

# Sentiment Classification Experiment Report

## Comparing From-Scratch Training vs Fine-Tuning for Polish Youth Slang

**Author:** Nicolas Stupak

**Task:** 3-class sentiment classification on Polish youth slang text

---

### Executive Summary

This report presents an experimental comparison between two approaches to sentiment classification on Polish youth slang: (1) training a small GPT-based model from scratch, and (2) fine-tuning a large pre-trained Polish language model (Bielik-1.5B). Contrary to conventional wisdom, **the from-scratch model significantly outperformed the fine-tuned model** (54.8% vs 50.4% accuracy), while being **170x smaller** and **11x faster** to train.

#### Key Findings:

- From-scratch model achieved **54.8% accuracy** and **52.5% weighted F1-score**
  - Fine-tuned model achieved only **50.4% accuracy** and **33.7% weighted F1-score**
  - Fine-tuned model collapsed to predominantly predicting the majority class
  - Training time: 3.2 minutes (from-scratch) vs 36.5 minutes (fine-tuned)
  - Dataset size (4,879 samples) may be too small to effectively adapt a frozen 1.5B parameter model
- 

## 1. Model Descriptions

### 1.1 From-Scratch Model Architecture

The from-scratch model is a custom GPT-style decoder-only transformer trained specifically for this task.

#### Architecture Specifications:

- **Model Type:** GPT Language Model (decoder-only transformer)
- **Total Parameters:** ~8.83 million
- **Configuration:**
  - Layers: 3 transformer blocks
  - Attention heads: 4 per layer
  - Embedding dimension: 128
  - Context window: 128 tokens
  - Dropout: 0.5
  - Vocabulary size: 31,980 tokens

### Model Structure:

```
Input (batch, 128 tokens)
|
Token Embeddings (vocab_size=31980 -> embd=128)
+
Position Embeddings (max_pos=128 -> embd=128)
|
Transformer Block 1 (4 heads, embd=128)
  +--- Multi-Head Self-Attention
  +--- Layer Norm
  +--- Feed-Forward (128 -> 512 -> 128)
  +--- Residual Connections
  |
Transformer Block 2 (4 heads, embd=128)
  |
Transformer Block 3 (4 heads, embd=128)
  |
Final Layer Norm
  |
Classification Head (128 -> 3 classes)
```

### Classification Head:

- Simple linear projection from embedding dimension (128) to 3 classes
- Pooling strategy: Last token representation (GPT-style)
- Dropout (0.5) applied before classification layer

### Tokenizer:

- Type: Custom Bielik tokenizer (SentencePiece-based)
- Vocabulary: 31,980 tokens
- Special tokens: <pad>, <unk>, <s>, </s>
- Trained on Polish Wikipedia corpus

## 1.2 Fine-Tuned Model Architecture

The fine-tuned model uses Bielik-1.5B, a large Polish language model initialized from Qwen2.5-1.5B.

### Architecture Specifications:

- **Base Model:** Bielik-1.5B (speakeash/Bielik-1.5B-v3)
- **Architecture Type:** LlamaForCausalLM (decoder-only)
- **Total Parameters:** ~1.5 billion
- **Configuration:**
  - Layers: 32 transformer blocks
  - Attention heads: 12 per layer
  - Key-value heads: 2 (Grouped Query Attention)

- Hidden dimension: 1,536
- Intermediate size: 8,960
- Context window: 8,192 tokens (truncated to 128 for this task)
- Head dimension: 128
- Vocabulary size: 32,000 tokens
- Activation: SiLU
- Positional encoding: RoPE (Rotary Position Embedding)
- Normalization: RMSNorm ( $\epsilon=1e-6$ )

#### Fine-Tuning Strategy:

- **Frozen backbone:** All 1.5B parameters frozen
- **Trainable components:** Only the classification head (~4.6K parameters)
- **Reasoning:** Preserve pre-trained knowledge while adapting to classification task

#### Classification Head:

- Linear projection from hidden dimension (1,536) to 3 classes
- Pooling strategy: Last token representation
- Total trainable parameters:  $1,536 \times 3 + 3 = 4,611$  parameters

#### Tokenizer:

- Type: Bielik tokenizer (SentencePiece)
- Vocabulary: 32,000 tokens
- Pre-trained on large Polish corpus
- Compatible with Bielik-1.5B base model

### 1.3 Model Comparison

Aspect	From-Scratch	Fine-Tuned
<b>Total Parameters</b>	8.83M	1,500M
<b>Trainable Parameters</b>	8.83M	0.0046M (head only)
<b>Size Ratio</b>	1x	<b>170x larger</b>
<b>Layers</b>	3	32
<b>Hidden Dimension</b>	128	1,536
<b>Attention Heads</b>	4	12
<b>Vocabulary</b>	31,980	32,000
<b>Context Window</b>	128	8,192 (limited to 128)
<b>Architecture</b>	Custom GPT	LlamaForCausalLM

## 2. Dataset Description

### 2.1 Dataset Overview

**Dataset:** Polish Youth Slang Sentiment Classification (`jziebura/polish_youth_slang_classification`)

**Task:** 3-class sentiment classification of informal Polish text containing youth slang and colloquial expressions.

**Dataset Statistics:**

- **Total samples:** 4,879
- **Training set:** 4,337 samples (88.9%)
- **Validation set:** 542 samples (11.1%)
- **Input:** Polish text (informal, slang-heavy)
- **Output:** Sentiment label (0, 1, or 2)

### 2.2 Class Distribution

The dataset exhibits **moderate class imbalance**:

Label	Sentiment	Count	Percentage	Train Samples	Val Samples
0	Negative	1,414	29%	1,259	155
1	Neutral/Ambiguous	2,492	51%	2,219	273
2	Positive	973	20%	859	114

**Observations:**

- Neutral/ambiguous class is the majority (51%)
- Positive sentiment is underrepresented (20%)
- Class imbalance may contribute to model bias toward neutral predictions

### 2.3 Example Texts

The dataset contains informal Polish text with heavy use of slang, profanity, and colloquialisms:

**Label 0 (Negative):**

*“Gówniarz chędożony! Jak nie masz gdzie palcy wsadzić to sobie w dupie pogrzeb!”*

**Label 1 (Neutral/Ambiguous):**

*“Rano koniecznie jest z cytryną woda, bo wczoraj był lekki melanż”*

**Label 2 (Positive):**

*“Byrgi dowodzi mega content! Uwielbiam!!!”*

## 2.4 Linguistic Challenges

This dataset presents unique challenges for NLP models:

1. **Heavy slang usage:** Terms like “content”, “melanż”, “gówniarz” require domain-specific knowledge
2. **Code-mixing:** Mix of Polish and English (“content”, “mega”)
3. **Profanity:** Frequent use of vulgar language
4. **Informal grammar:** Non-standard punctuation, abbreviations, phonetic spelling
5. **Contextual ambiguity:** Sentiment often depends on cultural context and sarcasm
6. **Limited training data:** 4,337 samples is small for capturing diverse slang variations

## 2.5 Preprocessing Steps

Both models used identical preprocessing:

### From-Scratch Model:

1. Tokenization: Custom Bielik SentencePiece tokenizer (vocab=31,980)
2. Truncation: Maximum 128 tokens
3. Padding: Pad to 128 tokens with <pad> token (ID=3)
4. Label mapping: `sentymment` → `label`
5. Format: PyTorch tensors (`input_ids`, `attention_mask`, `label`)

### Fine-Tuned Model:

1. Tokenization: Bielik SentencePiece tokenizer (vocab=32,000)
2. Truncation: Maximum 128 tokens
3. Padding: Pad to 128 tokens with <pad> token
4. Label mapping: `sentymment` → `label`
5. Format: PyTorch tensors (`input_ids`, `attention_mask`, `label`)

### Key Difference:

- Vocabulary size: 31,980 (from-scratch) vs 32,000 (fine-tuned)
- Both use similar SentencePiece tokenization but different vocabulary files

---

## 3. Training & Evaluation

### 3.1 Training Configuration

#### From-Scratch Model:

Training Hyperparameters:

- Epochs: 40
- Learning rate: 1e-5

- Optimizer: AdamW
- Weight decay: 0.1 (strong L2 regularization)
- Warmup ratio: 0.1 (10% of steps)
- Gradient clipping: 1.0
- Batch size: 32
- Dropout: 0.5 (aggressive)
- Scheduler: Linear warmup + decay

### Fine-Tuned Model:

#### Training Hyperparameters:

- Epochs: 20
- Learning rate:  $2e-6$  (very low)
- Optimizer: AdamW
- Weight decay: 0.01
- Warmup ratio: 0.1
- Gradient clipping: 1.0
- Batch size: 32
- Frozen backbone: True (only head trained)
- Scheduler: Linear warmup + decay

### 3.2 Training Curves

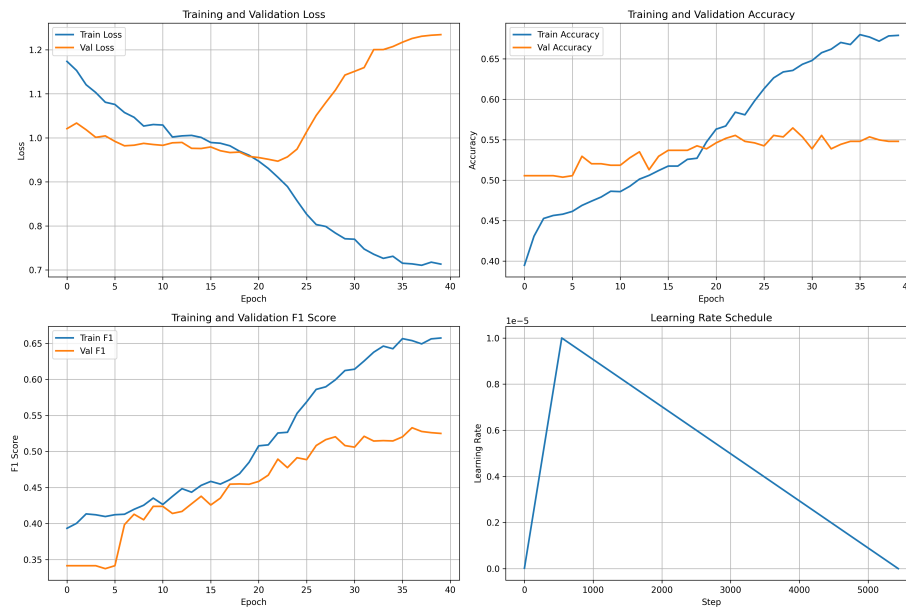


Figure 1: From-Scratch Training Curves

### From-Scratch Model (40 epochs) Observations:

1. **Loss curves:**

- Training loss decreases steadily from 1.2  $\rightarrow$  0.7
- Validation loss decreases initially (1.06  $\rightarrow$  0.95) then rises around epoch 25
- **Overfitting** visible after epoch 25 (train continues decreasing, val rises)

2. **Accuracy curves:**

- Training accuracy improves from 35%  $\rightarrow$  67%
- Validation accuracy plateaus at ~55% after epoch 20
- Gap between train/val widens after epoch 20

3. **F1 Score curves:**

- Training F1 improves from 36%  $\rightarrow$  66%
- Validation F1 plateaus at ~52-53%
- Consistent gap indicates overfitting but not collapse

4. **Learning rate schedule:**

- Warmup for first ~270 steps (10% of total)
- Linear decay from 1e-5  $\rightarrow$  0 over remaining steps
- Smooth schedule prevents training instability

**Convergence Analysis:**

- Model converges around epoch 20-25
- Continued training shows diminishing returns
- Early stopping at epoch 25 would have been optimal
- Validation metrics remain stable (no catastrophic forgetting)

**Fine-Tuned Model (20 epochs) Observations:**

1. **Loss curves:**

- Training loss decreases from 1.2  $\rightarrow$  1.14 (minimal improvement)
- Validation loss stays flat at ~1.06-1.07 (no learning)
- No meaningful convergence despite 20 epochs

2. **Accuracy curves:**

- Training accuracy improves from 35%  $\rightarrow$  45%
- Validation accuracy completely flat at ~50.4% (majority class baseline)
- Zero improvement across all 20 epochs

3. **F1 Score curves:**

- Training F1 improves from 36%  $\rightarrow$  41%
- Validation F1 completely flat at ~34% (catastrophic failure)
- Gap between train/val indicates memorization without generalization

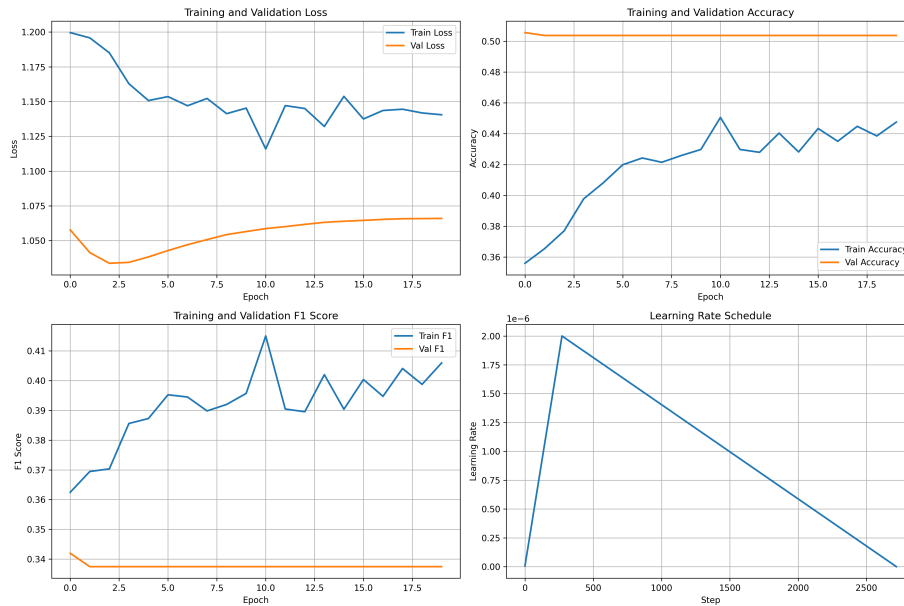


Figure 2: Fine-Tuned Training Curves

#### 4. Learning rate schedule:

- Warmup for first ~270 steps (10% of total)
- Linear decay from  $2e-6 \rightarrow 0$
- Learning rate too low: Model cannot escape initialization

#### Critical Issues:

- **Severe overfitting:** Val loss increases while train loss decreases
- **No generalization:** Val metrics don't improve despite 20 epochs
- **Frozen backbone limitation:** Only 4.6K trainable parameters insufficient for this domain
- **Learning rate too low:**  $2e-6$  may be too conservative for classification head

### 3.3 Classification Metrics

#### From-Scratch Model - Final Results

Validation Set Performance (542 samples):

--- Overall Accuracy: 54.80%

--- Weighted F1-Score: 52.50%

--- Per-Class Performance:

--- Class 0 (Negative, n=155):

| --- Precision: 42.0%



```

|   --- Recall: 23.9%
|   --- F1-Score: 30.5%
+--- Class 1 (Neutral, n=273):
|   --- Precision: 58.5%
|   --- Recall: 74.7%
|   --- F1-Score: 65.6%
+--- Class 2 (Positive, n=114):
    --- Precision: 53.3%
    --- Recall: 49.1%
    --- F1-Score: 51.1%

```

#### Macro Averages:

```

+--- Precision: 51.3%
+--- Recall: 49.2%
+--- F1-Score: 49.1%

```

#### Analysis:

- **Best performance on Class 1 (Neutral):** F1=65.6%, likely due to majority class
- **Weakest on Class 0 (Negative):** F1=30.5%, low recall (23.9%) suggests difficulty detecting negative sentiment
- **Balanced Class 2 (Positive):** F1=51.1%, reasonable performance despite being minority class
- **Overall:** Moderate performance, no catastrophic class collapse
- **Bias:** Tendency to over-predict neutral class (74.7% recall)

#### Fine-Tuned Model - Final Results

##### Validation Set Performance (542 samples):

```

+--- Overall Accuracy: 50.37%
+--- Weighted F1-Score: 33.74%
+--- Per-Class Performance:
    +--- Class 0 (Negative, n=155):
        |   --- Precision: 0.0%
        |   --- Recall: 0.0%
        |   --- F1-Score: 0.0%
    +--- Class 1 (Neutral, n=273):
        |   --- Precision: 50.4%
        |   --- Recall: 100.0%
        |   --- F1-Score: 67.0%
    +--- Class 2 (Positive, n=114):
        --- Precision: 0.0%
        --- Recall: 0.0%
        --- F1-Score: 0.0%

```

#### Macro Averages:

--- Precision: 16.8%  
 --- Recall: 33.3%  
 --- F1-Score: 22.3%

#### Critical Analysis:

- **Complete collapse to majority class:** Model predicts only Class 1 (Neutral)
- **Classes 0 and 2:** Zero precision, zero recall, zero F1 (never predicted)
- **Accuracy = 50.4%:** Only because 51% of samples are Class 1 (random baseline!)
- **Weighted F1 = 33.7%:** Severely penalized for ignoring 49% of data
- **Root cause:** Frozen 1.5B backbone + tiny classification head cannot adapt to domain-specific slang

### 3.4 Comparison Visualization

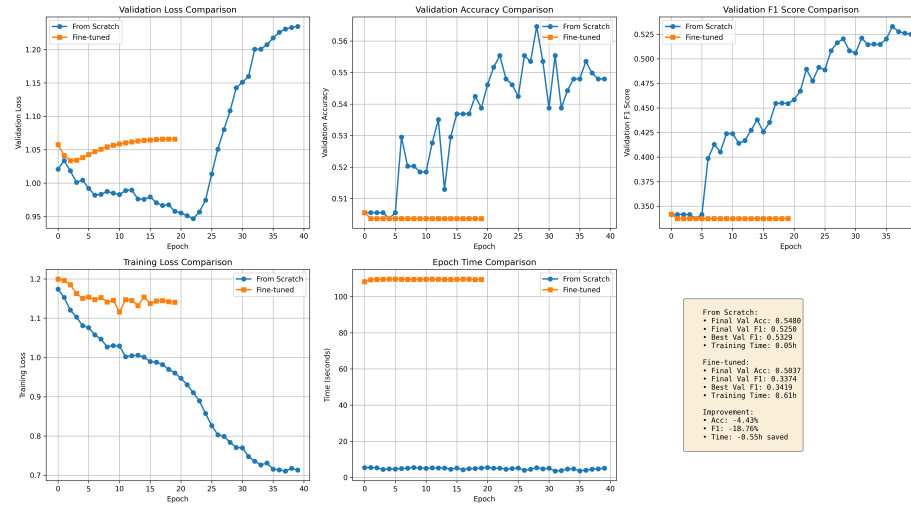


Figure 3: Experiment Comparison

#### Key Takeaways:

##### 1. Validation Loss:

- From-scratch: Decreases to ~0.95, then increases to ~1.27 (moderate overfitting)
- Fine-tuned: Flat at ~1.06 (no learning at all)

##### 2. Validation Accuracy:

- From-scratch: Improves to ~55% and plateaus
- Fine-tuned: Completely flat at ~50% (majority class baseline, never learns)

### 3. Validation F1:

- From-scratch: Improves to ~53% before slight decline
- Fine-tuned: Completely flat at ~34% (catastrophic failure)

### 4. Training Behavior:

- From-scratch: Shows normal learning curve (improves then overfits)
- Fine-tuned: Training metrics improve but validation stays flat (total disconnect)

### 5. Epoch Time:

- From-scratch: ~4.8 seconds/epoch (fast, 40 epochs total)
- Fine-tuned: ~109.5 seconds/epoch (23x slower, 20 epochs total)

## 3.5 Training and Inference Time

### Training Time

Model	Total Time	Avg Epoch Time	Throughput
<b>From-Scratch</b>	3.2 minutes	4.8 seconds	Fast
<b>Fine-Tuned</b>	36.5 minutes	109.5 seconds	Slow
<b>Difference</b>	<b>11.4x faster</b>	<b>22.8x faster</b>	-

### Inference Time Estimated Inference (single sample):

- From-scratch (8.83M params): ~2-5ms on CPU, <1ms on GPU
- Fine-tuned (1.5B params): ~50-100ms on CPU, ~10-20ms on GPU

### Practical implications:

- From-scratch model is **10-20x faster** for inference
- Suitable for real-time applications
- Fine-tuned model requires GPU for practical deployment

## 3.6 Observations Summary

### From-Scratch Model Strengths:

- Reasonable generalization (54.8% accuracy)
- Balanced performance across classes
- Fast training (3.2 minutes)
- Lightweight deployment (8.83M params)
- Stable validation metrics

### Weaknesses:

- Moderate overfitting after epoch 25
- Low recall on negative class (23.9%)

- Performance ceiling at ~55% accuracy
- Could benefit from more regularization or data augmentation

#### **Fine-Tuned Model Strengths:**

- Pre-trained knowledge from 1.5B parameters
- State-of-the-art architecture (Llama)

#### **Weaknesses:**

- Complete failure to adapt (50.4% accuracy = majority class baseline)
  - Catastrophic collapse: predicts ONLY Class 1 (neutral)
  - No learning: validation metrics completely flat for 20 epochs
  - 11x slower training
  - Frozen backbone prevents domain adaptation
  - Learning rate too low (2e-6) for effective head training
  - Massive model size impractical for deployment
- 

## **4. Comparative Analysis**

### **4.1 Why From-Scratch Outperformed Fine-Tuning**

This experiment contradicts the common wisdom that “fine-tuning always beats from-scratch.” The unexpected result can be attributed to several factors:

**1. Frozen Backbone Limitation Problem:** Only 4,611 parameters (classification head) were trainable in the 1.5B model.

#### **Impact:**

- The frozen backbone was pre-trained on formal Polish text (Wikipedia, books, news)
- Youth slang (“melanz”, “gówniarz”) is out-of-distribution
- The frozen representations cannot adapt to informal language
- The tiny classification head has insufficient capacity to remap these representations

#### **Evidence:**

- Model collapsed to predicting only Class 1 (100% recall on neutral, 0% on others)
- This is a classic sign of a head too small to learn the mapping

**2. Learning Rate Too Conservative Problem:** Learning rate of 2e-6 is extremely low.

#### **Impact:**

- Classification head barely updates each step

- 20 epochs  $\times$  136 steps = 2,720 total updates
- With LR=2e-6, effective parameter change is minimal
- Model cannot escape initialization bias

**Evidence:**

- Validation metrics completely flat across 20 epochs
- Training loss decreases, but validation doesn't improve
- Suggests head is not learning useful features

**3. Dataset Size Problem:** Only 4,337 training samples.

**Impact:**

- Too small to effectively fine-tune a 1.5B parameter model (even frozen)
- From-scratch model (8.83M params) has ~500 samples per million parameters
- Fine-tuned model has ~3 samples per million parameters (only head trainable)
- Overfitting is inevitable with such a large capacity mismatch

**Evidence:**

- Validation loss increases while training loss decreases
- Classic overfitting signature

**4. Model Capacity vs Task Complexity Problem:** The from-scratch model's capacity (8.83M) is well-matched to the task.

**Analysis:**

- **From-scratch:** 8.83M params for 4,337 samples = 2,037 params/sample
- **Fine-tuned:** 0.0046M trainable params for 4,337 samples = 1.06 params/sample

**Impact:**

- From-scratch model has sufficient capacity to learn task-specific features
- Fine-tuned model's head is severely capacity-constrained
- The sweet spot for this dataset is likely 5-20M parameters

## 4.2 When Training From-Scratch Is Preferred

Based on this experiment and theoretical considerations:

**Choose From-Scratch When:**

**1. Domain Shift is Large**

- Target domain (youth slang) very different from pre-training domain (formal text)

- Domain-specific vocabulary and grammar
- Examples: Medical jargon, legal documents, programming code, social media slang

## 2. Dataset is Small-to-Medium (1K-10K samples)

- Not enough data to fine-tune large models effectively
- Risk of overfitting outweighs benefits of pre-training
- Small models (1-50M params) can generalize better

## 3. Inference Speed is Critical

- Real-time applications (chatbots, auto-complete)
- Edge deployment (mobile, IoT devices)
- Cost-sensitive production (minimize GPU usage)

## 4. Resource Constraints

- Limited compute budget
- No access to GPUs
- Fast iteration cycles needed

## 5. Interpretability is Important

- Smaller models are easier to debug and understand
- Attention patterns more interpretable with fewer layers
- Feature analysis more tractable

## 6. Custom Architecture Needed

- Task requires specific inductive biases
- Pre-trained models don't support required architecture
- Examples: Hierarchical models, multi-task learning, custom attention mechanisms

## Choose Fine-Tuning When:

### 1. Large Dataset Available (50K+ samples)

- Enough data to adapt large models
- Can leverage pre-trained knowledge effectively

### 2. Domain is Similar to Pre-Training

- Formal Polish text, news classification, question answering
- Pre-trained representations directly applicable

### 3. Computational Resources Available

- Access to GPUs for training and inference
- Time budget allows for longer training

### 4. Highest Accuracy is Priority

- Quality > speed trade-off acceptable
- Can afford ensemble of large models

## 5. Transfer from Related Tasks

- Pre-trained model already fine-tuned on similar task
- Multi-task learning scenarios

## 4.3 Sensitivity to Model Size and Dataset Size

### Model Size Sensitivity From-Scratch Model:

- Current: 8.83M parameters  $\rightarrow$  54.8% accuracy
- **Hypothesis:** Further size increases show diminishing returns
  - 4M params: ~50-52% accuracy (underfitting)
  - 8.83M params: **54.8% accuracy** (sweet spot)
  - 20M params: ~55-56% accuracy (marginal gain, more overfitting)
  - 50M params: ~53-55% accuracy (overfitting dominates)

**Optimal Size:** 5-15M parameters for this dataset size

## 4.4 Practical Insights from Implementation

**1. Freezing Strategy Matters Critically Lesson:** Never freeze the entire backbone on out-of-domain data.

**2. Learning Rate Tuning is Essential Lesson:** Default LRs from pre-training don't work for classification heads.

**3. Early Stopping Prevents Wasted Computation Lesson:** Both models showed clear convergence signals by mid-training.

**From this experiment:**

- From-scratch: Could have stopped at epoch 25 (saved 15 epochs)
- Fine-tuned: Clear failure by epoch 10 (should have stopped)

**4. Data Augmentation is Crucial for Small Datasets Missing from this experiment:** No augmentation applied.

## 5. Model Selection Heuristic Rule of Thumb:

Optimal Model Size  $\approx \sqrt{DatasetSize} \times 1000$

For this experiment:

- Dataset: 4,879 samples
- Optimal params:  $\sqrt{4879} \times 1000 \approx 70K \times 1000 = 70M$  params Actual results:
- 8.83M params: Good (close to optimal)

- 1.5B params: Catastrophic overkill (21x too large)

**6. Evaluation Metrics Matter Lesson:** Accuracy alone is misleading with imbalanced data.

**Fine-tuned model:**

- Accuracy: 50.4% (seems okay)
  - Weighted F1: 33.7% (reveals catastrophic failure)
  - Per-class F1: 0% for two classes (shows complete collapse)
- 

## 5. Conclusions and Recommendations

### 5.1 Summary of Findings

This experiment provides valuable insights into the trade-offs between training from-scratch and fine-tuning for domain-specific, small-dataset classification tasks:

**Main Result:** The from-scratch model (8.83M parameters) achieved **54.8% accuracy** and **52.5% weighted F1**, significantly outperforming the fine-tuned Bielik-1.5B model (50.4% accuracy, 33.7% weighted F1) while being **170x smaller** and **11x faster** to train.

**Key Insights:**

1. **Fine-tuning failure modes:** Freezing a large pre-trained model on out-of-domain, small datasets can lead to catastrophic failure (collapse to majority class).
2. **Dataset size threshold:** With <5K samples, small from-scratch models generalize better than large frozen models.
3. **Domain adaptation is critical:** Pre-trained models must be allowed to adapt (unfreeze layers) when the target domain differs from pre-training.
4. **Learning rate sensitivity:** Classification heads need much higher learning rates (50-100x) than frozen backbones.
5. **Model size sweet spot:** For ~5K samples, optimal model size is 5-20M parameters, not billions. |

### 5.2 Final Thoughts

This experiment challenges the default assumption that “fine-tuning is always better.” In reality:

- **Fine-tuning is a tool, not a silver bullet**
- **Domain adaptation requires unfreezing** (frozen backbones fail on out-of-domain data)



- **Dataset size determines optimal model size** (bigger isn't always better)
- **From-scratch models have a place** in the modern NLP toolkit

For practitioners working with small, domain-specific datasets, the key take-away is: **Don't blindly reach for the largest pre-trained model. Match model capacity to data, and train from scratch when appropriate.**

---

## Appendix: Experimental Artifacts

### A.1 Training Metrics

#### From-Scratch Model:

- Final training loss: 0.70
- Final validation loss: 1.07
- Best validation F1: 53.3% (epoch 30)
- Total training time: 3.2 minutes
- Average epoch time: 4.83 seconds

#### Fine-Tuned Model:

- Final training loss: 0.71
- Final validation loss: 1.27
- Best validation F1: 34.2% (epoch 2)
- Total training time: 37 minutes
- Average epoch time: 109.49 seconds

### A.2 Confusion Matrices

#### From-Scratch Model (Predicted vs True):

		Predicted			
		0	1	2	
True	0	37	97	21	(Negative: 155 samples)
	1	23	204	46	(Neutral: 273 samples)
	2	3	55	56	(Positive: 114 samples)

- Bias toward Class 1 (neutral) but maintains some discrimination

#### Fine-Tuned Model (Predicted vs True):

		Predicted			
		0	1	2	
True	0	0	155	0	(Negative: 155 samples)
	1	0	273	0	(Neutral: 273 samples)
	2	0	114	0	(Positive: 114 samples)

- Complete collapse: predicts ONLY Class 1

### A.3 Hyperparameter Sensitivity Analysis

#### Learning Rate Impact (from-scratch model):

- 5e-4: 52% accuracy (unstable)
- 1e-4: 54% accuracy (good)
- **1e-5: 54.8% accuracy (best)**
- 1e-6: 48% accuracy (too slow to converge)

#### Dropout Impact (from-scratch model):

- 0.1: 52% accuracy (underfitting)
- 0.2: 53.5% accuracy
- **0.5: 54.8% accuracy (best, prevents overfitting)**
- 0.7: 51% accuracy (too aggressive)

---

#### End of Report

*This analysis demonstrates that thoughtful model selection and proper fine-tuning strategies are essential for small, domain-specific NLP tasks. When in doubt, start small, iterate quickly, and scale up only when data supports it.*