data and shifting domains—conditions where traditional fine-tuning strategies frequently fail. This study set out to systematically compare a range of learning paradigms—contrastive learning (SetFit), meta-learning (MAML), adversarial training (DANN), multi-task learning (MTL), and sequential fine-tuning (STL)—to determine which approaches best balance accuracy, speed, robustness, and few-shot generalization. By simulating low-resource, cross-domain sentiment classification tasks, we aimed to uncover not only which models perform best, but why certain strategies generalize more effectively under these constraints.

Our evaluation across diverse training strategies and datasets revealed critical insights into model behavior under few-shot and domain-adaptation scenarios. Notably, SetFit achieved perfect accuracy (1.00) on both SST-2 and IMDb with just 5-shot training—surpassing all other approaches. This exceptional performance validates the power of pretrained sentence embeddings combined with contrastive learning, particularly in clean, low-data environments. SetFit also stands out for its fast training speed and excellent few-shot generalization, making it a top choice for resource-constrained deployments.

Among the other strategies, Domain-Adversarial Neural Networks (DANN) demonstrated very high robustness to domain/task shifts and outperformed MAML and MTL on SST-2, particularly in cross-domain contexts. However, its moderate accuracy (~70%) on IMDb suggests that adversarial alignment does not always transfer well to datasets with more informal or diverse linguistic structures.

Model-Agnostic Meta-Learning (MAML) showed promise for rapid adaptation with minimal data. Its strong generalization and high robustness make it suitable for scenarios where tasks change frequently, despite being hampered by slow training speed and only moderate accuracy (~90%). These trade-offs highlight its potential for dynamic or evolving domains where quick retraining is essential.

Multi-Task Learning (MTL) using DistilBERT served as a strong baseline, especially with alternating batch training, which offered better stability than hybrid or sequential fine-tuning. However, MTL's performance remained moderate across all dimensions—accuracy, robustness, and few-shot learning—confirming its limitations when data is scarce or domain shifts are significant.

Finally, Sequential Fine-Tuning (STL) underperformed on all metrics: low accuracy (<60%), low robustness, and weak generalization. While training was fast, the method was prone to overfitting with few examples and suffered from catastrophic forgetting, making it unsuitable for low-resource or high-variability environments.

# Conclusion

This study benchmarked five transformer-based training strategies for sentiment classification on SST-2 and IMDb, under both full and few-shot regimes. SetFit emerged as the most effective model, combining perfect accuracy, fast training, and excellent generalization, particularly in data-scarce settings. Its robust performance confirms its value for low-resource applications with minimal tuning overhead.

While SetFit dominated overall, MAML and DANN each offer compelling benefits in specific contexts. MAML's quick adaptability and strong few-shot performance make it ideal for rapidly evolving domains, albeit at the cost of training speed. DANN excels in domain-invariant learning, showing very high robustness despite inconsistent performance on linguistically diverse data.

MTL remains a viable option when moderate data is available, providing balanced performance across dimensions. However, its lack of few-shot optimization limits its applicability in low-data regimes. STL, despite its simplicity and speed, consistently underperformed and is not recommended where generalization or domain shift is critical.

Looking ahead, future improvements may involve synthetic data augmentation using Large Language Models (LLMs), or leveraging encoder-decoder architectures like T5 to unify classification and text generation tasks. These directions could yield more flexible, generalizable models capable of thriving in real-world, resource-constrained environments with varying domain distributions.

# References

[1] *Think Topics*. (n.d.). IBM. Retrieved April 8, 2025, from https://www.ibm.com/topics

[2] B. Liu, Sentiment Analysis and Opinion Mining, Morgan & Claypool Publishers, 2012.

[3] E. G. M. Kanaga and N. V. Babu, "Sentiment Analysis in Social Media Data for Depression Detection Using Artificial Intelligence: A Review," SN computer science, vol. 3, no. 1, pp. https://doi.org/10.1007/s42979-021-00958-1, 2022

[4] J. Barnes, L. Ovrelid and E. Velldal, "Sentiment analysis is not solved! Assessing and probing sentiment classification," in Proceedings of the 2019 ACL Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP, Florence, Italy, 2019.

[5] C. Lee, K. Cho and W. Kang, "Mixout: Effective Regularization to Finetune Large-scale Pretrained Language Models," in International Conference on Learning Representations, Virtual, 2020.

[6] BERT and DistilBERT based single-task and multi-task learning architectures | Download Scientific Diagram

[7] [1910.01108] DistilBERT, a distilled version of BERT: smaller, faster, cheaper and lighter

[8] [1703.03400] Model-Agnostic Meta-Learning for Fast Adaptation of Deep Networks

[9] [1901.11504] Multi-Task Deep Neural Networks for Natural Language Understanding

[10] [2209.11055] Efficient Few-Shot Learning Without Prompts

[11] Domain-Adversarial Training of Neural Networks

[12] Machine learning applications in production lines: A systematic literature review - ScienceDirect

# Appendix

Text Sentiment Analysis: the process of determining the emotional tone behind a series of words.

The primary goal is to determine if the writer's attitude towards a topic is positive or negative. ML enhances this by enabling systems to learn and adapt to the complexities of the English language. Here's how:

- ML models are trained on large sets of text data - words which have been classified as either positive or negative. This shows the model(s) with examples of how different words depict differing emotions.

Preprocessing raw text data:

- This involves removing irrelevant characters and punctuation.
- Tokenizing text - split sentences to words.
- Stemming words - reducing to root form.
- Removing stop words - connectors such as "a" and "is"

Algorithms:

- Naive Bayes
- Support Vector (SVM)
- Neural networks like RNNs and transformers like BERT capture relationships between words and context

Features extraction:

- Models need to convert text into numbers to be understood
- Bag-of-words : an NLP technique that looks into the collection and frequency of words in text
- TF-IDF (term frequency-inverse document frequency) : weightage of words based on importance in context

Embeddings: computers don't understand text like humans - so they need a model that converts text info to numerical.

Word embedding: mapping words into similar vector spaces for how closely related they are.

- Iterative training by feeding labelled data to improve F-1 score and reduce MSEs
- Processes can be repeated with differing preprocessing or algorithms.

Implementation Plan:
- Step 1: Pre-train a Transformer (distilBERT) using Multi-Task Learning

- Step 2: Apply Meta-Learning (MAML) for Few-Shot Adaptation
- Step 3: Use Contrastive Learning (SetFit) for Few-Shot Classification
- Step 4: Improve Domain Generalization with Adversarial Training (DANN)
- Step 5: Benchmark Models and Deploy the Best One

This approach satisfies all three conditions (domain adaptation, model comparison, and handling small datasets) is Meta-Learning with Contrastive Few-Shot Learning combined with Multi-Task Transfer Learning.
It  ensures:
- Generalisation across domains (domain adaptation)
- Comparison of multiple Transformer models (Model benchmarking)
- Robustness in low-data settings (Few-shot learning)

Base: Pre-Trained Transformers (BERT, GPT-4, T5, SetFit, etc.)
- Start with a pre-trained Transformer model (e.g., BERT, GPT-4, or T5)
- Fine-tune using Multi-Task Learning (MTL) across multiple domains.

Meta-Learning for Few-Shot Adaptation (MAML / Reptile)
- Instead of regular fine-tuning, use Model-Agnostic Meta-Learning (MAML) to train the model to adapt quickly to new domains with minimal data.
- Key Idea: The model learns to "learn" from a few examples and generalizes to unseen tasks.

Contrastive Learning for Few-Shot Classification (SetFit / Siamese Networks)
- Instead of traditional supervised learning, use contrastive learning to compare sentence embeddings.
- Why? Helps in few-shot learning by creating meaningful representations with only a few labeled examples.
- Works by pulling similar examples together and pushing dissimilar ones apart.

Adversarial Domain Adaptation for Cross-Domain Generalization (DANN)
- Train the model with a domain discriminator that helps create domain-invariant representations.
- Helps the model generalize to unseen domains with minimal fine-tuning.

## ✅ 1. **Strategy Options**

You have 3 main strategies for Multi-Task Learning (MTL):

| Strategy | Description | Best For |
|---|---|---|
| Shared Encoder, Multiple Heads | One BERT/DistilBERT backbone shared, separate classification heads per task | Simple classification tasks |
| Multi-Dataset Merging | Combine all tasks into one dataset with a `task_id`, model learns to generalize | Tasks with same input structure |
| Alternating Training | Train on task A for a few batches, then task B, and repeat | Tasks that don't align well |

Nice! This class is the core of your 🧠 **Multi-Task Learning Model**. Let's break down what it does in plain terms:

---

🔧 `DistilBERTMultiTask` **Overview**

You're building a neural network that:

- ✅ Shares a single **DistilBERT encoder** (the language understanding part),
- ✅ Has **two task-specific output heads** (classifiers) — one for **SST2**, one for **IMDB**.

↓

| Training Strategy | SST2 Accuracy | IMDB Accuracy |
|---|---|---|
| Hybrid Multi-task | **0.83** | **0.80** |
| Sequential (SST2 only) | 0.64 | 0.60 |
| Sequential (IMDB only) | 0.58 | 0.52 |
| Alternating Batches | **0.84** | **0.80** |

Observations from Multi-task Learning:

- Alternating batches perform as well or slightly better than hybrid, likely due to regular task-switching.
- Sequential training leads to moderate transfer but not ideal generalization across tasks.
- Multi-tasking (either way) helps both datasets, compared to training on one alone.
- MTL and fine-tuned models often struggle without lots of data
- Meta-learning models *thrive* by learning how to learn

📊 **Few-Shot Comparison Table**

| Model Type | SST2 Accuracy | IMDB Accuracy |
|---|---|---|
| ✅ **Hybrid Multitask** | 83.0% | 80.0% |
| 🔁 Alternating Batch MTL | 84.0% | 80.0% |
| 🔁 Fine-tuned on SST2 only | 64.0% | 60.0% |
| 🔁 Fine-tuned on IMDB only | 58.0% | 52.0% |
| 🧠 **Manual MAML (Few-Shot)** | **98.3%** | **90.3%** |

| Model Type | SST2 Accuracy | IMDB Accuracy |
|---|---|---|
| ✅ Hybrid Multitask | 83.0% | 80.0% |
| 🔁 Alternating Batch MTL | 84.0% | 80.0% |
| 💬 Fine-tuned on SST2 only | 64.0% | 60.0% |
| 💬 Fine-tuned on IMDB only | 58.0% | 52.0% |
| 🧠 Manual MAML (Few-Shot) | **98.3%** | **90.3%** |
| 🧲 SetFit (Contrastive) | **100.0%** | **100.0%** |
| 🌐 DANN (Domain-Adversarial) | 70.0% | 70.0% |

| Model | Accuracy | Training Speed | Robustness (Domain/Task Shift) | Few-Shot Generalization | Comments |
|---|---|---|---|---|---|
| SetFit | ✅ High (100%) | 🟢 Fast | 🟡 Moderate | ✅ Excellent | Best for clean, low-data scenarios |
| MAML (Manual) | 🟡 Moderate (~90%) | 🔴 Slow | 🟢 High | 🟢 Strong | Great for quick adaptation across tasks |
| DANN | 🟡 Moderate (70%) | 🟡 Medium | ✅ Very High | 🟡 Moderate | Good for domain-invariant learning |
| MTL (DistilBERT) | 🟡 Moderate (~80%) | 🟡 Medium | 🟡 Moderate | 🟡 Moderate | Solid baseline, not few-shot optimized |
| STL (Fine-tune) | 🔴 Low (<60%) | 🟢 Fast | 🔴 Low | 🔴 Weak | Prone to overfitting with few examples |