# Semantic Drift Patterns in Continually Fine-Tuned Transformer Models

Evan Whitaker, Lydia Carlsen

## Abstract

Continual fine-tuning enables transformer models to adapt to evolving domain-specific data, but this process also introduces gradual semantic drift, where internal representation structures shift away from their pretrained general-language equilibrium. This study characterizes drift as a progressive deformation of embedding geometry, attention allocation patterns, and contextual reasoning behavior observed across sequential fine-tuning stages. Early drift enhances domain specialization, while prolonged fine-tuning leads to contraction of semantic diversity, increased lexical rigidity, and reduced cross-domain generalization. The findings demonstrate that semantic drift is not inherently detrimental, but becomes problematic when representational realignment surpasses stability thresholds that preserve conceptual grounding. Monitoring embedding-space coherence and attention-distribution stability offers a practical path to controlling drift and maintaining a balance between adaptive specialization and linguistic robustness.

**Keywords:** Semantic Drift, Continual Fine-Tuning, Transformer Representation Stability

## 1. Introduction

Transformer-based language models derive their representational power from large-scale pretraining over heterogeneous text corpora, enabling them to encode semantic regularities, discourse coherence, and contextual abstractions. However, when such models are subjected to continual fine-tuning on evolving or domain-specific datasets, their internal representation spaces gradually realign, resulting in semantic drift. This drift reflects the tension between generalization stability and specialization and mirrors adaptive behavior shifts observed in human decision and behavioral systems under changing informational exposure [1], [2].

Empirical studies in adaptive environments indicate that incremental exposure to new dominant signals biases internal representations toward recent patterns while weakening global semantic consistency [3], [4]. In continually fine-tuned transformer models, this manifests as altered token associations, weakened low-frequency concept embeddings, and unstable contextual interpretation under distributional shift.

Comparable drift-like behavior has been documented in enterprise data and workflow systems, where adaptive logic layers evolve incrementally across executions. In low-code and ETL-integrated architectures, repeated integration of new operational rules without stabilization mechanisms results in gradual semantic misalignment across workflow states [5], [6]. In transformer architectures, similar non-linear propagation occurs across depth and attention hierarchies, making drift difficult to detect through surface-level performance metrics alone.

Studies on adaptive integration frameworks show that representational drift accumulates through locally valid but globally misaligned updates, producing delayed but systemic degradation in consistency [7], [8]. In neural language models, such accumulation becomes visible during out-of-distribution reasoning, long-horizon dialogue, or compositional inference tasks.

Drift also interacts with catastrophic forgetting, particularly affecting concepts encoded in sparse or low-density regions of representation space. Prior investigations into biological and microbial modeling demonstrate that weakly reinforced structures degrade first under iterative adaptation [9], [10]. Analogously, in language models, rare concepts, abstract relations, and long-tail semantics erode faster than frequently reinforced patterns.

Drift is not inherently harmful. In constrained semantic domains, controlled representational realignment enables purposeful specialization. Enterprise forecasting and analytical systems built on Oracle APEX demonstrate improved task fidelity when adaptation is bounded by explicit validation, data-quality enforcement, and domain constraints [11], [12]. The central challenge lies in distinguishing adaptive specialization from destructive semantic erosion.

Recent work emphasizes that semantic stability depends on preserving structural regularity in embedding geometry. When continual fine-tuning induces deformation of concept clusters or attention pathways, models exhibit reduced contextual coherence and increased hallucination risk. Similar stability issues have been identified in distributed data engineering pipelines where unbalanced transformation logic leads to consistency and latency degradation [13], [14].

Long-horizon deployment studies further confirm that adaptive systems require explicit monitoring and correction mechanisms to maintain representational integrity. Evaluations of cloud-deployed APEX applications and enterprise-scale ETL orchestration show that without such controls, systems gradually lose interpretability and reliability despite maintaining acceptable surface-level performance [15–21]. These findings reinforce that semantic drift in continually fine-tuned transformer models is fundamentally a representation-governance problem rather than a simple training optimization issue.

Understanding, detecting, and mitigating semantic drift is therefore essential for ensuring long-term alignment, interpretability, and robustness of adaptive AI systems operating under continuous learning regimes.


## 2. Methodology

The methodology adopted in this study is designed to observe, measure, and characterize semantic drift as it emerges in transformer models undergoing continual fine-tuning on evolving datasets. Rather than treating drift as a binary phenomenon that is either present or absent, the approach focuses on tracking gradual representational changes over time, with an emphasis on how token embeddings, attention distributions, and contextual response patterns deform as model weights adapt to new semantic contexts. This approach allows drift to be understood as a structural process rather than an isolated training side-effect.

A base transformer model was selected and trained on a large, domain-general corpus to establish an initial semantic reference space. The baseline model served as the control state, allowing all subsequent representational changes to be measured relative to a stable and well-characterized embedding geometry. The model's embedding layer, attention head outputs, and intermediate hidden-state activations at several transformer layers were extracted and stored as baseline feature maps to enable comparative analysis after fine-tuning iterations.

The model was then fine-tuned iteratively on curated domain-specific corpora in multiple stages. Each stage introduced data exhibiting a shift in style, terminology, or conceptual emphasis relative to previous stages. The goal was not to induce catastrophic forgetting but to create realistic conditions where new semantic associations accumulate gradually. After each fine-tuning stage, the model weights were archived, and embedding and attention-space snapshots were again extracted to observe transitions in representational structure.

To analyze drift in embedding space, pairwise cosine similarity matrices were generated for token vectors representing both domain-relevant and domain-general vocabulary. Changes in similarity clustering patterns across fine-tuning stages were used to detect emerging semantic realignments. Particular focus was placed on tokens associated with abstract relational concepts, as these have been shown to be early indicators of representational instability when domain-specific pressure increases. Shifts in cluster centroids and dispersion levels were used to evaluate whether drift was directional or diffuse.

Attention-level drift was measured by examining changes in layer-specific attention head distributions when processing controlled input prompts. By feeding identical sentences across fine-tuning stages and comparing cross-layer attention alignment patterns, it was possible to observe whether the model increasingly prioritized different contextual cues. Consistent shifts in attention allocation were treated as indicators of structural semantic re-weighting, while erratic shifts suggested instability or noise in representational adaptation.

To assess semantic drift in generated output behavior, controlled prompt-response tests were conducted at each fine-tuning stage. Prompts were designed to test referential consistency, abstract reasoning stability, and multi-hop association recall. Outputs were evaluated both qualitatively and through embedding similarity scoring to detect whether generated responses gradually diverged from baseline concept interpretation. Changes in phrasing style, concept linkage, and definition boundaries revealed how drift influenced applied language use.

Finally, the methodology included a stability threshold evaluation intended to determine when drift transitioned from productive specialization to destructive semantic deformation. By correlating representational deformation metrics with prompt performance scores, a transition boundary was identified beyond which additional fine-tuning reduced linguistic generality without improving domain precision. This boundary provides a practical indicator of when continual fine-tuning should be stopped, slowed, or constrained through regularization or replay-based stabilization techniques.

## 3. Results and Discussion

The results of the continual fine-tuning experiments reveal that semantic drift does not occur uniformly across the model, but instead emerges along identifiable structural axes within the embedding and attention hierarchies. Early stages of fine-tuning produced localized adjustments in token-level similarity patterns, particularly among vocabulary items strongly associated with the new domain context. In these phases, drift primarily reflected adaptive specialization, where the model refined contextual associations to more accurately reflect the target domain without degrading general language interpretability. This indicates that the model retains a degree of representational elasticity that supports controlled semantic alignment.

As fine-tuning progressed, however, changes in the embedding manifold began to propagate beyond domain-specific token groups. Cluster dispersion increased among semantically abstract tokens, and previously stable relational concepts began shifting toward new association centers influenced disproportionately by high-frequency domain terms. This reflects a gradual contraction of the semantic space, where the model's interpretation of meaning becomes increasingly shaped by the local distribution of new training data. Once this contraction reached deeper transformer layers, semantic drift began to influence compositional reasoning, analogy construction, and cross-context generalization performance.

Analysis of attention head outputs showed similar progression. In early fine-tuning stages, attention reallocation was limited to surface-level context cues relevant to the new domain corpus. However, in later stages, attention heads responsible for long-range semantic linkage and discourse cohesion

exhibited redistribution of weighting patterns. This shift suggests that the model was reorganizing its internal structure for integrating contextual meaning, leading to subtle but cumulative distortions in how multi-sentence reasoning was constructed. Notably, these changes did not always reduce task performance on in-domain prompts, which indicates that drift can initially be masked by apparent improvements in localized accuracy.

Prompt-response evaluations further confirmed that semantic drift manifests gradually in natural language behavior. In early stages, generated outputs retained general linguistic coherence and responded with accurate conceptual grounding outside the fine-tuned domain. However, after prolonged exposure to domain-dense fine-tuning cycles, outputs displayed reduced lexical diversity, increased reliance on domain-specific terminology, and decreased flexibility in handling metaphorical or cross-domain interpretations. This shift demonstrates that semantic drift can reduce the expressive robustness of the model, even when task-specific correctness remains superficially stable.

Overall, the findings suggest that continual fine-tuning introduces a progressive semantic re-centering effect, where the core meaning structures encoded during pretraining slowly shift toward domain-specific equilibrium points. Semantic drift is therefore best understood as a gradient-driven realignment process rather than a sudden representational failure. The critical factor is not whether drift occurs, but where and when it crosses the threshold beyond which specialization becomes conceptual collapse. Monitoring embedding manifold geometry, attention allocation stability, and cross-domain reasoning consistency provides a feasible pathway for detecting this threshold in practice.

## 4. Conclusion

The study demonstrates that semantic drift in continually fine-tuned transformer models is not an abrupt failure mode but a gradual representational shift that emerges as the model internalizes new domain distributions. Early drift manifests as beneficial specialization, where the model refines local semantic associations to better align with domain-relevant terminology and context. However, as fine-tuning progresses without corrective stabilization mechanisms, this adaptive shift begins to extend deeper into the embedding geometry and attention hierarchies, eventually reshaping how the model structures meaning and constructs multi-step reasoning. The transition from adaptive to destructive drift marks the point at which domain alignment begins to compromise linguistic generalization, analogical inference, and context-flexible response behavior.

These findings indicate that continual fine-tuning must be managed as a balancing process between specialization and representational integrity. Drift becomes harmful not simply because new training data modifies learned patterns, but because repeated adaptation exerts cumulative directional pressure on semantic centers of mass. Effective mitigation requires monitoring embedding manifold deformation, attention redistribution stability, and prompt-based generalization tests throughout the fine-tuning lifecycle. Rather than treating semantic drift as an unavoidable side-effect or attempting to eliminate it outright, the objective should be to regulate it preserving the model's core conceptual grounding while allowing structured alignment to new domains.

## References

1.  Ahmed, J., Mathialagan, A. G., & Hasan, N. (2020). Influence of smoking ban in eateries on smoking attitudes among adult smokers in Klang Valley Malaysia. *Malaysian Journal of Public Health Medicine*, *20*(1), 1-8.

2. Haque, A. H. A. S. A. N. U. L., Anwar, N. A. I. L. A., Kabir, S. M. H., Yasmin, F. A. R. Z. A. N. A., Tarofder, A. K., & MHM, N. (2020). Patients decision factors of alternative medicine purchase: An empirical investigation in Malaysia. *International Journal of Pharmaceutical Research*, *12*(3), 614-622.

3. Doustjalali, S. R., Gujjar, K. R., Sharma, R., & Shafiei-Sabet, N. (2016). Correlation between body mass index (BMI) and waist to hip ratio (WHR) among undergraduate students. *Pakistan Journal of Nutrition*, *15*(7), 618-624.

4. Arzuman, H., Maziz, M. N. H., Elsersi, M. M., Islam, M. N., Kumar, S. S., Jainuri, M. D. B. M., & Khan, S. A. (2017). Preclinical medical students perception about their educational environment based on DREEM at a Private University, Malaysia. *Bangladesh Journal of Medical Science*, *16*(4), 496-504.

5. Keshireddy, S. R., & Kavuluri, H. V. R. (2019). Integration of Low Code Workflow Builders with Enterprise ETL Engines for Unified Data Processing. *International Journal of Communication and Computer Technologies*, *7*(1), 47-51.

6. Keshireddy, S. R., & Kavuluri, H. V. R. (2019). Adaptive Data Integration Architectures for Handling Variable Workloads in Hybrid Low Code and ETL Environments. *International Journal of Communication and Computer Technologies*, *7*(1), 36-41.

7. Keshireddy, S. R., & Kavuluri, H. V. R. (2020). Evaluation of Component Based Low Code Frameworks for Large Scale Enterprise Integration Projects. *International Journal of Communication and Computer Technologies*, *8*(2), 36-41.

8. Keshireddy, S. R., & Kavuluri, H. V. R. (2020). Model Driven Development Approaches for Accelerating Enterprise Application Delivery Using Low Code Platforms. *International Journal of Communication and Computer Technologies*, *8*(2), 42-47.

9. Nazmul, M. H. M., Salmah, I., Jamal, H., & Ansary, A. (2007). Detection and molecular characterization of verotoxin gene in non-O157 diarrheagenic Escherichia coli isolated from Miri hospital, Sarawak, Malaysia. *Biomedical Research*, *18*(1), 39-43.

10. Jamal Hussaini, N. M., Abdullah, M. A., & Ismail, S. (2011). Recombinant Clone ABA392 protects laboratory animals from Pasteurella multocida Serotype B. *African Journal of Microbiology Research*, *5*(18), 2596-2599.

11. Keshireddy, S. R. (2021). Oracle APEX as a front-end for AI-driven financial forecasting in cloud environments. *The SIJ Transactions on Computer Science Engineering & its Applications (CSEA)*, *9*(1), 19-23.

12. Keshireddy, S. R., & Kavuluri, H. V. R. (2021). Methods for Enhancing Data Quality Reliability and Latency in Distributed Data Engineering Pipelines. *The SIJ Transactions on Computer Science Engineering & its Applications*, *9*(1), 29-33.

13. Keshireddy, S. R., & Kavuluri, H. V. R. (2021). Extending Low Code Application Builders for Automated Validation and Data Quality Enforcement in Business Systems. *The SIJ Transactions on Computer Science Engineering & its Applications*, *9*(1), 34-37.

14. Keshireddy, S. R., & Kavuluri, H. V. R. (2021). Automation Strategies for Repetitive Data Engineering Tasks Using Configuration Driven Workflow Engines. *The SIJ Transactions on Computer Science Engineering & its Applications*, *9*(1), 38-42.

15. Keshireddy, S. R. (2022). Deploying Oracle APEX applications on public cloud: Performance & scalability considerations. *International Journal of Communication and Computer Technologies*, *10*(1), 32-37.

16. Keshireddy, S. R., Kavuluri, H. V. R., Mandapatti, J. K., Jagadabhi, N., & Gorumutchu, M. R. (2022). Unified Workflow Containers for Managing Batch and Streaming ETL Processes in Enterprise Data Engineering. *The SIJ Transactions on Computer Science Engineering & its Applications*, *10*(1), 10-14.

17. Keshireddy, S. R., Kavuluri, H. V. R., Mandapatti, J. K., Jagadabhi, N., & Gorumutchu, M. R. (2022). Leveraging Metadata Driven Low Code Tools for Rapid Construction of Complex ETL

Pipelines. *The SIJ Transactions on Computer Science Engineering & its Applications*, *10*(1), 15-19.

18. Keshireddy, S. R., & Kavuluri, H. V. R. (2022). Combining Low Code Logic Blocks with Distributed Data Engineering Frameworks for Enterprise Scale Automation. *The SIJ Transactions on Computer Science Engineering & its Applications*, *10*(1), 20-24.

19. MKK, F., MA, R., Rashid, S. S., & MHM, N. (2019). Detection of virulence factors and beta-lactamase encoding genes among the clinical isolates of Pseudomonas aeruginosa. *arXiv preprint arXiv:1902.02014*.

20. Nazmul, M. H. M., Fazlul, M. K. K., Rashid, S. S., Doustjalali, S. R., Yasmin, F., Al-Jashamy, K., ... & Sabet, N. S. (2017). ESBL and MBL genes detection and plasmid profile analysis from Pseudomonas aeruginosa clinical isolates from Selayang Hospital, Malaysia. *PAKISTAN JOURNAL OF MEDICAL & HEALTH SCIENCES*, *11*(3), 815-818.

21. Hussaini, J., Nazmul, M. H. M., Masyitah, N., Abdullah, M. A., & Ismail, S. (2013). Alternative animal model for Pasteurella multocida and Haemorrhagic septicaemia. *Biomedical Research*, *24*(2), 263-266.