



Felix Anda, Mark Scanlon, Nhien-An Le-Khac
School of Computer Science, University College Dublin, Ireland.
felix.andabasabe@ucdconnect.ie, {mark.scanlon, an.lekhac}@ucd.ie

Introduction

The exponential growth of data storage is a problem that must be tackled in a cyber realm era where criminal activities are perpetrated constantly. Data acquisition entails time consuming analysis that demands digital forensic experts. Ideally, seized devices should be explored immediately due to the potentially urgent need of evidence to prosecute a cybercrime that could be a matter of life or death [1], e.g., child exploitation and human trafficking. The shortage of digital forensic specialists and the eminent backlog of apprehended devices that haven't been processed is a burden [2]. However, it awakens the need to automate the analysis procedure of a digital investigation process model with both machine learning and computer vision techniques.

Automated machine learning and computer vision-based digital evidence classification and identification techniques could significantly ease the backlog and data deduplication methods in a cloud environment could assist substantially in the creation of a sifted "ready to analyse" image which resembles a sample of the entire disk drive and is proved to be a subset of the generated forensically sound image. The extraction of a model of the forensically sound image is compulsory due to the amount of computer power required to examine a device for probative data which is a problem that must be addressed when artificial intelligence takes place.

In Figure 1, the 3 different hashes corresponding to different stages can be seen when a whole disk is acquired (a), when a sample is taken (b) and the sample which is to be analysed (c).

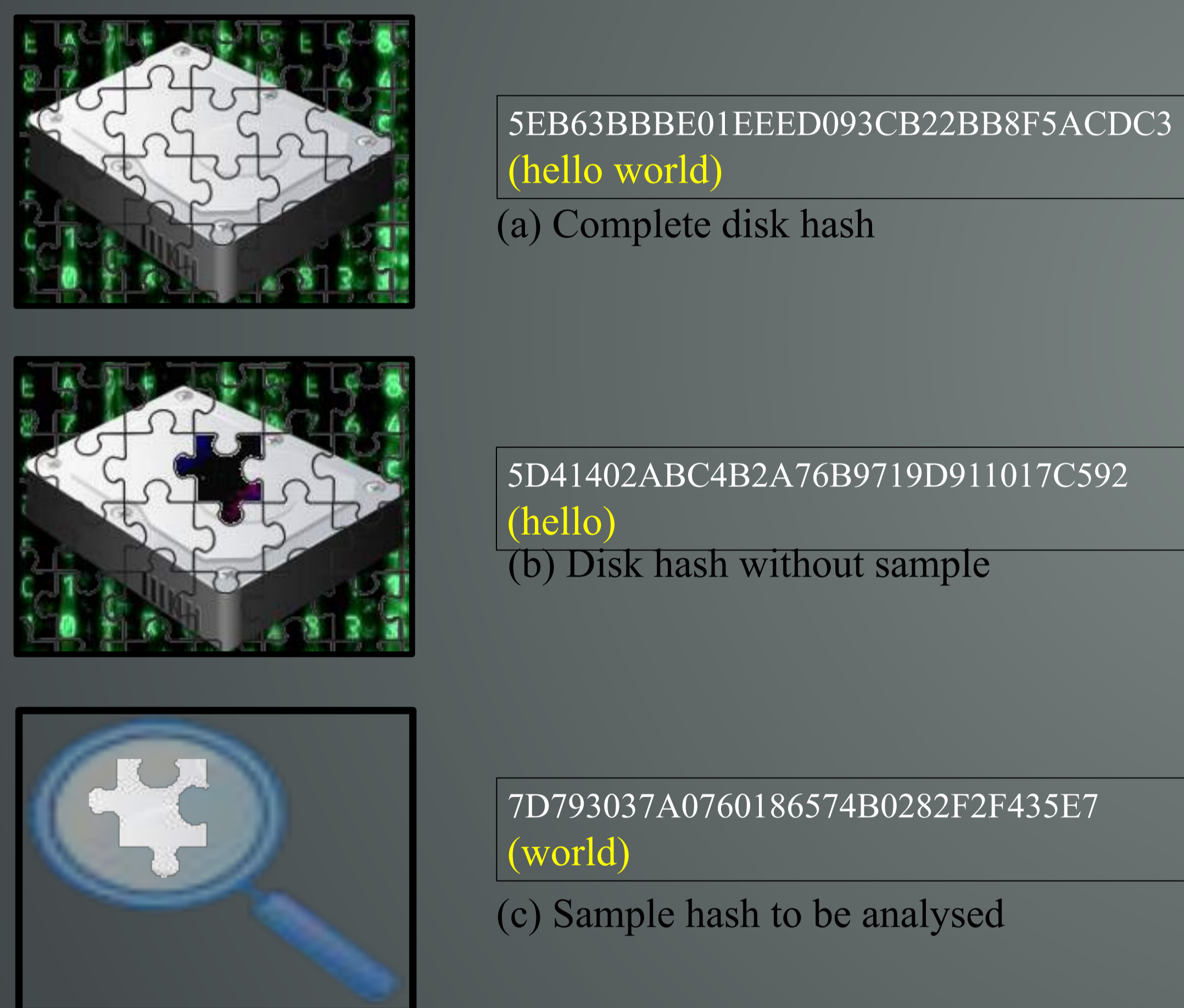


Figure 1. Example of MD5 Hashes for Disk with a Sample

Objectives of This Work

- Determine if the fusion of data deduplication, machine learning and computer vision techniques can assist a triage process model to alleviate the load for the forensic investigator and decrease processing times of seized devices in a digital forensic lab.
- Carry out a performance evaluation of computer vision online tools provided by Microsoft, Amazon, IBM and Kairos. Include offline frameworks like Caffe and Weka to gain a better understanding of the capabilities of both frameworks and how they interact with diverse pre-existing models and newly trained models.
- Train a model that can classify illegal content and predict the output of further incriminating data that hasn't been encountered previously .

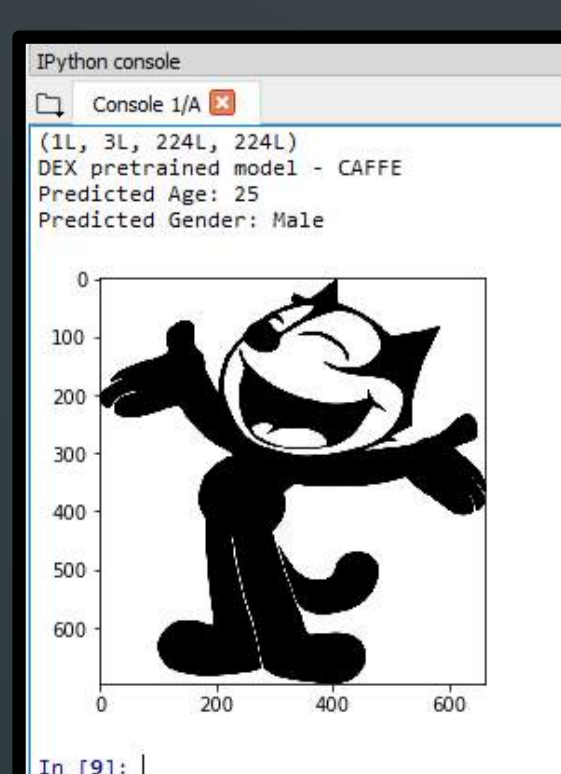


Figure 2. Age and Gender Prediction using Caffe [3]

Conclusions

Relevant classification to computer forensics is possible with the use of different pre-trained models that can tag images in a broad range. Labels include gender, ethnicity, age and adult content detection identification [3, 4, 5]. The mean absolute error (MAE) for age prediction has been reduced to approximately 3 years [6], which is highly encouraging research in this area due to its capability of outperforming age estimations accomplished by humans. In Figure 2 the reliability of machine learning tools is shown.

The use of digital forensics as a service (DFaaS) and cloud forensic data deduplication techniques enable the creation of a sample disk drive which is considered a subset of a forensically sound image and can be proved to be part of the original disk image. Using big data to store good know files and bad know files is valuable due to the ability to trigger faster a warning about implicating files that are being processed for the recreation of a copy [7]. An entire hard disk reconstruction is executed in less time on the server side than the whole duplication of a disk when reading the entire image and creating the bit to bit sequential copy [2]. The creation of a multithread software that works in parallel while the duplicate image is in progress allows the files to be analysed and classified for the ease of carrying out the investigation. Figure 3 depicts the cloud solution.

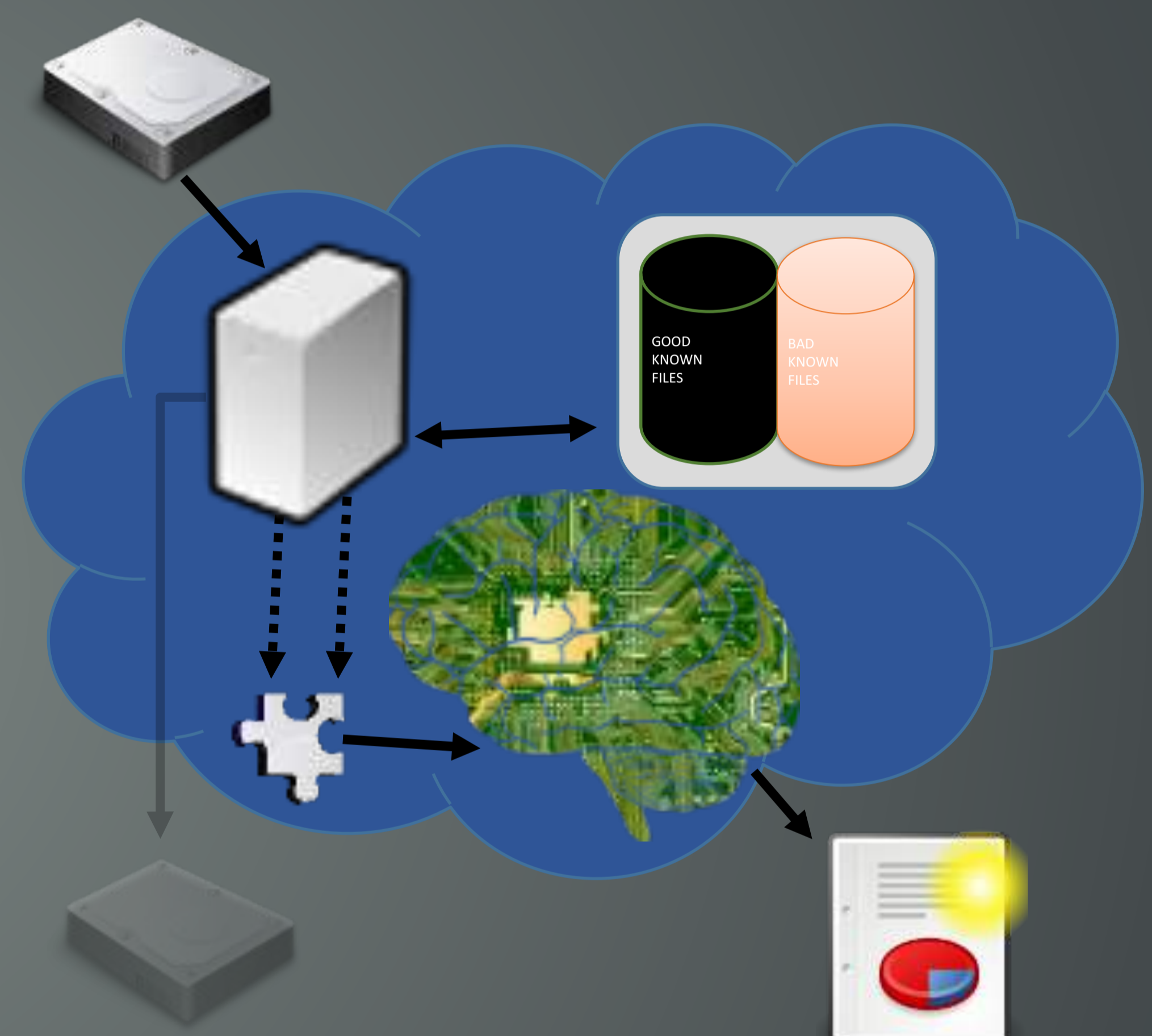


Figure 3. Cloud Solution with Deduplication & AI Techniques

References

1. Rogers, M. K., Goldman, J., Mislán, R., Wedge, T., & Debrota, S. (2006). *Computer Forensics Field Triage Process Model*. Proceedings of the Conference on Digital Forensics, Security and Law, , 27-40.
2. Scanlon M. *Battling the Digital Forensic Backlog through Data Deduplication*. In: Proceedings of the 6th IEEE International Conference on Innovative Computing Technologies (INTECH 2016). Dublin, Ireland: IEEE;2016.
3. R. Rothe, R. Timofte, and L. Van Gool, *DEX: Deep EXpectation of Apparent Age from a Single Image*, Proc. IEEE Int. Conf. Comput. Vis., vol. 2016–Febru, pp. 252–257, 2016.
4. B. Zhou, A. Lapedriza, J. Xiao, A. Torralba, and A. Oliva, *Learning Deep Features for Scene Recognition using Places Database*, Adv. Neural Inf. Process. Syst. 27, pp. 487–495, 2014.
5. Levi, G., & Hassner, T. (2015). *Age and gender classification using convolutional neural networks*. 2015 IEEE Conference on Computer Vision and Pattern Recognition Workshops (CVPRW), 34–42.
6. Xing, J., Li, K., Hu, W., Yuan, C., & Ling, H. (2017). *Diagnosing deep learning models for high accuracy age estimation from a single image*. Pattern Recognition, 66(November 2016), 106–116.
7. Carrier, Brian. *File system forensic analysis*. Addison-Wesley Professional, 2005.

Acknowledgements

I gratefully acknowledge the funding received towards my PhD from the UCD School of Computer Science.