# STREAM-VAE: Dual-Path Routing for Slow and Fast Dynamics in Vehicle Telemetry Anomaly Detection

Kadir-Kaan Özer*†, René Ebeling*, Markus Enzweiler†

*Abstract*—Automotive telemetry data exhibits slow drifts and fast spikes, often within the same sequence, making reliable anomaly detection challenging. Standard reconstruction-based methods, including sequence variational autoencoders (VAEs), use a single latent process and therefore mix heterogeneous time scales, which can smooth out spikes or inflate variances and weaken anomaly separation.

In this paper, we present STREAM-VAE, a variational autoencoder for anomaly detection in automotive telemetry time-series data. Our model uses a dual-path encoder to separate slow drift and fast spike signal dynamics, and a decoder that represents transient deviations separately from the normal operating pattern. STREAM-VAE is designed for deployment, producing stable anomaly scores across operating modes for both in-vehicle monitors and backend fleet analytics.

Experiments on an automotive telemetry dataset and the public SMD benchmark show that explicitly separating drift and spike dynamics improves robustness compared to strong forecasting, attention, graph, and VAE baselines.

*Index Terms*—automotive telemetry; anomaly detection; sensor data; generative models; variational autoencoder; intelligent vehicles

## I. INTRODUCTION

Modern intelligent vehicles stream high-rate telemetry signals from powertrain, chassis, ECUs, and body controllers. Detecting abnormal patterns in these signals is important for identifying emerging faults early and ensuring reliable vehicle operation, both in on-board monitoring and in backend fleet analytics. Yet several properties make this task difficult. The data are high dimensional and contain strong cross-sensor couplings. Vehicles operate in nonstationary modes that depend on traffic, environment, and driver behavior. Abnormal behavior may appear as brief transients caused by driver input or electrical disturbances, or as slow-evolving drifts driven by load or temperature changes.

Variational autoencoders (VAEs) are widely used in reconstruction-based anomaly detection, where a model is trained on normal data and high reconstruction error indicates abnormal behavior. They also provide a probabilistic notion of expected behavior through a learned likelihood [1]. In practice, time-series VAEs encode each temporal segment into a single latent trajectory, requiring the model to capture both high-frequency variations and low-frequency drifts within a shared representation. This often leads to over-smoothed
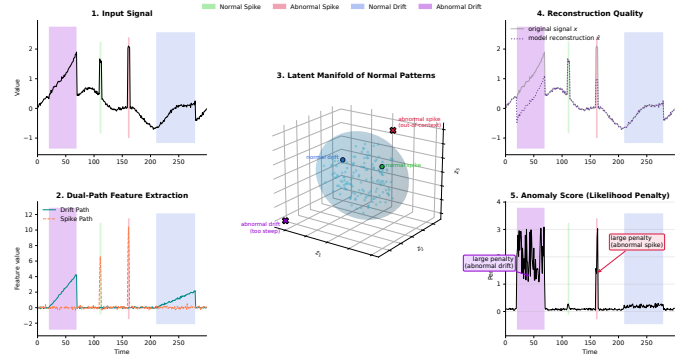


Fig. 1. Illustration of the core detection challenge: fast spikes and slow drifts distort different aspects of the signal. STREAM-VAE maps clean segments to a compact latent manifold, while anomalous patterns should fall outside this region.

spikes or inflated variances that weaken the separation between nominal and anomalous scores (see Fig. 1). For deployment in intelligent vehicles, this separation matters both for lightweight in-vehicle monitors and for high-volume backend fleet analytics, where stable thresholds and consistent behavior across operating modes are critical.

Prior work has addressed parts of this challenge. Smoothing regularizers stabilize reconstructions but can attenuate short events [2]. Attention-based encoders capture long-range structure [3]–[5], although they do not enforce any separation between fast and slow components and therefore mix heterogeneous time scales in a single latent representation. Automotive-oriented designs show how to prevent attention from bypassing the latent bottleneck [6], but they do not provide a mechanism for handling transient spikes separately from slow drifts. Posterior collapse remains a risk with expressive decoders, motivating feedback-based KL control [7]. Outside the VAE family, strong baselines such as TFT-Residual [8], Anomaly Transformer [9], and GDN [10] perform well on specific data regimes. However, they typically rely on residual thresholds or per-series tuning, and they do not directly address the need for consistent modeling of mixed time-scale behavior. In summary, existing methods lack an explicit mechanism to separate fast and slow dynamics and to explain transient spikes without widening the nominal tail.

Our main contribution in this paper is STREAM-VAE, which stands for **S**pike **T**rend **R**outing with **E**vent Residual

*Mercedes-Benz AG, Germany.
†Institute for Intelligent Systems, Esslingen University of Applied Sciences, Germany.

Attention and Mixture of Experts **VAE**. The model introduces a dual-path architecture that separates slow drift and fast spike feature dynamics in the encoder through two individual attention paths [11]. The decoder combines a per-feature mixture of experts (MoE) with an event-residual block that pairs a residual connection [12] with soft-thresholding [13] so that transient deviations can be represented without widening the nominal likelihood tail. These design choices directly address the two limitations outlined above: the absence of explicit slow-fast separation and the difficulty of representing brief spikes without inflating nominal tails. The architecture is intended to support both in-vehicle monitoring, where models must remain compact and predictable, and backend fleet-level analysis, where stable tail behavior enables consistent thresholding across vehicles. See Fig. 2 for an overview.

## II. RELATED WORK

Work on time-series anomaly detection in vehicle telemetry spans three lines that motivate our choices, especially in the context of automotive sensor streams with mixed time scales and frequent mode changes.

First, classical and forecasting-residual approaches treat anomalies as departures from a predictive signal. Isolation Forest isolates points via random partitions and remains competitive when feature scales are well behaved [14]. Sequence forecasters such as the Temporal Fusion Transformer with a residual scoring head ("TFT-Residual") fit rich conditional dynamics and flag large residuals [8], but residuals can overreact to mode switches (e.g., gear changes, drive-cycle transitions between urban and highway segments, or sudden load steps), which are common in vehicle telemetry unless operating mode changes are modeled explicitly.

A second line of work emphasizes self-attention and structure-aware discrepancy. Anomaly Transformer scores each timestamp by how atypical its associations are relative to learned patterns, yielding sharp localization on periodic or quasi-periodic series [9]. Graph-based methods such as GDN encode inter-sensor relations and measure deviations from graph-conditioned expectations, which is effective when the sensor topology is stable [10], although automotive subsystems often exhibit context-dependent couplings.

A third complementary line of research builds probabilistic generative models that return calibrated likelihoods. Omni-Anomaly uses a stochastic recurrent VAE with flow-based posteriors [15], [16]. Li et al. [2] add smoothing to stabilize reconstructions. MA-VAE integrates attention while protecting the latent bottleneck [6]. Models utilizing variational self-attention [4], Wasserstein distance–based latent anomaly scores [17], and sparse/structured VAEs such as VASP [18] further improve optimization and latent selectivity. These models provide calibrated scores, but a single latent process is typically forced to explain both driver-induced spikes and slow environmental drifts, which is problematic for in-vehicle monitoring where the two have distinct diagnostic meaning. For brevity, we will refer to [4] as VS-VAE, [17] as W-VAE (Wasserstein-similarity scoring), and [2] as SIS-VAE.

Across these model families, progress has improved context modeling, graph structure, expressivity, and likelihood calibration. Yet two frictions remain when deploying to intelligent vehicles: (i) a single latent path is often required to explain both spikes and drifts, and (ii) many systems depend on per-series thresholds, which complicates fleet-wide calibration.

STREAM-VAE maintains a probabilistic likelihood but separates fast and slow content in a way that aligns with both actuator-level dynamics and operating mode shifts. Additionally, it decodes with a per-feature mixture of experts plus a soft-thresholded event path so that transients do not broaden the nominal tail.

## III. MODEL ARCHITECTURE

Modern vehicle telemetry exhibits two types of changes that matter for detection: very short, high-amplitude *spikes* (e.g., pedal jabs or brief current surges) and slower *drifts* (e.g., load or temperature trends and operating mode switches). If a single latent stream is forced to explain both, spikes are often over-smoothed or the decoder inflates variance to cover drifts, in both cases narrowing the gap between nominal and anomalous scores. STREAM-VAE addresses this by (i) routing fast and slow content through two encoder paths and (ii) decoding with a per-feature mixture of experts that absorbs benign mode changes, plus a gated event residual that explains sparse transients *without* broadening the nominal likelihood tail. Fig. 2 shows the full architecture; we now describe each component following the left-to-right flow of the diagram.

**Input** $X$**.** We process standardized windows of length $T$ with $F$ features. The encoder produces a per-timestep latent mean and variance logit in $D$ dimensions. Logits are mapped to positive values by softplus and clipped for numerical stability. To keep the latent representation informative even when the decoder is expressive, we add a small learnable linear projection of the raw input into the posterior heads. This projection is initialized to zero and gated by a sigmoid so that it influences the posterior only if it proves useful during training.

**BI-LSTM Encoder.** A light two-layer BI-LSTM trunk [19] maps the input window to contextual representations for each timestep. These features feed the posterior heads, which produce the latent sequence $Z$. The same encoder features also feed the two attention branches described below.

**Latent** $Z$**.** From the latent sequence $Z$ we form the first difference $\Delta Z$ and a smoothed sequence $\mathrm{EMA}(Z)$. Here $\mathrm{EMA}(\cdot)$ denotes an exponential moving average applied along the time dimension with a learnable smoothing factor. The difference $\Delta Z$ later drives the event residual in the decoder, and $\mathrm{EMA}(Z)$ supplies the value stream for the slow branch so that gradual trends are modeled directly.

**Drift Features.** Let $H_{\mathrm{E}} = \mathrm{Enc}(X)$ denote the encoder features and $Z$ the corresponding latent sequence. To emphasize slow evolution, we apply a (separately learned) EMA to obtain a slowly varying baseline $H_{\mathrm{slow}} = \mathrm{EMA}(H_{\mathrm{E}})$ and then take a first difference along time to obtain $\Delta H_{\mathrm{slow}}$. Linear projections of $\Delta H_{\mathrm{slow}}$ form the queries and keys for the drift
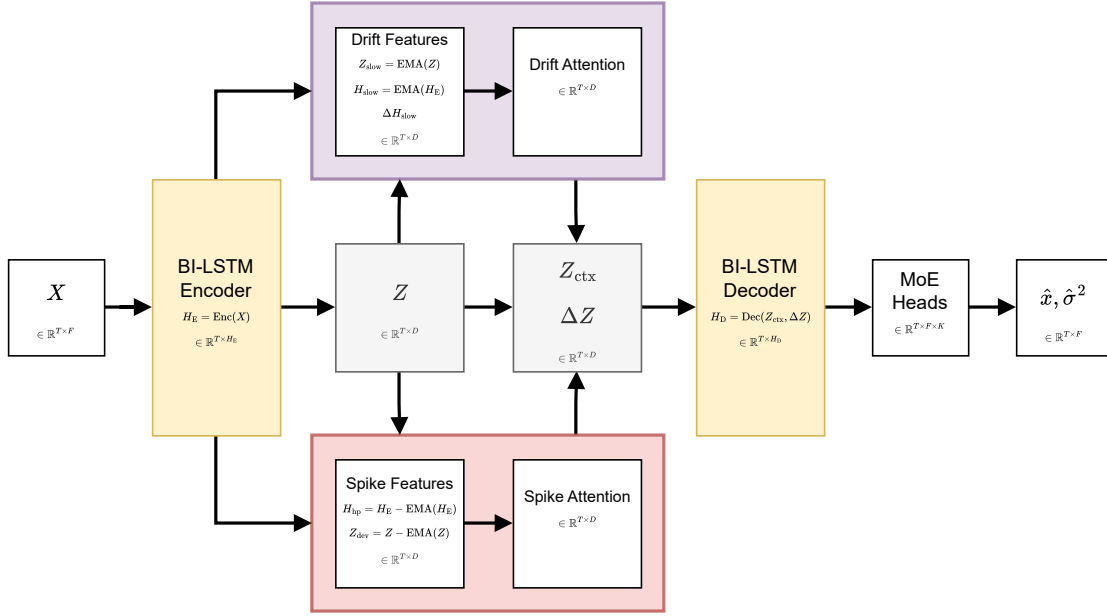
Fig. 2. Overview of the STREAM-VAE architecture with dual drift (top, purple) and spike (bottom, red) paths. Each block is annotated with the shape of its output tensor. The BI-LSTM Encoder (yellow) maps the input sequence $X \in \mathbb{R}^{T \times F}$ to encoder states $H_E \in \mathbb{R}^{T \times H_E}$ and latent states $Z \in \mathbb{R}^{T \times D}$. Slow EMA-based features define the drift path, while high-pass residual features define the spike path; both produce attention outputs in $\mathbb{R}^{T \times D}$. A gated fusion yields the latent context $Z_{\mathrm{ctx}} \in \mathbb{R}^{T \times D}$, and the first difference $\Delta Z \in \mathbb{R}^{T \times D}$ carries transient information. The BI-LSTM Decoder maps these to hidden states $H_D \in \mathbb{R}^{T \times H_D}$, which drive per-feature mixture-of-experts (MoE) heads that output Gaussian reconstruction parameters $(\hat{x}, \hat{\sigma}^2) \in \mathbb{R}^{T \times F}$ used for anomaly scoring.

attention branch, and the values come from the smoothed latent sequence $Z_{\mathrm{slow}} = \mathrm{EMA}(Z)$ so that this branch follows gradual trends rather than noise.

**Spike Features.** To expose brief and localized deviations, we compute a high-pass residual $H_{\mathrm{hp}} = H_E - \mathrm{EMA}(H_E)$ on the encoder features using a learnable EMA baseline. Linear projections of $H_{\mathrm{hp}}$ provide the queries and keys for the spike attention branch, while the values use the corresponding deviation in latent space $Z_{\mathrm{dev}} = Z - \mathrm{EMA}(Z)$.

**Drift Attention.** Multi-head attention uses queries, keys, and values from the Drift Features. This produces a $D$-dimensional representation that focuses on slow and persistent context. Queries and keys are $\ell_2$-normalized so that attention matching is scale-invariant. Grouped or multi-query attention (GQA) is used to reduce key–value memory cost without changing model behavior [20].

**Spike Attention.** A parallel branch uses queries, keys, and values from the Spike Features to target brief and localized transients. Fig. 3 illustrates the distinct attention patterns learned by the drift and spike attention pathways.

**Decoder Input ($Z_{\mathrm{ctx}}$, $\Delta Z$).** The decoder uses two latent signals. A sigmoid gate blends the drift and spike outputs at each timestep, and a light residual block refines the mixture to form the final encoder context $Z_{\mathrm{ctx}}$, which is fed to the BI-LSTM decoder. In parallel, we take the first difference of the latent sequence to obtain $\Delta Z$, which is used for the event-residual.

**BI-LSTM Decoder.** The decoder is a two-layer BI-LSTM that maps $Z_{\mathrm{ctx}}$ to hidden states used to parameterize the reconstruction mean and variance. Brief transients are modeled additively in the mean through an event residual driven by the first difference $\Delta Z$. A linear map produces a per-feature residual, which is soft-thresholded with a per-feature shrinkage parameter so that small activations are suppressed. A global gate and per-feature gains control the magnitude of this residual, and an RMS-based scaling keeps it commensurate with the MoE base mean. The variance head is a single linear layer followed by a softplus and clipping, and is shared across experts so that variance does not absorb transient spikes.

**MoE Heads.** For each timestep and feature, the decoder uses its current hidden state to produce non-negative softmax weights over the $K$ experts. From that same hidden state, it also generates $K$ expert means (via a low-rank factorization), and forms a base mean by mixing them with the softmax weights. This routing lets the model follow different nominal operating modes by shifting weight across experts instead of broadening the likelihood tails, consistent with classic mixture-of-experts ideas [21], [22]. In practice, only a few experts (e.g., $K=2$–$4$) are sufficient.

**Output ($\hat{x}$, $\hat{\sigma}^2$).** The final decoder mean $\hat{x}$ is the sum of the base mean from the MoE heads and the gated, soft-thresholded event residual. The variance $\hat{\sigma}^2$ is produced by a shared variance head. The anomaly score for each window is the Gaussian negative log-likelihood, where larger values indicate more anomalous behavior. Labels are aligned to the
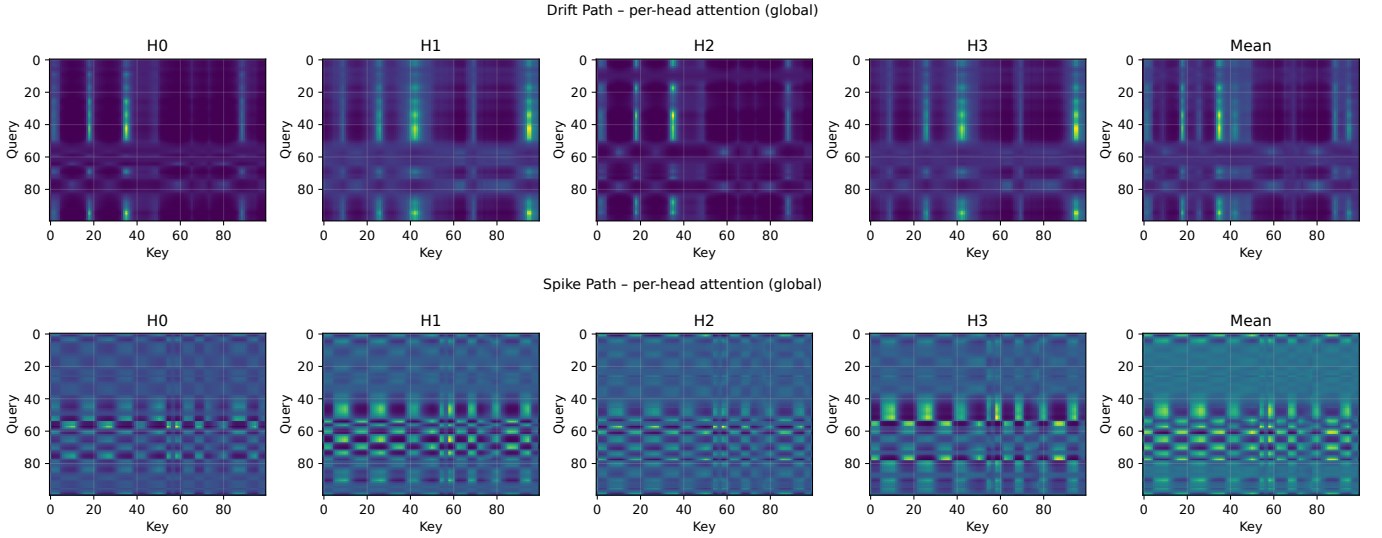
Fig. 3. Dual-path global attention (per head). **Drift path (top).** Multiple heads show bright vertical bands at consistent key indices, so many queries align to the same columns. This produces a few global reference anchors and a low-rank, anchor-based alignment that reflects slow and persistent context. **Spike path (bottom).** Heads display dense, oscillatory checkerboard patterns and narrow query-centric stripes, which indicate short-lag, high-frequency coupling typical of localized transients. The mean maps remain structured in both rows, which shows that heads specialize consistently and that the two paths split cleanly: The drift path captures global, low-frequency drift, and the spike path captures localized, high-frequency spikes.

end of each window so that scores and labels refer to the same temporal point.

**Training Objective and KL Control.** Training minimizes a Gaussian reconstruction term, a KL term with a feedback-controlled coefficient $\beta$ to keep the latent informative, and two light regularizers: an $\ell_1$ penalty on the event residual to encourage sparsity and a weak entropy target on the MoE gates to prevent expert collapse without forcing uniformity:

$$
\begin{aligned}
\mathcal{L} = & -\log p_\theta(X \mid Z) \, + \, \beta \, D_{\mathrm{KL}}\big(q_\phi(Z \mid X) \, \| \, \mathcal{N}(0, I)\big) \\
& + \, \lambda \, \|r\|_1 \, + \, \eta \, \big(H^\star - H\big)_+ .
\end{aligned}
\tag{1}
$$

Here $H$ is the average entropy of the MoE gates and $H^\star$ is a mid-level target equal to half of $\log K$. We use $(x)_+ = \max(x, 0)$ to denote the positive part. The coefficient $\beta$ is updated online with a proportional controller and EMA smoothing, following controllable VAE methods [7]. Optimization uses Adam with gradient-norm clipping [23].

**Scoring and Thresholding.** Thresholds are calibrated *once* per time series on *normal* training windows via Peaks-Over-Threshold (POT) [24], which models exceedances with a Generalized Pareto Distribution (GPD). To make this reliable during model selection (see Sec. IV), we penalize heavy nominal tails using a normalized upper-quantile width and a GPD-shape penalty evaluated on validation normals, using robust scale estimators to avoid outlier bias [25].

**Implementation Notes.** EMA memory factors for the drift and spike baselines are learned and initialized empirically at 0.9, a conventional choice that provides a smooth yet responsive starting point for separating slow and fast components, e.g., [23]. All smoothing and differencing operations are applied to the specific signal used in each branch,

following the routing shown in Fig. 2. Attention heads are split evenly across drift and spike branches and can optionally use GQA [20]. MoE expert means are implemented with a low-rank parameterization, and the event-residual thresholds are softplus-parameterized and initialized to small values. The shared variance head uses clipping for numerical stability. A feedback controller maintains the KL term near its target value and removes the need for manual schedule tuning [7].

## IV. EXPERIMENTAL SETUP

We evaluate STREAM-VAE and strong baselines on a proprietary automotive telemetry dataset and the public Server Machine Dataset (SMD) [15]. Each dataset provides a train/test split per entity, and anomalies are labeled.

**Baselines.** We compare STREAM-VAE against a diverse set of strong baselines that are representative of current practice in time-series anomaly detection. The set includes sparse and structured VAEs (VASP [18], VS-VAE [4], W-VAE [17], SIS-VAE [2], MA-VAE [6]), a stochastic recurrent VAE with flow-based posteriors (OmniAnomaly [15]), a graph-based detector (GDN [10]), an attention-based forecaster with residual scoring (TFT-Residual [8]), a structure-aware attention model (Anomaly Transformer [9]), and a classical tree ensemble (Isolation Forest [14]).

**Automotive Dataset.** Our automotive dataset consists of clean in-vehicle test-drive telemetry augmented with expert-designed synthetic anomalies that match field distributions and remain physically consistent. The dataset has been acquired in different locations and scenarios from a test vehicle fleet over a period of several days. A controlled injection simulator introduces short, contiguous perturbations affecting one or more correlated features within their valid operating

TABLE I
DATASET AND FEATURE STATISTICS FOR THE AUTOMOTIVE DATASET

**Dataset Summary**

| | |
|---|---|
| Total Records | 40,000 |
| Features (see below) | 8 |
| Normal Instances | 36,426 |
| Anomalous Instances | 3,574 |
| Anomaly Rate (%) | 8.94 |
| Completeness (%) | 100.00 |
| Outlier Percentage (%) | 14.46 |
| Range Coverage (%) | 100.00 |

**Feature-Level Statistics**

| Feature | Mean | Std | Skewness |
|---|---|---|---|
| Brake Torque | 189.530 | 658.052 | 5.13 |
| Battery Current | -1.774 | 20.977 | 2.25 |
| Accelerator Pedal Position (Raw) | 11.366 | 22.548 | 6.44 |
| Accelerator Pedal Position (OBD) | 13.338 | 12.631 | 2.86 |
| Steering Wheel Angular Speed | 0.034 | 37.294 | -0.54 |
| Vehicle Speed | 37.247 | 49.637 | 2.79 |
| Wheel RPM (Front Left) | 286.966 | 359.348 | 1.41 |
| Wheel RPM (Rear Left) | 288.820 | 366.703 | 1.44 |

ranges. Each anomaly follows one of six canonical fault types (*spikes*, *drifts*, *level shifts*, *variance jumps*, *flatlines*, *correlation breaks*), with magnitudes and durations sampled from type-specific calibrated ranges. During training, 30% of the nominal data are held out for validation. Corpus- and feature-level summaries are given in Table I.

**SMD (Server Machine Dataset).** SMD is a widely used benchmark in general-purpose time-series anomaly detection but is not commonly used in automotive research. We include it to evaluate generalization outside the vehicle domain and to compare STREAM-VAE against established baselines under a standardized, publicly reproducible protocol. SMD contains 28 entities (server nodes, referred to as *machines*) with multivariate metrics and labeled anomalies. Although the domain differs, it shares structural challenges with vehicle telemetry: mixed time scales, multi-sensor coupling, and infrequent transients. For SMD we report macro averages across all entities.

**Preprocessing.** STREAM-VAE processes standardized windows of length $T=100$ with $F$ features ($F=8$ for the automotive dataset; SMD varies per entity). The encoder outputs per-timestep latent means and variance logits in $D=64$ dimensions; these values define the default configuration across all experiments unless otherwise noted.

All methods share identical segmentation and normalization. Each feature is z-scored using statistics from non-anomalous (nominal) training data only. Training uses 90% window overlap (stride 10), and evaluation uses 100% window overlap (stride 1). Labels are assigned to the window end to align scores and labels.

**Thresholding.** Thresholds are calibrated *once* per entity using Peaks-Over-Threshold (POT) on nominal train scores and then held fixed. For the automotive dataset we use $(q,p)=(10^{-3},10^{-4})$. For SMD we use group-specific $q$ values (machine-1: $10^{-3}$; machine-2: $2.5\times10^{-3}$; machine-3: $5\times10^{-3}$; all with $p=10^{-4}$). If exceedances are too sparse,

we fall back to the empirical $(1-\alpha)$ upper quantile without labels [24], [26].

**Scoring.** All models are mapped to a scalar anomaly score per window to make them comparable. VAE-style models that output $(\boldsymbol{\mu}, \boldsymbol{\sigma}^2)$ use the Gaussian negative log-likelihood; all others use mean squared reconstruction or forecast error, oriented so that larger scores indicate more anomalous behavior.

**Evaluation Metrics.** We evaluate all models using four standard anomaly-detection metrics. AUC-PR and AUC-ROC provide threshold-free assessments of ranking quality. Point-Adjusted F1 (PA-F1) is computed at the fixed POT threshold and treats any detection within an anomaly window as a correct hit [27]. We also report Oracle PA-F1, which selects the best achievable PA-F1 over all possible thresholds and serves as a diagnostic upper bound [15].

**Hyperparameter Optimization.** For hyperparameter optimization, we allocate 20 Optuna [28] trials (Tree-structured Parzen Estimator + MedianPruner) per model. The selected configuration is retrained with early stopping (up to 30 epochs; batch size 50), and results are reported as mean±standard deviation over 5 seeds. To align model selection with stable POT calibration, we minimize a label-free objective on validation normals:

$$J = \tilde{L}_{\text{val}} + \lambda \, \text{NUQ}_q + \gamma \, \xi_+, \tag{2}$$

where $\tilde{L}_{\text{val}}$ is the per-dimension normalized validation loss, $\text{NUQ}_q$ is a robust upper-quantile width ($q=0.995$), and $\xi_+$ is the positive part of the GPD shape parameter estimated from exceedances. Minimizing $J$ discourages heavy nominal score tails that destabilize POT and removes the need for manual KL tuning. We fix $(\lambda,\gamma)=(0.01,0.05)$ for all datasets. For SMD we use a two-phase protocol: hyperparameter search on machine-1-1, machine-2-1, machine-3-1, then retrain and evaluate on all 28 entities with per-entity thresholds as above.

## V. RESULTS

Our experimental results are given in Table II. On the automotive dataset, STREAM-VAE achieves the highest Oracle PA-F1 (0.857), PA-F1 (0.794), AUC-PR (0.532), and a competitive AUC-ROC (0.755), indicating both strong score ordering and stable detection at a fixed threshold. On SMD, where quasi-periodic structure benefits association-based models, Anomaly Transformer obtains the best AUC-PR (0.462) and AUC-ROC (0.881), while STREAM-VAE remains competitive in the threshold-free metrics and achieves the highest Oracle PA-F1 (0.935), suggesting that its scores remain well ordered even when a single global threshold is not perfectly tuned.

On the automotive data, the performance gain of STREAM-VAE over VAE baselines can be traced to three effects. First, separating fast spikes and slow drifts in the encoder prevents sharp events from being oversmoothed by slow context. Second, per-feature mixture of experts in the decoder allows benign operating mode changes to be absorbed by expert reweighting instead of by inflating variance. Third, the soft-thresholded event residual explains sparse transients additively

TABLE II

PERFORMANCE OF THE PROPOSED STREAM-VAE MODEL AND BASELINES ON THE AUTOMOTIVE AND SMD DATASETS. MEAN ± STANDARD DEVIATION OVER 5 SEEDS. **BOLD** INDICATES THE BEST-PERFORMING METHOD FOR EACH DATASET.

| Model | Automotive | | | | SMD | | | |
|---|---|---|---|---|---|---|---|---|
| | Oracle PA-F1 | PA-F1 | AUC-PR | AUC-ROC | Oracle PA-F1 | PA-F1 | AUC-PR | AUC-ROC |
| Isolation Forest [14] | 0.568 ± 0.008 | 0.512 ± 0.021 | 0.133 ± 0.003 | 0.610 ± 0.002 | 0.861 ± 0.143 | **0.552 ± 0.339** | 0.280 ± 0.220 | 0.764 ± 0.121 |
| VASP [18] | 0.641 ± 0.007 | 0.445 ± 0.000 | 0.120 ± 0.001 | 0.599 ± 0.004 | 0.837 ± 0.180 | 0.367 ± 0.348 | 0.313 ± 0.237 | 0.737 ± 0.133 |
| Anomaly Transformer [9] | 0.698 ± 0.022 | 0.380 ± 0.100 | 0.132 ± 0.014 | 0.602 ± 0.024 | 0.871 ± 0.137 | 0.399 ± 0.319 | **0.462 ± 0.272** | **0.881 ± 0.107** |
| W-VAE [17] | 0.701 ± 0.017 | 0.386 ± 0.091 | 0.197 ± 0.008 | 0.633 ± 0.012 | 0.890 ± 0.119 | 0.427 ± 0.362 | 0.399 ± 0.263 | 0.800 ± 0.132 |
| VS-VAE [4] | 0.735 ± 0.016 | 0.466 ± 0.021 | 0.258 ± 0.024 | 0.704 ± 0.012 | 0.845 ± 0.149 | 0.471 ± 0.364 | 0.336 ± 0.268 | 0.773 ± 0.153 |
| OmniAnomaly [15] | 0.756 ± 0.009 | 0.663 ± 0.027 | 0.346 ± 0.008 | 0.686 ± 0.010 | 0.843 ± 0.193 | 0.400 ± 0.375 | 0.332 ± 0.270 | 0.744 ± 0.173 |
| MA-VAE [6] | 0.808 ± 0.037 | 0.723 ± 0.097 | 0.451 ± 0.027 | 0.743 ± 0.015 | 0.810 ± 0.176 | 0.423 ± 0.336 | 0.338 ± 0.216 | 0.790 ± 0.131 |
| SIS-VAE [2] | 0.824 ± 0.014 | 0.701 ± 0.047 | 0.380 ± 0.030 | 0.753 ± 0.017 | 0.717 ± 0.277 | 0.316 ± 0.310 | 0.226 ± 0.212 | 0.658 ± 0.170 |
| GDN [10] | 0.825 ± 0.009 | 0.760 ± 0.007 | 0.498 ± 0.011 | 0.744 ± 0.003 | 0.930 ± 0.088 | 0.496 ± 0.359 | 0.451 ± 0.263 | 0.833 ± 0.130 |
| TFT-Residual [8] | 0.830 ± 0.004 | 0.781 ± 0.022 | 0.479 ± 0.019 | 0.750 ± 0.014 | 0.906 ± 0.114 | 0.442 ± 0.356 | 0.438 ± 0.253 | 0.831 ± 0.111 |
| STREAM-VAE (ours) | **0.857 ± 0.024** | **0.794 ± 0.026** | **0.532 ± 0.030** | **0.755 ± 0.027** | **0.935 ± 0.087** | 0.493 ± 0.374 | 0.430 ± 0.260 | 0.812 ± 0.132 |

in the mean, which sharpens point-adjusted F1 by reducing fragmented detections and keeps nominal tails tight enough for stable POT calibration.

The behavior of the baselines is consistent with their design goals. TFT-Residual performs well when anomalies are departures from locally predictable trends, but it is sensitive to unmodeled operating mode switches. GDN leverages inter-sensor structure and is strong on correlation anomalies, yet its performance can degrade when couplings change with context. SIS-VAE stabilizes reconstructions and helps on drifts but tends to attenuate sharp spikes. MA-VAE adds long-range context but does not explicitly separate time scales, so variance can grow to cover fast deviations. OmniAnomaly's expressive posterior improves ranking but may handle heterogeneity mainly through scale changes, which complicates operation at a fixed threshold. VS-VAE and W-VAE capture global structure yet lack explicit routing of operating modes or transients. VASP's sparsity promotes salient spikes but can under-represent slow trends. Anomaly Transformer excels on SMD, where periodic anchors repeat, but underperforms on the automotive data where associations depend strongly on driving context. Isolation Forest remains competitive at high recall for gross outliers but lacks multivariate temporal modeling, which limits its performance on structured automotive signals. Our ablation study in the following section confirms

this interpretation: performance degrades systematically when individual architectural components are removed, indicating that the improvements of STREAM-VAE indeed come from its architectural design.

Finally, we evaluate computational efficiency on the automotive dataset (Table III) for an in-vehicle application. All runtimes were measured on a single Apple MacBook Pro M3 Max and refer to processing the full dataset. For context, STREAM-VAE processes about 408 windows/s in Python (2.45 ms per window, $T$=100), which corresponds to real-time operation up to roughly 400 Hz. This is well above the 10–100 Hz update rates of typical automotive signals.

To reflect practical on-board in-vehicle requirements, the table reports recall at a fixed false-positive rate of 1%. This metric is not a replacement for the threshold-free results shown earlier in Table II, but an additional operating-point analysis chosen because in-vehicle monitors must trigger only rarely while still detecting meaningful anomalies. At this operating point, STREAM-VAE matches the top recall of TFT-Residual. Its training and inference times fall in the mid-range of learned detectors: not as heavy as SIS-VAE, OmniAnomaly, or GDN, but naturally slower than very lightweight methods such as Isolation Forest or VASP. Anomaly Transformer is efficient but performs poorly at low false-positive rates on this dataset.

Overall, STREAM-VAE combines state-of-the-art detection performance with computational cost comparable to other modern deep anomaly detectors. Its runtime is sufficiently low for real-time in-vehicle monitoring while remaining well-suited for large-scale offline fleet analysis.

TABLE III

AUTOMOTIVE DATASET EFFICIENCY METRICS. MEAN ± STD. OVER 5 SEEDS. TRAINING AND INFERENCE TIMES ARE GIVEN FOR THE WHOLE DATASET. BEST PER METRIC IN **BOLD**.

| Model | Recall@1% | Train Time (s) | Inference Time (s) |
|---|---|---|---|
| VASP [18] | 0.012 ± 0.000 | 42.6 ± 0.1 | 52.5 ± 0.5 |
| Anomaly Transformer [9] | 0.023 ± 0.011 | 69.7 ± 10.8 | 56.4 ± 1.5 |
| Isolation Forest [14] | 0.024 ± 0.010 | **0.7 ± 0.1** | **52.1 ± 4.0** |
| W-VAE [17] | 0.049 ± 0.005 | 79.5 ± 0.7 | 58.7 ± 0.5 |
| VS-VAE [4] | 0.050 ± 0.009 | 102.5 ± 2.5 | 62.6 ± 0.8 |
| OmniAnomaly [15] | 0.094 ± 0.003 | 605.5 ± 0.8 | 135.2 ± 0.5 |
| SIS-VAE [2] | 0.085 ± 0.004 | 1081.8 ± 3.3 | 130.8 ± 0.5 |
| MA-VAE [6] | 0.102 ± 0.004 | 320.3 ± 13.8 | 119.6 ± 4.2 |
| GDN [10] | 0.109 ± 0.001 | 456.0 ± 4.9 | 130.4 ± 1.1 |
| TFT-Residual [8] | **0.110 ± 0.000** | 169.1 ± 4.4 | 79.3 ± 0.9 |
| STREAM-VAE (ours) | **0.110 ± 0.000** | 262.3 ± 18.0 | 97.9 ± 5.4 |

## VI. ABLATION STUDIES

Table IV summarizes the effect of removing or simplifying individual components of STREAM-VAE. The full model leads in both threshold-free separation and Oracle Point-Adjusted detection. Removing the event residual lowers both Point-Adjusted F1 and AUC-PR while producing a small increase in AUC-ROC. This indicates that the residual improves precision-oriented ranking near operating thresholds while slightly smoothing the extremes of the score distribution that influence ROC. Disabling the per-feature MoE increases false

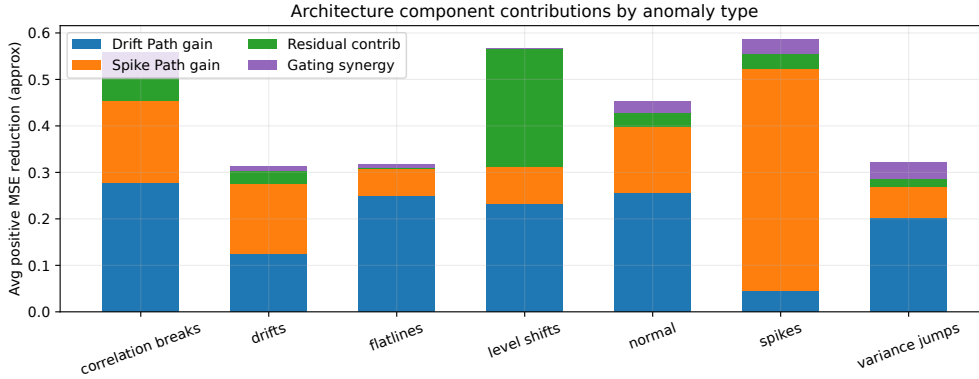| Variant | Oracle PA-F1 | PA-F1 | AUC-PR | AUC-ROC |
|---|---|---|---|---|
| STREAM-VAE (Full) | **0.792 ± 0.028** | **0.631 ± 0.018** | **0.309 ± 0.024** | 0.664 ± 0.011 |
| w/o MoE Decoder | 0.785 ± 0.037 | 0.619 ± 0.014 | 0.275 ± 0.034 | 0.658 ± 0.010 |
| Small Latent (32D) | 0.784 ± 0.019 | 0.623 ± 0.019 | 0.304 ± 0.024 | 0.665 ± 0.011 |
| Small MoE (k=4) | 0.783 ± 0.016 | 0.617 ± 0.008 | 0.301 ± 0.024 | 0.664 ± 0.012 |
| w/o Input Injection | 0.781 ± 0.039 | 0.612 ± 0.022 | 0.303 ± 0.036 | 0.661 ± 0.008 |
| Simple FFN (×1) | 0.781 ± 0.021 | 0.604 ± 0.018 | 0.300 ± 0.029 | 0.654 ± 0.006 |
| Concat Merge (No Gate) | 0.780 ± 0.034 | 0.614 ± 0.020 | 0.280 ± 0.023 | 0.656 ± 0.006 |
| Drift-Only Attention | 0.770 ± 0.034 | 0.611 ± 0.022 | 0.297 ± 0.023 | 0.656 ± 0.015 |
| w/o Event Residuals | 0.774 ± 0.032 | 0.614 ± 0.028 | 0.303 ± 0.021 | **0.666 ± 0.014** |
| Spike-Only Attention | 0.774 ± 0.016 | 0.622 ± 0.014 | 0.254 ± 0.012 | 0.642 ± 0.008 |
| w/o Input Injection + MoE | 0.770 ± 0.021 | 0.623 ± 0.016 | 0.289 ± 0.035 | 0.661 ± 0.007 |
| w/o Attention | 0.753 ± 0.018 | 0.601 ± 0.013 | 0.279 ± 0.022 | 0.658 ± 0.007 |



Fig. 4. Component-wise contribution to reconstruction error by anomaly type.

alarms during benign operating mode shifts, supporting its role in explaining away mode changes through expert routing. Collapsing the dual-path encoder to a single path or removing attention reduces score ordering quality across thresholds even if a single operating point can look acceptable; explicit separation of slow drift and fast spike information prevents variance inflation and drift-spike blurring. Ungated merges, either for injecting the residual or for combining the drift and spike pathways, tend to overfire and destabilize calibration, while the gated input injection keeps latents informative so that attention and MoE can operate effectively. Reduced-capacity variants, such as those with a smaller latent dimension, fewer experts, or a shallower decoder FFN, track the full model closely, which indicates that the gains come from the architectural design rather than parameter count.

To relate these behaviors to how the architecture reacts to specific anomaly types, we quantify component-wise contributions using a controlled decoder analysis that matches the metric in Fig. 4. The drift path and spike path correspond to the attention shown in Fig. 2, where the drift path exhibits stable global anchors and the spike path displays localized, high-frequency structure. For each window we encode once and freeze the latent sequence, then reconstruct the window with different component subsets, such as the drift path only, the spike path only, or no residual. We compare the recon-

struction error (mean squared error, MSE) of these restricted variants to that of the full model and record only positive increases in MSE. This measures how much each component helps the model reduce reconstruction error when present. Fig. 4 reports the average effects across anomaly categories. Spikes are mainly handled by the spike path, which matches the short-lag, high-frequency coupling seen in the spike-path attention maps. Drifts benefit from all components: the spike path and drift path provide most of the gain, with only a small residual contribution, consistent with the anchor-based attention patterns in Fig. 2. Level shifts are led by the residual and the drift path. Variance jumps are also drift path-led with moderate support from the spike path. Flatlines rely mainly on the drift path with a moderate spike path contribution. Correlation breaks show strong, shared improvements from both paths, with the drift path slightly ahead. On normal windows the contributions are small and distributed, indicating that the mixture cleans reconstructions without overusing the residual.

## VII. CONCLUSION

In this paper, we presented STREAM-VAE, a time-scale-aware variational autoencoder for anomaly detection in vehicle telemetry. The model centers on a dual-path encoder that separates slow drift from fast spike dynamics, supported by a lightweight decoder design that models transient deviations

without confusing them with normal operating behavior. Experiments on automotive data and a public benchmark show that this separation of time scales improves robustness while keeping computation suitable for both in-vehicle deployment and backend fleet analysis. STREAM-VAE therefore provides a compact and practical basis for reliable telemetry anomaly detection in intelligent vehicles. A current limitation is that the method relies on per-entity calibration that may shift across vehicles and environments. This suggests future work on more transferable calibration.

## REFERENCES

[1] D. P. Kingma and M. Welling, "Auto-encoding variational Bayes," in *2nd International Conference on Learning Representations, ICLR 2014, Banff, AB, Canada, April 14-16, 2014, Conference Track Proceedings*, Y. Bengio and Y. LeCun, Eds., 2014.

[2] L. Li, J. Yan, H. Wang, and Y. Jin, "Anomaly detection of time series with smoothness-inducing sequential variational auto-encoder," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 32, no. 3, pp. 1177–1191, 2021.

[3] Y. Zhao, X. Zhang, Z. Shang, and Z. Cao, "DA-LSTM-VAE: Dual-stage attention-based LSTM-VAE for KPI anomaly detection," *Entropy*, vol. 24, no. 11, 2022.

[4] J. Pereira and M. Silveira, "Unsupervised anomaly detection in energy time series data using variational recurrent autoencoders with attention," in *2018 17th IEEE International Conference on Machine Learning and Applications (ICMLA)*, 2018, pp. 1275–1282.

[5] H. Bahuleyan, L. Mou, O. Vechtomova, and P. Poupart, "Variational attention for sequence-to-sequence models," in *Proceedings of the 27th International Conference on Computational Linguistics*, E. M. Bender, L. Derczynski, and P. Isabelle, Eds. Santa Fe, New Mexico, USA: Association for Computational Linguistics, Aug. 2018, pp. 1672–1682.

[6] L. Correia, J. Goos, P. Klein, T. Bäck, and A. Kononova, "MA-VAE: Multi-head attention-based variational autoencoder approach for anomaly detection in multivariate time-series applied to automotive endurance powertrain testing," in *Proceedings of the 15th International Joint Conference on Computational Intelligence - NCTA*, INSTICC. SciTePress, 2023, pp. 407–418.

[7] H. Shao, S. Yao, D. Sun, A. Zhang, S. Liu, D. Liu, J. Wang, and T. Abdelzaher, "Controlvae: controllable variational autoencoder," in *Proceedings of the 37th International Conference on Machine Learning*, ser. ICML'20. JMLR.org, 2020.

[8] B. Lim, S. O. Arık, N. Loeff, and T. Pfister, "Temporal fusion transformers for interpretable multi-horizon time series forecasting," *International Journal of Forecasting*, vol. 37, no. 4, pp. 1748–1764, 2021.

[9] J. Xu, H. Wu, J. Wang, and M. Long, "Anomaly transformer: Time series anomaly detection with association discrepancy," in *International Conference on Learning Representations*, 2022.

[10] A. Deng and B. Hooi, "Graph neural network-based anomaly detection in multivariate time series," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 35, no. 5, 2021, pp. 4027–4035.

[11] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, "Attention is all you need," in *Proceedings of the 31st International Conference on Neural Information Processing Systems*, 2017, p. 6000–6010.

[12] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016, pp. 770–778.

[13] D. Donoho, "De-noising by soft-thresholding," *IEEE Transactions on Information Theory*, vol. 41, no. 3, pp. 613–627, 1995.

[14] F. T. Liu, K. M. Ting, and Z.-H. Zhou, "Isolation forest," in *2008 Eighth IEEE International Conference on Data Mining*, 2008, pp. 413–422.

[15] Y. Su, Y. Zhao, C. Niu, R. Liu, W. Sun, and D. Pei, "Robust anomaly detection for multivariate time series through stochastic recurrent neural network," in *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, ser. KDD '19. New York, NY, USA: Association for Computing Machinery, 2019, p. 2828–2837.

[16] D. J. Rezende and S. Mohamed, "Variational inference with normalizing flows," in *Proceedings of the 32nd International Conference on International Conference on Machine Learning - Volume 37*, ser. ICML'15. JMLR.org, 2015, p. 1530–1538.

[17] J. Pereira and M. Silveira, "Unsupervised representation learning and anomaly detection in ECG sequences," *Int. J. Data Min. Bioinformatics*, vol. 22, no. 4, p. 389–407, Jan. 2019.

[18] J. von Schleinitz, M. Graf, W. Trutschnig, and A. Schröder, "VASP: An autoencoder-based approach for multivariate anomaly detection and robust time series prediction with application in motorsport," *Engineering Applications of Artificial Intelligence*, vol. 104, p. 104354, 2021.

[19] A. Graves and J. Schmidhuber, "Framewise phoneme classification with bidirectional LSTM networks," in *Proceedings. 2005 IEEE International Joint Conference on Neural Networks, 2005.*, vol. 4, 2005, pp. 2047–2052 vol. 4.

[20] J. Ainslie, J. Lee-Thorp, M. de Jong, Y. Zemlyanskiy, F. Lebrón, and S. K. Sanghai, "GQA: Training generalized multi-query transformer models from multi-head checkpoints," *ArXiv*, vol. abs/2305.13245, 2023.

[21] R. A. Jacobs, M. I. Jordan, S. J. Nowlan, and G. E. Hinton, "Adaptive mixtures of local experts," *Neural Computation*, vol. 3, no. 1, pp. 79–87, 1991.

[22] M. Jordan and R. Jacobs, "Hierarchical mixtures of experts and the EM algorithm," in *Proceedings of 1993 International Conference on Neural Networks (IJCNN-93-Nagoya, Japan)*, vol. 2, 1993, pp. 1339–1344 vol.2.

[23] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization." in *ICLR*, 2015.

[24] S. Coles, *An Introduction to Statistical Modeling of Extreme Values*. Springer London, 2001.

[25] P. J. Rousseeuw and C. Croux, "Alternatives to the median absolute deviation," *Journal of the American Statistical Association*, vol. 88, no. 424, pp. 1273–1283, 1993.

[26] A. Siffer, P.-A. Fouque, A. Termier, and C. Largouet, "Anomaly detection in streams with extreme value theory," in *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, ser. KDD '17. New York, NY, USA: Association for Computing Machinery, 2017, p. 1067–1075.

[27] H. Xu, W. Chen, N. Zhao, Z. Li, J. Bu, Z. Li, Y. Liu, Y. Zhao, D. Pei, Y. Feng, J. Chen, Z. Wang, and H. Qiao, "Unsupervised anomaly detection via variational auto-encoder for seasonal KPIs in web applications," in *Proceedings of the 2018 World Wide Web Conference*, ser. WWW '18. Republic and Canton of Geneva, CHE: International World Wide Web Conferences Steering Committee, 2018, p. 187–196.

[28] T. Akiba, S. Sano, T. Yanase, T. Ohta, and M. Koyama, "Optuna: A next-generation hyperparameter optimization framework," in *The 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, 2019, pp. 2623–2631.