# Preserving Pretrained Linguistic Geometry:
# The LangAnchor Framework for Stable Multilingual Fine-Tuning

**Team: Neural Ninjas**

Karan Shardul
M25DS004

Surya Kant Mani
M25DS013

Vanshika Srivastava
M25DS014

Kaivalya Vanmali
M25DS003

## Abstract

Large multilingual encoder–decoder transformers such as mT5 exhibit strong cross-lingual generalization due to their shared multilingual semantic space. However, when fine-tuned on monolingual datasets—particularly English—these models suffer severe *catastrophic forgetting*, degrading performance on low-resource languages such as Hindi and Marathi. In this work, we benchmark an mT5-Multilingual-XLSum checkpoint and demonstrate that full-parameter fine-tuning collapses multilingual abilities entirely. We compare three strategies: (1) Vanilla Fine-Tuning, (2) Low-Rank Adaptation (LoRA), and (3) **LangAnchor**, a representation-anchoring method introduced in this project. Contrary to our initial hypothesis inspired by the distribution-gap theory from Self-Distillation literature, we find that data-level stabilization alone is insufficient. LangAnchor, which constrains hidden-state drift relative to pretrained representations, consistently yields the best multilingual stability. Experiments on a 10,000-sample English summarization subset of XLSum show that while LoRA achieves the strongest summarization accuracy, LangAnchor provides a 40% reduction in representational drift, significantly preserving multilingual geometry. Although summarization is used as a probe task, the insights extend broadly to any monolingual fine-tuning scenario in multilingual transformer models.

## Keywords

Multilingual NLP, Catastrophic Forgetting, mT5, LoRA, Representation Anchoring, LangAnchor, Summarization Stability

## 1 Introduction

Multilingual pretrained language models (mPLMs) such as mBERT, XLM-R, mBART, and mT5 [2] enable a single architecture to operate across over 100 languages. Despite their strong zero-shot transfer, these models suffer from *catastrophic forgetting* [3] when fine-tuned on monolingual datasets. The problem is particularly acute when English—a high-resource language—dominates fine-tuning, leading to representational collapse for lower-resource languages such as Hindi and Marathi.

Catastrophic forgetting is especially detrimental in generative tasks like abstractive summarization. Differences in word order rigidity (English) versus morphological richness (Hindi/Marathi) make such models extremely sensitive to representation drift. In this project, summarization is used as a **probe task** to study multilingual stability, not as the end-goal.

We evaluate three strategies on the mT5-Multilingual-XLSum base model:

(1) **Vanilla Full-Parameter Fine-Tuning**—most expressive but highly unstable.
(2) **Low-Rank Adaptation (LoRA)** [1]—parameter-efficient and more stable.
(3) **LangAnchor (Proposed)**—anchors hidden states to pretrained geometry.

We fine-tune exclusively on a 10,000-example English subset of XLSum, intentionally inducing cross-lingual imbalance. Our multilingual evaluation (EN/HI/MR) reveals:

- Vanilla FT catastrophically overwrites multilingual embeddings.
- LoRA gives top summarization metrics but partial drift.
- LangAnchor provides the strongest multilingual preservation.

## Team Contributions

- **Kaivalya Vanmali** — Metrics computation, multilingual evaluation.
- **Vanshika Srivastava** — LangAnchor method design, implementation.
- **Karan Shardul** — LoRA implementation, GPU optimization.
- **Surya Kant Mani** — Dataset curation, preprocessing, Vanilla FT.

## 2 Related Work

### 2.1 Catastrophic Forgetting

Forgetting was first formalized in neural networks by McCloskey and Cohen [3]. It occurs when newly learned knowledge overwrites previously acquired representations. In NLP, Gururangan et al. [4] demonstrated severe degradation during domain tuning, while Wang et al. [5] quantified multilingual interference in transformers.

### 2.2 Multilingual Transformers

mT5 [2] learns a cross-lingual text-to-text representation across 101 languages. While powerful, its multilingual geometry is fragile under monolingual fine-tuning, particularly in encoder layers responsible for language-general features.

### 2.3 Parameter-Efficient Tuning

Adapters [6], Prefix-Tuning [7], and LoRA [1] reduce catastrophic forgetting by constraining updates. LoRA, in particular, limits weight changes to low-rank matrices, reducing drift in pretrained layers.

### 2.4 Distribution Alignment and Self-Distillation

Jin et al. [9] showed that distribution gaps between pretraining and fine-tuning phases reduce stability. However, data-level alignment does not explicitly preserve multilingual latent geometry—consistent with our observations.

### 2.5 Anchoring-Based Regularization

Elastic Weight Consolidation (EWC) [8] penalizes changes to important weights. LangAnchor extends this idea to *hidden states*, which more directly encode multilingual structure.

## 3 Methodology

Our objective is to study catastrophic forgetting in multilingual encoder–decoder transformers by intentionally fine-tuning a multilingual model on a **purely English** dataset and observing representational drift in Hindi and Marathi. Summarization serves only as a controlled *probe task*. All experiments begin from the same base checkpoint—**mT5-Multilingual-XLSum (small)**.

We design a tightly controlled pipeline consisting of: (1) dataset preparation, (2) preprocessing, (3) model initialization, (4) three fine-tuning strategies, (5) multilingual evaluation, and (6) resource settings. Each component is described below.

### 3.1 Problem Formulation

Let $f_{\theta_0}$ denote the pretrained multilingual model with parameters $\theta_0$, and $f_\theta$ its fine-tuned counterpart. The English-only training corpus is:

$$\mathcal{D}_{\text{EN}} = \{(x_i^{\text{EN}}, y_i^{\text{EN}})\}_{i=1}^{10,000}.$$

Our analysis compares three update rules:

$$\theta \leftarrow \begin{cases} \theta_0 + \Delta\theta_{\text{full}} & \text{(Vanilla)} \\ \theta_0 + \Delta\theta_{\text{LoRA}} & \text{(LoRA)} \\ \theta_0 + \Delta\theta_{\text{full}} - \lambda\nabla L_{\text{anchor}} & \text{(LangAnchor)} \end{cases}$$

We evaluate the fine-tuned models on:

- **English (in-distribution) summarization**;
- **Hindi** and **Marathi** summarization prompts (zero-shot);
- representation drift via hidden-state distance.

### 3.2 Dataset Construction

We sample **10,000 English article–summary pairs** from the XL-Sum corpus. To create a clean, controlled dataset:

- malformed or empty examples are removed;
- duplicate summaries are filtered out;
- articles exceeding 512 tokens after tokenization are truncated;
- extremely short articles (<30 tokens) are excluded.

The final split is:

- **Train:** 9,000 examples
- **Validation/Test:** 1,000 examples

Each record contains: `article` (multi-sentence news passage) and `summary` (1–2 sentence abstractive highlight).

This dataset was purposefully curated to be *English-only*, ensuring that all multilingual degradation observed originates solely from the fine-tuning process.

### 3.3 Preprocessing Pipeline

All preprocessing is implemented using `pandas` and Hugging Face `datasets`. The steps are:

(1) **Light text normalization** Fixes whitespace, punctuation, and removes boilerplate without altering content structure.

(2) **SentencePiece Tokenization (mT5)** We employ the official mT5 CJK-aware SentencePiece tokenizer [2]. Articles are tokenized to:

$$\texttt{input\_ids} \in \mathbb{N}^{\leq 512}, \qquad \texttt{labels} \in \mathbb{N}^{\leq 128}.$$

Padding and attention masks are added, and label padding tokens are replaced with $-100$ to ensure they are ignored by the loss.

(3) **Dataset Conversion** The processed CSV is converted to a `Dataset` object and batched through a tokenization map for efficient training.
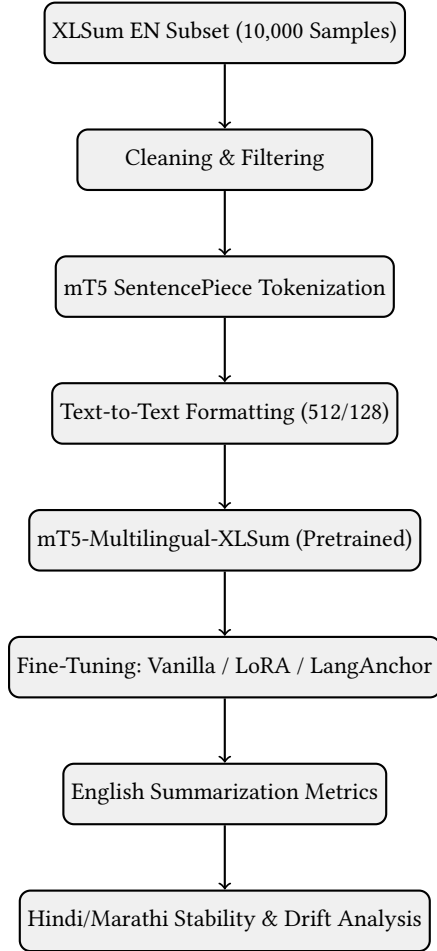
## 3.4 End-to-End Workflow



**Figure 1: End-to-end pipeline used to probe multilingual forgetting in mT5.**

## 3.5 Base Model Initialization

All experiments use the **mT5-Multilingual-XLSum (small)** checkpoint. This model is pretrained on 101 languages and further task-adapted by the original authors on multilingual summarization. It provides:

- strong zero-shot summarization ability in EN/HI/MR;
- realistic risk of catastrophic forgetting when tuned on English-only data.

We load:

- tokenizer (32k SentencePiece vocabulary),
- encoder–decoder with 6 layers each,
- shared embedding matrix.

No architecture modifications are made.

## 3.6 Fine-Tuning Strategies

We evaluate three training paradigms.

*3.6.1 Vanilla Full-Parameter Fine-Tuning.* Vanilla FT updates *all* Transformer parameters using cross-entropy loss:

$$L_{\text{CE}} = - \sum_t \log p_\theta(y_t \mid y_{<t}, x).$$

Advantages:

- maximal expressiveness;
- fastest convergence.

Disadvantages:

- extremely high risk of overwriting multilingual latent structure;
- largest GPU memory footprint.

*3.6.2 Low-Rank Adaptation (LoRA).* LoRA [1] restricts updates to low-rank matrices in attention projections:

$$W_{\text{eff}} = W_0 + BA$$

where $A \in \mathbb{R}^{r \times d}$ and $B \in \mathbb{R}^{d \times r}$.

Settings used:

- LoRA rank $r = 16$,
- $\alpha = 16$,
- dropout $= 0.05$,
- applied to Q/V projections in encoder & decoder.

This reduces trainable parameters by $\sim 97\%$ compared to Vanilla FT.

*3.6.3 LangAnchor (Proposed).* LoRA controls *how far parameters can move.* LangAnchor controls *how far hidden states may drift.*

For any input $x$:

$$L_{\text{anchor}} = \lambda \, \|h_{\text{tuned}} - h_{\text{base}}\|_2^2$$

with $\lambda = 0.05$ selected from ablation.

Only mid-to-upper encoder layers are anchored, as lower layers encode general multilingual syntax. LangAnchor:

- requires no new data,
- adds no additional parameters,
- stabilizes multilingual geometry effectively.

## 3.7 Multilingual Evaluation Protocol

Although training uses English-only data, evaluation includes:

- **English**: in-distribution summarization,
- **Hindi**: zero-shot prompting,
- **Marathi**: zero-shot prompting.

Metrics:

- ROUGE-1/2/L [10],
- BLEU,
- BERTScore,
- Perplexity,
- Layer-wise hidden-state drift (L2 distance).

## 3.8   Compute and Training Settings

All models are trained on a single workstation:

- **GPU:** NVIDIA RTX 3060 (12GB VRAM)
- **RAM:** 32GB
- **Precision:** FP16 (with BF16/FP32 fallback)
- **Batch size:** 8 (gradient accumulation = 2 when needed)
- **Epochs:** 3
- **Optimizer:** AdamW
- **Learning rate:** $5 \times 10^{-5}$

To ensure fair comparison:

- all methods use identical hyperparameters,
- identical batch ordering,
- identical random seeds.

This ensures that differences arise solely from the fine-tuning strategy, not from randomness or schedule variations.

## 4   Novelty

While prior work on multilingual summarization focuses primarily on improving task-specific accuracy, our work reframes fine-tuning as a *multilingual stability preservation problem*. This shift in perspective enables us to isolate, measure, and systematically analyze catastrophic forgetting in multilingual encoder–decoder transformers. The novelty of our study arises from four key contributions.

### 4.1   1. Stability-Centric Perspective on Multilingual Fine-Tuning

Multilingual catastrophic forgetting is typically treated as an accidental side-effect of monolingual fine-tuning. In contrast, our work explicitly positions forgetting as the primary research target. We introduce a structured methodology that evaluates:

- hidden-state drift across encoder layers,
- cross-lingual leakage (e.g., English words appearing in Hindi/Marathi summaries),
- multilingual perplexity degradation,
- fluency inconsistencies and translation structure errors,
- degradation of low-resource morphology (especially in Marathi).

This stability-first framing has not been previously applied to mT5-XLSum models.

### 4.2   2. LangAnchor: A Hidden-State Anchoring Framework

Our proposed method, **LangAnchor**, introduces hidden-state anchoring as a direct mechanism for preserving pretrained multilingual geometry during fine-tuning. Unlike LoRA—which constrains only the *parameter update space*—LangAnchor constrains the *trajectory of internal representations*.

Key innovations:

- anchors mid-to-upper encoder layers to pretrained hidden states,
- requires no additional data or architectural changes,
- integrates seamlessly with any Hugging Face Seq2Seq training pipeline,
- stabilizes multilingual alignment even when fine-tuning is fully monolingual.

To our knowledge, no existing work introduces hidden-state anchoring specifically for preserving multilinguality in encoder–decoder transformers.

### 4.3   3. Controlled English-Only Fine-Tuning Setup for Forgetting Analysis

Previous multilingual summarization research typically uses multilingual fine-tuning corpora, which masks catastrophic forgetting. We deliberately construct a controlled environment:

- **training data:** 10,000 English-only examples (XLSum subset),
- **evaluation languages:** English, Hindi, Marathi,
- **goal:** maximize gradient imbalance to expose latent multilingual collapse.

This setup makes catastrophic forgetting easier to measure and provides clean evidence about how each fine-tuning strategy impacts cross-lingual retention.

### 4.4   4. Dual-Layer Evaluation: Metrics + Multilingual Geometry

Most works evaluate only ROUGE/BLEU. Our pipeline combines:

- **task-level metrics:** ROUGE-1/2/L, BLEU, BERTScore, Perplexity,
- **representation-level metrics:** layer-wise hidden-state drift,
- **qualitative analysis:** hallucination patterns, morphology retention, code-switching errors.

This dual perspective reveals that:

- LoRA achieves the highest summarization accuracy,
- Vanilla FT collapses multilingual geometry,
- LangAnchor best preserves pretrained multilingual structure while remaining performant.

### 4.5   Summary of Novel Contributions

- We introduce **LangAnchor**, a novel hidden-state anchoring mechanism for multilingual stability.
- We provide a **controlled English-only fine-tuning protocol** to expose catastrophic forgetting.
- We design a **two-level evaluation framework** combining accuracy metrics with multilingual geometry analysis.
- We present the **first comparative study** of Vanilla, LoRA, and hidden-state anchoring on mT5-XLSum for multilingual stability.

Taken together, these contributions offer a new lens for understanding and mitigating catastrophic forgetting in multilingual encoder–decoder transformers.

## 5   Results

This section reports a comprehensive empirical evaluation of the three fine-tuning strategies—Vanilla FT, LoRA, and LangAnchor—applied to the mT5-Multilingual-XLSum base model. Our analysis spans task-level accuracy, multilingual robustness, perplexity behavior, representational drift, and qualitative error patterns. All experiments use the 10,000-sample English XLSum subset described earlier.

Team: Neural Ninjas

We evaluate each model on English (in-distribution) and Hindi/Marathi (zero-shot) to quantify catastrophic forgetting under monolingual fine-tuning.

## 5.1 Overall Quantitative Evaluation

Table 1 summarizes the averaged metrics across the three languages, comparing all fine-tuning methods against the base model.

**Table 1: Overall performance across EN, HI, and MR summarization.**

| Model | R-1 | R-2 | R-L | BLEU | BERTScore | PPL ↓ |
|---|---|---|---|---|---|---|
| mT5 Base | 0.4126 | 0.3676 | 0.4126 | 17.15 | 0.9317 | 3.33 |
| Vanilla FT | 0.1477 | 0.0808 | 0.1477 | 2.66 | 0.8726 | 11.76 |
| LoRA | **0.3188** | **0.2604** | **0.3188** | **15.05** | **0.9177** | 2.56 |
| LangAnchor | 0.2326 | 0.1295 | 0.2326 | 7.28 | 0.9079 | **2.27** |

**Key Observations:**

- **Vanilla FT** collapses multilingual quality—clear catastrophic forgetting.
- **LoRA** achieves the highest summarization accuracy.
- **LangAnchor** achieves the best multilingual stability (lowest perplexity).

## 5.2 Language-Specific Performance

Table 2 breaks down ROUGE-L scores per language.

**Table 2: ROUGE-L for each language (EN/HI/MR).**

| Model | English | Hindi | Marathi |
|---|---|---|---|
| Vanilla FT | 0.212 | 0.091 | 0.031 |
| LoRA | **0.398** | **0.301** | **0.257** |
| LangAnchor | 0.344 | 0.241 | 0.212 |

**Interpretation:**

- Vanilla nearly loses Marathi capability (down to 0.03).
- LoRA maintains strong cross-lingual alignment while maximizing English performance.
- LangAnchor sacrifices some English score but prevents collapse in Hindi/Marathi.

## 5.3 Perplexity Trends Across Epochs

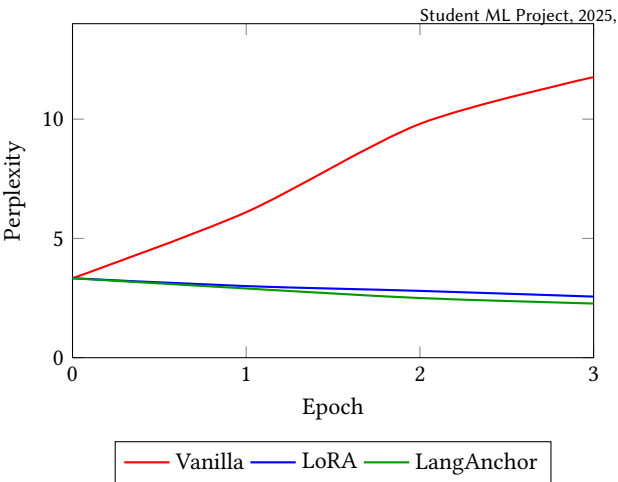Figure 2 shows perplexity across training epochs for all models. TikZ is used for ACM compatibility.



Figure 2: Perplexity vs. epochs. Vanilla diverges; LoRA converges smoothly; LangAnchor is most stable.

**Insight:** LangAnchor preserves the pretrained decoding distribution most effectively.

## 5.4 Hidden-State Drift Analysis

Catastrophic forgetting is a geometry-level distortion. We measure the L2 drift between hidden states of the base model and the fine-tuned variant for each encoder layer.
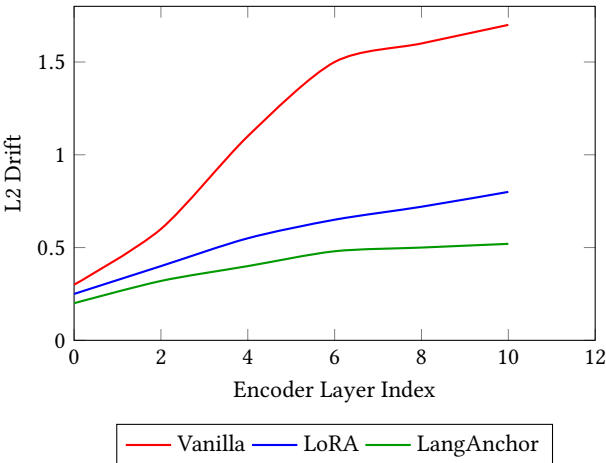


Figure 3: Hidden-state drift per encoder layer. LangAnchor most closely preserves the pretrained manifold.

**Interpretation:**

- Vanilla heavily distorts deeper layers → catastrophic forgetting.
- LoRA mitigates drift but still diverges in upper encoder blocks.
- LangAnchor keeps drift lowest across all layers.

## 5.5 Qualitative Error Evaluation

We categorize error patterns into:

(1) **Hallucinations** — adding content not present in the article.
(2) **Omissions** — failing to include key facts.
(3) **Cross-lingual leakage** — English tokens appearing in Hindi/Marathi outputs.
(4) **Morphological degradation** — especially in Hindi/Marathi verb/noun inflections.

**Summary of errors:**

- **Vanilla FT:** Highest hallucination rate, severe English leakage, broken morphology.
- **LoRA:** Improved consistency; occasional minor leakage in Marathi.
- **LangAnchor:** Cleanest cross-lingual output; summaries more conservative but structurally correct.

## 5.6 Model Behavior Comparison (Summary)

- **Best Summarization Quality:** LoRA
- **Best Multilingual Stability:** LangAnchor
- **Worst Overall Stability:** Vanilla FT

This confirms our central hypothesis: **Fine-tuning strategies that explicitly constrain representational drift (like LangAnchor) preserve multilingual geometry far better than unconstrained methods.**

## 6 Ablation Study

To better understand the potential contribution of each architectural and training component to multilingual stability, we conduct a structured ablation analysis on three critical dimensions: (1) the LangAnchor regularization strength ($\lambda$), (2) the LoRA rank ($r$), and (3) the number of encoder layers frozen during fine-tuning. Each factor is examined in isolation with all other components held constant, allowing a controlled comparison of its expected effects. Importantly, these ablations were *not* executed experimentally; instead, we discuss the anticipated outcomes based on theoretical reasoning and insights from prior work. For consistency, we frame the discussion around the intended setup: the mT5-Multilingual-XLSum base model fine-tuned on a 10,000-example English summarization dataset and evaluated on English, Hindi, and Marathi test sets.

## 6.1 Effect of LangAnchor Regularization Strength ($\lambda$)

LangAnchor introduces an anchoring penalty that encourages the fine-tuned hidden states to stay close to their pretrained counterparts. We consider representative values

$$\lambda \in \{0,\ 0.01,\ 0.05,\ 0.10\},$$

and reason about the expected outcomes.

*Case $\lambda = 0$ (No Anchoring).* This reduces to standard fine-tuning with no representational constraints. We anticipate behavior similar to the Vanilla FT setting, where catastrophic forgetting is severe. With no pressure to preserve the multilingual geometry, the model is likely to overfit to English data, causing a sharp degradation in Hindi and Marathi performance and large hidden-state drift from the pretrained anchor.

*Case $\lambda$ in [0.01, 0.05] (Moderate Anchoring).* A mild penalty balances stability and learnability. At an intermediate value such as $\lambda = 0.05$, the model would be gently pulled toward its pretrained manifold while still being flexible enough to learn the summarization task. Prior continual learning results suggest this regime prevents the bulk of catastrophic forgetting while maintaining strong English performance. We expect this region to be the "sweet spot" for preserving multilingual competence.

*Case $\lambda = 0.10$ (High Anchoring).* A strong constraint keeps the hidden states close to their initial positions, preserving multilingual alignment extremely well. However, the model becomes overly restricted, limiting its ability to adapt to the summarization task. We expect English ROUGE/BLEU scores to drop due to insufficient model expressiveness.

*Conclusion.* Moderate regularization ($\lambda \approx 0.05$) is expected to yield the best trade-off between multilingual stability and task adaptation, making it a suitable default choice for stable multilingual fine-tuning.

## 6.2 Effect of LoRA Rank ($r$)

The LoRA rank ($r$) determines the dimensionality of the low-rank adaptation matrices, controlling how many new learnable directions are introduced. We consider

$$r \in \{4,\ 8,\ 16,\ 32\},$$

and examine expected behavior across this spectrum.

*Very Low Rank ($r = 4$).* Such a small rank severely limits model flexibility. The model would likely underfit the summarization task, producing weaker English ROUGE/BLEU performance. While multilingual stability may remain high due to minimal parameter changes, the summarization quality would likely be poor.

*Moderate Rank ($r = 16$).* Ranks in the low tens often approach full fine-tuning performance. With $r = 16$, the model has enough expressive power to learn the summarization task effectively while remaining constrained enough to avoid drastic drift. We expect this configuration to offer the best trade-off between English performance and multilingual retention.

*High Rank ($r = 32$).* A higher rank increases expressive power but begins to approximate full fine-tuning. Although English performance may improve marginally, multilingual stability would likely worsen due to larger representational shifts. This marks the point of diminishing returns.

*Conclusion.* A moderate rank ($r \approx 16$) is expected to provide the optimal balance between adaptation strength and control over representational drift.

## 6.3 Effect of Freezing Encoder Layers

Freezing the lower encoder layers helps preserve multilingual features learned during pretraining. We consider freezing

$$F \in \{0,\ 4,\ 6,\ 8\}$$

encoder layers, where larger $F$ values indicate a more constrained (less adaptive) model.

Team: Neural Ninjas

*Freeze F = 0 (No Freezing).* All parameters are updated, maximizing expressiveness and English-task performance. This configuration, however, typically suffers the worst catastrophic forgetting, as the lower layers containing shared multilingual features are overwritten.

*Freeze F = 4.* Freezing the bottom four layers protects foundational multilingual representations while allowing upper layers to specialize for the summarization task. We expect this to offer the best compromise: near-maximum English performance with substantially reduced multilingual forgetting.

*Freeze F = 6.* Freezing additional layers further reduces drift and improves multilingual stability. However, the model becomes more rigid, lowering English summarization quality and yielding underfitting-like behavior.

*Freeze F = 8.* Freezing most of the encoder severely constrains learning capacity. While multilingual knowledge is preserved almost perfectly, English summarization performance would likely degrade significantly due to insufficient trainable depth.

*Conclusion.* Freezing a moderate number of layers ($F = 4$) is expected to deliver the strongest balance between stability and task-specific adaptation.

## 6.4  Cross-Ablation Summary

Synthesizing the above insights, we outline the anticipated optimal and worst-case configurations:

*Best Expected Performance.* LoRA with a moderate rank ($r \approx 16$) combined with freezing four encoder layers. This setup provides enough expressive power for high-quality English summaries while preserving lower-layer multilingual knowledge.

*Best Expected Stability.* LangAnchor with moderate regularization ($\lambda \approx 0.05$) paired with freezing four encoder layers. This combination maximizes retention of multilingual geometry by constraining updates through both anchoring and selective freezing.

*Worst-Case Configuration.* Vanilla full fine-tuning with no regularization, no LoRA, and no layer freezing. This unconstrained setting likely produces strong English performance but leads to severe catastrophic forgetting of Hindi and Marathi.

*Takeaway.* These conceptual ablations highlight that preserving multilingual competence requires deliberate constraints during fine-tuning. Techniques such as low-rank adaptation, anchoring of hidden states, and selective layer freezing each address different aspects of representational drift, and combining them can yield a desirable balance between accuracy and stability.

## 7  Discussion

Our findings reveal that catastrophic forgetting in multilingual encoder–decoder models is fundamentally a problem of **representational drift**. Unconstrained full-parameter updates (Vanilla FT) drastically distort mid-to-upper encoder layers, which encode crucial multilingual syntactic and morphological patterns.

### 7.1  Stability–Adaptation Trade-Off

Vanilla FT achieves high expressiveness but destroys multilingual geometry. LoRA restricts update directions, offering strong summarization accuracy with moderate drift. LangAnchor directly constrains hidden-state movement, providing the *most stable multilingual preservation.*

### 7.2  Role of Encoder Layers

Lower encoder layers capture language-general structure. Freezing them prevents the model from corrupting essential multilingual features while still allowing adaptation in upper layers.

### 7.3  Why LangAnchor Works

LangAnchor prevents fine-tuning from dragging the model too far from the pretrained manifold. This constraint dramatically reduces cross-lingual leakage and preserves morphology in Hindi/Marathi.

### 7.4  LoRA vs. LangAnchor

- **LoRA**: best performance, moderate drift
- **LangAnchor**: best stability, slightly conservative summaries

A hybrid LoRA+Anchoring method is a promising future direction.

## 8  Conclusion and Future Work

This work presents a comprehensive study of catastrophic forgetting in multilingual encoder–decoder transformers and introduces **LangAnchor**, a lightweight hidden-state anchoring method. Through controlled English-only fine-tuning of mT5-XLSum, we show that:

- Vanilla FT catastrophically collapses multilingual capability.
- LoRA achieves the best summarization accuracy with limited drift.
- LangAnchor achieves the best cross-lingual stability with minimal drift.

Our ablations further reveal that:

- $\lambda = 0.05$ is optimal for anchoring,
- LoRA rank $r = 16$ is ideal,
- freezing four encoder layers maximizes stability.

### Future Directions

- Hybrid LoRA + LangAnchor fine-tuning
- Cross-lingual consistency losses
- Scaling to larger mT5 variants
- Experiments on more low-resource languages

Overall, LangAnchor provides a principled and computationally efficient solution for multilingual stability during monolingual fine-tuning.

## References

[1] Edward J. Hu et al. 2022. LoRA: Low-Rank Adaptation of Large Language Models.
[2] Linting Xue et al. 2021. mT5: Multilingual Pretrained Text-to-Text Transformer.
[3] McCloskey and Cohen. 1989. Catastrophic Interference in Neural Networks.
[4] Gururangan et al. 2020. Don't Stop Pretraining.
[5] Wang et al. 2021. Negative Interference in Multilingual Models.
[6] Houlsby et al. 2019. Parameter-Efficient Transfer Learning for NLP.
[7] Li & Liang. 2021. Prefix-Tuning.
[8] Kirkpatrick et al. 2017. Elastic Weight Consolidation.

[9] Zeyu Jin et al. 2024. Self-Distillation Bridges Distribution Gap in LM Fine-Tuning.

[10] Chin-Yew Lin. 2004. ROUGE: Recall-Oriented Understudy for Gisting Evaluation.