

On the utility of coarse-grained mobile network telemetry in detecting performance degradation

Georgios Patounas*, Andres Gonzalez†, Ahmed Elmokashfi*

Abstract—This paper investigates whether passive coarse-grained Radio Access Network measurements can help inferring end users’ experience. To this end, we combine network-side logs with active end to end measurements from a number of probes that act as end users. The combined measurements are used to train a machine learning classifier that uses network side logs to infer whether a particular hour is problematic. We find that coarse-grained network telemetry can be a reliable indicator of performance degradation. Our methodology allows for the automatic detection of problems and root cause analysis through examination of the dominant performance indicators.

I. INTRODUCTION

Mobile Broadband (MBB) networks are becoming the connectivity mean of choice for a large number of users and services. The fifth generation of mobile networks (5G), that is currently being rolled out, promises transforming mobile networks from providing basic voice and data services to flexibly catering to a wide array of industries. These in turn have a diverse set of requirements (e.g. low latency, high reliability, high throughput and low energy consumption) [1]. The high reliance on MBB networks comes with high expectations on the users’ side who need consistent, predictable and reliable network performance as well as coverage almost everywhere. These expectations mean that mobile operators need to build robust approaches to Service Assurance (SA).

The SA process will vary depending on the underlying service. For example, critical services with extreme requirements (e.g. ultra-reliable low latency services like smart grid communications) may require a mixture of continuous monitoring, preemptive resource allocation and swift response to performance degradation. Coarser monitoring and relatively longer response cycles can be adequate for other services like fixed broadband access. While the details of SA architectures are yet to be standardized, what is clear is that it will involve a combination of high and low frequency telemetry. Furthermore, it will likely build on telemetry available in today’s networks. Motivated by this, we investigate whether existing coarse-grained base-station telemetry can help illuminate performance degradation impacting end users.

Base-stations keep track of a range of performance and availability counters, referred to as cell Key Performance Indicators (KPIs). Examples include traffic volumes and the

number of dropped data sessions. These KPIs are coarse-grained aggregates, often reported every hour and are typically used to characterize a base-station’s performance over a long timescale. For instance, a base-station with a sustained level of dropped sessions can be flagged for reconfiguration. Such decisions are made based on past experiences and expert opinion and in absence of direct correlation with users’ Quality of Experience (QoE). The way cell KPIs are leveraged today relegates them to the province of long-term network planning.

This paper investigates whether cell KPIs can be leveraged differently. We approach this by pairing hourly cell KPIs from a commercial mobile operator with fine-grained End-to-End (e2e) measurements from a set of dedicated probes that act as end users. The e2e measurements capture users’ QoE and can thus be used to identify periods with degraded performance. This allows us to label the collected cell KPIs as anomalous or not, which we use to train a Machine Learning (ML) classifier to assess whether changes in KPIs values can reveal increases in packet loss rate or jumps in latency for end users. Our goal is not to provide an approach for detecting poor performance in real-time (i.e. SA for critical services), since this will require more frequent measurements. Instead, we provide a mean for flagging poor performance for services with relatively looser SA expectations like fixed and mobile broadband access.

We find that although coarse-granular, cell KPIs can indeed illuminate performance degradation affecting end users. Our classifier achieves accuracy close to 90% in flagging problem hours, providing interpretable results that point to the root causes of poor performance. It is *timely* and *light* imposing minimal measurement overhead. We believe that our findings can help mobile operators exploiting cell KPIs for more than long-term network planning and sporadic troubleshooting. Thus adding an important input to the process of building suitable SA architectures for current and future mobile networks.

In summary, this paper makes the following contributions:

- 1) We present a novel study that combines both e2e active measurements and network-side cell KPIs in order to improve the inference of end-user performance degradation.
- 2) We evaluate a range of Machine Learning algorithms and approaches to thresholding for inferring end-user performance degradation, based on coarse grained cell KPIs.
- 3) We evaluate the importance of specific KPIs for a range of scenarios and find that a solution with low overhead and high performance is possible.
- 4) We highlight the role of handovers and provide a methodology to allow identification of their root cause.

* G. Patounas and A. Elmokashfi are with the Simula Metropolitan CDE, Oslo, Norway (e-mail: {gpatounas, ahmed}@simula.no). † A. Gonzalez is with Telenor Research, Oslo, Norway (e-mail: andres.gonzalez@telenor.com).

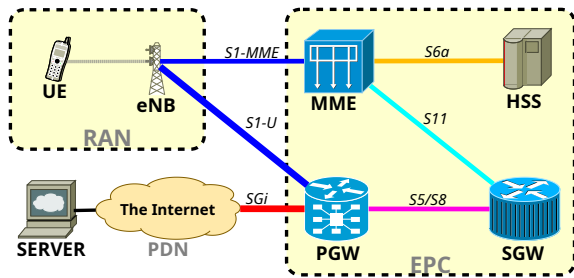


Fig. 1: Architecture of LTE MBB and measurements setup components.

II. BACKGROUND AND MEASUREMENT INFRASTRUCTURE

This section gives an overview of the architecture of MBB networks and presents the infrastructure that was employed to conduct the measurements used in this study.

A. MBB Network Architecture

In this paper, we use measurements from a mobile operator providing 2G, 3G, Long Term Evolution (LTE) and soon 5G services to users nationwide, totaling almost three million subscriptions in 2019. Fig. 1 shows the main components of a typical LTE MBB network divided into the Evolved Universal Terrestrial Radio Access Network (E-UTRAN) and Evolved Packet Core (EPC). The E-UTRAN consists of User Equipment (UE) and evolved Node B (eNB) base-stations, while the EPC consists of several gateways and supporting systems. Out of the main components, the Mobility Management Entity (MME) handles all control plane functions for subscriber and session management, whereas Serving Gateway (SGW) and Packet Delivery Network Gateway (PGW) are the packet data service termination points towards E-UTRAN and packet data network, respectively. For the remainder of the paper we refer to E-UTRAN and EPC as simply the Radio Access Network (RAN) and Core Network (CN).

B. Measurement setup

We leverage the NorNet edge infrastructure [2], which is a testbed dedicated to measurements and experimentation on operational mobile networks. It consists of several hundreds of stationary and mobile, geographically distributed measurement nodes and a well-provisioned server-side infrastructure for the collection of measurements.

A single measurement node is a single-board computer running a standard Linux distribution. Each node is connected to one or more MBB network. 3G/LTE Commercial Off-The-Shelf (COTS) UE is used to connect to the MBB networks. Modems are configured to prefer LTE Radio Access Technology (RAT) when available and fall back to 3G and 2G otherwise. The nodes are capable of running different experiments, which are performed against measurement servers that are part of the measurement setup back-end (Fig. 1). Software running on the nodes also collects various metadata attributes (RAT, signal strength, etc.) from the modems and periodically sends them to the back-end.

For the purposes of this study, the focus is placed on stationary nodes with constant LTE connectivity. The nodes connect to a number of LTE base-stations concentrated in the metropolitan area of Oslo as well as its outskirts and three cities in the south, west and north of Oslo, at distances up to 340 kilometers as the crow flies. The server-side infrastructure is located in Oslo, providing a mix of short and long range links to the measurement nodes through the MBB back-haul and the Internet.

III. EXPERIMENTAL METHODOLOGY

Next, we describe the measurements we use in this paper.

A. Active Measurements

Each measurement node sends a 20-byte UDP packet every second to an Echo Server (ES) using all available MBB connections. The ES is a part of the testbed back-end. Each UDP packet sent by a node is recorded into a database along with a timestamp when the packet was sent and an incremental sequence number that is part of the packet's payload. If a reply from the ES arrives within 60 seconds, the **Round-Trip Time (RTT)** is recorded. Otherwise, the packet is considered lost. To measure the overall **loss rate**, we aggregate the measurements into 5-minute bins. Each bin contains the number of sent and received packets. In addition, we collect node status measurements. These include the current serving cell, radio access technology and signal quality, collected every minute.

B. Passive Measurements

One of the fundamental tools for mobile operators to track the quality provided by their cells are the cell KPIs. In order to have a standard reference and metrics, the ETSI together with the 3GPP released the specification documents [3], [4], where the concepts and requirements of the KPIs for E-UTRAN are presented. Based on the ETSI and 3GPP work, those KPIs can be classified in five different testing groups. *Accessibility* KPIs evaluate the potential difficulties that a user has in order to get the service. These measurements assist the network operator with information about the accessibility provided to their customers. *Retainability* KPIs can be used to estimate the ability of the network to retain an established service for a user. They can indicate whether the end user has been interrupted during the use of a service or not. *Availability* is a simple and illustrative set of KPIs that measures the percentage of time that the cell is considered available to end users. For the analysis performed throughout this paper, we have combined the accessibility, retainability and availability KPIs under the umbrella of *stability*. *Mobility* KPIs measure how the system behaves during handover procedures. Finally, *Integrity* KPIs are in charge of showing how the E-UTRAN impacts the service quality provided to an end-user. A sample of basic metrics from each group is given in table I. Our dataset comes with 270 unique KPIs, which are themselves only a subset of what is provided by the eNB vendors.

KPIs group	Metric
Stability	Call, Radio Resource Control (RRC) and Radio Access Bearer (E-RAB) establishment success rate RRC re-establishment failures, E-RAB abnormal release rate, UE context release success rate
Mobility	Intra and inter-frequency handover success rate, S1 handover success rate
Integrity	Up-link peak user throughput, downlink peak user throughput, downlink latency and transport block error rate

TABLE I: Sample of measurements collected by base-stations.

IV. DATA PRE-PROCESSING

The measurements and logging described in Sec. III, provide us with three datasets that need to be pre-processed and fused. These are: 1) *Node status*, i.e. timestamp, serving cell, 2) *Node performance*, i.e. timestamp, packet loss, RTT and 3) *Cell KPIs*. In this section, we go through the basic pre-processing that is necessary before we can dive into analysis.

A. Filtering

The first two datasets include logs and active measurements for 26 nodes. For this work, we have selected a subset of 10 stationary nodes that for the duration of the measurement period were connected exclusively through LTE.

The third dataset provided by the operator includes KPIs from 53 cells, 23 of which served the selected nodes. This dataset comes with a number of artifacts due to the fact that the full suite of defined KPIs is not punctually collected by all cells. It is common that a cell has not activated logging of certain KPIs or they are not applicable to its state. In order to unify the dataset, we filter KPIs with sparse, absent or incorrect values which leaves us with 88. We then manually examine them and remove ones that do not apply to our use-case (e.g. LTE to 3G handovers) to arrive at 76 KPIs that can be split according to the categorization introduced in Sec. III-B to 41 characterizing availability, 28 traffic and 7 mobility.

B. Aggregating and Fusing the Datasets

Due to the different aggregation granularity utilized by the testbed and the operator, the timestamps need to be set on a time-step equal to the lowest common denominator. In our datasets, this is the 1-hour time-step, dictated by the cell KPIs.

First, the node status (metadata) which is given every minute, needs to be matched and fused with node performance which is available as an average over 5-minute bins. Due to infrequent events of loss of mobile connectivity, there are cases where the node status is not available for a portion of the 5-minute bin. Such bins amount to 0.3% of our dataset and are discarded to avoid inconsistencies and possible loss of handover events. Following this, the new 5-minute bins of node status and performance, are matched and fused with the cell KPIs. The final fused dataset contains 2314 hours of observations.

Threshold	Value
Packet loss rate	0.003 / 0.006 / 0.01
Round-Trip Time	75th percentile / 85th percentile
Violation frequency	2 bins / 5 bins

TABLE II: Range of evaluated thresholds (all combinations).

C. Setting Thresholds and Labeling

As a starting point, we need a definition for performance degradation. To this end, we use the active measurements (i.e. packet loss and RTT) that are carried by the nodes as a performance yardstick.

We define thresholds for packet loss based on the performance of all nodes combined, as a universal value of close to zero packet loss is desired for any node under nominal operation. The packet loss thresholds are applied to each 5-minute bin and expressed as a percentage of seconds with lost packets. On the other hand, thresholds for RTT are individually calculated for each of the nodes, to account for persistent characteristics (e.g. location, cell equipment) that alter the expected baseline of performance. The RTT thresholds are applied to each 5-minute bin and expressed as a percentile of observed RTT for each node.

To produce the final label of an observation, we introduce a threshold on the frequency of violations per hour. The frequency threshold is calculated for the number of 5-minute bins that violate the thresholds for packet loss or RTT (e.g. an hour that includes 1 bin that exceeds its packet loss threshold, 1 that exceeds its RTT threshold and 1 that exceeds both, produces frequency = 4). The complete dataset obtains a number of labels according to a wide range of thresholds that we define based on the methodology described here and listed in table II.

Fig. 2 presents an example of 1 day of measurements. The top panel shows two KPIs that capture the change in the number of users being served by the eNB that is serving the measurement node, while the bottom panel shows the packet loss and RTT measured by the node. The RTT increases as the number of users connected to the eNB increases, while packet loss does not seem to correlate with the number of connected users. This simple example illustrates that cell KPIs may be useful in inferring performance deprecations. However, there is a need for methods that scale beyond basic visual inspection and pairwise correlations.

V. CLASSIFICATION METHODS AND SENSITIVITY ANALYSIS

The problem of inferring whether a particular hour is problematic (i.e. according to the thresholds in Sec. IV-C) or not is essentially a classification problem. Furthermore, the sheer number of cell KPIs means that simple pairwise correlations between end to end measurements and cell KPIs may prove intractable and not easy to map to a final label. Instead, we explore whether a ML-based classifier can overcome these limitations.

In the following, we discuss various choices that we make to prepare our data before applying ML. We also compare

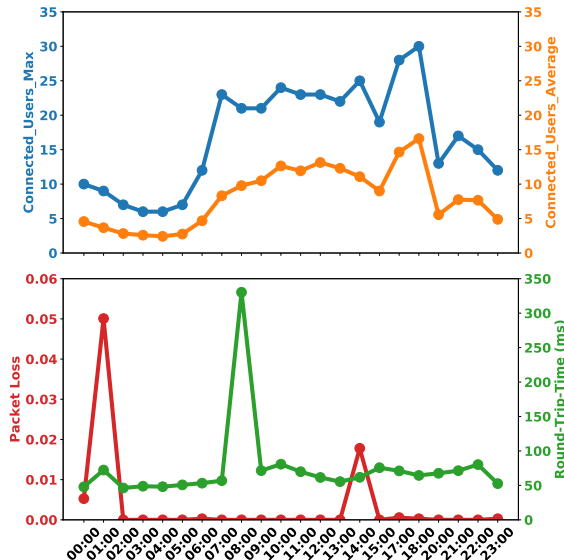


Fig. 2: A day’s sample of KPIs & performance measurements.

the performance of various ML classifiers and evaluate their sensitivity to our performance identification thresholds.

A. Handovers

In cases where one or more handovers are performed over the course of an hour, a node will be connected to multiple cells. Therefore KPIs from all the cells that a node has been connected to over the hour are collected. Out of the 2314 hours that our final dataset contains, 677 involve one or more handover events.

Handovers, by nature, introduce additional delay and an elevated risk of packet loss [5]. Therefore, we examine the performance of time-bins where the node has been connected to a single cell throughout the hour separately from time-bins that include handovers. Splitting the dataset provides two benefits. First, the dataset containing the hours of single-cell operation remains unaffected by handovers, which are already a well established cause of performance degradation. Second, the dataset containing only the hours of handovers will establish a new baseline for user performance and the correlated KPIs. This allows the classifier to be trained in detecting causes of performance degradation other than the known handover.

B. Class Imbalance

The datasets are labeled according to the thresholds set on the active measurements (see Sec. IV-C), to separate normal hours from those with degraded performance. As should be expected, the normal hours outnumber those with anomalies making our data-set imbalanced. Table III lists the initial datasets for selected threshold levels (see Sec. V-D for details).

Before attempting to classify, the dataset needs to be balanced. Imbalanced datasets can affect accuracy by tricking the classification algorithms into placing more emphasis on the dominant class. In this way, while usual metrics may indicate good classification performance, the algorithm has in fact learned to ignore the under-represented class. There are

Dataset	Threshold level	N	A	Ratio
No-handovers	Low	1147	490	2.3:1
No-handovers	High	1367	270	5.1:1
Handovers	Low	519	158	3.3:1
Handovers	High	606	71	8.5:1

TABLE III: Ratio of normal (N) and anomalous (A) hours.

two approaches to balance classes. The first is under-sampling, where the observations belonging to the majority class are reduced to match the number of observations in the minority class. This can be achieved by simply selecting a subset of the observations or using more sophisticated techniques such as creating new synthetic observations that summarize the original. Under-sampling comes at the expense of loss of information, that can be significant for small or severely skewed datasets. The second is oversampling, where the minority class is augmented with synthetic samples to match the number of observations in the majority class. However, over-sampling introduces the risk of over-fitting. Fortunately, a number of over-sampling techniques that attempt to mitigate this exist as well as suitable evaluations that can be performed to ensure our solutions are not affected (see Sec.V-C). The results presented henceforth are thus based on oversampling with ADASYN [6].

C. Classification Techniques

With the thresholds defined, the features cleaned up and the dataset labeled and balanced, we turn to methods that can automatically classify our observations.

There is a range of well-known classification methods for binary classification. We evaluate four commonly used algorithms and their variations, namely Logistic Regression, Random Forest (RF), State Vector Machines (SVM) and the Naive Bayes classifier. [7]

Here we summarize the performance of the algorithms using the metrics of Accuracy (ACC), Positive Predictive Value (PPV) commonly known as precision and True Positive Rate (TPR) commonly known as recall. These metrics stem from the ratios of True Positive (TP), True Negative (TN), False Positive (FP), False Negative (FN) that are produced by the classifiers. Accuracy is the overall percentage of correct identifications relative to the size of the complete dataset. Precision describes how many of the observations that the algorithm selected are actual anomalies while recall describes how many of the total anomalies were pointed out. Random search of a wide range of hyper-parameters and k-fold cross validation is used throughout all our evaluations to ensure performance and generalization.

In Fig. 3 we examine the evaluation metrics for each of the classifiers over the range of thresholds defined in Sec. IV-C and for both the datasets of no-handovers and handovers. While all classifiers provide promising results and other interesting alternatives exist including k-Nearest Neighbors and Neural Networks, by this summary it is evident that considering our dataset, the RF outperforms the other methods in both overall performance and consistency. The use of bootstrap aggregating and random splitting of features allows RF to achieve high accuracy and low over-fitting that consistently outperforms

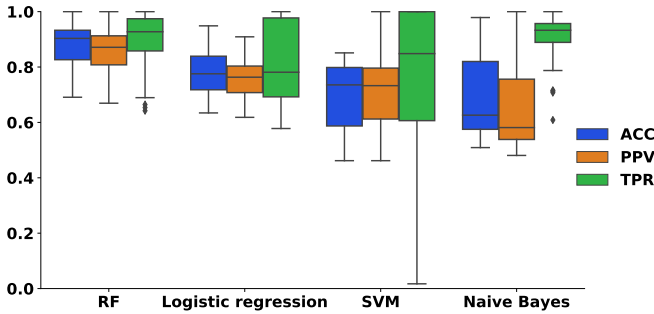


Fig. 3: Evaluation of classifier performance.

alternatives [8] without sacrificing interpretability. Given the results of our exploratory analysis, the relatively small size of the dataset and the simplicity of tuning RF, we defer evaluation of advanced techniques to future work. The results presented henceforth are thus based on RF utilizing 1000 decision trees.

D. Evaluation and Final Selection of Thresholds

Having selected an algorithm, we now move to evaluating the effect of the thresholds defined in Sec. IV-C on the performance of our method. We examine the same set of performance metrics that were presented in Sec. V-C.

In detail, Fig. 4a presents the results for packet loss rate. Looking back to table II, we fix the threshold for packet loss to 0.003, 0.006 (omitted for clarity) and 0.01 accordingly and perform multiple classification rounds with all other combinations of the defined datasets (i.e. no-handover, handover) and remaining thresholds (i.e. RTT, frequency). The resulting box plot isolates the effect that the packet loss threshold has on the performance of the classification algorithm. Likewise RTT is shown in Fig. 4b and frequency in Fig. 4c. On average, all combinations achieve high accuracy, above 0.85. Precision is trailing behind, denoting an increased number of FP. However, what is most valuable is that recall, on average, achieves very high values of ≈ 0.9 meaning that $\approx 90\%$ of anomalies can be consistently identified with most thresholds. It is also apparent that higher thresholds tend to produce improved performance. This can be attributed to their nature of capturing severe degradation that is more likely to be imprinted in the coarse KPIs. It is worth mentioning here that the lowest threshold, by design, captures all packet loss instances, no matter how small. The highest threshold is meant to capture, on average, half of the packet loss focusing on only the more severe instances.

Based on the results of threshold evaluation and adopting the perspective of an operator, we finally define two levels of thresholds to use henceforth, listed in table IV. The first (low) level is meant to flag benign performance degradation, that is indicative of a node or cell approaching the limits of its performance. These are cases that do not necessarily warrant immediate attention but may be interesting to examine or keep a closer watch on. These could also be considered an urgent issue in cases of a sudden surge of low threshold violations on multiple nodes in a cell. The second (high) level is reserved for severe performance degradation that

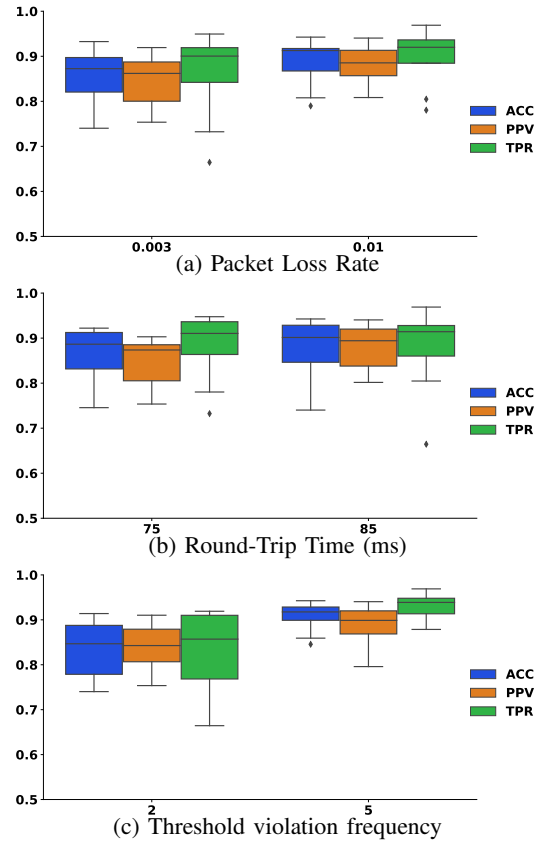


Fig. 4: Evaluation of thresholds.

Threshold lvl	Packet loss	RTT	Frequency
Low	0.003	75th percentile	2 occurrences
High	0.01	85th percentile	5 occurrences

TABLE IV: Final threshold levels selected.

warrants immediate attention. Three parameters define a level of thresholds: sensitivity to packet loss, sensitivity to RTT and sensitivity to frequency of violations.

VI. NETWORK SIDE DETECTION OF PERFORMANCE DEGRADATION

In this section we dive into the effectiveness of cell KPIs in inferring performance degradation experienced by end users. To this end, we use the insights gained in Sec. V, to train our classifier. Given the complete dataset, we feed the KPIs to the RF as features and training is performed based on the labels produced by the selected thresholds.

As explained in V-A, we have split our dataset to hours where a node was connected to a single cell and to hours with one or more handovers, meaning that the node has been connected to multiple cells. We examine each case below.

A. Hours with no Handovers

To compare the performance of the classifier for each of the thresholds, we examine the corresponding Receiver Operating Characteristic (ROC) curves and the F-score [9]. A ROC curve is commonly used to visualize the performance of binary classifiers over varying thresholds or configurations while the

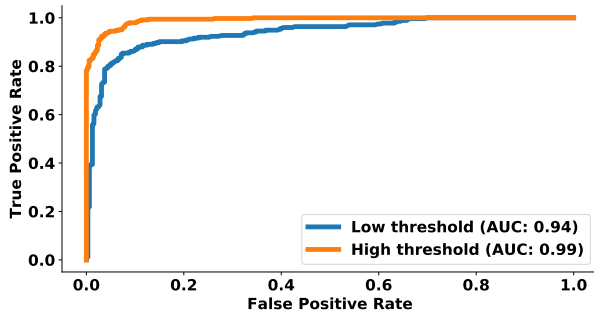


Fig. 5: ROC of no-handover dataset.

KPIs	Imp.	Cat.
HueR_Connected_Users_Max	0.090	traffic
HueR_Connected_Users_Average	0.082	traffic
HueR_Paging_Success_Rate_UU	0.057	stability
HueR_RRC_ReEstablishment_Setup_Attempt_All	0.048	traffic
HueR_RRC_Setup_Request_Signalling	0.046	traffic

TABLE V: Most important KPIs. No-handovers.

F-score considers the precision and recall of the classifier in a unified metric. As shown in Fig. 5 the TP rate is plotted against the FP rate meaning that an ideal classifier that identifies all anomalies correctly and does not produce any FP either, would produce an Area Under the Curve (AUC) = 1.

Both levels of thresholds perform well, achieving an AUC equal to 0.94 and F-score of 0.9 for the low level and AUC 0.99 and F-score of 0.94 for the high level. Translating this to exact numbers of correct and incorrect predictions, looking at the worst case of the low level thresholds: out of 274 normal hours in our testing dataset 39 (14%) were incorrectly thought to be anomalous (FP) while out of 321 anomalous hours 29 (9%) were not identified (FN). Both rates are relatively low. While we can not explain each of these misclassifications, we note that a FP may correspond to an actual problem that only impacts a fraction of connected users which may not involve the measurement node. Further, a FN may be caused by a performance issue that does not lie in the RAN (e.g. packet loss in the CN).

Having confirmed that we can successfully identify user-end performance degradation with high confidence using only the RAN KPIs we turn to feature selection. In studying the utility that each of the KPIs provides, we gain insights and provide a path to lowering the overhead of collection and analysis, by using a subset of appropriate KPIs to achieve comparable performance of identification. Table V lists the five top features in order of importance.

It is evident that during periods of stable connectivity, the performance is directly related to the overall load of the cell, with the top two KPIs characterizing the number of users connected to the cell. To understand how each feature provides the information needed for classification, we use the Kernel Density Estimation (KDE) shown in Fig. 6. Both KPIs exhibit a higher mean and more variability (i.e. longer tail) during anomalous hours. The patterns observed in the KPIs, with the distinct peaks and valleys, are the key pieces that a classification algorithm exploits to infer the existence of

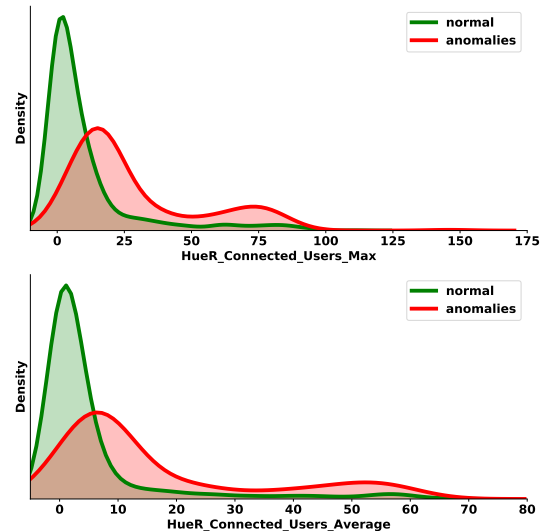


Fig. 6: KDE of most important KPIs. No-handovers.

anomalies. Here it is easy to observe that anomalies tend to happen in periods of elevated user activity.

It is also interesting to consider the categories of KPIs laid out in Sec.III-B as a whole and the contribution each of them made. On average, for both thresholds, traffic KPIs impact was by far the largest, contributing 58%. This is despite the fact that traffic makes up just 37% of the complete set of KPIs used. Stability KPIs contributed 30% while making up 54% of the complete set, making them the least significant group compared to its size. Mobility KPIs contributed 12% while making up 9% of the complete set.

B. Hours with Handovers

To examine the dataset containing hours with handovers some additional consideration is required. In each hour, a node will be connected to multiple cells, meaning that multiple sets of features will be relevant. A methodology is needed to select the set of features that will be associated with the observed end-user performance. This is simple enough during periods of normal end-user performance where every set of features, i.e. every cell that the user was connected to, can be considered as normal. However, when examining an anomalous period, which cell should be blamed is not clear. Once we take into account that we are only dealing with stationary nodes, i.e. the handovers examined are not initiated by physical movement out of a cell's coverage area, we can make some assumptions about the reasons that led to them.

The 3GPP standards dictate that a decision of handover depends mainly on the signal strength provided to the UE [10]. Once the signal strength by the serving cell is lower than a predetermined threshold and a better alternative is available, the handover is initiated. Such conditions are frequently met when a mobile UE is moving out of the coverage area of a cell. However, a stationary UE will tend to connect to a single cell that consistently provides the best performance due to e.g. physical proximity. Indeed, examining our nodes we notice that their majority favors a single "dominant" cell. A handover

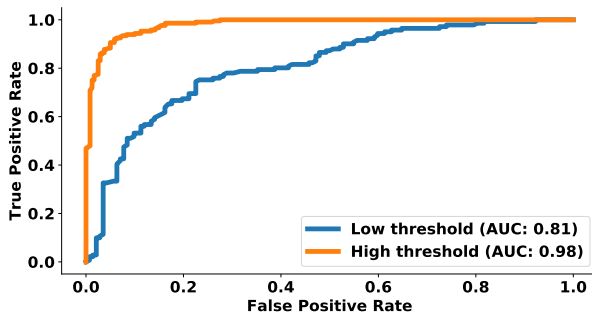


Fig. 7: ROC of handover dataset. All cells blamed.

KPIs	Imp.	Cat.
HueR_Intra_Freq_HO_Out_Attempt	0.073	mobility
HueR_eRAB_Setup_Attempt_Inc_HO	0.063	mobility
HueR_Inc_HO_Attempt	0.050	mobility
HueR_Paging_Attempt_UU	0.038	traffic
HueR_eRAB_Setup_Attempt_All	0.036	traffic

TABLE VI: Important KPIs. Handovers. All cells.

away from the dominant cell can then be an indication of performance degradation or complete cell unavailability.

Considering the above, we examine two approaches to labeling the cells. First, we naively blame all cells which a node was connected to during an anomalous hour. Second, we blame the previously identified dominant cell for each node. Next, we compare the performance of identification for each of the thresholds.

When **both cells are blamed**, while the high level of thresholds performs almost identically to the no-handover dataset, performance using the low level of thresholds has significantly deteriorated (Fig. 7). The low level threshold achieves an AUC of 0.81 and F-score of 0.74 while the high level achieves AUC of 0.98 and F-score of 0.93. We observe a drop in TP and increase in FP that can be attributed to the nature of this dataset which is made up of hours that are expected to, on average, experience worse performance due to the presence of handovers thus masking mild degradation.

Next we examine the most important KPIs for this approach. The top five results of feature selection are shown in table VI. It is immediately evident that a higher than normal incidence of handover related events, is a clear marker for performance degradation experienced by the users of the cells. The three most important features this time belong to the category of mobility KPIs.

On average for both levels of thresholds, the traffic KPIs impact was again the largest, contributing 54%. The stability KPIs contributed 28%. The mobility KPIs contributed 18% which makes them by far the largest contributor compared to the size of the set.

When the **dominant cell is blamed** both thresholds achieve very high levels of performance (Fig. 8). The low level threshold achieves an AUC of 0.98 and F-score of 0.92 while the high level achieves AUC of 0.99 and F-score of 0.97.

We once again examine the most important KPIs, shown in table VII. When focusing on the dominant cell only, which is

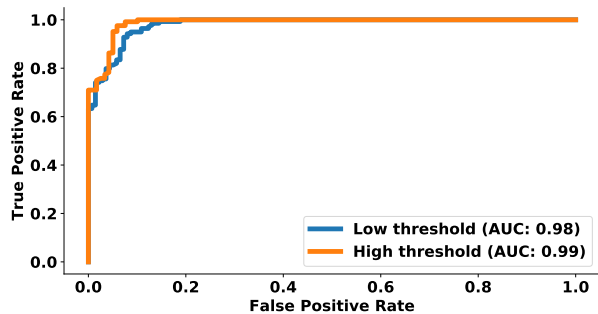


Fig. 8: ROC of handover dataset. Dominant cell blamed.

KPIs	Imp.	Cat.
HueR_Cell_Availability_Rate	0.100	stability
HueR_RRC_Setup_Request_Signalling	0.076	traffic
HueR_Connected_Users_Max	0.066	traffic
HueR_Used_UL_PRBs_Average	0.059	traffic
HueR_Connected_DRBs_All	0.052	traffic

TABLE VII: Important KPIs. Handovers. Dominant cell.

expected to be the initiator (i.e. in terms of being causative) of the handovers, the underlying problem is uncovered. Handover incidence is no longer important (as this is similar for other cells in the area experiencing the issue) but traffic and stability become the top KPIs. This shows that handovers are initiated due to the serving cell becoming unavailable or overwhelming traffic and the combination of this underlying problem and its solution (i.e. handover), strongly affects user-end performance.

On average for both levels of thresholds, the traffic KPIs impact was the largest, contributing 61%. The stability KPIs contributed 30%. The mobility KPIs contributed 9%. The distribution resembles the case of no handovers, where mobility plays a minor role and traffic comes first, both in absolute contribution and relative to the size of the set.

C. Classification Performance with a Reduced Set of KPIs

Here, we briefly examine whether feature selection can be effectively used to minimize the overhead of KPIs collection and processing while maintaining comparable success of identifying performance degradation. As seen in Fig. 9, the results are very promising. For the dataset of no handovers, using just 10 features we can achieve performance closely matching the one produced by the full set of 76. The same success is observed for the handover dataset (figures omitted for brevity).

VII. DISCUSSION

Our results are both unexpected and interesting. In a sense, they breathe new life into the utility of coarse-grained cell KPIs, indicating that network operators can exploit them to infer users' QoE. Next we briefly discuss the interpretation, implications and limitations of our findings.

Interpretation. We find that almost all episodes of degraded performance correlate with changes in cell KPIs. This agrees with the common-wisdom that most mobile network problems are RAN-related. Further, most degraded performance happens at times when eNBs experience an increase in connected users and traffic. We also find that stationary users handover root

causes lie with the original cell and related to cell availability and traffic volumes.

Implications. Inferring whether every hour is problematic or not can help operators reconfigure their networks in a timescale of an hour or two to alleviate problems before receiving users' complaints. Characterizing problematic hours can also help operators tune their configurations to react to future changes in demand in order to minimize such occurrences. As we mentioned before, this level of SA will suit services with best effort expectations (e.g. mobile broadband) and not critical services. We also note that the volume of involved measurements is fairly low (i.e. 10 values from each cell when using feature selection). Hence, the entire inference engine will impose minimal overhead on network resources and can be executed either in a centralized or distributed manner. Inference from different cells can be fed into another classifier to detect correlated outages and common hot spots. This promise, comes with the caveat that e2e measurements are necessary to train the classifier. While difficult today, operators are rolling out offerings that will result in massive stationary deployments e.g. fixed wireless access and stationary IoT sensors that could be leveraged for collecting e2e measurements.

Limitations. The used measurement nodes are fairly sparse (i.e. one node per cell). This may have resulted in the recorded 14% FP. Although this could be a weakness, we argue that the observed effect is limited. Our findings indicate that most problems (about 90%) are RAN-related, which is surprising since we may expect a slightly higher percentage of CN related problems. Investigating this further, we search for hours in our active measurements where several nodes report degraded performance as a proxy for problems beyond the RAN. We found three such hours and all of them were successfully flagged by our classifier. This hints that problems beyond the RAN may also manifest themselves in some cell KPIs like drop in paging and session establishment success rate. While the monitored operator may have had a limited number of CN-related problems, others may not. Hence, in addition to using cell KPIs operators may want to deploy other approaches for detecting correlated performance issues that do not impact the RAN (e.g. [11]). Another performance metric that directly affects the end-users' QoE is the available bandwidth of their connection. While the testbed measures the bandwidth of each connection, at present this is only done sparsely out of consideration for the shared nature of the testbed and the increased utilization that such measurements entail. Future work will seek to obtain frequent measurements of bandwidth for further study.

VIII. RELATED WORK

There is a plethora of work on monitoring and anomaly detection using cell KPIs in mobile networks. One direction foregoes user-side performance characterization, focusing on diagnosing large scale anomalies. [12] discusses autonomous alarm based detection, focusing on self-healing. [13] explores auto-diagnosing faults in 3G cells with the help of expert opinions and automatic thresholds. [14] evaluates supervised

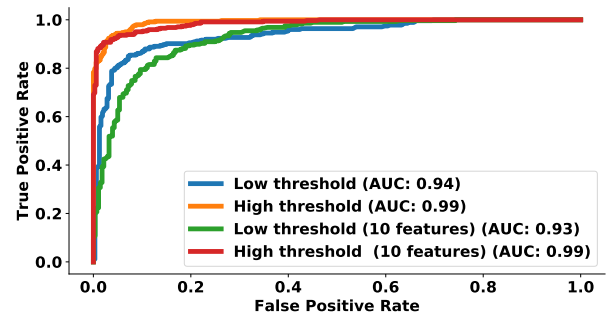


Fig. 9: ROC of no-handover dataset with feature selection.

and unsupervised techniques on real-time collected KPIs for the prediction and root-cause analysis of faults, demonstrating high accuracy for both approaches. [15] proposes a (mostly) unsupervised technique for automatic diagnosis of faults in LTE networks. Using cell KPIs, unsupervised clustering of cell states is performed and complemented by expert opinions for the final root cause analysis. While [16] introduces the aspect of user-end performance it does so post-mortem, through stored customer service calls and periodical crash logs. Another direction targets higher granularity aggregation of KPIs, with [17] predicting call drops by learning based on per-user datasets that present real-time data access and volume challenges. [18] proposes models to automate diagnosis of RAN problems, based on bearer-level aggregation of cell KPIs and bearer records, again facing challenging data volumes.

Our work begins with a strong foundation on a dataset uniquely combining active measurements using geographically distributed COTS UE connected to a national network, with passive measurements from the RAN provided by the operator. We present a methodology with a low barrier of entry and low overhead that can readily be implemented in any deployed network. Threshold-based labeling combines the alarm based systems commonly used by operators with a ML approach. Sourcing the thresholds from user-side KPIs, provides results proven to reflect conditions that directly affect the user QoE. The supervised method lends its high accuracy while forgoing the need for expert input in creating the labels. Finally, feature selection takes a step towards root cause analysis.

IX. CONCLUSIONS

This paper investigates whether coarse-grained RAN KPIs can be used to infer episodes of degraded performance impacting end users. To this end, we pair active measurements from a set of probes that act as end users, with operator side measurements, providing a complete view of e2e performance and network telemetry, on a commercial network. We identify periods of end-user performance degradation and train a classifier to detect them based on aggregated network KPIs. The results are extremely promising with over 90% accuracy in inferring problematic hours. Going further, we identify primary and auxiliary causes to the performance loss through KPIs selection. Our findings pave the way to easily accessible, low overhead Service Assurance.

X. ACKNOWLEDGMENTS

This work was supported by the Research Council of Norway grants 240850 and 209954. We thank our shepherd, Philippe Owezarski and the anonymous reviewers for their thoughtful feedback. We would like also to thank Džiugas Baltrunas for his inputs to the experiment design.

REFERENCES

- [1] 5GPPP Architecture Working Group, “View on 5G Architecture,” *White Paper*, 2017.
- [2] A. Kvalbein, D. Baltrūnas, K. Evensen, J. Xiang, A. Elmokashfi, and S. Ferlin-Oliveira, “The nomet edge platform for mobile broadband measurements,” *Computer Networks*, vol. 61, pp. 88–101, 2014.
- [3] ETSI Technical Specification (TS), “ETSI - TS 132 450. Key Performance Indicators (KPI) for Evolved Universal Terrestrial Radio Access Network (E-UTRAN): Definitions,” 2016-02.
- [4] —, “ETSI - TS 132 451. Key Performance Indicators (KPI) for Evolved Universal Terrestrial Radio Access Network (E-UTRAN): Requirements,” 2011-04.
- [5] K. Dimou, M. Wang, Y. Yang, M. Kazmi, A. Larmo, J. Pettersson, W. Muller, and Y. Timner, “Handover within 3gpp lte: Design principles and performance,” in *2009 IEEE 70th Vehicular Technology Conference Fall*. IEEE, 2009, pp. 1–5.
- [6] H. He, Y. Bai, E. A. Garcia, and S. Li, “Adasyn: Adaptive synthetic sampling approach for imbalanced learning,” in *2008 IEEE international joint conference on neural networks (IEEE world congress on computational intelligence)*. IEEE, 2008, pp. 1322–1328.
- [7] G. James, D. Witten, T. Hastie, and R. Tibshirani, *An introduction to statistical learning*. Springer, 2013, vol. 112.
- [8] L. Breiman, “Random forests,” *Machine learning*, vol. 45, no. 1, pp. 5–32, 2001.
- [9] M. Hossin and M. Sulaiman, “A review on evaluation metrics for data classification evaluations,” *International Journal of Data Mining & Knowledge Management Process*, vol. 5, no. 2, p. 1, 2015.
- [10] G. T. 36.300, “Evolved universal terrestrial radio access (e-utra) and evolved universal terrestrial radio access network (e-utran); overall description; stage 2,” 2011.
- [11] B. Nguyen, Z. Ge, J. Van der Merwe, H. Yan, and J. Yates, “Absence: Usage-based failure detection in mobile networks,” in *Proceedings of the 21st Annual International Conference on Mobile Computing and Networking*, 2015, pp. 464–476.
- [12] Y. Liu, J. Zhang, M. Jiang, D. Raymer, and J. Strassner, “A model-based approach to adding autonomic capabilities to network fault management system,” in *NOMS 2008-2008 IEEE Network Operations and Management Symposium*. IEEE, 2008, pp. 859–862.
- [13] R. Barco, V. Wille, L. Díez, and M. Toril, “Learning of model parameters for fault diagnosis in wireless networks,” *Wireless Networks*, vol. 16, no. 1, pp. 255–271, 2010.
- [14] W. Zhang, R. Ford, J. Cho, C. J. Zhang, Y. Zhang, and D. Raychaudhuri, “Self-organizing cellular radio access network with deep learning,” in *IEEE INFOCOM 2019-IEEE Conference on Computer Communications Workshops (INFOCOM WKSHPS)*. IEEE, 2019, pp. 429–434.
- [15] A. Gómez-Andrades, P. Muñoz, I. Serrano, and R. Barco, “Automatic root cause analysis for lte networks based on unsupervised techniques,” *IEEE Transactions on Vehicular Technology*, vol. 65, no. 4, pp. 2369–2386, 2015.
- [16] C.-Y. Hong, M. Caesar, N. Duffield, and J. Wang, “Tiresias: Online anomaly detection for hierarchical operational network data,” in *2012 IEEE 32nd International Conference on Distributed Computing Systems*. IEEE, 2012, pp. 173–182.
- [17] N. Theera-Ampornpant, S. Bagchi, K. R. Joshi, and R. K. Panta, “Using big data for more dependability: a cellular network tale,” in *Proceedings of the 9th Workshop on Hot Topics in Dependable Systems*, 2013, pp. 1–5.
- [18] A. P. Iyer, L. E. Li, and I. Stoica, “Automating Diagnosis of Cellular Radio Access Network Problems,” in *Proceedings of the 23rd Annual International Conference on Mobile Computing and Networking*. ACM, 2017, pp. 79–87.