

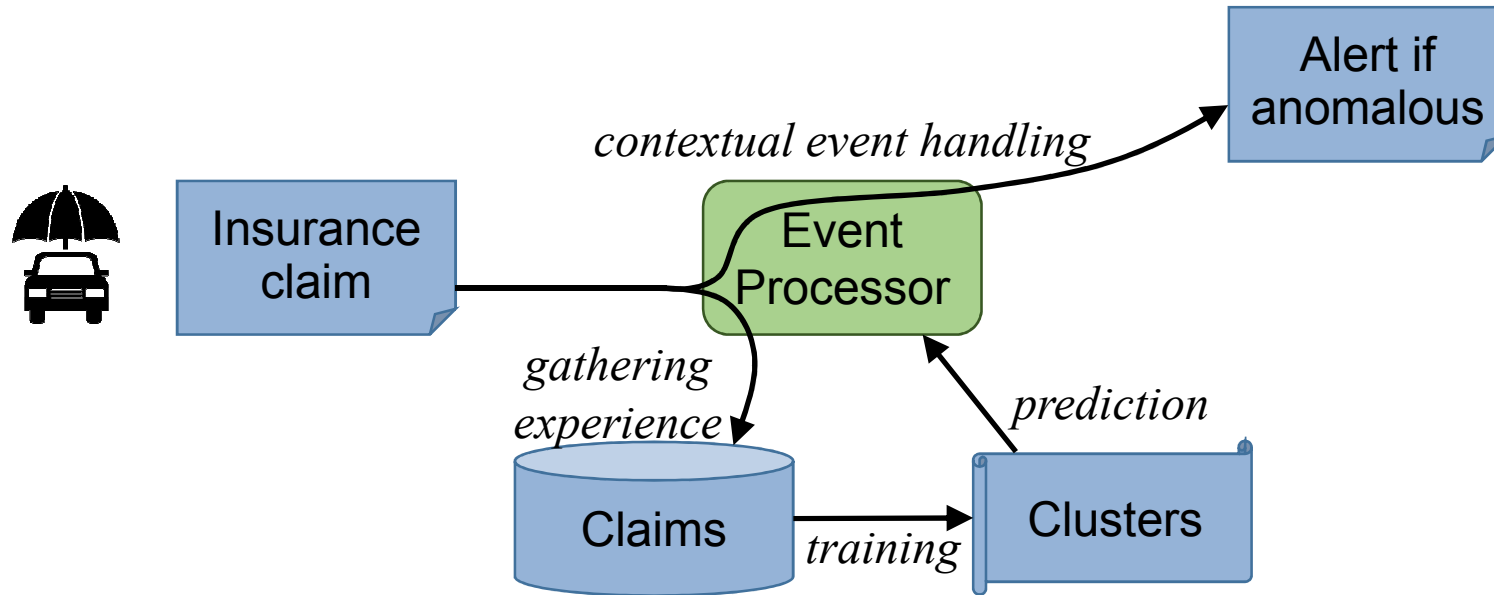
AQuA: Adaptive Quality Analytics

Conference on Distributed Event-Based Systems (DEBS), June 2016

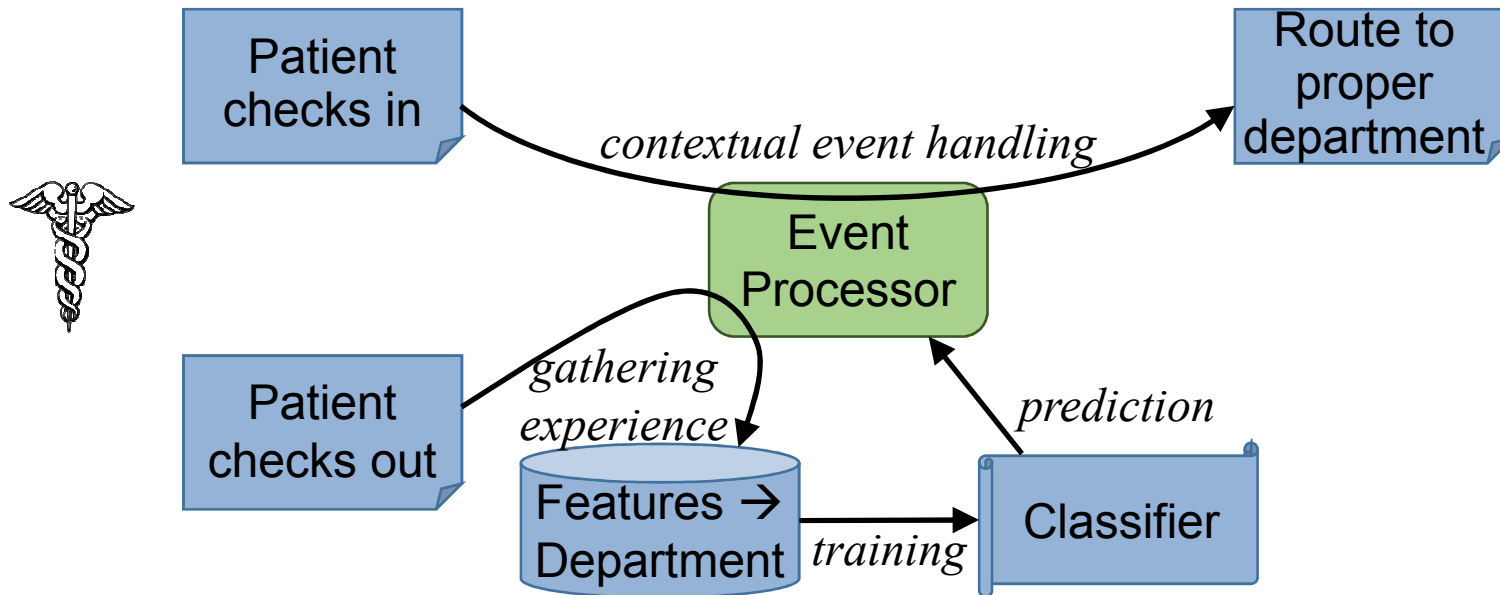
Wei Zhang, Martin Hirzel, and David Grove



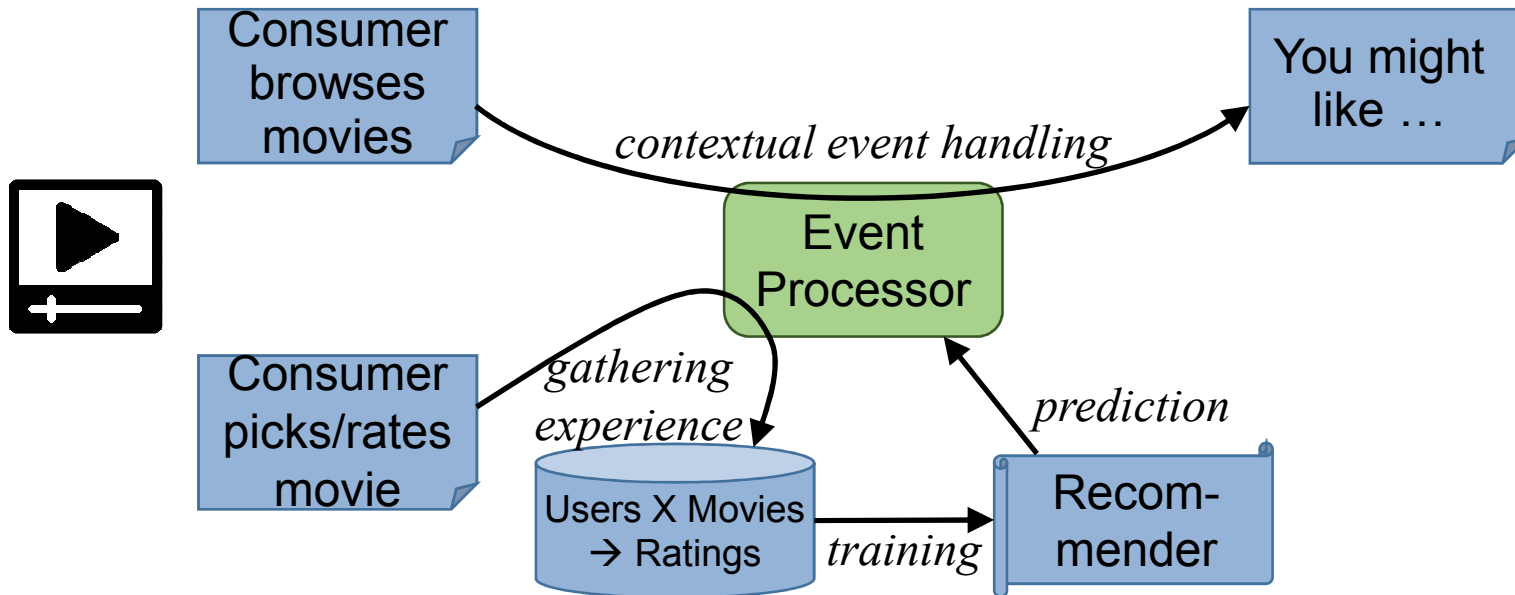
Unsupervised Learning Scenario: Clustering



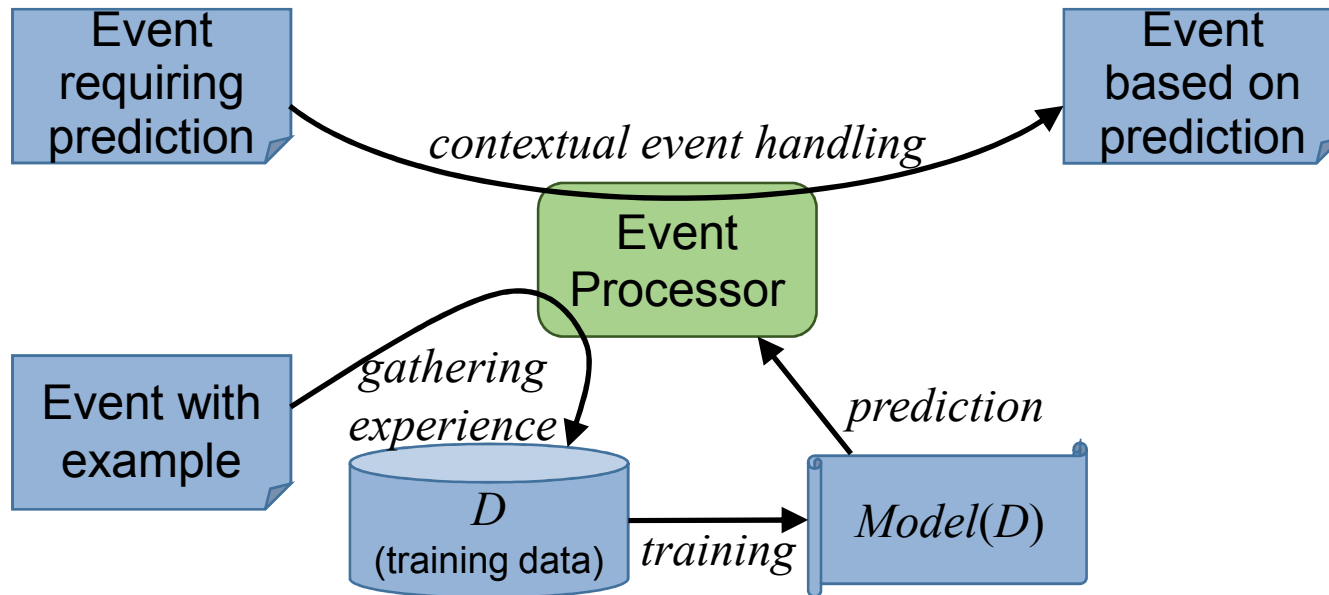
Supervised Learning Scenario: Classification



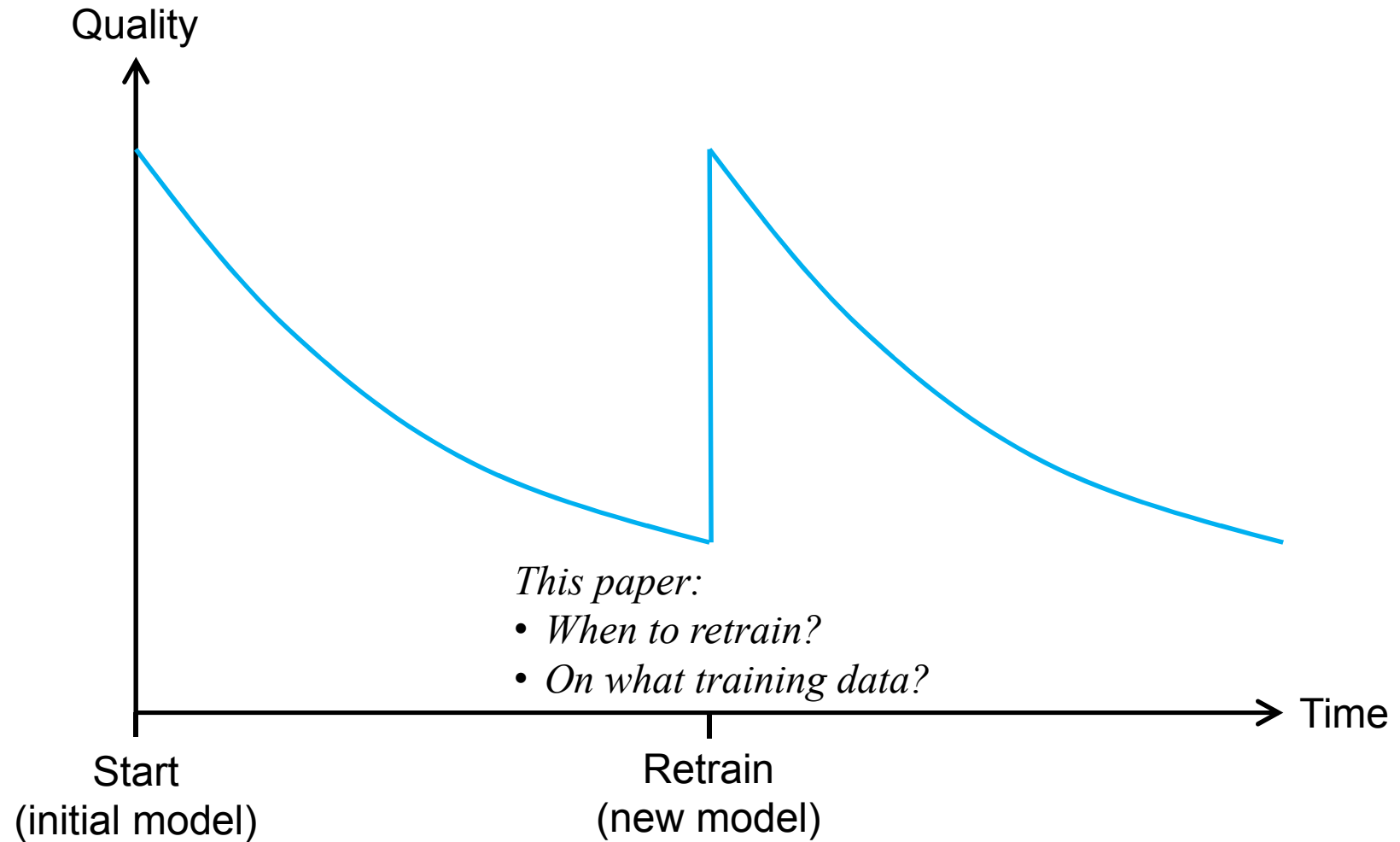
Supervised Learning Scenario: Collaborative Filtering



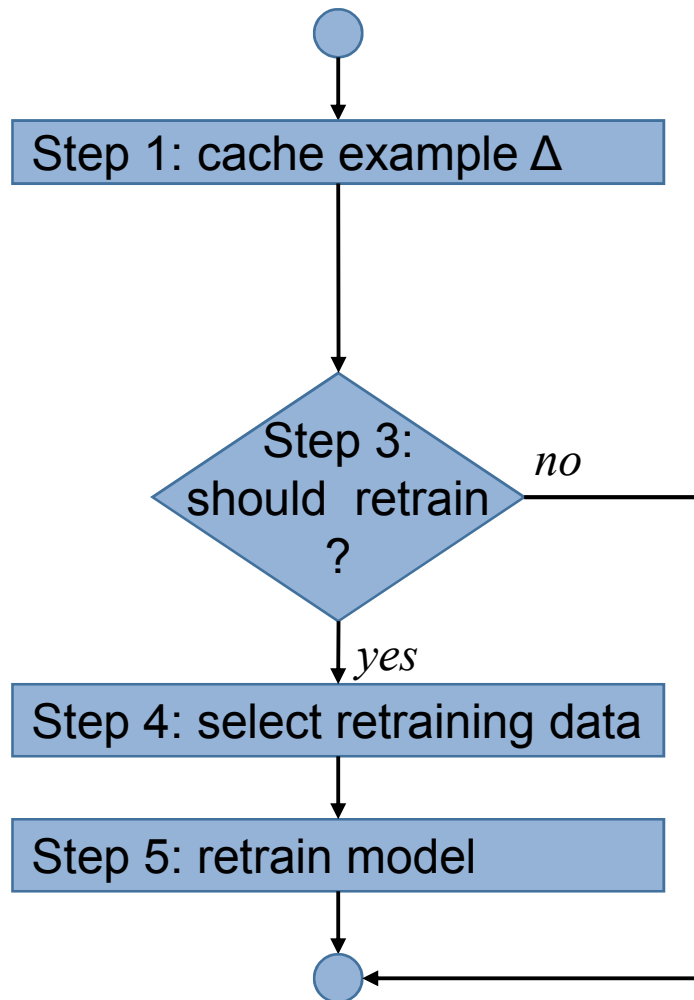
META: Middleware for Events, Transactions, and Analytics [IBMRD'16]



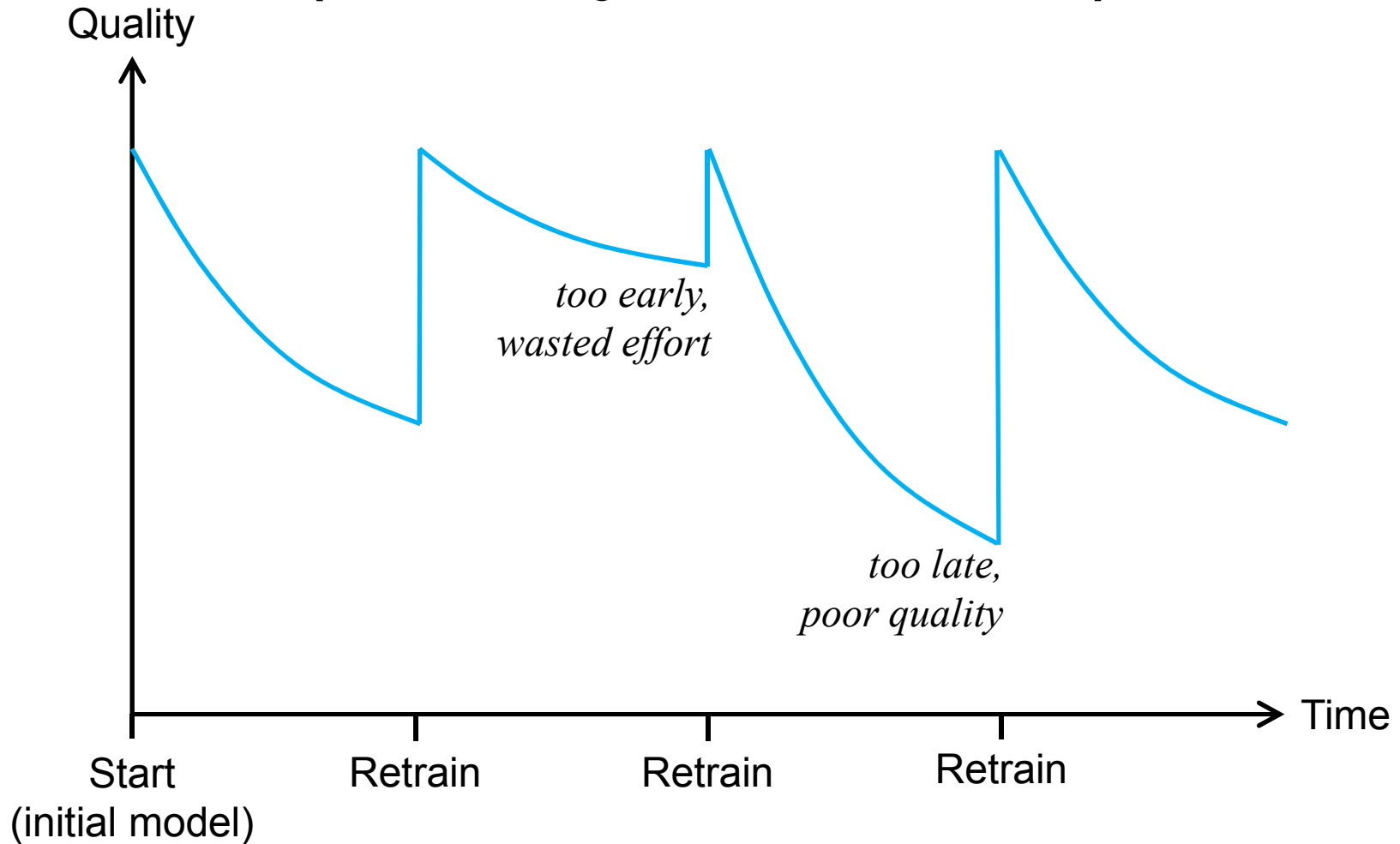
Model Drift



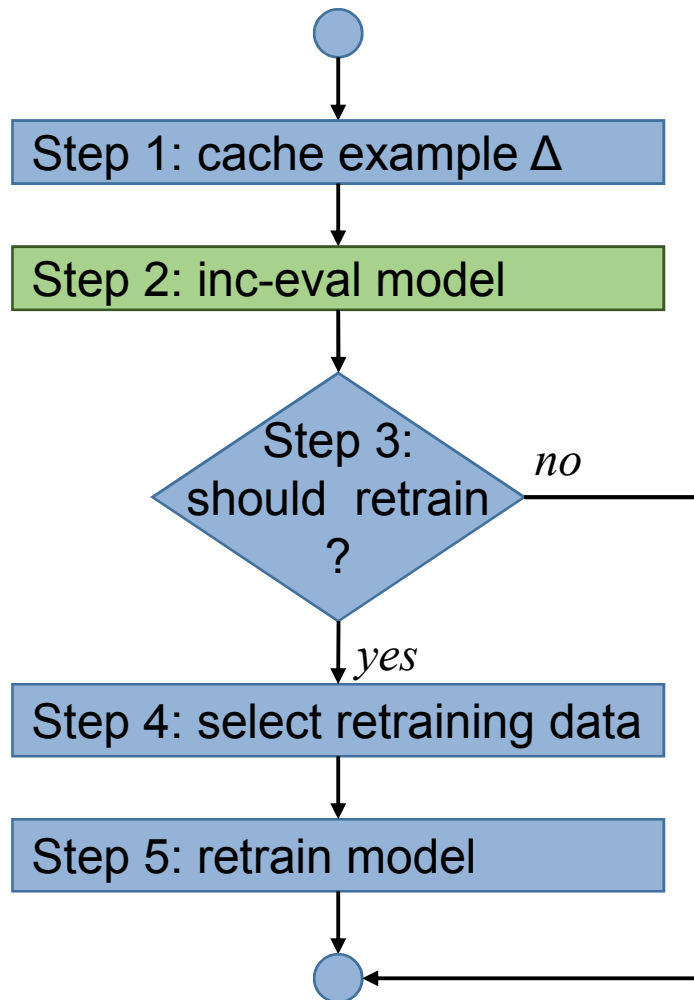
Retraining Workflow (Attempt 1)



◆ Fixed-Size Retraining Strategy (Quality-Oblivious)



Retraining Workflow (Attempt 2)



Training data D

Test data T

Evaluate error (inverse of quality):

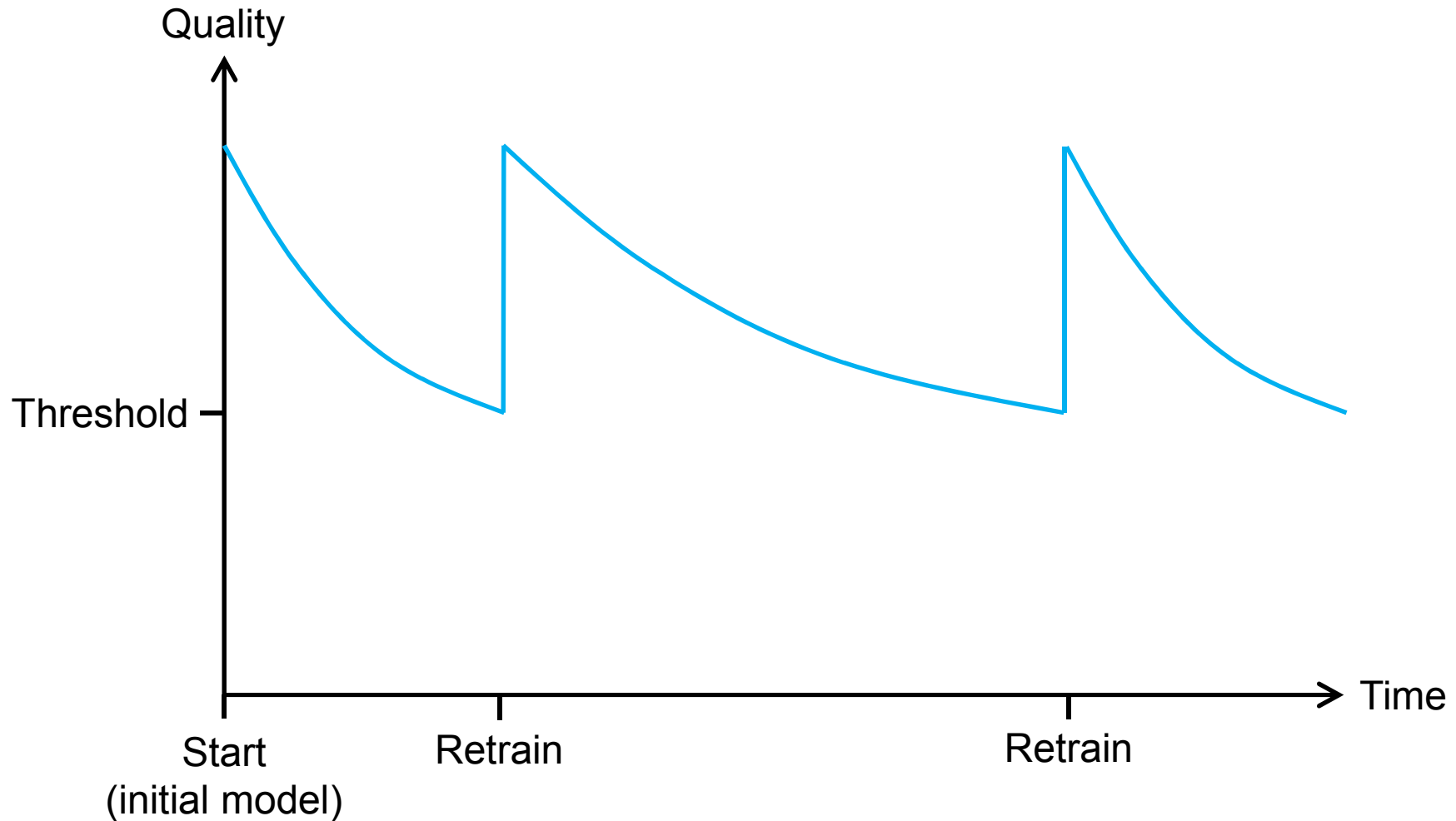
$$evaluate(Model(D), T) = \sqrt{\frac{\sum_{t \in T} error(t)^2}{|T|}}$$

New test data Δ

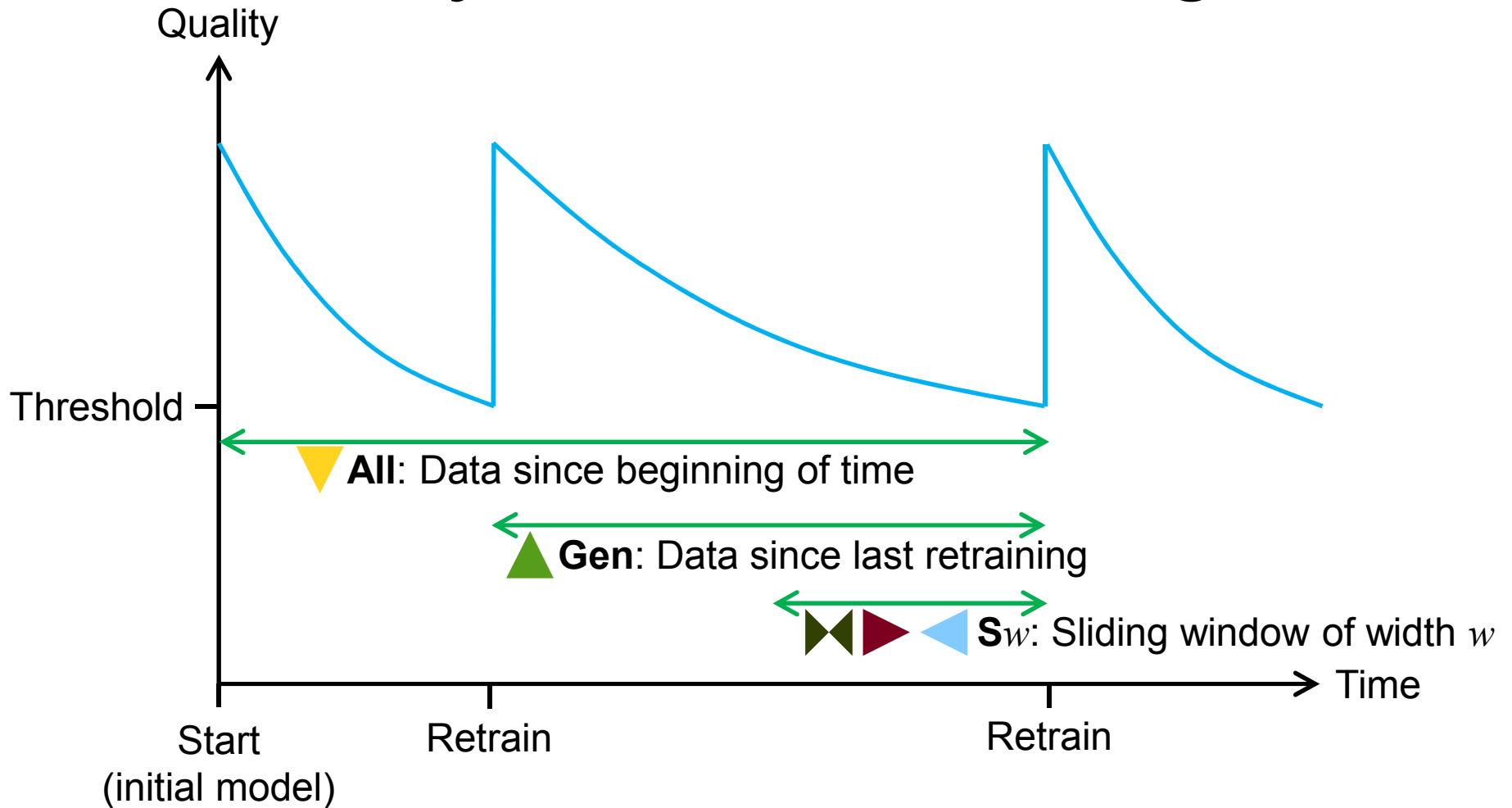
Incrementally update *sum* and *count*

$$evaluate(Model(D), T \pm \Delta) = \sqrt{sum/count}$$

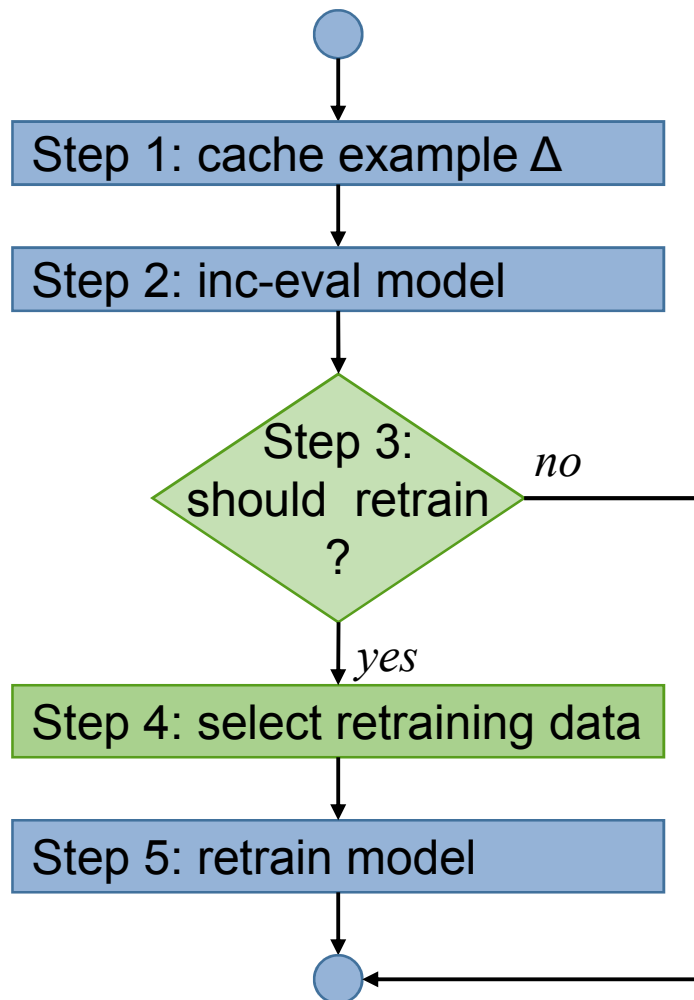
Retrain When Quality falls Below Threshold



Training Data Selection for Quality-Directed Strategies

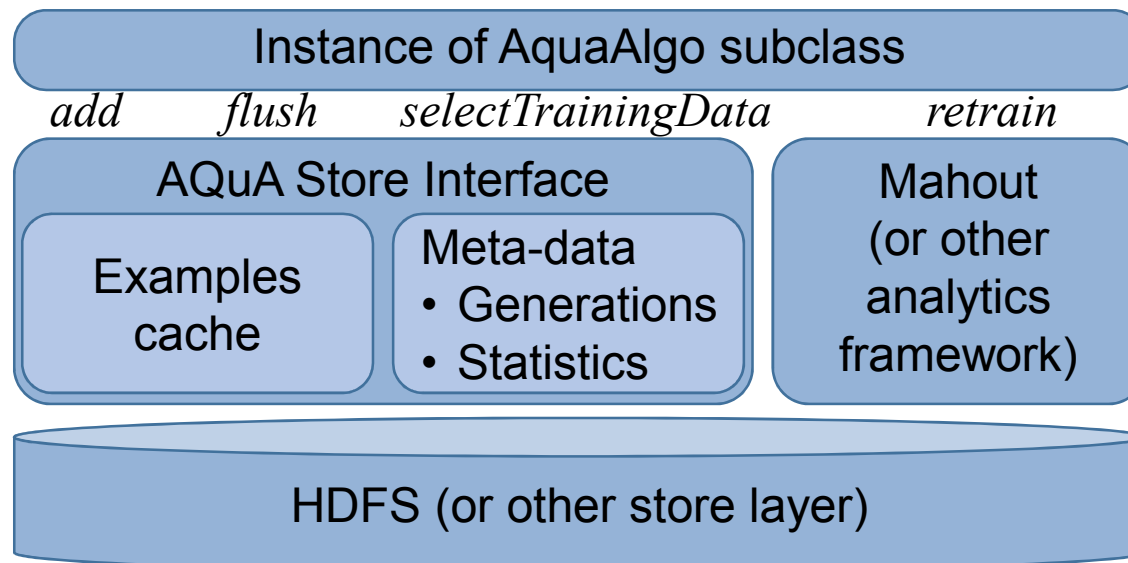


Retraining Strategies

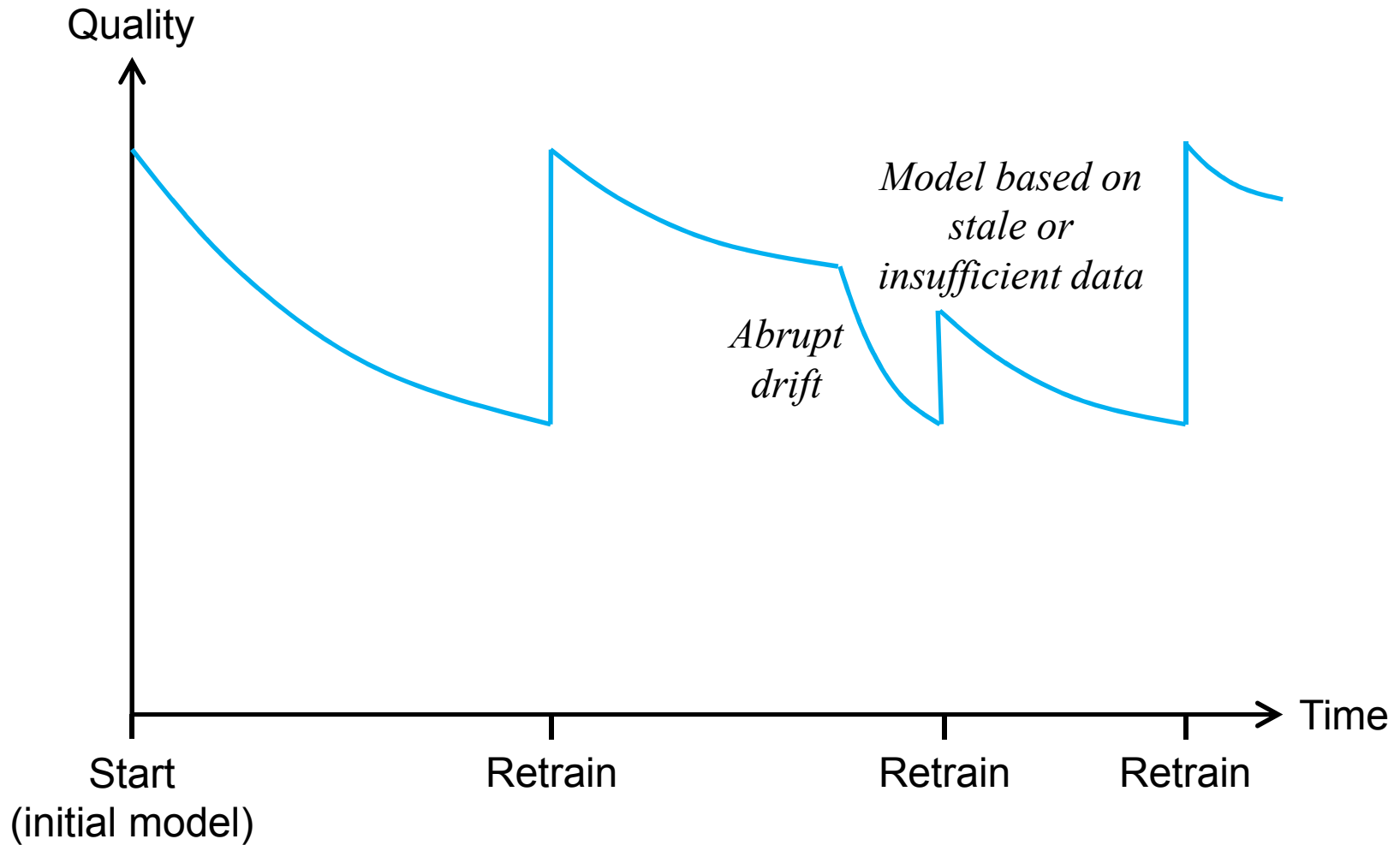


Quality-oblivious strategies	Quality-directed strategies
<ul style="list-style-type: none"> ■ Never: Never retrain model ◆ Fix: Fixed-size retrain interval 	<ul style="list-style-type: none"> ▼ All: Data since beginning of time ▲ Gen: Data since last retraining ◀ S_w: Sliding window of width w

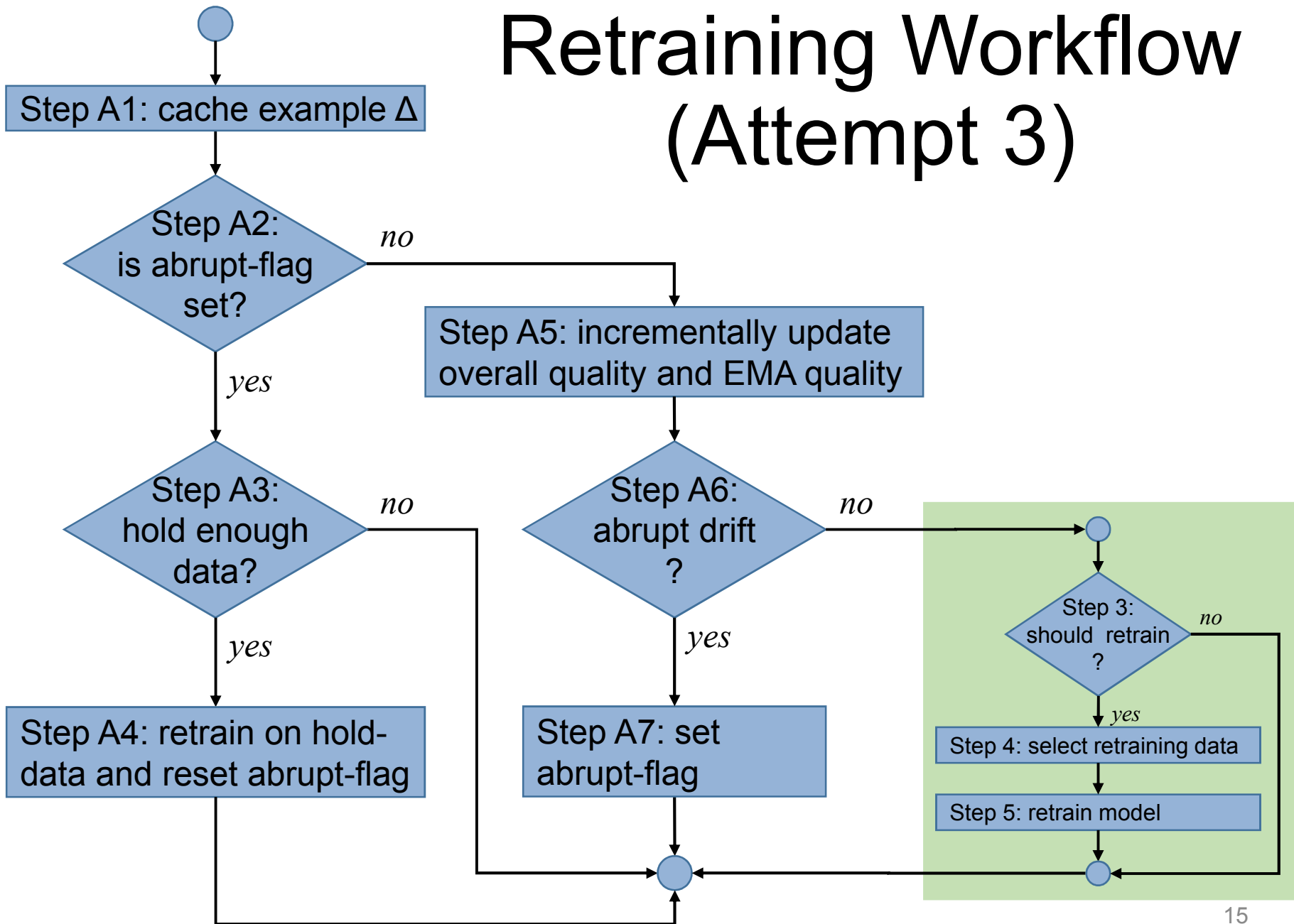
AQuA Store Interface



Gradual vs. Abrupt Model Drift

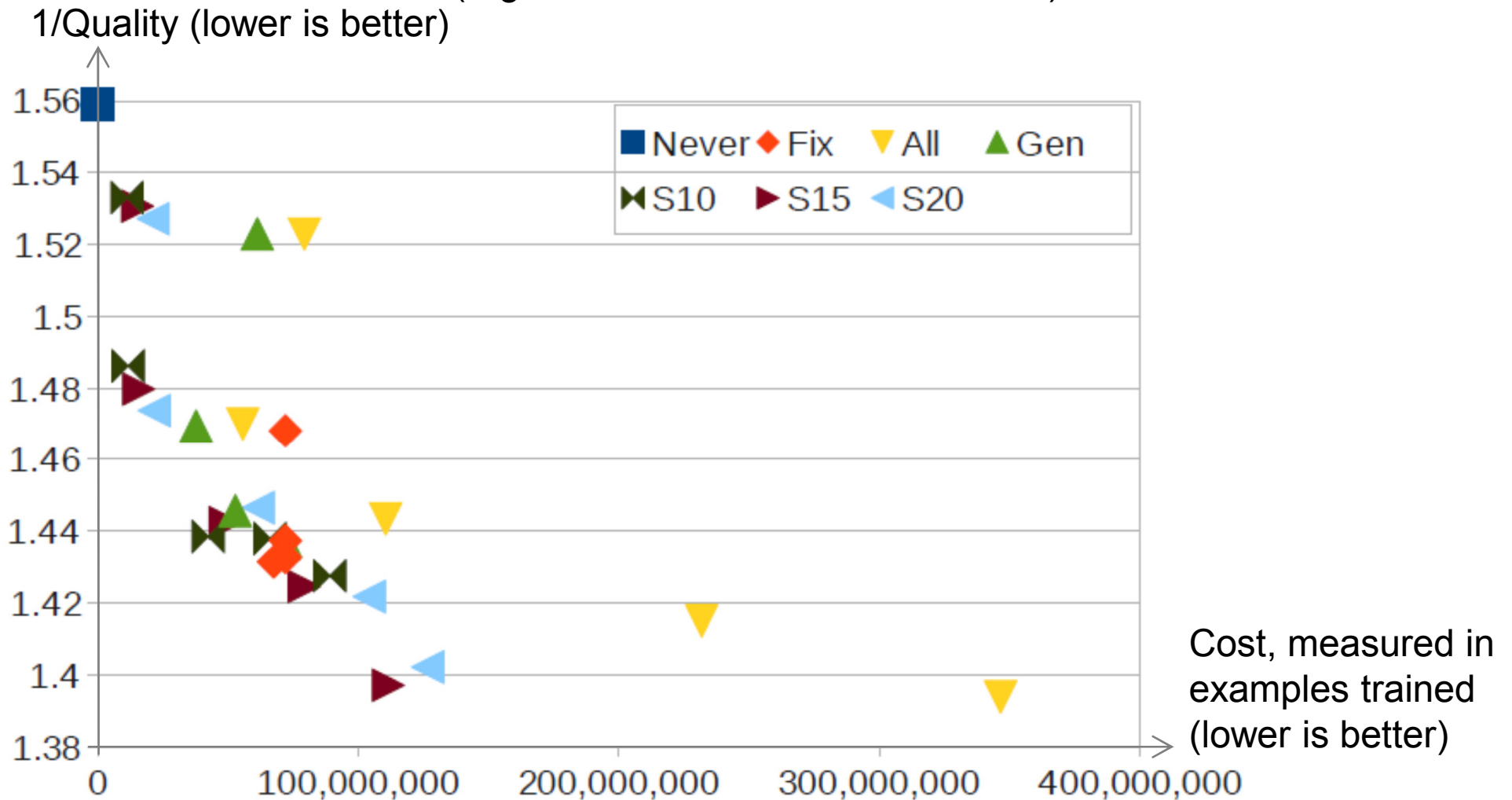


Retraining Workflow (Attempt 3)



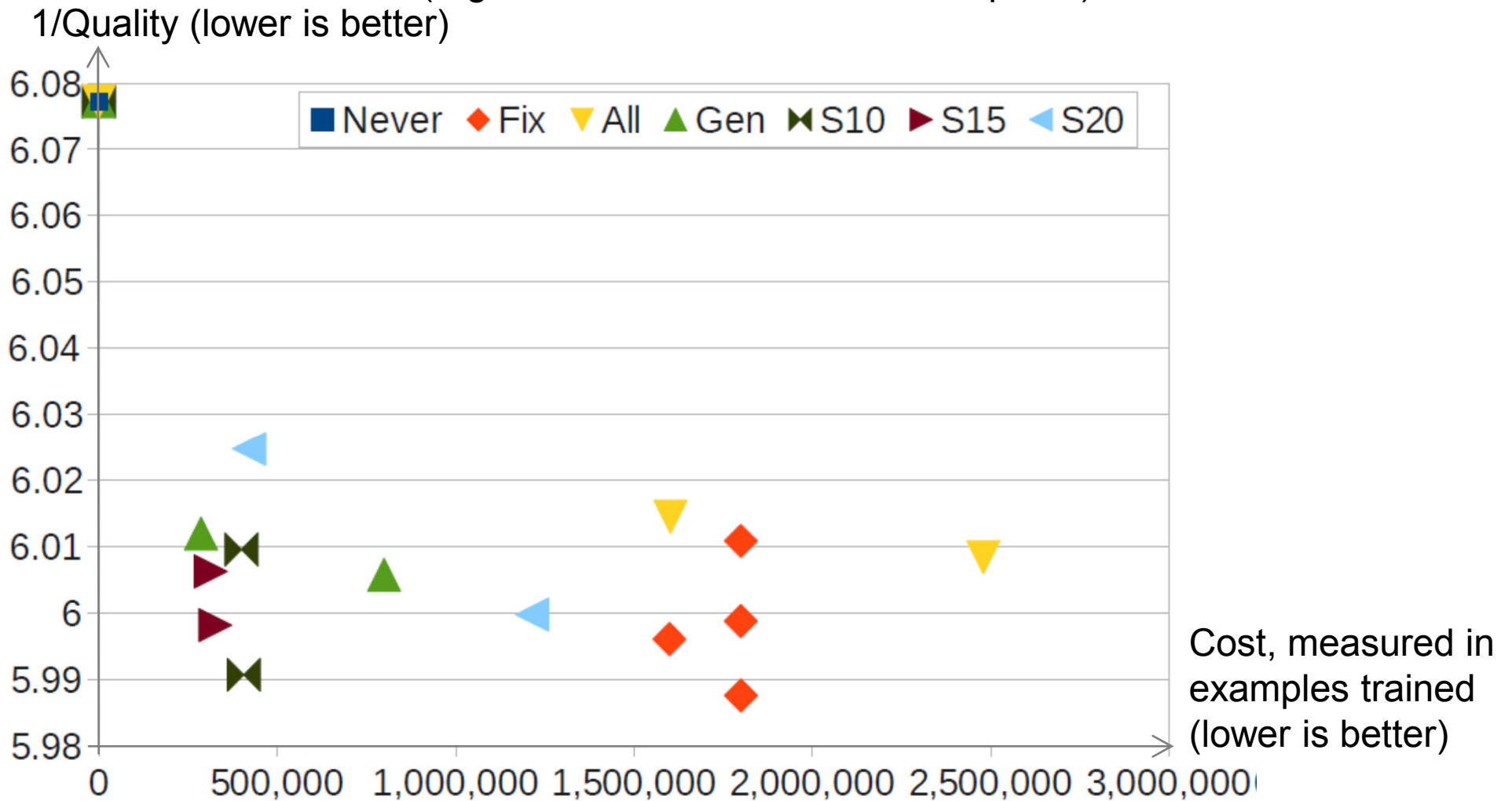
Performance: Collaborative Filtering

(Algorithm: ALSWR, dataset: Netflix)



Performance: Clustering

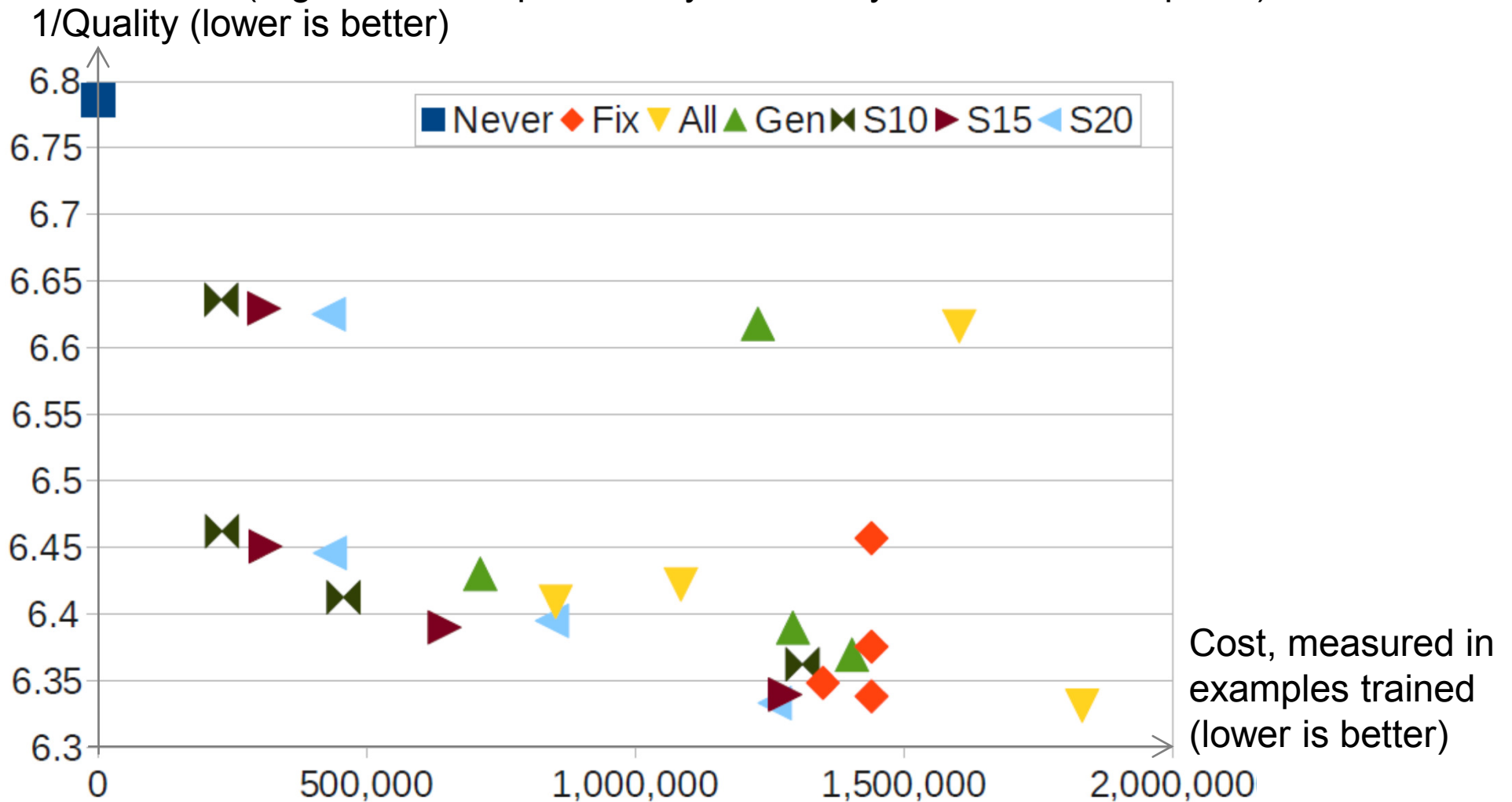
(Algorithm: KMeans, dataset: Wikipedia)



S_w (sliding window strategies) are on Pareto frontier

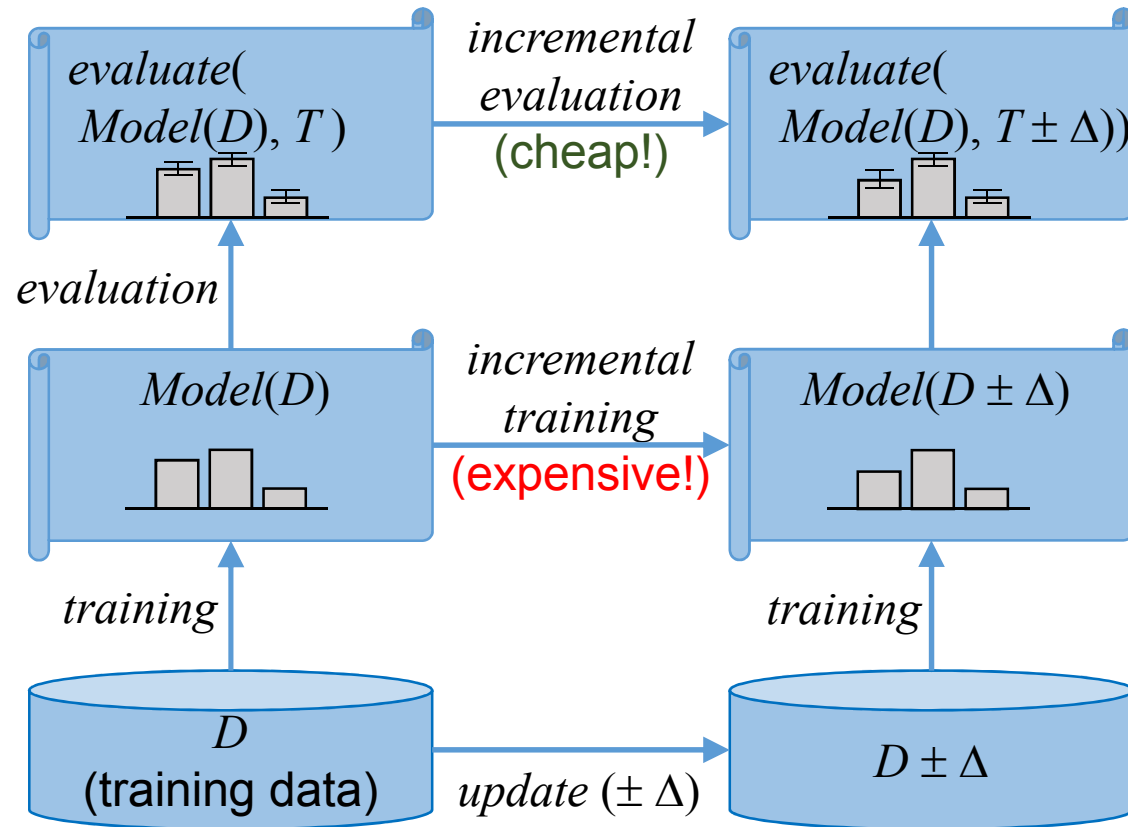
Performance: Classification

(Algorithm: Complementary Naïve Bayes, dataset: Wikipedia)

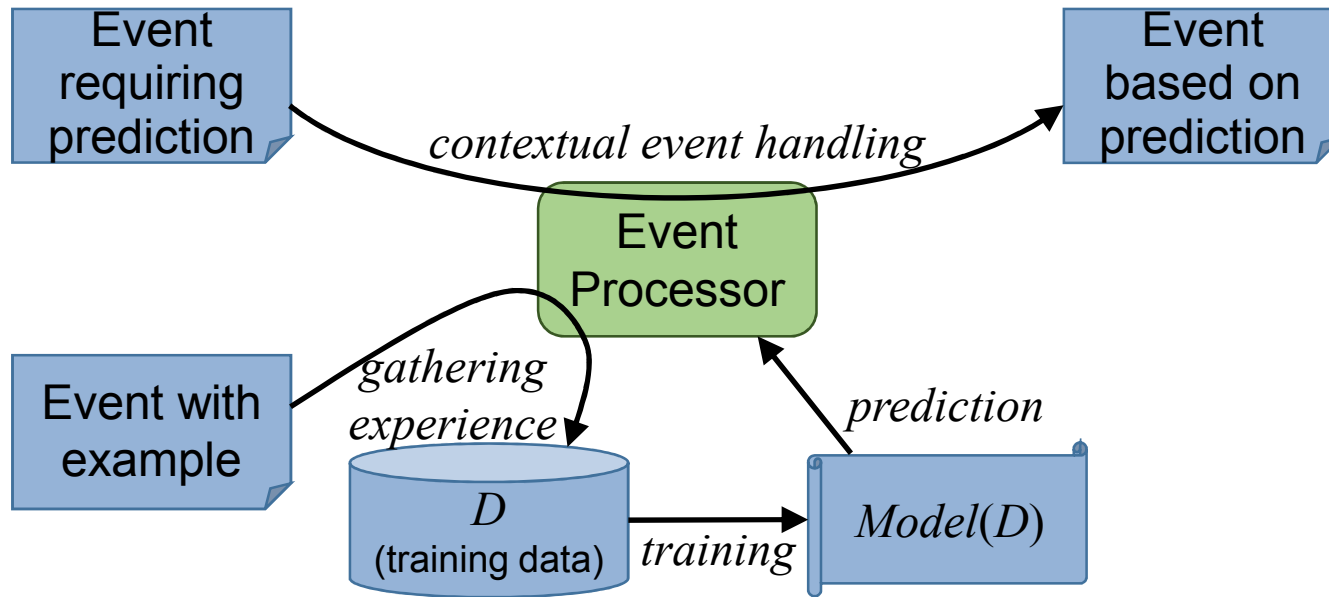


S_w (sliding window strategies) are on Pareto frontier

Related Work



Conclusions



- Incremental evaluation of model quality
- Quality-directed retraining
- Strategies for gradual and abrupt model drift
- Sliding window strategies are on Pareto frontier