# Paper Review- Modyn: Data-Centric Machine Learning Pipeline Orchestration

**Summary**: This paper introduces Modyn, an open-source orchestrator for ML pipelines on continuously growing datasets such as click logs and sensor streams, which also provides tooling for analyzing pipelines. Retraining on full data is costly and slow, so Modyn provides policy-driven retraining through two mechanisms: (i) triggering policies (deciding when to retrain) and (ii) selection policies (deciding which samples to use based on factors like importance/weight). Its pipeline is built around a supervisor (manages triggers and orchestrates flow), a selector (implements sample selection), storage (optimized for random sample retrieval at scale), a trainer (executes retraining), and an evaluator (measures performance). Key innovations include support for sample-level selection, drift-based (drift: an estimate of how different the current data is from the historical data) triggers that work without labels, and system optimizations like partitioning, parallelism, and prefetching. Experiments on datasets like Yearbook, CGLM landmarks, Criteo 1TB, and Kaggle arXiv, shows that Modyn achieves near full-data accuracy with only 50% of training data, adapts retraining to drift efficiently, and delivers 71-98% of local throughput despite the complexity of random sample selection.

**Evaluation**: The problem Modyn addresses is highly significant: in real-world ML deployments, data streams grow endlessly, and distribution shifts are inevitable, making full retraining costly and causing stale models if ignored. Modyn's solution is both valid and novel, offering an end-to-end open-source orchestrator (plus analyser) that supports sample-level data selection, policy-driven triggers, and a declarative pipeline abstraction. A key strength is the use of composite models for fair pipeline comparison and efficient sample-level data access. The experimental evaluation is detailed, spanning multiple modalities (images, text, recommender logs) and combining ML-centric metrics (accuracy vs. selection/triggering policy) with system-centric ones (throughput and overhead vs local baselines). Results clearly show near full-data accuracy at reduced cost and throughput close to system ideal. The paper is clear and well-structured, with figures/graphs to explain workflows/results in a clear manner.

## Main takeaways
1. Smart sample selection can nearly match full-data accuracy at half the cost, showing that data-centric retraining is practical and effective.
2. Different data selection strategies excel under different types of distribution shift (example: uncertainty sampling handles covariate shift, while RS2 or DLIS work well for prior-probability shift), showing that there is no single solution for all.
3. Systems-level optimizations like partitioning, parallelism, and prefetching allow Modyn to deliver close-to-local throughput, proving that sample-level selection is feasible at scale.

## Strengths
1. Modyn balances performance and ease-of-use by implementing critical components in C++ while exposing pluggable policies in Python, allowing researchers to add new models and strategies easily without worrying about low-level systems details.
2. Provides a modular pipeline (supervisor, selector, storage, trainer, evaluator) with clear interfaces, allowing data selection and triggering policies to be swapped or extended independently of the system.
3. Modyn is open-source and comes with an interactive dashboard for pipeline analysis and comparison, making it highly usable and extensible for researchers.

## Weaknesses
1. The choice of drift metrics and embedding spaces is not deeply explored, leaving uncertainty about how well drift detection generalizes across data types such as images, text, and tabular logs (unstructured data).
2. Evaluation focuses on controlled experiments: real-world deployment challenges such as delayed labels, noisy data, and scaling beyond a single server (like distributed environments) are underexplored.

## Discussion topics
1. How practical is drift-based triggering in real-world deployments where data is messy and embeddings may change over time? Could AutoDrift be reliably used in production?
2. Given Modyn reaches close to local throughput, what additional optimizations or architectural changes would be needed to scale it to distributed, multi-node training environments?