Abstract

Most existing failure detection algorithms rely on statistical methods, and very few use machine learning (ML). This paper explores the viability of ML in the field of failure detection: is it possible to implement an ML-based detector that achieves a satisfactory quality of service? We implement a prototype that uses a basic long short-term memory neural network algorithm, and study its behavior with real traces. Although ML model has comparatively longer computing time, our prototype performs well in terms of accuracy and detection time.

Index Terms: failure detection, machine learning.

Introduction

Distributed systems should provide reliable and continuous service despite failures. In [1], Chandra and Toueg show that failure detection is the dominant factor in system unavailability, and introduce the notion of unreliable failure detector (FD) as a theoretical construct to extend the applicability of distributed algorithms such as consensus and atomic broadcast. An FD is an oracle which monitors a remote process *P*, and assesses in real time whether *P* is up or has crashed. The assessment is unreliable because an FD might provide incorrect information over limited periods of time.

Machine learning (ML) performs well upon analyzing extremely regular data, but network latency on a link can vary often and significantly over time, and failures or delays are transient events that occur irregularly. To achieve high accuracy rates consistently, an ML-based FD must constantly train over newly emerging data. Resource greediness can be prohibitive, since a crucial aspect of an FD is its output rate. Our main goal is to study the viability of FDs based on ML techniques; such FDs should require minimum resources and time during real-time training while maintaining high accuracy. In this paper, we present a preliminary approach that uses a basic long short-term memory neural network algorithm. After comparing its output with a baseline FD, we further optimize our model's performance by adjusting its parameters and structure. Our simulations based on real traces show that our model performs well in terms of accuracy and detection time, despite that the ML model incurs non-negligible computation time.

An FD is a process that monitors remote processes, and strives to estimate whether they're still up or

RELATED WORK

node in a distributed environment with 100% certainty. To make up for this, distributed systems can establish a network of mutual observation via *unreliable* FDs [1]. Such FDs require each node in the system to send periodical heartbeat messages to a monitor node to prove its liveliness. If a heartbeat message does not arrive before its expected arrival date, the monitor will suspect that the corresponding process has failed. A high performance FD is expected to generate a series of efficient but tolerant expected arrival dates, which are able to detect a failed process in time with minimum wrong suspicions of the healthy ones. Implementing an FD is a challenge, because it must contend with the unpredictable and asynchronous nature of network links while preserving its set properties in terms of completeness and accuracy [3]. Early

whether they've crashed. The authors of [2] prove formally that there is no way to determine the failure of a

approaches [4, 5] propose adaptive FDs that adjust the timeout delay tolerance dynamically; but they assume an unrealistic timing model with no bound on the delays. Chen et al. [6] overcome this issue by introducing quality of service metrics which quantifies the performance of an FD based on the speed of correct detection and the ratio of a false ones. Previous FDs are mostly based on statistical models. Although more recent studies introduce ML-driven methods for FD-related problems, their solutions are developed under different context. In [7], the authors use a stacked-LSTM model to classify potential anomaly events in cloud services by analyzing labeled static

sensor logs. They conclude that their S-LSTM approach has the ability to quickly learn from the historic patterns and adjust to unpredictable anomalous events. The authors of [8] propose three different ML methods for high-performance computing systems failure detection by analyzing informative hardware usage data. Their best solution, based on the Support Vector Machine method, achieves a precision of 90%. [9] acknowledges the significance of timestamp data in system logs, and passes both time and token sequence data into the model under the form of interpolated vectors. By using a CNN model, this approach reaches a 99.5% F1 score. The ML-driven methods mentioned above address a different, extended model where they might have more extensive information for error implication and require pre-labeled data to tag anomalies in the original datasets, addressing a classification problem. Because of this, they don't apply well to the classic FD problem

which is a regression problem since timestamps are the only retrievable information from the monitored sub-systems, yet this information is self-sufficient for detection training. Hence, in this paper we will focus on studying the viability of ML on the classic FD problem and compare the results with the statistic-driven FD models. III System Model In our distributed system setup, we consider a classic Push-FD model [10], where every node p sends

computes an estimate of the arrival date, also called as freshness point, τ_i of the next heartbeat h_i . If q receives h_i before τ_i , then it considers that p is up. Otherwise, q starts suspecting p of having crashed. The

 h_{i-1}

probability of a false positive detection [6].

possible detection time, highest possible accuracy.

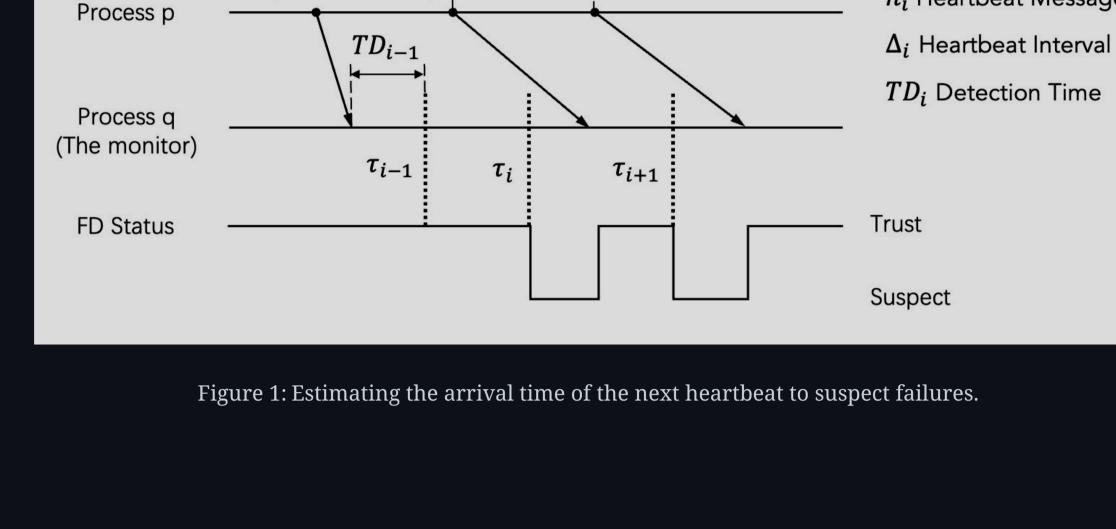
IV BASELINE

Training

time that elapses between the emission of h_i and τ_i is the detection time T_D . τ_i Freshness Point

 h_i Heartbeat Message

heartbeats to its counterparts q every Δ_i milliseconds (Figure 1). To assess the status of p (trust or suspect), q



The FD generates a *freshness point* for the monitored node's next heartbeat based on its previous behaviors. The freshness point is composed of two parts, an estimated arrival (EA) and a safety margin (α). The function

of a safety margin is to mitigate the impact of unexpected and erratic delays and thus to reduce the

We assume that our network follows a crash-stop model, where when a node fails, it will not come back alive. However, it can easily be extended to a crash-recovery model: The FD will consider a node that recovers as a new node. Restricting the model to crash-stop ensures that all nodes behave normally in the time period where incoming data is recorded. Therefore, the main goal of our algorithm is not to detect behavior anomalies, unlike [7, 8, 9]. Instead, our FD aims to learn the pattern of heartbeat signals sent by each node in real-time, and to determine an expected arrival time with a high quality of detection: shortest

probability of availability (P_A) , which is the ratio of the safe predictions over the total predictions. Making a safe prediction means the arrival time of the next heartbeat is earlier than the predicted time, and this is to avoid false-positive results where a node is wrongly suspected. Second, the detection time (T_D) is the difference between the forecast time stamp and the actual reception time stamp. This metric complements the P_A : it reduces the gap between the predicted time and the actual time, since a model can achieve a high P_A by always predicting absurdly long arrival times. Thus, both metrics balance each other: improving T_D aggressively increases the rate of false-positives, while doing so for P_A increases the risk of false-negatives. Our last metric is the *computation time* T_c . This is important because of the generally high computational cost of non-linear machine learning methods. An FD whose next prediction takes longer to compute than the next arrival date is pointless.

Three different metrics are used to assess FD performance. The first and the most significant one is the

Arrival EA) upon receiving a heartbeat message. Equation 1 shows how CFD computes its prediction as a statistical analysis of the n latest reception times. To reduce the ratio of false positives, CFD adds a constant safety margin (α) to EA. $EA_{k+1} \approx \frac{1}{n} \left(\sum_{i=k-n}^{k} A_i - \Delta_i * i \right) + (k+1) * \Delta_i$ (1)

We use Chen's FD (CFD) [6] as our baseline for comparison. CFD predicts the next arrival time (Expected

Our ML-Based Failure Detector Design Similarly to CFD, our approach also bases its prediction on an analysis of the η latest receptions. We use the long short-term memory (LSTM) model: it can effectively retain important long-term information, and it can

quickly fit the data compared to traditional time series forecasting models [11]. Moreover, its fitting method

One of the weaknesses of Chen is the constant α . Setting α by default reduces the performance of CFD for

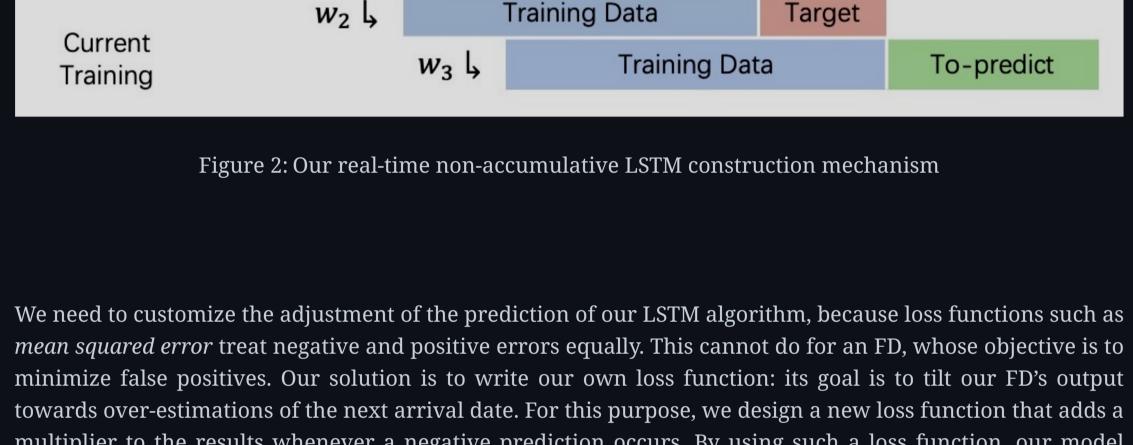
highly unstable traces. Besides, selecting the optimal value for α is a delicate task. It depends heavily on the

is non-linear, which makes its fitting potential stronger than ordinary linear models. As our goal is to make real-time predictions, we train our LSTM model real-time and non-accumulatively: we

only keep the η most recent data items as the training set, and always feed the model with the nearest fixedsize of data. Although we abandon data that falls out of the range, their impact are kept in the parameters. For each real-time training, the model retains the memory of the last training results, and will not restart from scratch. Since we don't store historical data for future training, their impact gradually decreases over time. This training method has two main advantages. First, it speeds up the training process and reduces its

computational cost by retaining the memory of the last training results. Second, it gradually reduces the

impact of the least recent historical data; this is essential to accommodate the high volatility of network link behaviours. Figure 2 illustrates the mechanism of our model. **Historical Data Current Data** Heartbeats **Training Data** Target Historical Training Data Target $w_1 \downarrow$



multiplier to the results whenever a negative prediction occurs. By using such a loss function, our model achieves a 95% accuracy rate on its predictions. To improve the PA further, we add a dynamic safety margin based on the ϵ most recent errors, unlike our

baseline model which uses a static one. A safety margin is used to help the model adapt to any unexpected

and unstable delay when it occurs, and such a delay is unbounded and hard to detect. So the goal is to adapt

the model as soon as such a series of delays occur, so having a dynamic array that keeps the ϵ most recent errors enables the model to adjust in time from the latest error. Altogether, our model computes its next freshness point τ with Equation 2. $\tau_{k+1} \approx LSTM.predict(k-\eta) + \frac{1}{\epsilon} \sum_{i=k-\epsilon}^{k} error_i$ **(2)**

size, and the epoch number. The training data set size η decides how much data to learn each time while realtime training; the batch sizes and epoch numbers influence the accuracy and timeliness of training from the structure and nature of training. Grid search is used to search for best combinations. Table I shows the top ten combinations of P_A among all since P_A is the most critical property for evaluation. The top three combinations share the same P_A and similar T_D , however, the second one has an outstanding T_C than the other two. As a result, we used a η of 500, a batch size of 64 and an epoch of 5.

 P_A

0.9957

0.9947

0.9945

batch size epoch

32

64

64

32

64

32

32

500

500

500

1000

1000

500

100

1000

In the LSTM model, three parameters are decided by exhaustive search: η (training data set size), the batch

 T_D

9.7832

11.1584

12.2609

9.2335

0.9938 11.3694 83.6161

0.9957 9.4416

0.9957 9.1369

 T_{C}

72.9630

48.1018

74.2897

150.0853

89.4486

116.3436

8.0131 45.8082

We assess our approach on top of real traces: heartbeat transmission logs collected over a week from 9 nodes on the PlanetLab network (http://www.planet-lab.org/). Nodes 1 through 9 generate a heartbeat message containing the sender ID and a sequence number every 100 milliseconds, and send it to node 0, which we consider as the monitoring node. Upon reception of a heartbeat, node 0 records its arrival time in the log associated with the sender.

We set the parameter values of Equations 1 and 2 as follows: the heartbeat arrivals window size n for Chen's

FD (CFD) is 1000; for our Machine Learning-based FD (MLFD), the heartbeat window size η is 500, and the

error window size ϵ is 10. We aim to compare the performance of MLFD with that of CFD. However, CFD's

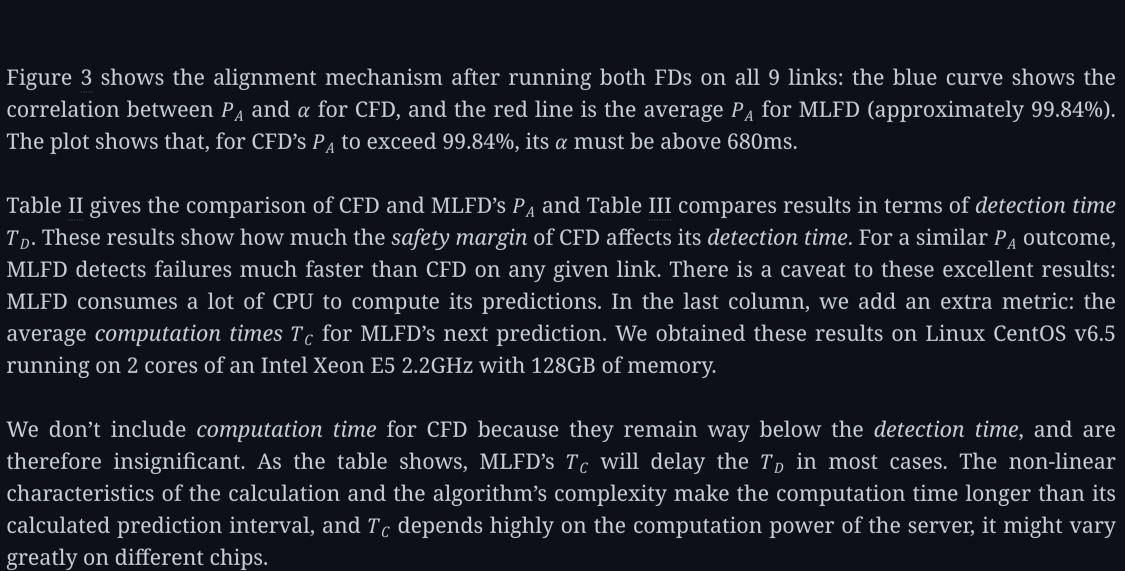
performance depends highly on the value of its constant safety margin α . To allow for a fair comparison, we

align the P_A values of both models before comparing their T_D ; this leads to a value of 680ms for α .

Chen's Model PA Average

Figure 3: Influence of CFD's safety margin on its accuracy.

LSTM Model PA Average



7 0.989449 0.998765 0.989559 8 0.995549 0.998366 0.995732 9 0.999055 0.998333 0.998362 0.998365 Average TABLE II: P_A comparison between CFD and MLFD

12.926

96.087

21.827

148.897

12.410

20.400

26.823

Chen's P_A (680ms) MLFD's P_A Chen's P_A (690ms)

0.999127

0.997104

0.998439

0.996781

0.999165

0.998457

CFD T_D (680ms) MLFD T_D CFD T_D (690ms) MLFD T_C

689.093

520.860

686.351

594.397

690.019

688.190

682.343

123.847

124.529

124.682

124.836

125.254

125.199

125.754

Link

6

Link

679.093

510.860

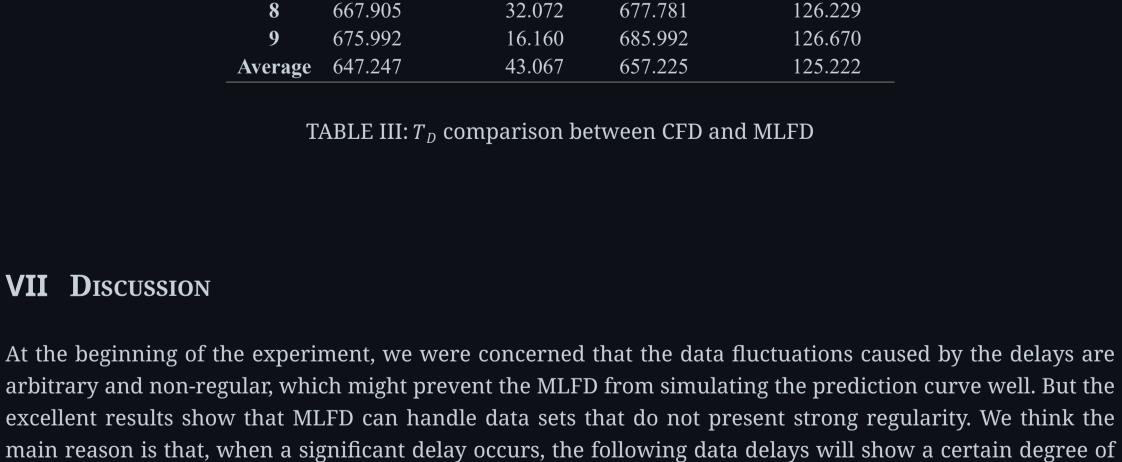
676.351

584.397

680.019

678.190

672.419



consistency. Although MLFD cannot predict a sudden and major delay variation, it can sensitively capture

obviously a major challenge for the application of machine learning in the field of failure detection.

the change and adjust the model parameters to adapt to the delays if they occur in bursts. And once the burst ends, the MLFD will adapt back to the normal trend. However, the heavy calculation cost does indeed weigh on the performance of our MLFD. On average, the T_c of the MLFD is three times as long as its T_D , which seriously affects the performance of the MLFD. This is

times for a similar probability of availability, but at a significant computation cost.

Nevertheless, we remain optimistic about our approach. One reason is that, despite the computational cost issue, the accuracy and detection time of our MLFD is still significantly better than that of CFD. The other reason is that the processing power of the server at our disposal to run the MLFD is relatively low. Running our model on a more powerful server would shorten the calculation time significantly. **VIII** Conclusion In this paper, we present a preliminary study about the feasibility of failure detector implementations based

on machine learning algorithms. Our results suggest that ML may be a viable approach: our LSTM-based FD

implements a unilateral penalty loss function, dynamic *safety margin*, and it trains on the 500 most recent

message receptions in real-time. Upon comparison with Chen's, our FD achieves much shorter detection

Our results illustrate the potential of machine learning models in this field, and we hope it will elicit further

research. Our next step is to refine our prototype, and to test it against more aggressive FDs with dynamic safety margins [3]. We also intend to pursue our work with a focus on the trade-off between probability of availability and computation time. Evolved machine learning models, such as the transformer's attention model, seem like strong candidates for better performance. The advent of DPUs [12] also opens another promising avenue of research for our work: FDs are obvious candidates for network function virtualization, and SmartNICs have the power to support advanced implementations such as our MLFD. References

T. D. Chandra and S. Toueg, "Unreliable Failure Detectors for Reliable Distributed [1] [2]

[11]

- Systems," *Journal of the ACM*, vol. 43(2), p. 225–267, 1996. M. J. Fischer, N. A. Lynch, and M. S. Paterson, "Impossibility of distributed consensus with one faulty process," *Journal of the ACM*, vol. 32, no. 2, p. 374–382, 04 1985.
- M. Bertier, O. Marin, and P. Sens, "Implementation and performance evaluation of [3] an adaptable failure detector," in Proc. of the Int. Conf. on Dependable Systems and Networks, 2002, pp. 354–363.
- C. Fetzer, M. Raynal, and F. Tronel, "An adaptive failure detection protocol," in [4] Pacific Rim Int. Symposium on Dependable Computing, 2001, pp. 146–153. [5] I. Sotoma and E. R. M. Madeira, "DPCP (Discard Past Consider Present)-a novel approach to adaptive fault detection in distributed systems," in *Proc. 8th IEEE* Workshop on Future Trends of Distributed Computing Systems, 11 2001, pp. 76–82.
- W. Chen, S. Toueg, and M. K. Aguilera, "On the quality of service of failure [6] detectors," in *Proc. of the Int. Conf. on Dependable Systems and Networks*, vol. 51, no. 1, 01 2002, p. 13–32. R. Vinayakumar, K. P. Soman, and P. Poornachandran, "Long short-term memory
- [7] based operation log anomaly detection," in 2017 International Conference on Advances in Computing, Communications and Informatics (ICACCI), 2017, pp. 236– 242.
- K. Yamnual, P. Phunchongharn, and T. Achalakul, "Failure detection through [8] monitoring of the scientific distributed system," in 2017 International Conference on Applied System Innovation (ICASI), 2017, pp. 568–571.
- Y. Huangfu, S. Habibi, and A. Wassyng, "System failure detection using deep [9] learning models integrating timestamps with nonuniform intervals," IEEE Access, vol. 10, pp. 17629–17640, 2022. P. Felber, X. Defago, R. Guerraoui, and P. Oser, "Failure detectors as first class [10] objects," in Proceedings of the International Symposium on Distributed Objects and

Applications, 1999, pp. 132–141.

NY, USA, 2021, p. 772–787.

Copyright

predictable through time-window approaches," in *Proc. of the Int. Conf. on Artificial Neural Networks*, 2001, pp. 669–676. Y. Qiu, J. Xing, K.-F. Hsu, Q. Kang, M. Liu, S. Narayana, and A. Chen, "Automated [12] smartnic offloading insights for network functions," in *Proceedings of the ACM* SIGOPS 28th Symposium on Operating Systems Principles, ser. SOSP '21, New York,

D. Gers, Felix A.and Eck and J. Schmidhuber, "Applying LSTM to time series

View original Conversion Report aruv report (OK) on arXiv an issue

Generated on Thu Mar 14 02:52:15 2024 by LATEXML

Privacy Policy