



<http://www.diva-portal.org>

Postprint

This is the accepted version of a paper presented at *2017 IEEE 86th Vehicular Technology Conference (VTC Fall), Toronto, Canada, 24–27 September 2017*.

Citation for the original published paper:

Garcia, J., Alfredsson, S., Brunström, A., Beckman, C. (2017)  
Train Velocity and Data Throughput: A Large Scale LTE Cellular Measurements Study  
In: *Proceedings of the 2017 IEEE 86th Vehicular Technology Conference (VTC Fall)*  
(pp. 1-6). New York: IEEE  
IEEE Vehicular Technology Conference VTC

N.B. When citing this work, cite the original published paper.

Permanent link to this version:

<http://urn.kb.se/resolve?urn=urn:nbn:se:kau:diva-65251>

# Train Velocity and Data Throughput - A Large Scale LTE Cellular Measurements Study

Johan Garcia, Stefan Alfredsson, Anna Brunstrom

Department of Mathematics and Computer Science

Karlstad University, Karlstad, Sweden

Email: {johan.garcia, stefan.alfredsson, anna.brunstrom}@kau.se

Claes Beckman

Center for Wireless Systems, Wireless@KTH

KTH Royal Institute of Technology, Stockholm, Sweden

Email: claesb@kth.se

**Abstract**—Train-mounted aggregation routers that provide WiFi access to train passengers and bundle external communication over multiple cellular modems/links is an efficient way of providing communication services on trains. However, the characteristics of such systems have received limited attention in the literature. In this paper we address this gap by examining the communication characteristics of such systems based on a large data set gathered over six months from an operational Swedish railway system. We focus our examination on the relationship between per link throughput and train velocity. Using Levenberg-Marquardt non-linear regression a noticeable critical point is observed for an RS-SINR of around 12 dB. At this point the impact of increased train velocity on per link throughput changes from being negative to becoming positive. Using a machine learning approach we also explore the relative importance of several observed metrics in relation to per link throughput.

**Index Terms**—Cellular networks, LTE, 4G, Trains

## I. INTRODUCTION

As the usage of smart-phones and tablets continues to rise, computer communications has become a more integrated part of everyday life for a large majority of the population. The cellular infrastructure caters for a large fraction of the communication needs of this user base, something that is especially true when users are mobile. One important mode of transportation is railway travel, and users increasingly expect to be able to communicate effortlessly while traveling on the train. However, the radio conditions inside a train, and a potentially large number of simultaneous users that move at high speed, create considerable challenges for the cellular infrastructure. Providing a good Quality-of-Experience for a large number of users inside a train puts considerable demands on the density of cell tower placement along the train track. Instead of dimensioning the cellular network to handle a large and fast-moving group of users, an alternative approach is to improve the communication of users on trains by employing train-mounted aggregation routers. The users on the train access these via WiFi, and the routers then provide caching services and bundle external connections through attached modems connected to external antennas which are optimized

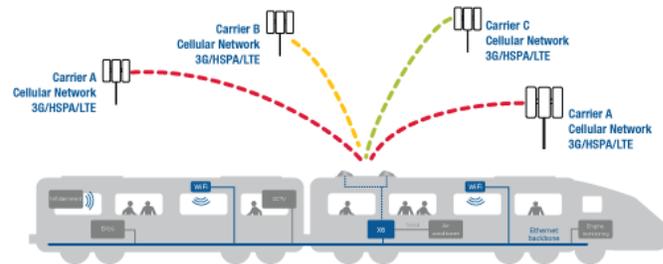


Figure 1: Train-based cellular router structure [courtesy of Icomera AB]

for train mounting. An illustration of such a system is provided in Figure 1.

The system connects to multiple operators and can aggregate the traffic over multiple cellular links. As far as the authors know, the communication characteristics of operational systems on trains have not been previously examined at scale in the literature. In this paper we present a study based on a large data set from such a system that is in operation in the Swedish railway system. Data is collected from over 90 trains completing over 7000 journeys between four major cities in Sweden.

This paper provides two main contributions. First we provide an examination of the relationship between train velocity and observable throughput performance. Here, a smaller subset of the data with high load is used to mitigate for dependency on aggregate demand. An unexpected critical point is observed, and Levenberg-Marquardt non-linear regression is used to locate it at an RS-SINR of around 12 dB. Below this point the change in throughput with increased velocity is negative, and above it becomes positive. Secondly, we examine the relative importance of the measured metrics in relation to per link throughput performance. Using a random-forest predictor importance approach we find that no metric is strongly dominant, and that there appears to be collinearity between the different metrics.

The remainder of this paper is organized as follows. The next section discuss related work, while Section 3 elaborates on data collection procedures, and provides an overview of the data set. Section 4 examines the impact of train velocity and elaborates on the machine learning based relative metric importance examination, followed by the conclusions in Section 5.

## II. RELATED WORK

There are a number of studies related to cellular mobile communication at high vehicular speeds, for railway trains in particular, and studies on characterization of LTE metrics. The feasibility of LTE on trains have been studied for example by [2], [4], [8], [10], [11], [12], [13]. These studies are mainly focused on radio level issues or the general feasibility of using LTE for train communication; Masson et al. [10] perform a survey on the technical solutions available to provide Internet access on high speed trains. They note that train access terminals “represents the best technical solution to optimize performance and throughput”. Train access terminals consist of a mobile router that provides local user access, for example via WiFi, and is connected via a cellular or satellite uplink. The system we study in this paper thus belongs to this technical approach. Mueller et al. [12] perform a simulation study to compare the performance of using a system based on train access terminals to direct communication with the UEs on the train. Using train access terminals is shown superior, although direct communication with UEs is not so far behind. In [8], Lin et al. do an analytical study of using a train access terminal with dual radios (3G/LTE), and show that this can improve handover performance for high-speed trains as there is an overlap coverage area of base stations, where one radio can do an early handover while the other radio holds on longer to the existing cell before handover. Merz et al. [11] study the radio link performance of LTE at high velocities based on two real life train traces, of two hours and 40 minutes each. SNR is shown to be the most important factor for reliable operation.

For HSPA+, Li et al. [7] reported measurements on high-speed trains which revealed severe TCP problems due to RTT spikes, drops and disconnections. In [9], Lutu et al. use NorNet Edge nodes on Norwegian trains. The nodes act as, and experience the performance of, end-users in the train. From these measurements they build an operator coverage map, where the country is divided in mosaic segments. Within each segment, packet loss and HTTP download performance is assessed for each of the operators. These studies both perform measurements of individual user terminals on the train. In contrast, our work investigates measurements of passengers traffic sent via an on-board WiFi router, and aggregated over multiple live mobile broadband links in the backhaul.

The effect of LTE metrics and the impact on higher layers are discussed in e.g. [5], where Elnashar and El-Saidny show LTE performance in a test-driving scenario. Here, they show a clear correlation between the SINR and the download throughput, for both field measurements and simulated results. The measurements also show the effect of path loss (typically due



Figure 2: Train lines used in evaluation

to increasing distance) on throughput. They also investigate handover time and the data interruption time, which gives a more detailed understanding of the channel and system interaction.

## III. DATA SET

The data used in this study was collected during six months in 2016. Train journeys performed along three different stretches of rail track were included in the data set: Stockholm - Göteborg, Stockholm - Malmö, and Stockholm - Karlstad. The geographical outline of these train lines are shown in Figure 2. As can be seen, the lines are sharing some sections, but are mostly non-overlapping. Data from over 7000 journeys are included in the data set. The data set contains 97 unique router ids, which can be expected to correspond to the number of unique train sets. Train-line details are provided in Table I. The number of unique cell ids seen by the modems are also recorded. As can be derived from the values in the table, the average number of observed cell ids per km of track is quite varied between the lines, with StoGbg having 11.5 cell ids per km whereas StoMal only has 8.2. This observation is mostly of cursory importance, however, as the geographical clustering of cells and their distance to the track is of more importance.

Data is collected by the on-board router at five second intervals. Data is collected for general metrics such as number of active devices, GPS positions and current velocity of the train, as well as radio-related metrics and performance metrics such as aggregated and per link throughput and ping delay.

Table I: Examined train lines

Line name	Track Length	Nr of journeys	Nr router ids	Nr of cell-ids	Avg velocity
StoGbg	485 km	2380	54	5589	103 km/h
StoMal	615 km	3900	36	5048	137 km/h
StoKsd	325 km	1458	52	2765	135 km/h
Overall		7738	97	11644	125 km/h

Additional characterization of these collected performance metrics is available in [6].

Each onboard router uses four Sierra Wireless MC7710 LTE modems, with two modems assigned to each of two cellular operators. The operators are in the following referred to as operator 1 and 2 (Op 1 and Op 2). The router monitors the individual communication conditions on each link and schedules traffic over links according to a proprietary scheme.

Before performing analysis, the data is adapted and filtered. Spurious zero values are reported for various metrics, possibly due to router reboots, modem failures or other causes. Journeys with such data are filtered out, along with those where values are outside the possible operating conditions.

#### IV. MEASUREMENT ANALYSIS

Various mobile network metrics [1] are collected by the user equipment, which here are the mobile broadband modems. The RSSI (Received Signal Strength Indicator) represents the average total received power in the measurement bandwidth and is reported in dB. The RSSI metric includes all signals received over the measured frequency band. Even if the received signal is strong, it also includes noise, which can drown the data-carrying signal. An alternate metric is the RS-SINR, which measures the signal to interference and noise for the reference signal. Further indicating the signal quality, the RSRP (Reference Signal Received Power) metric represents the average of received power of a reference signal, and is reported in dBm. To get a generic comparable metric, the RSRP and RSSI are used to calculate the RSRQ (Reference Signal Received Quality) as the ratio  $N \cdot \text{RSRP} / \text{RSSI}$  over the  $N$  resource blocks used in the RSSI measurement. Also relevant for throughput is the link bandwidth. In LTE, this is up to 20 MHz carrier bandwidth. Another factor is the LTE carrier frequency, typically at the 800, 900, 1800 or 2600 MHz bands. Higher frequencies result in shorter transmission range, i.e., resulting in the signal strength and the capacity dropping off earlier as user equipment move further from the base station, compared to a lower carrier frequency.

##### A. Relationship between train velocity and per link throughput

This section examines to what extent the velocity of the train is related to variations in observed per link downlink throughput. As the data set is collected in a purely passive way, the degree of control over the data collection conditions is limited. A challenge in this data set is that the offered load is varying. Observed performance is not always bounded by the system throughput but by offered load. For per-link analysis, the operation of the link scheduling is also a factor.

To minimize the risk that the offered load is insufficient to load the network, we in the following only consider measurements which have more than 150 active devices. This corresponds to approximately the 90th percentile of all measurements. Further, the analysis only uses data for the one link out of the four concurrent links that reports the highest throughput at each measurement time, to minimize link scheduling artifacts. The data is further filtered so that only

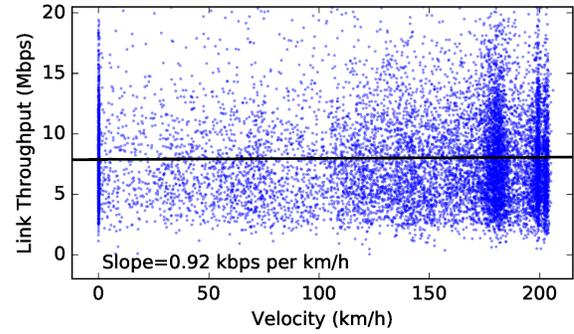


Figure 4: Scatter plots between link throughput and velocity for one train

measurements that have an RS-SINR reading are considered, which is only reported by the modem for a fraction of the measurements. To allow analysis to be performed with sufficient data, trains which now have less than 10000 measurements are filtered away.

After filtering, 514437 measurements from 23 trains remain, with between 50077 and 10560 measurements per train. A scatter plot of throughput over velocity values for the train with 50077 measurements is provided in Figure 4, along with a linear regression to receive an indication of any potential trend for throughput as velocity increases. Also shown in the figure is the slope coefficient for the linear fit, which here is weakly positive at 0.92 kbps per km/h. However as 13 of the trains have negative slopes and the median slope is -0.35, no clear trend for the relationship of train velocity and link throughput is discernible when looking at all data on a per train basis.

As a result of the large number of measurement points collected, it is however possible to further analyze the data for potential interaction effects. In particular, the RS-SINR was found interesting. Using the same measurement set for one train as Figure 4, but binning the measurements into different intervals of RS-SINR values allows a more detailed examination as shown in Figure 3. The center of each 3 dB RS-SINR interval is indicated in each subfigure. Each subfigure also includes a line with a linear fit to the data and the 95% confidence interval of the fit. Looking at the placement of the linear fit line along the y axis it is consistent with expectations. The average throughput increases from ca 5 Mbps to ca 10 Mbps as the RS-SINR increase. However, an intriguing trend is observable in the slope of the linear fit. It is seen that for low RS-SINR values the slope shows a negative relationship between velocity and throughput. For the 13 dB bin this however changes, and the slope shows an increasing positive trend as the RS-SINR increases.

While the data shown in Figure 3 represents only one train, a similar behavior of the slope is consistently observed also for the other trains. This is shown in Figure 5, which shows the slope coefficient over different RS-SINR bins for each of

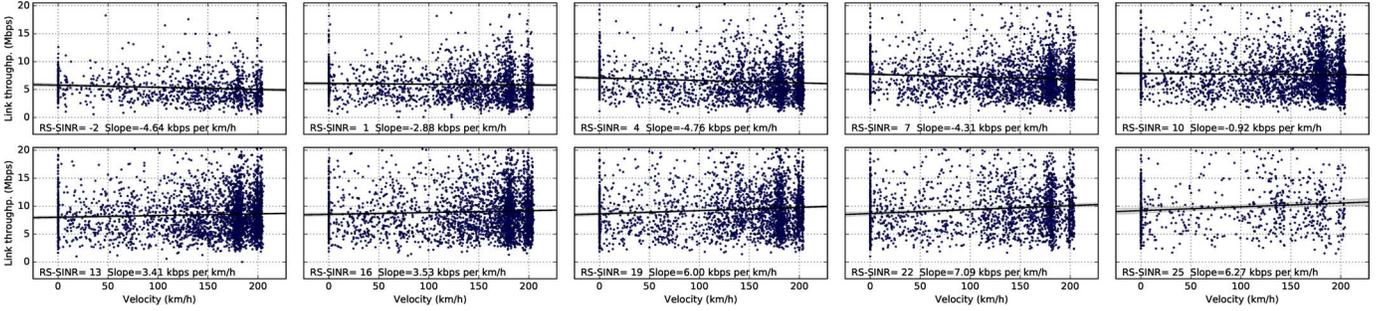


Figure 3: Scatter plots between link throughput and velocity per RS-SINR bin of 3 dB for one train

Table II: Relative metric importance for per link throughput

All metrics		One metric removed		Two metrics removed		Three metrics removed	
Metric	Weight	Metric	Weight	Metric	Weight	Metric	Weight
Velocity	0.30	Active devices	0.27	RS-SINR	0.31	RSRP	0.44
Active devices	0.18	RSRP	0.25	RSRP	0.27	RSSI	0.22
RSRP	0.16	RS-SINR	0.20	RSSI	0.19	RSRQ	0.16
RS-SINR	0.16	RSRQ	0.13	RSRQ	0.14	LTE Rx channel	0.12
RSRQ	0.09	RSSI	0.12	LTE Rx channel	0.05	Operator	0.03

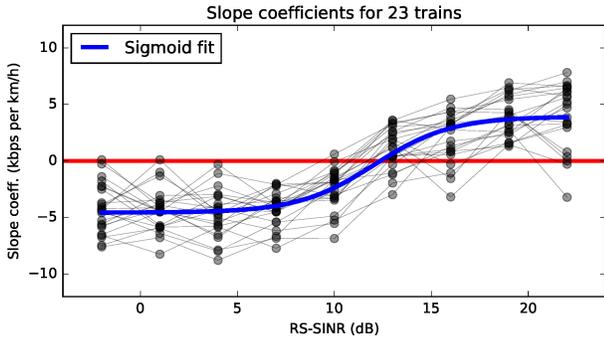


Figure 5: Slope coefficient for link throughput versus velocity over RS-SINR bins computed for 23 trains

the 23 trains which had a sufficient number of measurements. Also included in the figure is a fit of a sigmoid function in the form of a generalized logistic function:

$$S(R) = \frac{S_s - S_0}{1 + e^{-R_s*(R-R_0)}} + S_0 \quad (1)$$

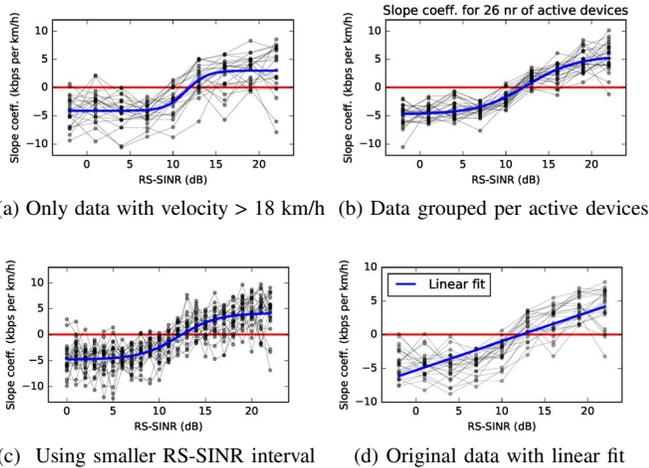
where  $S_x$  corresponds to slope related parameters, and  $R_x$  to RS-SINR related parameters. Fitting of the sigmoid model parameters  $S_0, S_s, R_0, R_s$  is done using non-linear Levenberg-Marquardt least-squares regression. The fitted model crosses the zero point of the slope coefficient at an RS-SINR slightly above 12 dB. This behavior is further examined in the next subsection.

### B. Relative metric importance for per link throughput

A further exploration of possible relationships between per link throughput and other collected metrics can also be done. As seen in Figure 5, there appears to be nonlinear interaction effects between metrics, and the relative similarity of several radio-related metrics likely results in collinearity and possibly

there are also multicollinearity. These issues along with the large spread observed in the measured values creates challenges when considering multiple linear regression as a means to learn structure from the data. Instead, a Machine Learning (ML) based approach is used to explore the relative importance of the observed metrics in relation to link throughput. There are several possible ML techniques that could be used for such task, and here we use random forests. Random forest regression is an ensemble-based ML technique that uses a large set of decision trees together with majority voting to obtain the predicted value. In this context, the random forest approach has several advantages as it is more robust, does not need metric value normalization, and provides metric importance data. The purpose here is not to build a ML model for prediction, but rather to use the relative importance values reported during the model training phase to provide insights on relationships between the metrics.

Using the same data set as in the previous subsection, and performing several iterative ML runs, yields the results shown in Table II. The leftmost subtable shows the overall metric ranking when all metrics were considered. The metric weights are normalized so that the sum of all weights is always one. In the leftmost subtable it can be observed that velocity has the highest weight, followed by number of active devices. This matches well to observations in the previous subsection. Here it can be noted that the random forest approach does not make any assumptions about the nature of the relationship of a particular metric and the target metric regarding linearity, etc. The interaction aspect of RS-SINR and velocity seems to not confuse the random forest metric ranking. To examine how a restriction of the velocity interval affects the behavior observed in Figure 5, a new set of calculations were performed where 45299 measurements in the velocity region 0-18 km/h were removed. In particular, this removes all measurements



(a) Only data with velocity > 18 km/h (b) Data grouped per active devices  
(c) Using smaller RS-SINR interval (d) Original data with linear fit

Figure 6: Variations on the computation of slope coefficient for link throughput versus velocity over RS-SINR bins

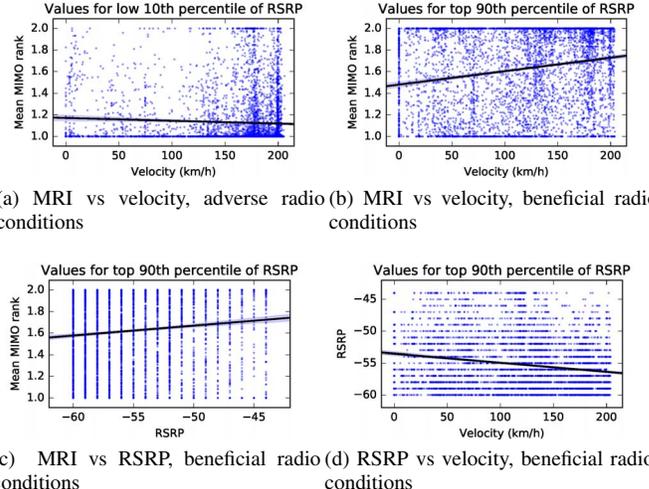
when the train is stationary at stations. The results are visible in Figure 6a. The behavior of the slope coefficients are still similar to Figure 5, although the spread is now larger.

The second subtable in Table II shows the metric ranking when the velocity metric is no longer included in the ML training data. Now the number of active devices is the highest weighted metric, although RSRP is very close. This seems to suggest that although the measurement data was selected to have the top decile with regards to number of active devices, there may be some remaining effect of the variations in offered load on the measurements. To explore the potential impact this may have, a new set of calculations was performed, this time with the data blocked according to the number of active devices instead of blocking per train. This data set comprised of 454230 measurements for 26 different active device numbers in the range 151-176 which had at least 10000 measurements each. Figure 6b shows that the slope coefficient behavior is similar also in this case, but the spread is now smaller as compared to Figure 5.

In the third subtable the active devices metric is removed, and it can be seen that the order of the remaining metrics have now changed so that RS-SINR is now before RSRP. To examine the dependence of the slope characteristics on the RS-SINR interval size, computations were also performed with tighter RS-SINR intervals. Instead of dividing the RS-SINR range into intervals 3 dB wide, the calculations shown in Figure 6c uses an interval of 1 dB. Again the trend is similar.

A computation was also performed to compare the fit of the sigmoid regression in Figure 5 versus the fit of a linear regression, which is shown in Figure 6d. For the sigmoid fit the Residual Sum of Squares (RSS) is 870, whereas the illustrated linear fit gives a worse RSS at 1148.

To summarize, the measurements indicate that the relationship between train velocity and per link throughput is dependent on the radio conditions, here in the form of RS-SINR, and that the nature of this relationship changes from



(a) MRI vs velocity, adverse radio conditions (b) MRI vs velocity, beneficial radio conditions  
(c) MRI vs RSRP, beneficial radio conditions (d) RSRP vs velocity, beneficial radio conditions

Figure 7: MIMO rank indication (MRI) relationships in second data set

negative to positive at a critical point of ca 12dB. This behavior is robustly present over a range of test cases. The reasons for this observed behavior is currently not conclusively understood although we suspect that MIMO-related effects could be at play.

### C. MIMO Rank Indication Examination

In addition to the results reported above we have also performed an evaluation of a second data set which had additional radio layer information. This data set could only be collected from one train equipped with updated hardware, which includes a Sierra MC7455 Modem. In particular, this setup allowed us to collect data on the mean MIMO rank indicator. The data set contains approximately 400,000 measurements. With the new information provided in this data set, further hints towards the underlying mechanisms can be explored. Here we focus on the MIMO rank indication (MRI). Although the evaluation here focuses on radio conditions rather than data throughput, similar filtering as before is used to select data points with the highest offered load. Thus, only data with more than 150 active devices are used in this examination. For these measurements we focus on the good and bad signal strength as captured by the lowest 10th percentile and the upper 90th percentile of RSRP values, respectively. This results in around 6000 data points in each percentile bin. The results for a regression of MRI over train velocity are shown in Figures 7a and 7b. While no strong trend is shown for the adverse radio conditions in Figure 7a, a strong positive trend is visible for the beneficial radio conditions in Figure 7b. Next, the relationship between reported RSRP values and the MRI is shown for the values of the 90th percentile in Figure 7c. A positive impact on MRI is seen as RSRP increases, which fits with the intuition. Finally, Figure 7d shows that the RSRP has a slight downward trend as train velocity increases. The trends observed by the regression in Figures 7c and 7d are

similar also when considering the complete data set rather than the deciles. Similar to the case for throughput versus velocity earlier, MRI versus velocity (Figures 7a and 7b) is almost flat when instead averaged over all radio conditions in the complete data set.

We hypothesize that the use of spatial multiplexing as reflected in the higher MIMO rank indication are coupled to the conditions between the train and the base station. For a train traveling along the track, the average distance between the train and the base station can be expected to be in order of 2 km, and average tower height be around 40+ meters [3]. Here it should be noted that even though the antennas for most of the time are in line-of-sight, they are actually not operating in free space. The antennas are transmitting over a large and diffusively scattering plane which is likely to create additional uncorrelated paths between the train and the base station antennas. As the train increases the velocity, the signals it receives will thus become more uncorrelated and the richer the scattering channel will become. The reason for this not occurring at low SNR is that under such conditions beamforming are typically used by the network, rather than spatial multiplexing. As the direction of arrival estimation is crucial for performance under such circumstances, increased train velocity could now have a negative impact. Improved MIMO performance with increasing train velocity was also observed in an earlier measurement campaign related to [2], although these observations were at that point considered measurement artifacts as there were no similar previous observations reported in the literature.

## V. CONCLUSIONS

In this paper we provide an examination on the relationship between per link throughput and train velocity based on analysis of large scale operational data. A noticeable critical point is observed for an RS-SINR around 12 dB above which the relationship between train velocity and observed throughput changes from being negative to becoming positive. A machine learning guided extended evaluation showed the robustness of the observed RS-SINR interaction effect, and the critical point location, for several variations of the data set.

Further examination of a separate train data set indicates that this phenomena can potentially be related to interaction between the scattering channel as seen by the train, and the networks choice of transmission modes for different radio condition information fed back from the UE.

For future work it is possible to extend this initial examination of aggregated data with a per-cell-focused study using a subset of the data and less aggregation. Additionally, we hope to complement the current passive measurements with additional active measurements which could provide further insights into the observed behaviors.

## ACKNOWLEDGMENTS

The authors wish to thank Mats Karlsson and Constanze Deiters for assisting with data collection and processing.

## REFERENCES

- [1] 3GPP. Evolved Universal Terrestrial Radio Access (E-UTRA); Physical layer - Measurements. TS 36.214, 3rd Generation Partnership Project (3GPP).
- [2] M. Alasali and C. Beckman. LTE MIMO performance measurements on trains. In *European conference on Antennas and Propagation (EuCAP)*. IEEE, 2013.
- [3] M. Alasali, C. Beckman, and M. Karlsson. Providing internet to trains using mimo in lte networks. In *2014 International Conference on Connected Vehicles and Expo (ICCVE)*, pages 810–814, Nov 2014.
- [4] J. Calle-Sánchez et al. Long term evolution in high speed railway environments: feasibility and challenges. *Bell Labs Technical Journal*, 18(2):237–253, 2013.
- [5] A. Elnashar and M. A. El-Saidny. Looking at LTE in practice: A performance analysis of the LTE system based on field test results. *IEEE Vehicular Technology Magazine*, 8(3):81–92, 2013.
- [6] J. Garcia, S. Alfreðsson, and A. Brunstrom. Examining Cellular Access Systems on Trains: Measurements and Change Detection. In *IEEE/IFIP Workshop on Mobile Network Measurement(MNM)*, 2017.
- [7] L. Li et al. A measurement study on TCP behaviors in HSPA+ networks on high-speed rails. In *IEEE INFOCOM*, 2015.
- [8] Y.-B. Lin, S.-N. Yang, and C.-T. Wu. Improving handover and drop-off performance on high-speed trains with multi-RAT. *Trans. Intelligent Transportation Systems*, 15(6), 2014.
- [9] A. Lutu et al. The Good, the Bad and the Implications of Profiling Mobile Broadband Coverage. *Computer Networks*, 2016.
- [10] É. Masson, M. Berbineau, and S. Lefebvre. Broadband Internet access on board high speed trains, A technological survey. In *Communication Technologies for Vehicles*. Springer, 2015.
- [11] R. Merz et al. Performance of LTE in a high-velocity environment: A measurement study. In *Proc. All Things Cellular (ATC)*. ACM, 2014.
- [12] M. K. Müller, M. Tarantetz, and M. Rupp. Providing current and future cellular services to high speed trains. *Communications Magazine, IEEE*, 53(10):96–101, 2015.
- [13] J. Rodriguez-Pineiro et al. LTE downlink performance in high speed trains. In *IEEE Vehicular Technology Conference (VTC Spring)*, 2015.