

# GRID HTM: Hierarchical Temporal Memory for Anomaly Detection in Videos

Blinded for review

email

blinded

blinded

## ABSTRACT

The interest for video anomaly detection systems has gained traction for the past few years. The current approaches use deep learning for performing anomaly detection in videos, but this has multiple problems. For starters, deep learning in general has issues with noise, concept drift, explainability, and training data volumes. Additionally, anomaly detection in itself is a complex task and faces challenges such as unknownness, heterogeneity, and class imbalance. Anomaly detection in deep learning is therefore mainly constrained to generative models such as generative adversarial networks and autoencoders due to their unsupervised nature, but even they suffer from general deep learning issues and are hard to train properly. In this paper, we explore the capabilities of the hierarchical temporal memory (HTM) algorithm to perform anomaly detection in videos, as it has favorable properties such as noise tolerance and online learning which combats concept drift. We introduce a novel version of HTM, namely, Grid HTM, which is an HTM-based architecture specifically for anomaly detection in complex videos such as surveillance footage.

## CCS CONCEPTS

- Computing methodologies → Computer vision.

## KEYWORDS

HTM, deep learning, surveillance, anomaly detection

### ACM Reference Format:

Blinded for review. 2018. GRID HTM: Hierarchical Temporal Memory for Anomaly Detection in Videos. In *Proceedings of Make sure to enter the correct conference title from your rights confirmation email (Conference acronym 'XX)*. ACM, New York, NY, USA, 7 pages. <https://doi.org/XXXXXXX.XXXXXXX>

## 1 INTRODUCTION

As the global demand for security and automation increases, many seek to use video anomaly detection systems. In the US alone, the surveillance market is expected to reach \$23.60 Billion by 2027 [1]. Leveraging modern computer vision, modern anomaly detection systems play an important role in increasing monitoring efficiency and reducing the need for expensive live monitoring. Their use

---

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

Conference acronym 'XX, June 03–05, 2018, Woodstock, NY

© 2018 Association for Computing Machinery.

ACM ISBN 978-1-4503-XXXX-X/18/06...\$15.00

<https://doi.org/XXXXXXX.XXXXXXX>

cases can vary from detecting faulty products on an assembly line to detecting car accidents on a highway, and everything in between.

The most important component in video anomaly detection systems is the intelligence behind it. The intelligence ranges from simple on-board algorithms to advanced deep learning models, where the latter has experienced increased popularity in the past few years. Yet, despite the major progress within the field of deep learning, there are still many tasks where humans outperform models, especially in anomaly detection where the anomalies are often undefined. Deep learning approaches also perform poorly when dealing with noise and concept drift.

The cause for the discrepancy lies in the difference between how humans and machine learning algorithms represent data and learn. Most machine learning algorithms use a dense representation of the data and apply back-propagation in order to learn. Human learning happens in the neocortex, where evidence points to that the neocortex uses a sparse representation and performs Hebbian-style learning. For the latter, there is a growing field of machine learning dedicated to replicating the inner mechanics of the neocortex, namely Hierarchical Temporal Memory (HTM) theory [2]. This theory outlines its advantages over standard machine learning, such as noise-tolerance and the ability to adapt to changing data.

With the advantages of HTM and the rise of video anomaly detection in mind, a natural question one could pose is whether it is possible to apply HTM for anomaly detection in videos. Combined with a lack of related works, it is this very question that is the motivation behind this paper.

## 2 RELATED WORK

Anomaly detection is often defined as detecting data points that deviate from the general distribution [3]. Unlike most other problems in deep learning, anomaly detection deals with unpredictable and rare events which makes it hard to apply traditional deep learning for anomaly detection. A subset of anomaly detection is smart surveillance [4], which is the use of video analysis specifically in surveillance.

An issue for deep-learning models in general is that they are susceptible to noise in the dataset [5, 6], which leads to decreased classification accuracy and poor prediction results. Due to the nature of training deep learning models, they are also in most cases not self-supervised and therefore require constant tuning in order to stay effective on changing data. Not to mention that they require a lot of data before they can be considered effective, and that performance increases logarithmically based on the volume of training data [7]. Deep learning models also suffer from issues with out-of-distribution performance, where a model might perform great on the dataset it is tested on, but performs poorly when deployed in real life. This could be caused by selection bias in the dataset

or when there are differences in the causal structure between the training domain and the deployment domain [8].

Deep learning models generally suffer from a lack of explainability [9]. While it is known *how* the models make their decisions, their huge parametric spaces make it unfeasible to know *why* they make those predictions. Combined with the vast potential that deep learning offers in critical sectors such as medicine, makes approaches that offer explainability highly attractive.

HTM theory [2] introduces a machine learning algorithm which works on the same principles as the brain and therefore solves some of the issues that deep learning has. HTM is considered noise resistant and can perform online learning, meaning that it learns as it observes more data. HTM replicates the structure of the neocortex which is made up of cortical regions, which in turn are made of up of mini-columns and then neurons.

The data in an HTM is represented using a Sparse Distributed Representation (SDR), which is a sparse bit array. An encoder converts real world values into SDRs, and there are currently encoders for numbers, geospatial locations, categories, and dates. One of the difficulties with HTM is making it work on visual data, where creating a good encoder for visual data is still being researched [10, 11, 12].

The learning mechanism consists of two parts, the Spatial Pooler (SP) and the Temporal Memory (TM). The SP learns to extract semantically important information into output SDRs. The TM learns sequences of patterns of SDRs and forms a prediction in the form of a predictive SDR.

Research has shown that HTM is very capable of performing anomaly detection on low-dimensional data and is able to outperform other anomaly detection methods [13]. However, related works, such as Daylidyonok, Frolenkova, and Panov [12], show that HTM struggles with higher dimensional data. Therefore, a natural conclusion to make is that HTM should be applied differently, and that a new type of architecture using HTM should be explored for the purpose of video anomaly detection and surveillance.

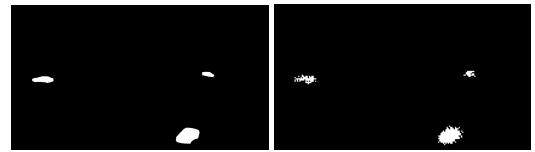
### 3 GRID HTM

This paper explores a new type of architecture, named Grid HTM, for anomaly detection in videos using HTM, and proposes to use segmentation techniques to simplify the data into an SDR-friendly format. These segmentation techniques could be anything from simple binary thresholding to deep learning instance segmentation. Even keypoint detectors such as ORB [14] could in theory be applied. When explaining Grid HTM, the examples will be taken from deep learning instance segmentation of cars on a video from the VIRAT [15] dataset. An example segmentation is shown in Figure 1. The idea is that the SP will learn to find an optimal general



**Figure 1:** Segmentation result of cars, which is suited to be used as an SDR. Original frame taken from VIRAT [15].

representation of cars. How general this representation is can be configured using the various SP parameters, but ideally they should be set so that different cars will be represented similarly while trucks and motorcycles will be represented differently. An example representation by the SP is shown in Figure 2.



**Figure 2:** The SDR (left) and its corresponding SP representation (right). Note that the SP is untrained.

The task of the TM will then be to learn the common patterns that the cars exhibit, their speed, shape, and positioning will be taken into account. Finally, the learning will be set so that new patterns are learned quickly, but forgotten slowly. This will allow the model to quickly learn the norm, even if there is little activity, while still reacting to anomalies. This requires that the input is stationary, in our example this means that the camera is not moving.

It is possible to split different segmentation classes into their respective SDRs. This will give the SP and the TM the ability to learn different things for each of the classes. For instance, if there are two classes "person" and "car", then the TM will learn that it is normal for objects belonging to "person" to be on the sidewalk, while objects belonging to "car" will be marked as anomalous when on the sidewalk.

Ideally, the architecture will have a calibration period spanning several days or weeks, during which the architecture is not performing any anomaly detection, but is just learning the patterns.

### 3.1 Improvements

Daylidyonok, Frolenkova, and Panov [12] was shown to not be particularly effective, therefore multiple improvements needed to be introduced to increase effectiveness. **Invariance.** One issue that becomes evident is the lack of invariance, due to the TM learning the global patterns. Using the example, it learns that it is normal for cars to drive along the road but only in the context of there being cars parked in the parking lot. It is instead desired that the TM learns that it is normal for cars to drive along the road, regardless of whether there are cars in the parking lot. We propose a solution based on dividing the encoder output into a grid and have a separate SP and TM for each cell in the grid. The anomaly scores of all the cells are then aggregated into a single anomaly score using an aggregation function.

**Aggregation Function.** Selecting the correct aggregation function is important because it affects the final anomaly output. For instance, it might be tempting to use the mean of all the anomaly scores as the aggregation function:

$$X : \{x \in \mathbb{R} : x \geq 0\}$$

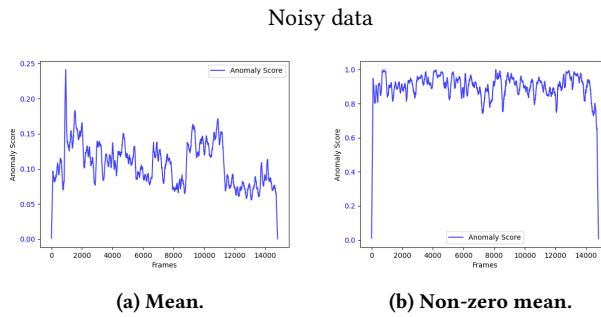
$$anomScore = \frac{\sum_{x \in X} x}{|X|}$$

However, this leads to problems with normalization, meaning that an overall anomaly score of 1 is hard to achieve due to many cells having a zero anomaly score. In fact, it becomes unclear what a high anomaly score is anymore. Using the mean also means that anomalies that take up a lot of space will be weighted higher than anomalies that take up a little space, which might not be desirable. To solve the aforementioned problem and if the data has little noise, a potential aggregation function could be the non-zero mean:

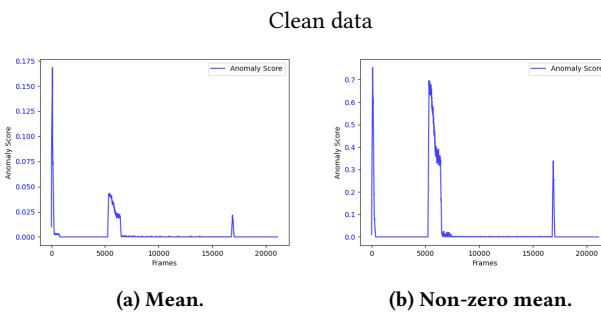
$$X : \{x \in \mathbb{R} : x > 0\}$$

$$\text{anomScore} = \begin{cases} \frac{\sum_{x \in X} x}{|X|} & \text{if } |X| > 0 \\ 0 & \text{otherwise} \end{cases}$$

Meaning that only the cells with a non-zero anomaly score, denoted  $X$ , will be contributing to the overall anomaly score which helps solve the aforementioned normalization and weighting problem. On the other hand, this will perform poorly when the architecture is exposed to noisy data which could lead to there always being one or more cells with a high anomaly score. Figure 3 illustrates



**Figure 3: How the two aggregation functions perform on the same noisy data.**



**Figure 4: How the two aggregation functions perform on the same clean data.**

the effect of an aggregation function for noisy data, where the non-zero mean is rendered useless due to the noise. On the other hand, Figure 4 shows how the non-zero mean gives a clearer anomaly

score when the data is clean. **Explainability.** Having the encoder output divided into a grid has the added benefit of introducing explainability into the model. By using Grid HTM it is now possible to determine where in the input an anomaly has occurred by simply observing which cell has a high anomaly score. It is also possible to estimate the number of predictions for each cell which can be used as a measure of certainty, where fewer predictions means higher certainty. Making it possible to measure certainty per cell creates a new source of information which can be used for explainability or robustness purposes. **Flexibility and Performance.** In addition, it is also possible to configure the SP and the TM in each cell independently, giving the architecture increased flexibility and to use a non-uniform grid, meaning that some cells can have different sizes. Last but not least, dividing the frame into smaller cells makes enables it to run each cell in parallel for increased performance. **Reviewing Encoder Rules.** A potential challenge with the grid approach is that the rules for creating a good encoder, may not be respected and therefore should be reviewed:

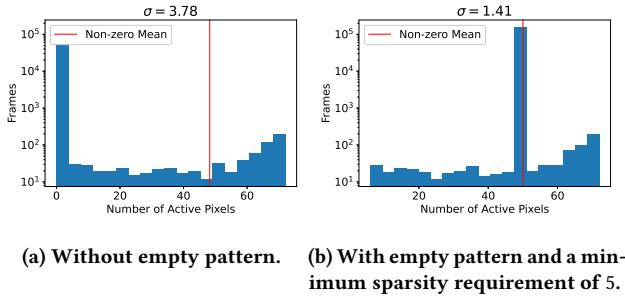
- **Semantically similar data should result in SDRs with overlapping active bits.** In this example, a car at one position will produce an SDR with a high amount of overlapping bits as another car at a similar position in the input image.
- **The same input should always produce the same SDR.** The segmentation model produces a deterministic output given the same input.
- **The output must have the same dimensionality (total number of bits) for all inputs.** The segmentation model output has a fixed dimensionality.
- **The output should have similar sparsity (similar amount of one-bits) for all inputs and have enough one-bits to handle noise and subsampling.** The segmentation model does not respect this. An example is that there can be no cars (zero active bits), one car ( $n$  active bits), or two cars ( $2n$  active bits), and that this will fluctuate over time.

The solution for the last rule is two-fold, and consists of imposing a soft upper bound and a hard lower bound for the number of active pixels within a cell. The purpose is to lower the variation of number of active pixels, while also containing enough semantic information for the HTM to work:

- Pick a cell size so that the distribution of number of active pixels is as tight as possible, while containing enough semantic information and also being small enough so that the desired invariance is achieved. The cell size acts as a soft upper bound for the possible number of active pixels.
- Create a pattern representing emptiness, where the number of active bits is similar to what can be expected on average when there are cars inside a cell. This acts as a hard lower bound for the number of active pixels.

There could be situations where a few pixels are active within a cell, which could happen when a car has just entered a cell, but this is acceptable as long as it does not affect the distribution too much. If it does affect the distribution, which can be the case with noisy data, then an improvement would be to add a minimum sparsity requirement before a cell is considered not empty, e.g. less than 5 active pixels means that the cell is empty. In the following example,

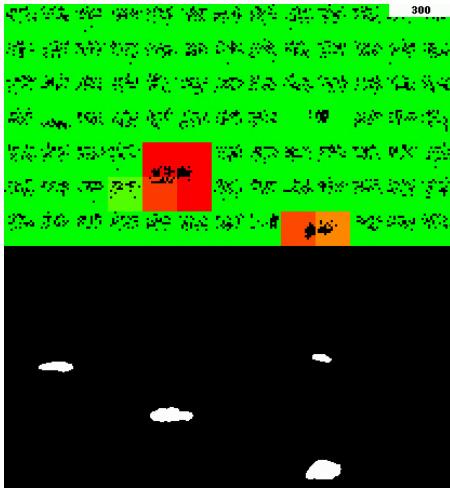
the number of active pixels within a cell centered in the video was used to build the distributions seen in Figure 5:



**Figure 5: Distribution of number of active pixels within a cell of size  $12 \times 12$ .**

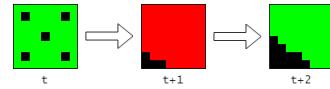
With a carefully selected empty pattern sparsity, **the standard deviation of active pixels was lowered from 3.78 to 1.41**. It is possible to automate this process by developing an algorithm which finds the optimal cell size and empty pattern sparsity which causes the least variation of number of active pixels per cell. This algorithm would run as a part of the calibration process.

The visual output resulting from these changes, which is an equally important output as the aggregated anomaly score, can be seen in Figure 6. Since there are now cells that are observing



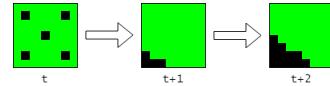
**Figure 6: Example Grid HTM output and the corresponding input. The color represents the anomaly score for each of the cells, where red means high anomaly score and green means zero anomaly score. Two of the cars are marked as anomalous because they are moving, which is something Grid HTM has not seen before during its 300 frame long lifetime.**

an empty pattern for a lot of the time in sparse data, boosting is recommended to be turned off, otherwise the SP output for the empty cells would change back and forth in order to adjust the active duty cycle.



**Figure 7: High anomaly score when an empty cell (represented with an empty pattern with a sparsity value of 5) changes to being not empty, as something enters the cell.**

**3.1.1 Stabilizing Anomaly Output.** Another issue with the grid based approach is when a car first comes into a cell. The TM in that cell has no way of knowing that a car is about to enter, since it does not see outside its own cell, and therefore the first frame that a car enters a cell will cause a high anomaly output. This is illustrated in Figure 7 where it can be observed that this effect causes the anomaly output to needlessly fluctuate. The band-aid solution is to ignore the anomaly score for the frame during which the cell goes from being empty to being not empty, which is illustrated in Figure 8. A



**Figure 8: In this case, the anomaly score is ignored (set to 0) for the frame in which the cell changes state from empty to not empty.**

more proper solution could be to allow the TM to grow synapses to the TMs in the neighboring cells, but this is not documented in any research papers and might also hinder invariance.

**3.1.2 Multistep Temporal Patterns.** Since the TM can only grow segments to cells that were active in the previous timestep, it will struggle to learn temporal patterns across multiple timesteps. This is especially evident in high framerate videos, where an object in motion has a similar representation at timestep  $t$  and  $t + 1$ , as an object standing still.

This could cause situations where an object that is supposed to be moving, suddenly stands still, yet the TM will not mark it as an anomaly due to it being stuck in a contextual loop. A contextual loop is when one of the predictions at  $t$  becomes true at  $t + 1$ , and then one of the predictions at  $t + 1$  is almost identical to the state at  $t$ , which becomes true if the object is not moving, causing the TM to enter the same state that it was in at  $t$ . A solution is to concatenate the past  $n$  SP outputs as input into the TM, which is made possible by keeping a buffer of past SP outputs and shifting its contents out as new SP outputs are inserted. This follows the core idea behind encoding time in addition to the data, which makes time act as a contextual anchor. However, in this case there are no timestamps that are suitable to be used as contextual anchors, so as a replacement, the past observations are encoded instead.

Concatenating past observations together will force the TM input, for when an object is in motion and when an object is still, to be unique. High framerate videos can benefit the most from this, and the effect will be more pronounced for higher values of  $n$ .

A potential side effect of introducing temporal patterns, is that because the TM is now exposed to multiple frames at once, it will be more tolerant to temporal noise. An example of temporal noise is

when an object disappears for a single frame due to falling below the classification threshold of the deep learning segmentation model encoder. The reason for the noise tolerance is that instead of the temporal noise making up the entire input for the TM, it now only makes up  $\frac{1}{n}$  of the TM input.

**3.1.3 Use Cases.** The most intuitive use case is to use Grid HTM for semi-active surveillance, where personnel only have to look at segments containing anomalies, leading to drastically increased efficiency. One example is making it possible to have an entire city be monitored by a few people. This is made possible by making it so that people only have to look at segments that Grid HTM has found anomalous, which is what drastically lowers the manpower requirement for active monitoring of the entire city.

## 4 EXPERIMENTAL DETAILS AND RESULTS

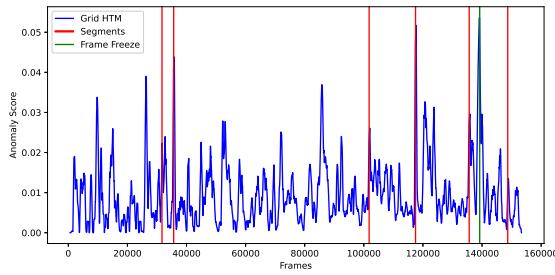
As stated earlier, one of the use cases of Grid HTM is anomaly detection in surveillance. This example will show how Grid HTM could perform. The video to be used is part of the VIRAT [15] video dataset, and was selected due to its long duration and stationary camera.

The downside is that the video does not contain any non-technical anomalies, but consists of technical anomalies in the form of several segments with sudden frame skips in between. There is also a synthetic anomaly introduced in the form of a frame repeat lasting a couple of seconds, essentially "freezing" time, in order to test whether Grid HTM is able to understand how objects should be moving in time.

In this experiment, a segmentation model which can extract classes into their respective SDRs is employed. Meaning that there could be an SDR for cars and an SDR for persons, that are then concatenated before being fed into the system.

The segmentation model used is PointRend [16] with a ResNet101 [17] backbone, pretrained on ImageNet [18], and implemented using PixelLib [19]. For the sake of simplicity, this experiment will focus only on the segmentation of cars.

While on the topic of segmentation, it is important to mention that the segmentation model is not perfect and that there are cases where objects are misclassified as well as cases where cars repeatedly go above and below the confidence threshold.

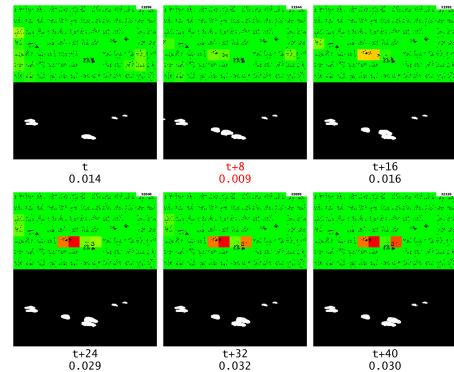


**Figure 9: Anomaly score output from Grid HTM.**

**4.0.1 Results.** We can see in Figure 9 that Grid HTM is detecting when segments begin and end, however it is not possible to use a

threshold value to isolate them, and they also have vastly different anomaly scores compared to each other. This is due to the way the aggregation function works, which means that the anomaly output is dependent on the physical size of the anomaly. It should also be noted that a moving average ( $n = 200$ ) was applied to smooth out the anomaly score output, otherwise the graph would be too noisy.

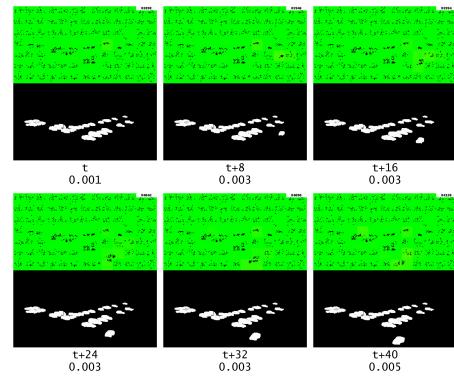
With the aggregation functions presented in this paper in mind, it is safe to conclude that looking at the anomaly score output is meaningless for complex data such as a surveillance video. This however does not mean that Grid HTM is completely useless, and this can be observed by looking at the visual output of Grid HTM. The visual output during which the first segment anomaly occurs can be seen in Figure 10. Here it is observed that Grid HTM correctly



**Figure 10: The first segment anomaly, which is marked with red text, and the corresponding changes detected by Grid HTM.**

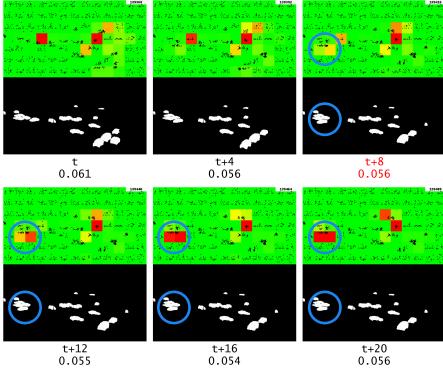
marks the sudden change of cars when the current segment ends and a new segment begins.

**4.0.2 Road.** In the original video, there is a road on which cars regularly drive. By observing the visual output, it becomes evident that after some time Grid HTM has mostly learned that behavior and does not report those moving cars as anomalies. This is shown in Figure 11.



**Figure 11: Visual output during when a car is driving along the main road.**

**4.0.3 Frame Repeat.** To prove that Grid HTM has learned that cars on the road should be moving, it is possible to look at the visual output during the period when the video is repeating the same frame and observe if the architecture marks the cars standing still on the road as anomalies. It can be observed in Figure 12 that the

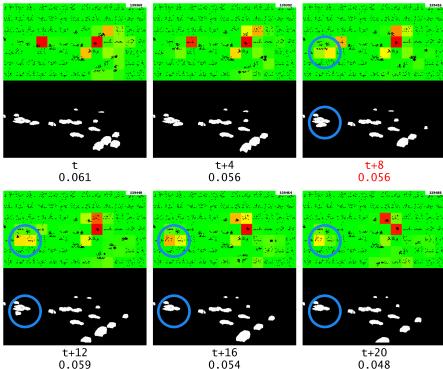


**Figure 12:** Anomaly output during the repeating frame, the start of the frame repeat is marked with red text. The blue circle highlights the object of interest.

cars along the main road are not marked as anomalies, but this could be attributed to the fact that there is a crossing there and that cars periodically have to stop at that point to let pedestrians cross.

On the other hand, when looking at the anomaly marked with a blue circle, the car on the road in the parking lot is marked as an anomaly that increases in severity as the time goes on during the frame repeat. The reason why that car causes an anomaly is because, unlike the cars on the main road, a car is rarely observed as standing still at that position.

To prove that the anomaly was actually directly caused by the repeating frame, and not just due to repeating the anomaly in time, it should be compared to the anomaly output if there was no repeating frame. It can be observed in Figure 13 that the anomaly

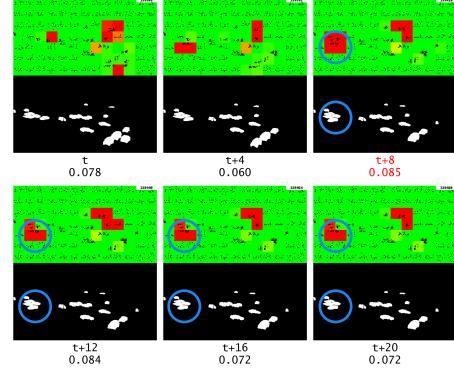


**Figure 13:** Anomaly output when there is no frame repeating, where it should have repeated is marked in red. The blue circle highlights the object of interest.

output is minor compared to when there was a repeating frame,

proving that the anomaly was indeed a product of the repeating frame and that Grid HTM was able to learn how objects should be moving in time.

Finally, it is interesting to look at how Grid HTM handles the repeating frames without multistep temporal patterns, which is shown in Figure 14. Unfortunately, simply disabling multistep tem-



**Figure 14:** Anomaly output during the repeating frame, the start of the frame repeat is marked with red text. The blue circle highlights the object of interest. This time without multistep temporal patterns.

poral patterns without adjusting the other TM parameters causes the same car to be marked as an anomaly before and during the frame repeat. In fact, as mentioned in Section 3.1.2, disabling multistep temporal patterns causes Grid HTM to be less noise tolerant which causes a lot more anomalies to be wrongly detected. This is evident in Figure 14, where a higher amount of severe anomalies can be observed compared to previous examples. This also highlights how sensitive HTM can be regarding parameters. The working code for Grid HTM and the parameters for the experiments conducted in for this paper can be found on GitHub.

## 5 CONCLUSION

We presented a novel method to perform anomaly detection in videos. Experiments showed that the proposed Grid HTM can be used for anomaly detection in complex videos such as surveillance footage. One of the most important future work would be to create a dataset with videos that are several days long and contain anomalies such as car accidents, jaywalking, and other similar anomalous behaviors. For Grid HTM, more time can be spent exploring other aggregation functions so that the aggregated anomaly score can be used more efficiently. Additionally, it would be a big benefit to create an algorithm which can decide the parameters for each cell during the calibration phase. It is also possible to improve explainability and robustness by implementing a measure of certainty for each cell.

Finally, experiments should be performed to validate the possibility of having the TM in each cell grow synapses to neighboring cells in order to solve the issue with unstable anomaly output, which was mentioned in Section 3.1.1.

## REFERENCES

- [1] Divyanshi Tewari. 2019. U.s. video surveillance market by component (solution, service, and connectivity technology), application (commercial, military & defense, infrastructure, residential, and others), and customer type (b2b and b2c): opportunity analysis and industry forecast, 2020–2027. (March 2019). <https://www.alliedmarketresearch.com/us-video-surveillance-market-A06741>.
- [2] J. Hawkins, S. Ahmad, S. Purdy, and A. Lavin. Biological and machine intelligence (bam). Initial online release 0.4, (2016). <https://numenta.com/resources/biological-and-machine-intelligence/>.
- [3] Guansong Pang, Chunhua Shen, Longbing Cao, and Anton Van Den Hengel. 2021. Deep learning for anomaly detection. *ACM Computing Surveys*, 54, 2, (April 2021), 1–38. ISSN: 1557-7341. doi: 10.1145/3439950.
- [4] Sijie Zhu, Chen Chen, and Waqas Sultani. 2020. Video anomaly detection for smart surveillance. (2020). arXiv: 2004.00222 [cs.CV].
- [5] Shivani Gupta and Atul Gupta. 2019. Dealing with noise problem in machine learning data-sets: a systematic review. *Procedia Computer Science*, 161, 466–474. The Fifth Information Systems International Conference, 23-24 July 2019, Surabaya, Indonesia. ISSN: 1877-0509. doi: <https://doi.org/10.1016/j.procs.2019.11.146>.
- [6] Dan Hendrycks and Thomas Dietterich. 2019. Benchmarking neural network robustness to common corruptions and perturbations. (2019). arXiv: 1903.12261 [cs.LG].
- [7] Chen Sun, Abhinav Shrivastava, Saurabh Singh, and Abhinav Gupta. 2017. Revisiting unreasonable effectiveness of data in deep learning era. In *Proceedings of the IEEE international conference on computer vision*, 843–852.
- [8] Alexander D'Amour, Katherine Heller, Dan Moldovan, Ben Adlam, Babak Alipanahi, Alex Beutel, Jonathan Deaton, Jacob Eisenstein, Matthew Hoffman, Farhad Hormozdiari, Neil Houlsby, Shaobo Hou, Ghassen Jerfel, Alan Karthikesalingam, Mario Lucic, Yian Ma, Cory McLean, Diana Mincu, and D. Sculley. 2020. Underspecification presents challenges for credibility in modern machine learning. (November 2020).
- [9] Alejandro Barredo Arrieta, Natalia Díaz-Rodríguez, Javier Del Ser, Adrien Bennetot, Siham Tabik, Alberto Barbado, Salvador Garcia, Sergio Gil-Lopez, Daniel Molina, Richard Benjamins, Raja Chatila, and Francisco Herrera. 2020. Explainable artificial intelligence (xai): concepts, taxonomies, opportunities and challenges toward responsible ai. *Information Fusion*, 58, 82–115. ISSN: 1566-2535. doi: <https://doi.org/10.1016/j.inffus.2019.12.012>.
- [10] Y. Zou, Y. Shi, Y. Wang, Y. Shu, Q. Yuan, and Y. Tian. 2018. Hierarchical temporal memory enhanced one-shot distance learning for action recognition. In *2018 IEEE International Conference on Multimedia and Expo (ICME)*, 1–6. doi: 10.1109/ICME.2018.8486447.
- [11] David McDougall (ctrl-z 9000-times). 2019. (September 2019). <https://github.com/htm-community/htm.core/issues/259#issuecomment-53333336>.
- [12] 2019. *Extended hierarchical temporal memory for motion anomaly detection: proceedings of the ninth annual meeting of the bica society*. (January 2019), 69–81. ISBN: 978-3-319-99315-7. doi: 10.1007/978-3-319-99316-4\_10.
- [13] Subutai Ahmad, Alexander Lavin, Scott Purdy, and Zuhra Agha. 2017. Unsupervised real-time anomaly detection for streaming data. *Neurocomputing*, 262, 134 –147. Online Real-Time Learning Strategies for Data Streams. ISSN: 0925-2312. doi: <https://doi.org/10.1016/j.neucom.2017.04.070>.
- [14] Ethan Rublee, Vincent Rabaud, Kurt Konolige, and Gary Bradski. 2011. Orb: an efficient alternative to sift or surf. In *2011 International Conference on Computer Vision*, 2564–2571. doi: 10.1109/ICCV.2011.6126544.
- [15] Sangmin Oh, Anthony Hoogs, Amitha Perera, Naresh Cuntoor, Chia-Chih Chen, Jong Taek Lee, Saurajit Mukherjee, J. K. Aggarwal, Hyungtae Lee, Larry Davis, Eran Swears, Xioyang Wang, Qiang Ji, Kishore Reddy, Mubarak Shah, Carl Vondrick, Hamed Pirsiavash, Deva Ramanan, Jenny Yuen, Antonio Torralba, Bi Song, Anesco Fong, Amit Roy-Chowdhury, and Mita Desai. 2011. A large-scale benchmark dataset for event recognition in surveillance video. In *CVPR 2011*, 3153–3160. doi: 10.1109/CVPR.2011.5995586.
- [16] Alexander Kirillov, Yuxin Wu, Kaiming He, and Ross Girshick. 2020. Pointrend: image segmentation as rendering. (2020). arXiv: 1912.08193 [cs.CV].
- [17] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2015. Deep residual learning for image recognition. (2015). arXiv: 1512.03385 [cs.CV].
- [18] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. 2009. Imagenet: a large-scale hierarchical image database. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*, 248–255. doi: 10.1109/CVPR.2009.5206848.
- [19] Ayoola Olafenwa. 2021. Simplifying object segmentation with pixellib library. In.