

# Adversarial Attacks on Deep Models for Financial Transaction Records

Ivan Fursov<sup>\*1</sup>, Matvey Morozov<sup>\*1</sup>, Nina Kaploukhaya<sup>1,2,4</sup>, Elizaveta Kovtun<sup>1</sup>, Rodrigo Rivera-Castro<sup>1</sup>, Gleb Gusev<sup>†3,4</sup>, Dmitry Babaev<sup>3</sup>, Ivan Kireev<sup>3</sup>, Alexey Zaytsev<sup>\* ‡ 1</sup>, Evgeny Burnaev<sup>1</sup>

<sup>1</sup>Skoltech, Moscow, Russia

<sup>2</sup>IITP RAS, Moscow, Russia

<sup>3</sup>Sber AI lab, Moscow, Russia

<sup>4</sup>MIPT, Moscow, Russia

## Abstract

Machine learning models using transaction records as inputs are popular among financial institutions. The most efficient models use deep-learning architectures similar to those in the NLP community, posing a challenge due to their tremendous number of parameters and limited robustness. In particular, deep-learning models are vulnerable to adversarial attacks: a little change in the input harms the model’s output.

In this work, we examine adversarial attacks on transaction records data and defences from these attacks. The transaction records data have a different structure than the canonical NLP or time series data, as neighbouring records are less connected than words in sentences, and each record consists of both discrete merchant code and continuous transaction amount. We consider a black-box attack scenario, where the attack doesn’t know the true decision model, and pay special attention to adding transaction tokens to the end of a sequence. These limitations provide more realistic scenario, previously unexplored in NLP world.

The proposed adversarial attacks and the respective defences demonstrate remarkable performance using relevant datasets from the financial industry. Our results show that a couple of generated transactions are sufficient to fool a deep-learning model. Further, we improve model robustness via adversarial training or separate adversarial examples detection. This work shows that embedding protection from adversarial attacks improves model robustness, allowing a wider adoption of deep models for transaction records in banking and finance.

## 1 Introduction

Adversarial attacks are a fundamental problem in deep learning. Alone a small targeted perturbation in the machine learning model’s input results in a wrong pre-

diction by the model [11]. These type of attacks are pervasive and present in various domains such as computer vision, natural language processing, event sequence processing, and graphs [4].

In this study, we consider a specific application domain of data models based on transaction records similar to Figure 1. This type of data arises in the financial industry, a sector that sees a vertiginous adoption of deep-learning models with a large number of parameters [3, 2]. Analysts use them to detect a credit default or fraud. As a result, deep-learning models on sequences of transactions stemming from a client or a group of clients are a natural combination. However, given their vast number of parameters, these models can be vulnerable to adversarial attacks [11]. Hence, risk mitigation related to the possibility of adversarial attacks becomes paramount. Despite the great attention of both academy and industry to adversarial attacks in machine learning, there are still no papers studying this problem in the finance domain. We must make these models robust to such attacks if we want to see them widely adopted by the industry.



Figure 1: Example of a transaction records data sequence: Each transaction has information about its MCC (Merchant Category Codes) describing the purchase category and purchase amount. Data can also include transaction location and time values.

Transaction records data consists of sequences, making it possible to transfer some techniques from the domain of natural language processing [30]. However, there

\*Equal contribution

†E-mail: gleb57@gmail.com

‡E-mail: a.zaytsev@skoltech.ru

are several peculiarities related to the generation of adversarial sequences of transactions:

1. Unlike neighbouring words in natural language, there may be no logical connection between neighbouring transactions in a sequence.
2. An attacker cannot insert a token at an arbitrary place in a sequence of transactions in most cases. Often, the attacker can only add new transactions to the end of a sequence.
3. The concept of semantic similarity is underdeveloped when we compare it to what we see in the NLP world [29].
4. Transaction data is complex. Besides MCC, which is categorical, we also have the transaction amount, and can also have its location, timestamp, and other fields. It is particularly an essential question if adversarial transactions' amount affects an attack's success.

To make an attack scenario realistic, we consider inference-time black-box attacks, i.e. an attacker has access only to a substitute model different from the target attacked model and trained using a different dataset and attacks a model, when it is used in production. We compare this scenario to a more dangerous but less realistic white-box attack. This setting is close to works such as [9] for NLP and [22] for CV models.

## 1.1 Novelty

This work addresses three research questions: (1) The vulnerability of finance models to different attack methods in a realistic scenario, (2) the effectiveness of adversarial training, (3) the role of the amount of adversarial transactions. We propose several black-box adversarial attacks for transaction records data for realistic scenarios and provide defences from these attacks accounting for the above challenges. Our black-box attacks change or *add* only a couple of tokens with minimal money amounts for each token to fool the model. We adopt the idea of loss gradients and a substitute model to create effective adversarial attacks using a separately trained generative model for transactions models. The scheme of our attack is in Figure 2.

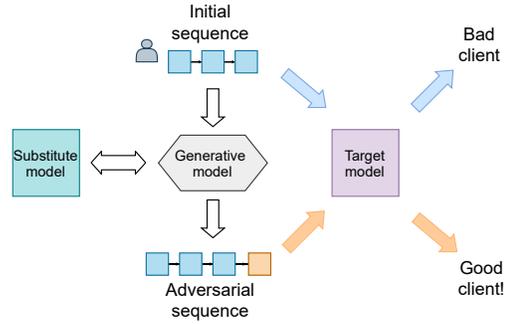


Figure 2: Proposed attack scheme: we consider a realistic black box scenario, when a malicious model user has access to a separate substitute model during an attack on a target model. To fool a target model, he adds new transactions to the end of a sequence making additional purchases. To select new sequence, the attacker adopts a conditioned generative model.

Our defences make most of these attacks ineffective besides an attack based on sampling from a language model. Banks and other financial actors can use our results to improve their fraud detection and scoring systems.

The main claims and structure of the paper are as follows:

- We address an important and practical research questions: the vulnerability of deep finance models to adversarial attacks in the most realistic black-box scenarios, and the effectiveness of adversarial learning;
- We develop black-box adversarial attacks for transaction records data from the financial industry and defences from these attacks in the light of the above challenges. Our approach adopts the idea of loss gradients and a substitute model to create effective adversarial attacks based on a generative model that takes into account mixed discrete-continuous structure of bank transactions data. [subsection 3.1](#) describes our attacks and defences.
- Our black-box attacks change or add only a couple of tokens with small monetary values for each token to fool the model. We show that in [subsection 5.2](#).
- Our defences render most of these attacks ineffective except for an attack based on sampling from a language model. We present it in [subsection 5.4](#).
- We conduct experiments on relevant datasets from the industry and provide an investigation on the effectiveness of such attacks and how one can defend her model from such attacks. Banks can use the obtained results to improve their fraud detection and scoring systems. We provide a discussion on our attacks and defences in [subsection 5.6](#).

## 2 Related work

Adversarial attacks exist for most deep-learning models in domains as wide as, [25], computer vision [1, 13], natural language [30, 23], and graph data [21]. A successful attack generates adversarial examples that are close to an initial input but result in a misclassification by a deep-learning model [14]. [24, 20] provides a comprehensive review. For sequential data and especially for NLP data, there are also numerous advances in recent years. [30, 23] discuss them. Similarly, the work of [28] emphasizes adversarial attacks in financial applications.

Our study uses transaction records data. As a consequence, we generate adversarial examples for discrete sequences of tokens. For NLP data, [19] is one of the first approaches to deal with this challenge. The authors consider white-box attacks for an RNN model. In our work, we focus on black-box attacks. We consider this scenario more realistic than the white-box one. Along the same lines [10] also consider black-box attacks for discrete sequences. They present an algorithm, DeepWordBug, to identify the most relevant tokens for the classifier whose modifications lead to wrong predictions.

Popular methods for generating adversarial sequence use the gradient of the loss function to change the initial sequences. The first practical and fast algorithm based on the gradient method is the Fast Gradient Sign Method (FGSM) [11]. One attack inspired by the FGSM is the high-performing BERT-attack [15]. The attack uses a pre-trained language model BERT to replace crucial tokens with semantically consistent ones.

The existence of adversarial attacks implies the necessity of developing defence methods. [26] provides an overview, for example. One of the most popular approaches for defence is an adversarial training advocated in [11, 12]. The idea of adversarial training is to expand a training sample with adversarial examples equipped with correct labels and fine-tune a deep-learning model with an expanded training sample. Similarly, in [27], the authors improve a fraud detection system by adding fraudulent transactions in the model’s training set. Another approach is adversarial detection. In this case, we train a separate detector model. [18] describe the detector of adversarial examples as a supplementary network to the initial neural network. Meanwhile [7] provide a different approach to detecting adversarial examples based on Bayesian theory. They detect adversarial examples using a Bayesian neural network and estimating the uncertainty of the input data.

Adversarial attacks and defences from them are of crucial importance in financial applications [28]. The literature of adversarial attacks on transaction records includes [8, 9]. However, the authors do not consider the peculiarities of transaction records data and apply general approaches to adversarial attacks on discrete sequence data. Also, they pay little attention to efficient defences from such attacks and how realistically one can obtain such sequences.

There are no comprehensive investigations on the reliability of deep-learning models aimed at processing se-

quential data from transaction records in the literature. It stems from the lack of attention to the transaction records data’s peculiarities and the limited number of problems and datasets.

This study fills this void by being the first encompassing work on adversarial attacks on transaction records data and corresponding defences for diverse datasets while taking data peculiarities into account. According to our knowledge, no one before explored attacks with additions of tokens to the end and training specific generative model for financial transactions data that is a sequence of a mix of discrete and continuous variables. Moreover, as the data is interpretable, it is vital to consider, what factors affect vulnerability of models to adversarial attacks.

## 3 Methods

**Blackbox vs whitebox scenario** We focus on the most realistic blackbox scenario [22]. The attack in this scenario does not have access to the original model. Instead, the attacker can use scores for a substitute model trained with a different subset of data. A substitute model can have different hyperparameters or even a different architecture replicating real-life situations. The attacker has little knowledge about the model’s real architecture or parameters but can collect a training data sample to train a model.

### 3.1 Attack methods

We consider realistic attack scenarios in the financial sector. For this reason, our attacks resemble existing approaches in other communities. We adapt them to the specificity of transaction records data. We categorize the algorithms by attack method into two types. In the first case, an attack edits an original sequence of transactions by **replacing existing tokens**. We present the scheme of such attack in Figure 3. For the second case, an attack **concatenates new adversarial transactions to the end** of the original sequence. For such attacks, we use the *Concat* prefix in the name. The scheme of such attack is in Figure 4.

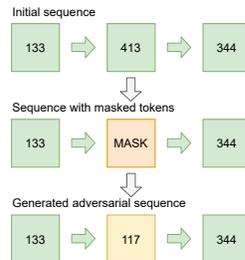


Figure 3: Example of an attack with token replacements where we select tokens to replace and then find the most effective tokens for a replacement to obtain an adversarial sequence

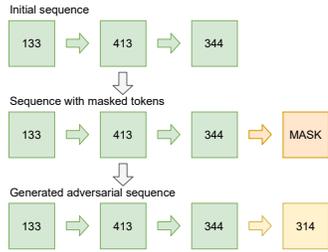


Figure 4: Example of an attack with token additions where we add the most adversarial tokens to the end of an adversarial sequence

**Replacement and concatenation attacks Sampling Fool (SF)** [9]. It uses a trained Masked Language Model (MLM) [5], to sample from a categorical distribution new tokens to replace random masked tokens. Then it selects the most adversarial example among a fixed number of generated sequences. We generate 100 sequences in our experiments.

**FGSM (Fast Gradient Sign Method)** [16]. The attack selects a random token in a sequence and uses the Fast Gradient Sign Method to perturb its embedding vector. The algorithm then finds the closest vector in an embedding matrix and replaces the chosen token with the one that corresponds to the identified vector. We show the FGSM attack in an embedded space of tokens in Figure 5. If we want to replace several tokens, we do it greedily by replacing one token at a time. We select the amounts from a categorical distribution given the overall limit for these attacks if not specified otherwise.

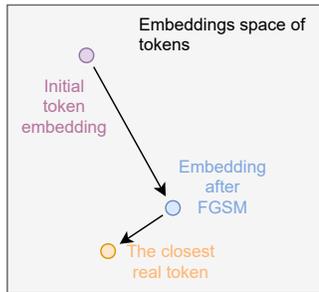


Figure 5: FGSM attack scheme in an embedded space of tokens

**Concat FGSM.** This algorithm extends the FGSM idea. Here, we add  $k$  random transactions to the original sequence and run the FGSM algorithm only on the added transactions. As a result, the edit-distance between the adversarial example and the original will be exactly  $k$ . We use two variations of this approach. They are Concat FGSM, [seq], and Concat FGSM, [sim]. In the first option, sequential, we add tokens one by one at the end of the sequence. In the second, simultaneous, we start by adding a random number of mask tokens and then replacing them by FGSM simultaneously taking into account the context. We show the difference between both approaches in Figure 6.

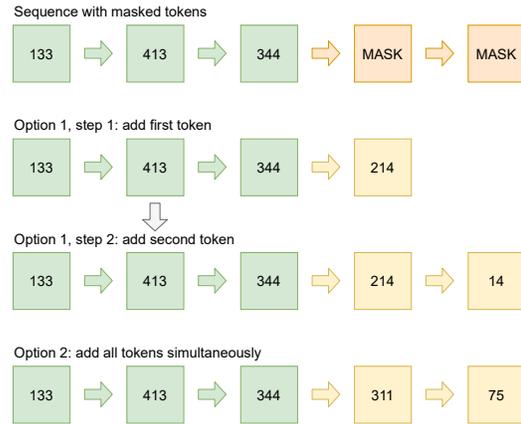


Figure 6: Difference between sequential (Option 1) and simultaneous (Option 2) Concat FGSM attacks

**LM FGSM.** To increase the plausibility of generated adversarial examples, LM FGSM uses a trained autoregressive language model to estimate which transaction is the most appropriate in a given context. The algorithm uses FGSM to find the closest token with the perplexity less than a predefined threshold  $\tau$ . Thus, LM FGSM increase the plausibility of generated adversarial examples. A similar idea appears in [6] for NLP data.

**Concat Sampling Fool (Concat SF).** This method replicates the idea of Concat FGSM and adds  $k$  random transaction tokens to the original sequence. Then, using the pre-trained BERT language model, we receive the vocabulary probabilities for each added token. As a result, the algorithm samples from the categorical distribution and chooses the best adversarial example using a given temperature value.

**Sequential Concat Sampling Fool (Seq Concat SF).** The difference between this algorithm and the usual Concat SF is that this algorithm, instead of immediately adding  $k$  tokens to the original sequence, adds tokens one at a time. We choose the token at each step that reduces the model’s probability score for the target (correct) class the most.

### 3.2 Defense methods

We use two approaches to defend from adversarial attacks with ideas from the general defences literature [23]. They are adversarial training and adversarial detection.

**Adversarial training.** The adversarial training objective is to increase the model robustness by adding adversarial examples into the training set. The most popular and common approach is adding the adversarial examples with correct labels to the training set, so the trained model correctly predicts the label of future adversarial examples after fine-tuning. Other popular adversarial training methods generate adversarial examples at each training iteration [17]. From the obtained examples, we calculate the gradients of the model and do a backpropagation.

**Adversarial detection.** This method consists of training an additional model, a discriminator, which solves a binary classification model addressing whether

a present sequence is real or generated by an adversarial attack. The trained discriminator model can detect adversarial examples and prevent possible harm from them.

## 4 Data

We consider sequences of transaction records stemming from clients at financial institutions. For each client, we have a sequence of transactions. Each transaction is a triplet consisting of three elements. The first one is a discrete code to identify the type of transaction (Merchant Category Code). Among others, it can be a drug store, ATM or grocery store. The second corresponds to a transaction amount, and the third is a timestamp for the transaction time. For the transaction codes and the formatting of the data, we take [2] as an inspiration.

Moreover, we consider diverse real-world datasets for our numerical experiments taken from three major financial institutions. In all cases, the input consists of transactions from individual clients. For each client, we have an ordered sequence of transactions. Each entry contains a category identifier such as a bookstore, ATM, drugstore, and more as well as a monetary value. Further, the first and second sample’s target output, *Age 1* and *Age 2*, is the client’s age bin with four different bins in total. For the third sample, *Leaving*, the goal is to predict a label corresponding to a potential customer churn. The fourth sample, *Scoring*, contains a client default for consumer credit as a target. We present the main statistics for both the data and the models we use in Table 1. For our critical experiments, we provide results for all datasets. In other cases, due to space constraints, we focus on one or two datasets, as results are similar if not specified otherwise.

We train the target model on 65% of the original data. Similarly, we use for the substitute model the rest of the data, 35%. Additional implementation details are in Appendix.

## 5 Experiments

We provide the code, experiments, running examples and links to the transaction datasets on a public online repository\*.

### 5.1 Metrics

To measure the quality of attacks, we find an adversarial example for each example from the test sample. We use the obtained adversarial examples to calculate a set of metrics and assess the attacks’ quality. We then average over all considered examples. As we compare multiple adversarial attacks, we provide a unified notation to facilitate their understanding. We consider the following elements in our notation:

- $C^t$  — The target classifier that we want to ”deceive” during the attack;
- $\mathbf{x}_i$  — An original example from the test sample.
- $\tilde{\mathbf{x}}_i$  — An adversarial example that we obtain during the attack for an example  $\mathbf{x}_i$ .

- $C^t(\mathbf{x}_i)$  – The prediction of the original classifier for the object  $\mathbf{x}_i$ , a class label.
- $C^t(\tilde{\mathbf{x}}_i)$  – The prediction of the original classifier for the adversarial example  $\tilde{\mathbf{x}}_i$ , a class label.

With this notation, we can define a list of four metrics to evaluate our data against them.

1. **Word Error Rate (WER)**. Before we can create an adversarial attack, we must change the initial sequence. Our change can be either by inserting, deleting, or replacing a token in some position in the original sequence. In the *WER* calculation, we treat any change to the sequence made by insertion, deletion, or replacement as 1. Therefore, we consider the adversarial sequence perfect if  $WER = 1$  and the target classifier output has changed.
2. **Adversarial Accuracy (AA)**. AA is the rate of examples that the target model classifies correctly after an adversarial attack. We use the following formula to calculate it:  $AA(A) = \frac{1}{N} \sum_{i=1}^N \mathbb{1}\{C^t(\mathbf{x}_i) = C^t(\tilde{\mathbf{x}}_i)\}$ .
3. **Probability Difference (PD)**. The PD demonstrates how our model’s confidence in the response changes after an adversarial attack:  $PD(A) = \frac{1}{N} \sum_{i=1}^N C_{\mathbb{P}}^t(\mathbf{x}_i) - C_{\mathbb{P}}^t(\tilde{\mathbf{x}}_i)$ , where  $C_{\mathbb{P}}^t(\mathbf{x}_i)$  is the confidence in the response for the original example  $\mathbf{x}_i$ , and  $C_{\mathbb{P}}^t(\tilde{\mathbf{x}}_i)$  is the confidence in the response for the adversarial example  $\tilde{\mathbf{x}}_i$ .
4. **NAD (Normalized Accuracy Drop)** NAD is a probability drop, normalized on WER. For a classification task, we calculate the Normalized Accuracy Drop in the following way:

$$NAD(A) = \frac{1}{N} \sum_{i=1}^N \frac{\mathbb{1}\{C^t(\mathbf{x}_i) \neq C^t(\tilde{\mathbf{x}}_i)\}}{WER(\mathbf{x}_i, \tilde{\mathbf{x}}_i)},$$

NAD reflects a trade-off between the number of changed tokens and model quality on adversarial examples.

### 5.2 Overall attack quality

We compare the performance of six attacks. They are Sampling Fool (SF), Concat Sampling Fool (Concat SF), Sequential Concat Sampling Fool (Seq Concat SF), FGSM, Concat FGSM, Concat LM FGSM and we describe them in subsection 3.1. The results are in Table 2 for the GRU target and the substitute model architectures. Similarly, in Table 3, we have them for the sake of comparison for the white-box scenario, when we use the target model as the substitute for the most promising and representative attacks.

We observe that all adversarial attacks have a high quality. After we change on average no more than 2 tokens, we can fool a model with a probability of at least

\*[https://github.com/fursovia/adversarial\\_sber](https://github.com/fursovia/adversarial_sber)

Dataset		Age 1	Age 2	Client leaving	Scoring
#classes		4	4	2	2
Mean sequence length		86.0	42.5	48.4	268.0
Max sequence length		148	148	148	800
Number of transactions		2 943 885	3 951 115	490 513	3 117 026
Number of clients		29 973	6886	5000	332 216
Number of code tokens		187	409	344	400
Models' accuracy for the test data					
LSTM	Substitute	0.537	0.675	0.653	0.825
	Target	0.551	0.668	0.663	0.881
GRU	Substitute	0.548	0.675	0.644	0.833
	Target	0.562	0.663	0.672	0.862
CNN	Substitute	0.553	0.674	0.660	0.872
	Target	0.564	0.674	0.670	0.903

Table 1: Datasets statistics and models' accuracy for the test data

0.8. This probability makes the attacked model useless, while the number of added or modified transactions is small compared to the total number of transactions in the sequence judging by Word Error Rate. We can completely break the model for some attacks by adding only one token at the end. Attacks, where we add tokens at the end, perform comparably to attacks, where we include tokens in the middle of a sequence. So, these more realistic attacks are also powerful for money transactions data. The quality of attacks is comparable to the quality of the Greedy baseline attack. It is the Greedy attack based on a brute force selection of tokens for insertion or editing. Thus, the attack provides close to the best achievable performance in our black-box scenario. However, FGSM-based attacks provide better performance scores than SamplingFool-based ones due to the random search for adversarial examples in the second case. It can be useful to unite these approaches to create a generative model that can generate sequences of transaction records that are both realistic and adversarial [8]. Also, more realistic concatenation of tokens to the end of a sequence results in lower performance scores. In sections below we consider the most successful and the most representative attacks SF, Concat SF, FGSM, and Concat FGSM.

### 5.3 Dependence on the architecture

A malicious user can lack knowledge regarding the specific architecture used for data processing. Hence, we want to address how the attack's quality changes if the architectures of both the attacked target and the substitute model used for the attack differ. For this, we present the results when the attack targets an LSTM model and a CNN the substitute model in Table 4. Comparing these results to the previous section results, we see that the attack's quality in the case of models of different nature deteriorates markedly. However, in both cases, the Concat FGSM attack works reasonably by adding strongly adversarial tokens at the end of the sequence. We conjecture that this is because RNN-based models pay more attention to the tokens in the end, while a CNN model is position-agnostic regarding a particular

token. So, Concat FGSM attacks are successful, even if a substitute model's architecture is different from that of a target model.

## 5.4 Defenses from adversarial attacks

### 5.4.1 Adversarial training

Figure 7 presents the adversarial training results for one dataset. The target and substitute models are GRU. In the figure, we average the results from 10 runs and present mean values as solid curves and a mean  $\pm$  with two standard deviations as a shaded area. For most attacks, the Adversarial Accuracy quickly increases when we compare it against the overall accuracy presented as the dashed line. The results for other datasets are similar. For most attacks, we see that it is enough to add about 15000 adversarial examples to the training sample and fine-tune model using an expanded training sample to make model robust against a particular adversarial attack. However, after adversarial training, a SamplingFool attack is possible despite being worse on our attack quality metrics.

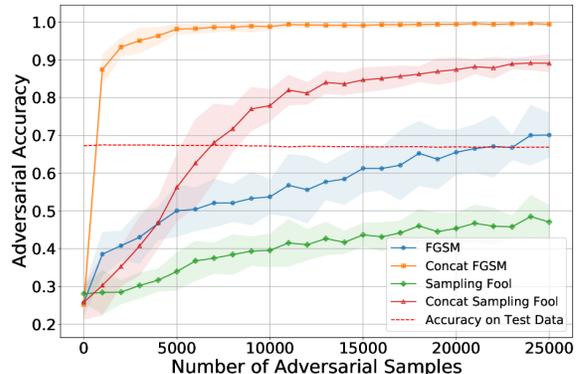


Figure 7: Adversarial Accuracy metric vs number of adversarial examples added to the training set for dataset Age 2.

Attack	NAD $\uparrow$	WER $\downarrow$	AA $\downarrow$	PD $\uparrow$
Age 1 (Accuracy 0.562)				
SF	0.09	12.87	<b>0.45</b>	0.07
Concat SF	0.26	2	0.49	0.06
Seq Concat SF	0.28	2	<b>0.45</b>	0.09
FGSM	<b>0.45</b>	1.79	0.49	0.04
Concat FGSM	0.13	4	0.47	0.07
LM FGSM	0.44	<b>1.78</b>	0.50	0.04
Age 2 (Accuracy 0.663)				
SF	0.45	5.18	0.29	0.07
Concat SF	0.38	<b>2</b>	0.25	0.07
Seq Concat SF	0.39	<b>2</b>	<b>0.12</b>	0.09
FGSM	<b>0.71</b>	3.46	0.21	0.08
Concat FGSM	0.2	4.0	0.19	0.07
LM FGSM	<b>0.71</b>	3.45	0.20	0.08
Client leaving (Accuracy 0.672)				
SF	0.14	14.78	0.69	0.00
Concat SF	0.19	<b>2</b>	0.62	0.00
Seq Concat SF	0.23	<b>2</b>	0.55	0.03
FGSM	0.23	8.69	0.47	0.10
Concat FGSM	0.14	4	0.46	0.09
LM FGSM	<b>0.24</b>	8.84	<b>0.44</b>	0.11
Scoring (Accuracy 0.86)				
SF	0.15	<b>5.74</b>	0.78	0.00
Concat SF	0.05	6	0.73	0.05
FGSM	<b>0.27</b>	8.43	0.67	0.10
Concat FGSM	0.10	6	<b>0.44</b>	0.22
LM FGSM	<b>0.27</b>	7.76	0.67	0.09

Table 2: Summary of the effectiveness of black-box attacks with GRU as the target and the substitute models architectures for different datasets with initial accuracy values in brackets.  $\uparrow$  marks metrics we want to maximize,  $\downarrow$  marks metrics we want to minimize.

#### 5.4.2 Adversarial detection

Another approach for defending against adversarial attacks is the adversarial sequences detection. Carrying out this type of attack entails that we train a separate classifier detecting a sequence from an original data distribution or stemming from an adversarial attack. If this classifier is of high quality, we can easily protect our model from adversarial attacks by discarding potentially adversarial sequences.

As the detector classifier, we use a GRU model. We train it for 10 epochs. We notice that increasing the number of epochs does not lead to improving the detector’s quality.

The classification data is balanced, so we calculate the accuracy for evaluating the quality of adversarial detection. These results for different attacks are in [Table 5](#).

We can detect with accuracy greater than 0.9 the adversarial examples generated by most of the attacks. However, adversarial sequences from the Sampling Fool attack are more difficult to detect. Therefore, we postulate that we can effectively repel many attacks with the

Attack	NAD $\uparrow$	WER $\downarrow$	AA $\downarrow$	PD $\uparrow$
Age 1				
SF	0.13	13.18	0.07	0.23
Concat SF	0.38	<b>2</b>	0.23	0.24
FGSM	<b>0.70</b>	2.51	<b>0.00</b>	0.29
Concat FGSM	0.25	4	<b>0.00</b>	0.48
LM FGSM	0.68	2.68	0.01	0.27
Age 2				
SF	0.56	4.81	0.22	0.09
Concat SF	0.42	<b>2</b>	0.15	0.12
FGSM	0.79	2.40	<b>0.01</b>	0.15
Concat FGSM	0.24	4	0.05	0.25
LM FGSM	<b>0.84</b>	2.13	0.03	0.14
Client leaving				
SF	0.29	14.37	0.50	0.09
Concat SF	0.23	<b>2</b>	0.53	0.06
FGSM	<b>0.42</b>	6.51	<b>0.09</b>	0.15
Concat FGSM	0.16	4	0.35	0.18
LM FGSM	<b>0.42</b>	6.37	0.11	0.15
Scoring				
SF	0.15	<b>5.74</b>	0.79	0.15
Concat SF	0.03	6	0.80	0.22
FGSM	<b>0.47</b>	8.43	<b>0.02</b>	0.28
Concat FGSM	0.14	6	0.14	0.43
LM FGSM	0.46	7.76	0.04	0.28

Table 3: Summary of the effectiveness of the most promising white-box attacks with GRU as the target model architecture

help of a trained detector.

### 5.5 Dependence on additional amount

In addition to a discrete token, each transaction has an associated money amount. A strong dependence on an attack’s performance based on using additional monetary amounts can be a show stopper, as a malicious user tries to minimize the amount of money spent on additional transactions.

We experiment on the Age 1 dataset by measuring the Concat FGSM attack’s quality with varying constraints starting from few constraints and going to the case where there is practically no limit.

In [Table 6](#), we present the results. For this table, we consider the Age 1 dataset and add one token at the end of a transaction sequence, so WER equals 1 for all attacks. They demonstrate no difference in the attack’s quality concerning the considered amount. Therefore, we can fool the model using a quite limited amount of monetary funds. We elucidate that the model takes most of its information from the transaction token. As a result, attack strategies should focus on it.

Attack	NAD $\uparrow$	WER $\downarrow$	AA $\downarrow$	PD $\uparrow$
Age 1				
SF	0.10	13.75	0.45	0.12
Concat SF	0.30	2	0.40	0.14
FGSM	0.40	3.81	0.48	0.10
Concat FGSM	0.18	4	0.30	0.20
LM FGSM	0.43	3.53	0.45	0.10
Age 2				
SF	0.21	14.92	0.78	-0.10
Concat SF	0.09	2	0.82	-0.81
FGSM	0.21	8.90	0.78	-0.13
Concat FGSM	0.05	4	0.80	-0.10
LM FGSM	0.58	3.37	0.38	0.02
Client leaving				
SF	0.09	12.17	0.46	0.03
Concat SF	0.25	2	0.50	0.04
FGSM	0.46	6.67	0.30	0.10
Concat FGSM	0.14	4	0.43	0.08
LM FGSM	0.21	12.53	0.55	0.07
Scoring				
SF	0.10	6.31	0.84	0.05
Concat SF	0.06	6	0.65	0.18
FGSM	0.17	13.92	0.76	0.12
Concat FGSM	0.06	6	0.64	0.18

Table 4: Black-box attack effectiveness with LSTM as the target model and CNN as the substitute model

Attack	Age 1	Age 2	Client leaving	Scoring
SF	<b>0.500</b>	<b>0.623</b>	<b>0.661</b>	<b>0.422</b>
Concat SF	0.998	0.986	0.967	0.988
FGSM	0.493	0.946	0.800	0.953
Concat FGSM	1.000	0.995	0.991	1.000

Table 5: Accuracy of adversarial examples detection for different types of attacks and datasets

## 5.6 Closer look at generated sequences

Natural language data allow for a straightforward interpretation by experts. A human can assess how realistic an adversarial sentence is. For transaction records data, we can manually inspect any generated adversarial sequences and verify their realism level.

The histograms of the initial distribution of tokens and tokens inserted by the Sampling Fool attack are in Figure 8. We notice that most of the inserted tokens occur in the original sequences with similar frequencies, so the generated sequences are realistic from this point of view. However, some of the inserted tokens do not belong to the history of the client’s transactions. Nevertheless, they are in the training set of the model. This constellation also occurs in the case of the histogram of the Concat Sampling Fool.

We can observe different types of distributions for the Concat FGSM in Figure 9. For this attack, only a limited number of tokens has a significant effect on the model prediction. The same phenomenon happens for the cases of the FGSM attacks. So, we expect, that while these attacks are adequate given our metrics, it

remains easy to detect them or fine-tune the model to make it resistant to this type of attacks.

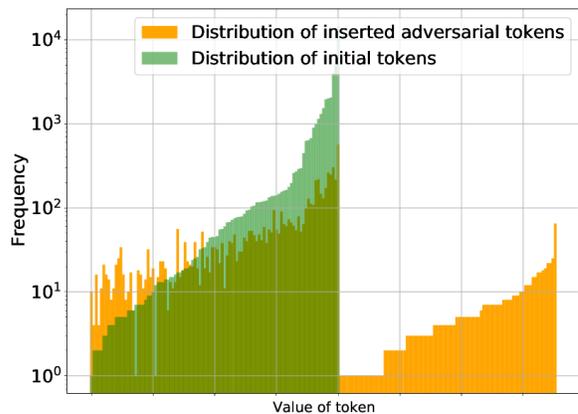


Figure 8: The histograms of tokens from the initial transactions are in green and tokens inserted by the Sampling Fool attack in orange for the Age 1 dataset

Amount limit	AA	PD
300	0.09	0.23
500	0.12	0.20
1000	0.2	0.17
3000	0.05	0.27
5000	0.01	0.31
10000	0.01	0.35
100000	0.04	0.31

Table 6: Overview of the Concat FGSM attack quality for different amounts of tokens for Age 1 dataset

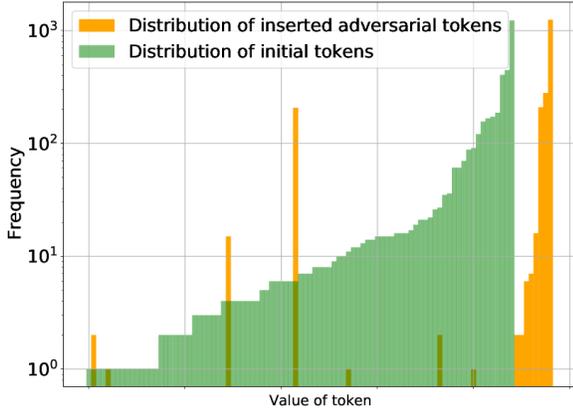


Figure 9: The histograms of tokens from the initial transactions are in green and tokens inserted by Concat FGSM attack in orange for the Age 1 dataset

We can summarize numerically the diversity of tokens added to the sequences through the considered attacks. Three of our metrics are a diversity rate, a repetition rate, and a perplexity. For the first case, **the diversity rate** represents the number of unique adversarial tokens divided by the size of the vocabulary from which we sample these tokens. High diversity values suggest that an attacker uses a vast number of unique tokens and the generated adversarial sequences are harder to distinguish from sequences present in the original data. In contrast, low diversity values suggest that an attack inserts the limited number of specific tokens. For the second case, **the repetition rate** equals the number of tokens added by an attack appearing in an original sequence divided by the total number of added adversarial tokens. High repetition rate means that an attack adds tokens that are already present in a sequence, thus the generated sequence is more realistic, while low repetition rate means that an attacker tries to insert new types of tokens, and that can be unrealistic. Finally, **Perplexity** is the inverse normalized probability density that the specified sequence  $\mathbf{x} = (x_1, x_2, \dots, x_n)$  is derived from the specified distribution. We define it as  $PP(\mathbf{x}) = \mathbb{P}(x_1, x_2, \dots, x_n)^{-\frac{1}{n}} = \sqrt[n]{\frac{1}{\mathbb{P}(x_1, x_2, \dots, x_n)}}$ . We calculate the density using a trained language model. The lower the perplexity, the more plausible the sequence is from the language model’s point of view.

We observe the obtained diversity rate, repetition rate, and perplexity in Table 7 for Age 1 dataset. For other datasets, the results are similar. As the baselines we perform uniform sampling of tokens from the training set (Unif. Rand.) and sampling w.r.t. their frequencies of the occurrence (Distr. Rand.). We see that SF and Concat SF’s attacks for the language model generate more realistic adversarial sequences based on the considered criterion, while FGSM-based approaches and the Greedy attack do not pay attention to how realistic the generated sequences are.

In addition to the added tokens’ realism, we consider the added amounts’ plausibility for those tokens. As we show, the transaction’s value has little effect on the attack’s success, so we only need certainty that the monetary value is realistic and not excessively large. We can find the histograms for the original and the Sampling Fool adversarial sequences in Figure 10. We see that the adversarial sequences meet both of the specified requirements. For other datasets and attacks, the figures are similar.

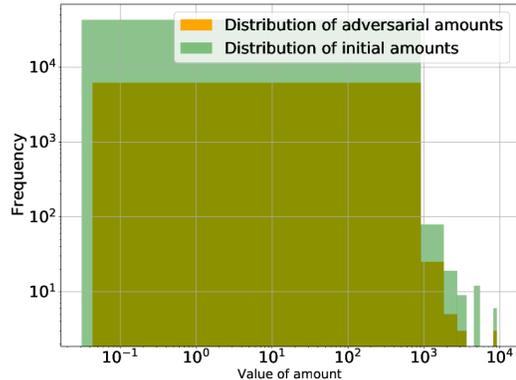


Figure 10: The histograms of amounts for original transactions and amounts for transactions inserted by the Sampling Fool attack for a balanced Age 1 dataset of adversarial and initial normal sequences

## 6 Conclusion

We consider an essential topic of adversarial attacks on machine-learning models for a highly relevant use case, transactions record data, and present defences against such attacks.

Data	Attack	Diversity $\uparrow$ rate	Repetition $\uparrow$ rate	Perplexity $\downarrow$
Age 1	Unif. Rand.	0.99	0.47	7.49
	Distr. Rand.	0.37	0.99	3.02
	SF	0.62	<b>0.80</b>	<b>4.28</b>
	Concat SF	<b>0.77</b>	0.58	6.01
	FGSM	0.15	0.57	4.71
	Concat FGSM	0.08	0.11	17.76

Table 7: Diversity rate, repetition rates and perplexity of the adversarial sequences generated by the attacks

We find out that alone by inserting a couple of transactions in a sequence of transactions, it is possible to fool a model. Even small model accuracy drops can lead to significant financial losses to an organization, while for most datasets we observe a severe model quality decrease. Moreover, attack are the most effective for important borderline decision cases. However, we can easily repel most classic attacks. It is straightforward to detect an adversarial sequence or to fine-tune a model to process adversarial sequences correctly.

Attacks stemming from generative data models represent the most promising ones. It is challenging to defend against them, and they show high efficiency. During our experiments, a manual assessment of the adversarial sequences show that they are realistic, and thus experts must be careful whenever deploying deep-learning models for processing transaction records data. Nevertheless, we still can detect adversarial sequences in most realistic scenarios, when an attack appends new transactions only to the end of a sequence.

We expect, that our result shades a light on how attacks and defences should look like for models based on transactions record data. They can be extended to more general event sequences data.

## References

- [1] Naveed Akhtar and Ajmal Mian. 2018. Threat of adversarial attacks on deep learning in computer vision: a survey. *IEEE Access* 6 (2018), 14410–14430.
- [2] D. Babaev, I. Kireev, N. Ovsov, M. Ivanova, G. Gusev, and A. Tuzhilin. 2020. Event sequence metric learning. *arXiv:2002.08232* (2020).
- [3] Dmitrii Babaev, Maxim Savchenko, Alexander Tuzhilin, and Dmitrii Umerenkov. 2019. ET-RNN: Applying Deep Learning to Credit Loan Applications. In *25th ACM SIGKDD*. 2183–2190.
- [4] Anirban Chakraborty, Manaar Alam, Vishal Dey, Anupam Chattopadhyay, and Debdeep Mukhopadhyay. 2018. Adversarial attacks and defences: A survey. *preprint arXiv:1810.00069* (2018).
- [5] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. BERT: pre-training of deep bidirectional transformers for language understanding. *preprint arXiv:1810.04805* (2018).
- [6] J. Ebrahim, A. Rao, D. Lowd, , and D. Dou. 2018. Hotflip: white-box adversarial examples for text classification. *arXiv:1712.06751* (2018).
- [7] R. Feinman, R. R. Curtin, S. Shintre, and A. B. Gardner. 2017. Detecting adversarial samples from artifacts. *preprint arXiv:1703.00410* (2017).
- [8] I. Fursov, A. Zaytsev, N. Kluchnikov, A. Kravchenko, and E. Burnaev. 2020. Differentiable language model adversarial attacks on categorical sequence classifiers. *arXiv:2006.11078* (2020).
- [9] Ivan Fursov, Alexey Zaytsev, Nikita Kluchnikov, Andrey Kravchenko, and Evgeny Burnaev. 2020. Gradient-based adversarial attacks on categorical sequence models via traversing an embedded world. *arXiv:2003.04173* (2020).
- [10] Ji Gao, Jack Lanchantin, Mary Lou Soffa, and Yanjun Qi. 2018. Black-box generation of adversarial text sequences to evade deep learning classifiers. In *IEEE Security and Privacy Workshops*. IEEE, 50–56.
- [11] Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. 2014. Explaining and harnessing adversarial examples. *preprint arXiv:1412.6572* (2014).
- [12] R. Huang, B. Xu, D. Schuurmans, and C. Szepesvári. 2015. Learning with a strong adversary. *preprint :1511.03034* (2015).
- [13] Valentin Khruikov and Ivan Oseledets. 2018. Art of singular vectors and universal adversarial perturbations. In *IEEE CVPR*. 8562–8570.
- [14] Alexey Kurakin, Ian Goodfellow, and Samy Bengio. 2016. Adversarial examples in the physical world. *preprint arXiv:1607.02533* (2016).
- [15] Linyang L., Ruotian M., Qipeng G., Xiangyang X., and Xipeng Q. 2020. BERT-attack: adversarial attack against BERT using BERT. *arXiv:2004.09984* (2020).
- [16] Bin Liang, Hongcheng Li, Miaoqiang Su, Pan Bian, Xirong Li, and Wenchang Shi. 2017. Deep text classification can be fooled. *preprint arXiv:1704.08006* (2017).

- [17] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. 2017. Towards deep learning models resistant to adversarial attacks. *preprint arXiv:1706.06083* (2017).
- [18] J. H. Metzen, T. Genewein, V. Fischer, and B. Bischoff. 2017. On detecting adversarial perturbations. *ICLR* (2017).
- [19] Nicolas Papernot, Patrick McDaniel, Ananthram Swami, and Richard Harang. 2016. Crafting adversarial input sequences for recurrent neural networks. In *MILCOM 2016-2016 IEEE Military Communications Conference*. IEEE, 49–54.
- [20] Nikolaos Pitropakis, Emmanouil Panaousis, Thanassis Giannetsos, Eleftherios Anastasiadis, and George Loukas. 2019. A taxonomy and survey of attacks against machine learning. *Computer Science Review* 34 (2019), 100199.
- [21] Lichao Sun, Ji Wang, Philip S Yu, and Bo Li. 2018. Adversarial attack and defense on graph data: a survey. *preprint arXiv:1812.10528* (2018).
- [22] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow, and Rob Fergus. 2013. Intriguing properties of neural networks. *preprint arXiv:1312.6199* (2013).
- [23] W. Wang, B. Tang, R. Wang, L. Wang, and A. Ye. 2019. Towards a Robust Deep Neural Network in Texts: A Survey. *preprint arXiv:1902.07285* (2019).
- [24] Han X., Yao M., Haochen L., Debayan D., Hui L., Jiliang T., and Anil K.J. 2019. Adversarial attacks and defenses in images, graphs and text: a review. *arXiv:1909.08072* (2019).
- [25] Xiaoyong Yuan, Pan He, Qile Zhu, and Xiaolin Li. 2019. Adversarial examples: attacks and defenses for deep learning. *IEEE transactions on neural networks and learning systems* 30, 9 (2019), 2805–2824.
- [26] Xiaoyong Yuan, Qile Zhu Pan He, and Xiaolin Li. 2017. Adversarial examples: attacks and defenses for deep learning. *arXiv:1712.07107* (2017).
- [27] Mary Zeager, Aksheetha Sridhar, Nathan Fogal, Stephen Adams, Donald Brown, and Peter Beling. 2017. Adversarial learning in credit card fraud detection. 112–116. <https://doi.org/10.1109/SIEDS.2017.7937699>
- [28] Mary Frances Zeager, Aksheetha Sridhar, Nathan Fogal, Stephen Adams, Donald E Brown, and Peter A Beling. 2017. Adversarial learning in credit card fraud detection. In *2017 Systems and Information Engineering Design Symposium (SIEDS)*. IEEE, 112–116.
- [29] T. Zhang, V. Kishore, F. Wu, K.Q. Weinberger, and Y. Artzi. 2019. Bertscore: Evaluating text generation with bert. *preprint arXiv:1904.09675* (2019).
- [30] Wei Emma Zhang, Quan Z Sheng, Ahoud Alhazmi, and CHENLIANG LI. 2019. Adversarial attacks on deep learning models in natural language processing: a survey. *preprint arXiv:1901.06796* (2019).

# A Research Methods

## A.1 Implementation details

Here, we give an overview of the implementation details. We base our implementation on the AllenNlp library.

- Reading data.** To read the original transactions record data, and to transform it into the so-called *allennlp.data.fields* format, which we require to manipulate our data, we use the *Transactions Reader*, based on the *DatasetReader* from the AllenNlp library, with a pre-trained discretizer to transform monetary values and further work with these amounts in a discrete space. The discretizer is the *KBinsDiscretizer* with hyperparameters,  $n\text{ bins}=100$ ,  $\text{strategy}=\text{quantile}$ ,  $\text{encode}=\text{ordinal}$ .
- Classifiers.** In this work, we use a bidirectional GRU, an LSTM with 1 layer, a dropout with 0.1, a hidden size of 256 or a CNN with 64 filters and a ReLU activation layer. We use cross-entropy as the loss function. We train all models with a batch size 1024 for 50 epochs with an early stop if the loss does not decrease during the validation phase within 3 epochs. Further, we use the Adam optimizer with a step size 0.001.
- Attacks.** We conduct most of the experiments with attacks using black-box settings. For our experiments, Table 2 is the result of an RNN with a GRU block as an attacked model and an RNN with a GRU block as a substitute model that an attacker uses to assess the so-called *adversality* of the generated sequences. We generate **1000** adversarial sequences from the test set and average evaluated model performance metrics over them.

All hyperparameters and implementation details for our attacks are in subsection 3.1. We list them as follows:

- **Sampling Fool (SF).** The algorithm hyperparameters of the Sampling Fool attack include temperature and the number of samples. The temperature regulates how close the distribution of tokens is to be uniform. At high temperatures, attacks are more diverse but less effective. We set the temperature parameter to 2.0 and the number of samples to 200.

We use the Masked Language Model (MLM) [5], as our language model for the Sampling Fool attack. The MLM accidentally masks some tokens from the original sequence during training and tries to predict which tokens were instead of the masked ones. In our task, we input the model not only of the embeddings of transactions but also the embeddings of the amounts of these transactions, which are concatenated and fed to the encoder. We train the model by

maximizing the probability of the sequences.  
$$\max_{\theta} \log p_{\theta}(\mathbf{x}) = \sum_{t=1}^T \log p_{\theta}(x_t | \mathbf{x}_{<t}).$$

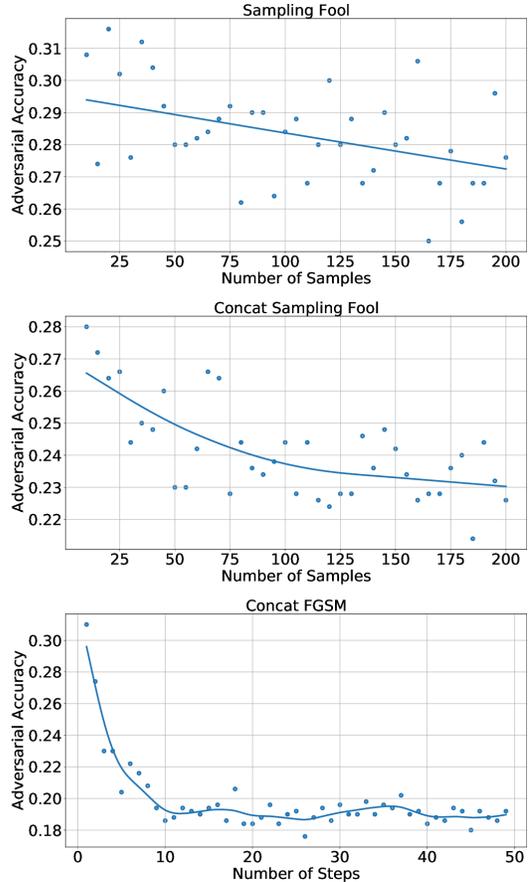


Figure 11: Dependence of the Adversarial Accuracy Metric on the number of samples for Sampling Fool and Concat Sampling Fool, and a number of iterations in the case of Concat FGSM. For this experiment, we used Age 2 dataset.

We use these models to sample the most plausible sequences from transactions with a constraint on the input amounts:

- **Concat SF, Sequential Concat SF.** The attacks have the same hyperparameters as SF with an additional parameter for the number of tokens, added in the attack. We set this parameter to 2.
- **FGSM.** We have the hyperparameters,  $\epsilon$  for the step size towards the antigradient, and  $n$  as the algorithm’s number of steps. We set  $\epsilon$  to 1.0 and number of steps to 30.
- **Concat FGSM.** The attack has the same hyperparameters as FGSM with an additional parameter for the number of tokens added in the attack. We set this parameter to 4.

Attack	NAD $\uparrow$	WER $\downarrow$	AA $\downarrow$	PD $\uparrow$
Age 1				
SF	0.08	12.79	0.48	0.08
Concat SF	0.46	2	0.09	0.26
FGSM	0.44	1.73	0.51	0.06
Concat FGSM	0.24	4	0.02	0.24
Client Leaving				
SF	0.17	14.71	0.67	0.01
Concat SF	0.18	2	0.63	0.01
FGSM	0.29	8.67	0.50	0.08
Concat FGSM	0.11	4	0.53	0.09
Scoring				
SF	0.12	6.16	0.83	-0.05
Concat SF	0.05	6	0.67	0.07
FGSM	0.23	8.72	0.67	0.08
Concat FGSM	0.08	6	0.50	0.19

Table 8: Attacking of model with CNN architecture using substitute model with GRU architecture

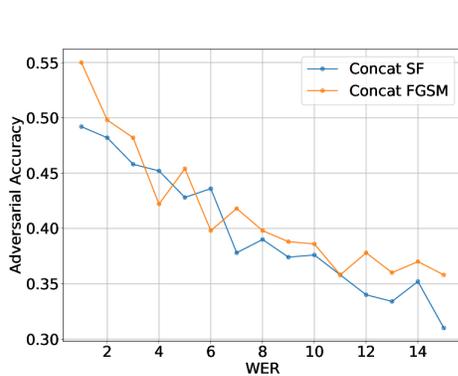


Figure 12: Dependence of the Adversarial Accuracy Metric on the number of added adversarial tokens (or WER) for Concat SF and Concat FGSM. Client leaving dataset is under consideration, GRU is used as the target and substitute models.

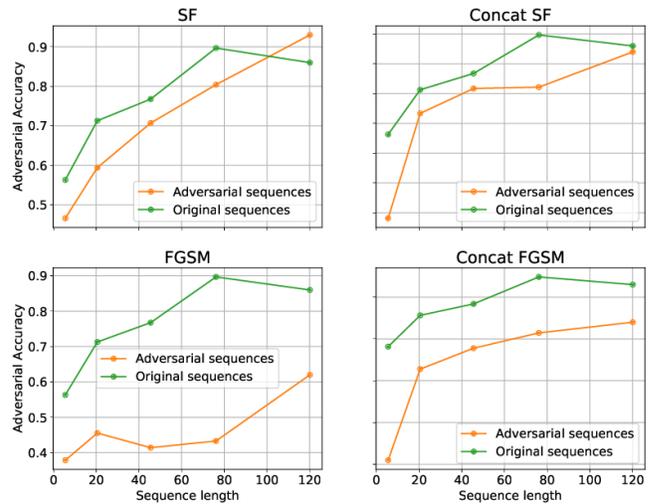


Figure 13: Dependence of the Adversarial Accuracy Metric on the length of the original sequences for Client leaving dataset. The target and substitute models are GRU. The green line corresponds to the accuracy of the target classifier for the original sequences and the orange line corresponds to the results on adversarial sequences.

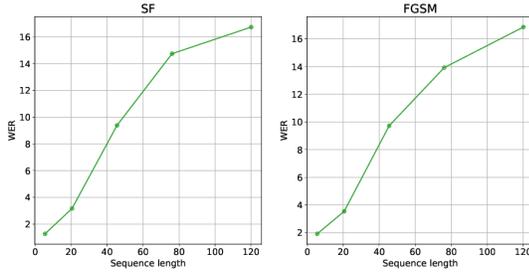


Figure 14: Dependence of WER on length of the original sequences for successful attacks. Client leaving dataset was used. The target and substitute models are GRU.

## A.2 Alternative model architectures

We perform additional experiments for different architectures of the target and substitute models. The results are in Table 8. We can attack the target model consisting of a CNN architecture with a substitute model entailing a GRU architecture, and the quality of some attacks will be still high. We are mainly considering when the target and substitute models are both a GRU architecture. However, most of our main claims hold when both models also have a CNN architecture.

## A.3 Varying hyperparameters

We select all parameters for the attacks by doing brute force on the grid. We examine the efficiency of the attacks by varying the number of sample hyperparameters. In Figure 11, we can observe the results for our considered attacks. An increase in the number of samples linearly impacts the scale attack time, but we observe

that it is sufficient to use about 100 to reach the attack quality plateau for all attacks.

## A.4 Influence of WER on attacks effectiveness

We evaluate the Concat SF and Concat FGSM attacks to understand how the number of appended adversarial tokens at the end affects the attacks' quality. In Figure 12, we depict the result of the dependence of Adversarial Accuracy on WER. From this experiment, we can see that we increase the attack's effectiveness by adding more tokens at the end of the sequence.

## A.5 Connection between attacks metrics and sequence lengths

We are also interested in the connection between the attacks' metrics and the length of the initial transaction sequences. In Figure 13, we assess the dependence of Adversarial Accuracy on sequence length. We divide all transaction records sequences into five groups according to their length. Each point corresponds to the center of one of five length intervals. Then, on each interval, we record the number of attacks where the initial label does not change. From our observations, we conclude that long sequences are more robust to the attacks, while short adversarial sequences are more vulnerable to the attacks because of the absence of long-term patterns in them. Also, we want to understand the relationship between the sequence lengths and the number of adversarial tokens required in the sequence to lead to misclassification. Our results are in Figure 14. We must add more adversarial tokens to the long sequences, so the target classifier generates a wrong prediction compared to the short sequences.