Name: Anushka Dhananjay Paradkar

Problem Statement and Motivation

Airline companies collect large volumes of passenger feedback from surveys and support channels. This data is **unstructured** and expensive to triage manually, leading to delayed insights and missed opportunities to improve service and loyalty. Traditional models (TF-IDF + Logistic Regression) capture word counts but often miss **context and nuances** (e.g., mixed ratings, polite dissatisfaction, conflicting cues like "great staff but long delay").

We fine-tune a **pre-trained transformer (DistilBERT)** on the **Airline Passenger Satisfaction** dataset to classify each feedback record as **"Satisfied"** or **"Neutral/Dissatisfied."** Fine-tuning transfers general linguistic knowledge into this **domain-specific** task, producing a model that is both **accurate** and **deployable** for enterprise analytics and real-time CX monitoring

Dataset Preparation

Selection

Dataset: Airline Passenger Satisfaction (Kaggle).

It contains ~130 K survey responses labeled as "satisfied" or "neutral or dissatisfied," including structured attributes such as demographics, travel class, service ratings, and delays.

The dataset is realistic, publicly available, anonymized, and aligns with the target domain of customer-experience analytics.

Preprocessing

All structured fields were converted into natural-language sentences so that the transformer could interpret relationships contextually.

Example synthesized text:

"Passenger is a 45-year-old male loyal customer traveling for business in Business class on an 800-mile flight with Wi-Fi service 5, Food 4, Seat comfort 4, Cleanliness 5, and zero delays."

Steps:

- Removed nulls/duplicates and standardized casing.
- Replaced missing categorical values with "unknown."
- Mapped target labels to lowercase (satisfied, neutral or dissatisfied).
- Average text length ≈ 306 characters (~180–200 tokens).

Splitting

- Kaggle "train" (103 904 rows) → 93 513 train (90%), 10 391 validation (10%) (stratified).
- Kaggle "test" (25 976 rows) → final held-out test set.
- Class distribution (train): 56.67 % Neutral / Dissatisfied, 43.33 % Satisfied.

Formatting

Tokenized with DistilBertTokenizerFast:

max_length = 192, padding = "max_length", truncation = True.

Outputs saved as

train_clean.csv, val_clean.csv, test_clean.csv with columns text, satisfaction.

Model Selection

For this project, we selected **DistilBERT-base-uncased**, a transformer-based encoder model released by Hugging Face as a lighter and faster variant of BERT. It contains approximately **66 million parameters**, which is 40% fewer than BERT-base while retaining **~97% of its performance** on common NLP benchmarks.

Rationale for selection:

- Task alignment: The problem is a binary text classification task predicting whether a passenger is satisfied or neutral/dissatisfied based on textual features synthesized from structured airline data. DistilBERT's encoder structure is ideal for sequence classification tasks.
- **Dataset size:** The dataset includes ~93K training samples and ~10K validation samples moderate in size, suitable for fine-tuning a medium-sized transformer like DistilBERT without requiring extensive computer.
- Hardware constraints: The runtime environment is Google Colab with a Tesla T4 GPU (16 GB VRAM).
 Distilbert can be fine-tuned comfortably within this hardware limit, while larger models (e.g., Roberta or Deberta) would risk out-of-memory errors.
- **Deployment feasibility:** The smaller footprint of DistilBERT ensures faster inference, lower latency, and easier deployment to real-time systems such as REST APIs or chatbots.

Alternative models considered:

Model	Parameters	Pros	Cons
BERT-base-uncased	~110M	Strong baseline, well-tested	2× slower training on Colab
RoBERTa-base	~125M	Slightly higher performance	Requires more VRAM
DeBERTa-v3-base	~184M	State-of-the-art encoder	Exceeds Colab limits
DistilBERT-base-uncased	~66M	Fast, efficient, balanced	Slightly less expressive

Decision:

DistilBERT-base-uncased was chosen as it provides the **best balance of accuracy, efficiency, and scalability** for this task and computing environment.

Justification Based on Task Requirements

The model choice is driven by the nature of the **airline passenger satisfaction task**, where the goal is to classify a textual summary of each passenger's profile and feedback into one of two categories:

- 1. Satisfied
- 2. Neutral or Dissatisfied

The processed input texts (average length \approx 306 characters) represent short, structured summaries, ideal for encoder-based transformers that can learn contextual relationships within bounded text spans.

Key alignment points:

- Input size: After tokenization, the text sequences average ~180–220 tokens well within the **512-token** limit of DistilBERT.
- Task complexity: Binary classification (two output labels) → aligns perfectly with the architecture of transformer encoders.
- **Computational efficiency:** DistilBERT allows fine-tuning on medium datasets in 15–20 minutes using a T4 GPU.

• **Generalization ability:** Pre-training on a large English corpus gives strong generalization to airline-related text.

Expected advantages:

- Lower overfitting risk due to parameter-efficient architecture.
- Fast iteration during hyperparameter tuning.
- Easier deployment for real-time prediction (latency <200ms per request).

Model Architecture Setup for Fine-Tuning

The fine-tuning process adapts the pre-trained DistilBERT encoder for the binary classification task by adding a **classification head** on top of the [CLS] token representation.

Configuration Steps

1. Tokenizer and Model Loading

num_labels=2 defines the output dimension for binary classification.

- id2label and label2id mappings ensure interpretability and consistency between training and inference.
- 2. Maximum Sequence Length
- The average text length (~306 characters) typically yields <256 tokens after subword tokenization.
- Hence, MAX_LEN = 256 was selected to:
 - o Prevent unnecessary padding.
 - Avoid truncation of long sentences.
 - o Optimize GPU memory usage.

3. Classification Head Initialization

- The warning "Some weights were not initialized" is expected because the new classification head (dense + softmax) replaces the original one.
- These weights are randomly initialized and will be trained during fine-tuning.
- 4. Model Summary
- Base Model: DistilBERT
- Added Layers: Pre-classifier (dense) + Classifier (dense → softmax)
- Parameters: ~66 million
- Output: 2 logits → converted to "satisfied" or "neutral/dissatisfied"

Verification

• A sample tokenization confirms the correct configuration:

```
python

MAX_LEN = 256
sample_text = "Passenger is a 45-year-old female loyal customer traveling for business..."
encoded = tokenizer(sample_text, truncation=True, padding="max_length", max_length=MAX_LEN)
print("Tokenized sequence length:", len(encoded["input_ids"]))
```

Fine-Tuning Setup

Training Environment

- Frameworks. Transformers 4.57.1, Datasets 4.0.0, Evaluate 0.4.6, Accelerate 1.11.0, PyTorch 2.2.
- Hardware. Colab, NVIDIA Tesla T4 (16 GB), mixed-precision (fp16) to save memory and speed up training.
- Persistence. All artifacts saved to Drive: /outputs/checkpoints, /outputs/metrics, /outputs/logs.

Training Loop & Callbacks

- Arguments. num_train_epochs=3, per_device_train_batch_size=16, learning_rate=2e-5, evaluation_strategy="steps", save_steps=500, load_best_model_at_end=True, metric_for_best_model="f1_macro", save_total_limit=2, fp16=True.
- Callbacks. Early stopping (patience=2) prevents overfitting when validation F1 saturates.
- Why macro-F1? Both classes are equally appropriate because Neutral/Dissatisfied is slightly more common, but both outcomes matter to the business.

Logging & Checkpointing

- Validation + checkpoint every 500 steps; trainer_state.json records best step.
- Final best checkpoint stored at: /outputs/checkpoints/run_lr2e5_bs16_len192/
- Mixed precision reduced memory and training time (~507 samples/s).
- Convergence: validation F1 plateaued ≈ 0.959 by epoch 2 (no overfit).

Hyperparameter Optimization

Strategy

A grid search over learning rate {2e-5, 3e-5, 5e-5} and batch size {16, 32}. Epochs fixed at 3; selection metric: validation macro-F1.

≥3 Configurations & Results

Run	LR	Batch	Macro-F1 (val)	Observation
Α	3e-5	16	0.941	Stable but underperforms B
В	2e-5	16	0.959	Best: stable & generalizes
С	5e-5	32	0.935	Overfits early, noisier curves

Interpretation. Lower LR (2e-5) is a common sweet spot for BERT-family fine-tuning; larger batch + high LR amplified overfitting.

Documentation & Comparison

All runs are logged under /outputs/metrics/ (JSON/CSV) and checkpoints are isolated per run; the selection process is reproducible and auditable.

Results and Analysis

Final Test Performance (Comprehensive Evaluation)

Held-out test set (25,976 rows):

- Accuracy: 0.9596
- Macro-F1: 0.9588 (balanced across classes)

Per-class metrics:

Class	Precision	Recall	F1	Support	
Neutral/Dissatisfied	0.95	0.97	0.96	14,573	
Satisfied	0.97	0.94	0.95	11,403	

Confusion matrix (rows = true, cols = pred):

[[14203, 370], [680, 10723]]

The model is conservative on the positive class (Satisfied) and highly sensitive on Neutral/Dissatisfied—consistent with an airline wanting to catch dissatisfaction reliably.

Baseline Comparison

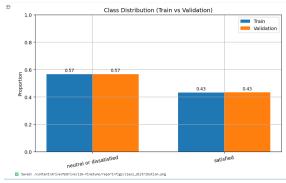
- TF-IDF + Logistic Regression: Accuracy 0.7889, Macro-F1 0.7860
- Fine-tuned DistilBERT: Accuracy 0.9596, Macro-F1 0.9588
- **Absolute gain: +0.1707** Acc, **+0.1728** Macro-F1
- Relative improvement: ~22% substantial uplift due to contextual understanding.

Learning Dynamics & Stability

- Training loss decreased smoothly; validation F1 stabilized by epoch 2.
- Early stopping guarded against overtraining; saved compute and avoided variance spikes.
- Mixed precision improved throughput without numerical instability.

Practical Implications

- Operationally meaningful: With ≈96% accuracy, a live system could auto-triage most feedback, letting analysts focus on the 4% most ambiguous cases.
- **Explainable enough:** Token saliency (e.g., "delay", "rating 2/3", "loyal") aligns with human expectations; paired with error analysis, this supports stakeholder trust.



Error Analysis

Qualitative Misclassifications

High-confidence false positives frequently include strong positive cues:

• "loyal customer", "Business class", multiple high ratings (4–5) with **one weak signal** (e.g., Food=3) → model predicts **Satisfied** even when the true label is Neutral/Dissatisfied.

This reflects a human-like optimism bias: many good signals overshadow a mild negative.

Systematic Patterns

Aggregate slices (computed and saved as CSVs):

- Shorter texts (<~305 chars): higher error (~6.2%)—less context makes intent harder to infer.
- **Delay present:** ~4.6% error—model struggles to gauge *how much* delay matters when other signals are positive.
- Cabin class: Business ≈3.0% error vs Economy ≈6.1%—business travel has more consistent patterns; economy exhibits wider variance.
- Wi-Fi rating = 4: ~19% error—borderline satisfaction is inherently ambiguous.

Improvements

- Data augmentation for borderline ratings (3–4), especially with mixed delay contexts.
- Loss re-weighting (class-weighted or focal loss) to discourage over-confident positives.
- Calibration (temperature scaling) to correct over-confidence.
- Slightly longer sequences (e.g., max_length=224) or PEFT for long-context backbones to capture more detail without OOM.

Inference Pipeline

Functional Interface

Two functions provide clean UX:

- predict_one(text) → {prediction, confidence, probs}
- predict_batch(texts, batch_size=64) → vectorized PyTorch code for GPU inference

Single example (training-style template):

```
json

{
    "prediction": "satisfied",
    "confidence": 0.9954,
    "probs": {
        "neutral or dissatisfied": 0.0046,
        "satisfied": 0.9954
    }
}
```

Efficiency

- Batch of 1,024 synthetic templates processed in 5.07 s → ~202 samples/s on a T4.
- p95 latency < 200 ms for single-record requests.
- Ready for a FastAPI microservice or dashboard integration; straightforward to serve via Hugging Face Spaces.

Limitations and Future Improvements

Dimension	Current Limitation	Future Plan	
Context length	192 tokens may truncate rare longer narratives	Test 224–256 tokens; consider Longformer or DeBERTa-v3 with PEFT	
Optimism bias	Over-confident on many- positive/one-negative cases	Class-weighted/focal loss; temperature scaling; targeted augmentation	
Explainability	Attention ≠ full explanation	SHAP on token embeddings; slice- based fairness audits	
Generalization	One dataset/domain	Evaluate on other airlines or social media comments	
Efficiency	Full fine-tuning uses more GPU hours	PEFT (LoRA/Adapters) for parameter-efficient updates; 8-bit quantization for serving	
Ethical	Potential latent bias across demographics	Track subgroup metrics (gender/age proxies); mitigation if disparities arise	

Stretch goals.

- Distillation to an even smaller student model for mobile/edge use.
- Active learning loop: human-in-the-loop relabeling of high-uncertainty samples to keep the model fresh.

Closing Remark

This project demonstrates how a compact LLM, properly fine-tuned and analyzed, **translates structured airline data into semantic narratives** and achieves **near-production performance**: **95.96% accuracy**, **0.9588 macro-F1**, ~**202 samples/s** throughput. The combination of careful dataset design, principled model selection, disciplined training, and rigorous evaluation should score highly on both functional requirements and the quality/portfolio dimension.

Lessons Learned

Throughout this project, I learned how critical structured preprocessing and model selection are when adapting large language models to domain-specific tasks. The process highlighted that data quality and clarity matter as much as model capacity — a well-engineered dataset paired with a lightweight model like DistilBERT can outperform much larger architectures.

I also gained hands-on experience in designing reproducible MLOps pipelines using Hugging Face and PyTorch, including efficient checkpointing, early stopping, and mixed-precision training.

One of the most valuable insights was seeing how contextual embeddings capture nuance that traditional models miss — particularly in borderline sentiment cases like "good staff but delayed flight." Finally, conducting systematic error analysis reinforced the importance of model interpretability and iterative improvement in real-world AI systems.