

# A Big Data Framework Approach in Healthcare Industry

K.Rajesh Khanna<sup>1</sup>, Dr.D.Suresh Babu<sup>2</sup>, Dr. Vaibhav Bansal<sup>3</sup>

<sup>1</sup> Research Scholar, Department of Computer science, OPJS University, Rajasthan.

<sup>2</sup> Professor & Head, Computer Science Department, Kakatiya Government College, Warangal.

<sup>3</sup> Associate Professor, Department of Computer science, OPJS University, Rajasthan.

**ABSTRACT:** It has provided tools to accumulate, manage, analyze, and assimilate large volumes of disparate, structured, and unstructured data produced by current healthcare systems. Big data analytics has been recently applied towards aiding the process of care delivery and disease exploration. However, the adoption rate and research development in this space is still hindered by some fundamental problems inherent within the big data paradigm. In this paper, we discuss some of these major challenges with a focus on three upcoming and promising areas of medical research: image, signal, and genomics based analytics.

**KEYWORDS-** Data Acquisition, Big data analytics, Data Storage and Retrieval.

## I. INTRODUCTION

Although there is already a huge amount of healthcare data around the world and while it is growing at an exponential rate, nearly all of the data is stored in individual silos. Data collected by a GP clinic or by a hospital is mostly kept within the boundaries of the healthcare provider. Moreover, data stored within a hospital is hardly ever integrated across multiple IT systems. For example, if we consider all the available data at a hospital from a single patient's perspective, information about the patient will exist in the EMR system, laboratory, imaging system and prescription databases. Information describing which doctors and nurses attended to the specific patient will also exist. However, in the vast majority of cases, every data source mentioned here is stored in separate silos. Thus deriving insights and therefore value from the aggregation of these data sets is not possible at this stage. It is also important to realize that in today's world a patient's medical data does not only reside within the boundaries of a healthcare provider. The medical insurance and pharmaceuticals industries also hold information about specific claims and the characteristics of prescribed drugs respectively.

Increasingly, patient-generated data from IoT devices such as fitness trackers, blood pressure monitors and weighing scales are also providing critical information about the day-to-day lifestyle characteristics of an individual. Insights derived from such data generated by the linking among EMR data, vital data, laboratory data, medication information, symptoms (to mention some of these) and their aggregation, even more with doctor notes, patient discharge letters, patient diaries, medical publications, namely linking structured with unstructured data, can be crucial to design coaching programmes that would help improve people's lifestyles and eventually reduce incidences of chronic disease, medication and hospitalization.

As the healthcare sector transitions from a volume to value-based care model, it is essential for different stakeholders to get a complete and accurate understanding of treatment trajectories of specific patient populations. The only way to achieve this is to be able to aggregate the disparate data sources not just within a single hospital's/GP clinic's IT infrastructure, but also across multiple healthcare providers, other healthcare players (e.g. insurance & pharma) and even consumer-generated data. Such unified data sets would benefit not only every player within the healthcare industry (thus allowing better quality care and access to healthcare at lower costs), but would also most importantly benefit the patient by providing first time right treatment, based on a sustainable pricing model.

However, achieving such a vision which involves the integration of such disparate healthcare datasets (in terms of data granularity, quality, type (e.g. ranging from free text, images, (streaming) sensor data to structured datasets) poses major legal, business and technical challenges from a data perspective, in terms of the volume, variety, veracity and velocity of the data sets. The only way to successfully address these challenges is to utilise Big Data technologies. "Big

data” has a wide range of definitions in health research<sup>1314</sup>. However, a viable definition of what big data means for health is the following: “Big data in health” encompasses high volume, high diversity biological, clinical, environmental, and lifestyle information collected from single individuals to large cohorts, in relation to their health and wellness status, at one or several time points. More general definition of Big Data, refers to “datasets whose size is beyond the ability of typical database software tools to capture, store, manage and analyse”. (McKinsey Global Institute). This definition puts the accent on the size/volume aspect but, as we described above, the dimensions are many: variety (handling with a multiplicity of types, sources and format), data veracity (related to the quality and validity of these data), and data velocity (availability in real time). In addition, there are other factors that should also be considered such as data trustworthiness, data protection, and privacy (due to the sensitivity of data managed). All these aspects lead to the need for new algorithms, techniques and approaches to handle these new challenges.

## II. SUGGESTED WORK

Streaming data analytics in healthcare may be defined as a scientific use of continuous waveform (signal varying in opposition to time) and related medical record statistics evolved via applied analytical disciplines (e.g., statistical, quantitative, contextual, cognitive, and predictive) to pressure choice making for patient care. The analytics workflow of actual-time streaming waveforms in scientific settings can be broadly defined the usage of Figure 1. Firstly, a platform for streaming data acquisition and ingestion is needed which has the bandwidth to address multiple waveforms at different fidelities. Integrating those dynamic waveform facts with static statistics from the EHR is a key component to provide situational and contextual attention for the analytics engine. Enriching the data ate up by way of analytics not most effective makes the device more strong, but also helps balance the sensitivity and specificity of the predictive analytics. The specifics of the signal processing will in large part depend upon the form of sickness cohort underneath investigation. A sort of signal processing mechanisms can be applied to extract a mess of target capabilities which can be then fed on by way of a pretrained

machinemastering version to supply an actionable insight. These actionable insights may want to both be diagnostic, predictive, or prescriptive. These insights may want to in addition be designed to cause other mechanisms together with alarms and notification to physicians.

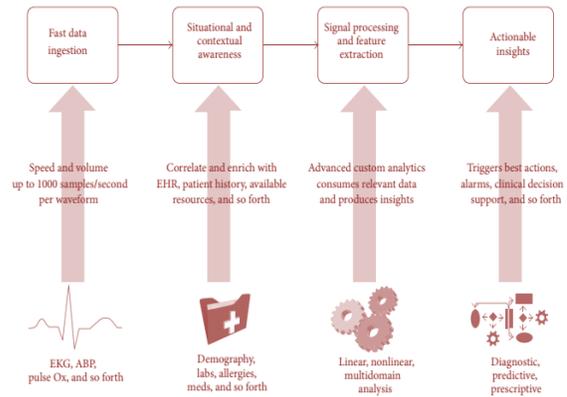


Figure 1: Generalized analytic workflow using streaming healthcare data

Harmonizing such non-stop waveform statistics with discrete records from different resources for locating necessary affected person statistics and conducting research in the direction of development of next technology diagnoses and remedies may be a daunting mission. For bed-side implementation of such systems in scientific environments, there are numerous technical concerns and requirements that want to be designed and implemented at device, analytic, and scientific ranges. The following subsections provide a top level view of different challenges and current processes within the development of monitoring systems that devour both excessive constancy waveform data and discrete data from noncontinuous assets.

**Data Acquisition.** Historically streaming data from continuous physiological signal acquisition devices was rarely stored. Even if the option to store this data were available, the length of these data captures was typically short and downloaded only using proprietary software and data formats provided by the device manufacturers. Although most major medical device manufacturers are now taking steps to provide interfaces to access live streaming data from their devices, such data in motion very quickly poses archetypal big data challenges. The fact that there are also governance challenges such as lack of data

protocols, lack of data standards, and data privacy issues is adding to this. On the other side there are many challenges within the healthcare systems such as network bandwidth, scalability, and cost that have stalled the widespread adoption of such streaming data collection. This has allowed way for system-wide projects which especially cater to medical research communities. Research community has interest in consuming data captured from live monitors for developing continuous monitoring technologies. There have been several indigenous and off-the-shelf efforts in developing and implementing systems that enable such data capture. There are also products being developed in the industry that facilitate device manufacturer agnostic data acquisition from patient monitors across healthcare systems.

**Data Storage and Retrieval.** With large volumes of streaming data and other patient information that can be gathered from clinical settings, sophisticated storage mechanisms of such data are imperative. Since storing and retrieving can be computational and time expensive, it is key to have a storage infrastructure that facilitates rapid data pull and commits based on analytic demands. With its capability to store and compute large volumes of data, usage of systems such as Hadoop, MapReduce, and MongoDB is becoming much more common with the healthcare research communities. MongoDB is a free cross-platform document-oriented database which eschews traditional table-based relational database. Typically each health system has its own custom relational database schemas and data models which inhibit interoperability of healthcare data for multi-institutional data sharing or research studies. Furthermore, given the nature of traditional databases integrating data of different types such as streaming waveforms.

This is where MongoDB and other document-based databases can provide high performance, high availability, and easy scalability for the healthcare data needs. Apache Hadoop is an open source framework that allows for the distributed processing of large datasets across clusters of computers using simple programming models. It is a highly scalable platform which provides a variety of computing modules such as MapReduce and Spark. For performing analytics on continuous

telemetry waveforms, a module like Spark is especially useful since it provides capabilities to ingest and compute on streaming data along with machine learning and graphing tools. Such technologies allow researchers to utilize data for both realtime as well as retrospective analysis, with the end goal to translate scientific discovery into applications for clinical settings in an effective manner.

#### IV. CONCLUSION

Big data analytics which leverages legions of disparate, structured, and unstructured data sources is going to play a vital role in how healthcare is practiced in the future. One can already see a spectrum of analytics being utilized, aiding in the decision making and performance of healthcare personnel and patients. Here we focused on three areas of interest: medical image analysis, physiological signal processing, and genomic data processing. The exponential growth of the volume of medical images forces computational scientists to come up with innovative solutions to process this large volume of data in tractable timescales. The trend of adoption of computational systems for physiological signal processing from both research and practicing medical professionals is growing steadily with the development of some very imaginative and incredible systems that help save lives.

#### REFERENCES

- [1] C. F. Mackenzie, P. Hu, A. Sen et al., "Automatic pre-hospital vital signs waveform and trend data capture fills quality management, triage and outcome prediction gaps," AMIA Annual Symposium Proceedings, vol. 2008, pp. 318–322, 2008.
- [2] M. Bodo, T. Settle, J. Royal, E. Lombardini, E. Sawyer, and S.W. Rothwell, "Multimodal noninvasive monitoring of soft tissue wound healing," Journal of Clinical Monitoring and Computing, vol. 27, no. 6, pp. 677–688, 2013.
- [3] P. Hu, S. M. Galvagno Jr., A. Sen et al., "Identification of dynamic prehospital changes with continuous vital signs acquisition," Air Medical Journal, vol. 33, no. 1, pp. 27–33, 2014.

- [4] D. Apiletti, E. Baralis, G. Bruno, and T. Cerquitelli, "Real-time analysis of physiological data to support medical applications," *IEEE Transactions on Information Technology in Biomedicine*, vol. 13, no. 3, pp. 313–321, 2009.
- [5] J. Chen, E. Dougherty, S. S. Demir, C. P. Friedman, C. S. Li, and S. Wong, "Grand challenges for multimodal bio-medical systems," *IEEE Circuits and Systems Magazine*, vol. 5, no. 2, pp. 46–52, 2005.
- [6] N. Menachemi, A. Chukmaitov, C. Saunders, and R. G. Brooks, "Hospital quality of care: does information technology matter? The relationship between information technology adoption and quality of care," *Health Care Management Review*, vol. 33, no. 1, pp. 51–59, 2008.
- [7] C. M. DesRoches, E. G. Campbell, S. R. Rao et al., "Electronic health records in ambulatory care—a national survey of physicians," *The New England Journal of Medicine*, vol. 359, no. 1, pp. 50–60, 2008.
- [8] J. S. McCullough, M. Casey, I. Moscovice, and S. Prasad, "The effect of health information technology on quality in U.S. hospitals," *Health Affairs*, vol. 29, no. 4, pp. 647–654, 2010.
- [9] J. M. Blum, H. Joo, H. Lee, and M. Saeed, "Design and implementation of a hospital wide waveform capture system," *Journal of Clinical Monitoring and Computing*, vol. 29, no. 3, pp. 359–362, 2015.
- [10] D. Freeman, "The future of patient monitoring," *Health Management Technology*, vol. 30, no. 12, article 26, 2009.
- [11] B. Muhsin and A. Sampath, "Systems and methods for storing, analyzing, retrieving and displaying streaming medical data," *US Patent 8,310,336*, 2012.
- [12] D. Malan, T. Fulford-Jones, M. Welsh, and S. Moulton, "Codeblue: an ad hoc sensor network infrastructure for emergency medical care," in *Proceedings of the International Workshop on Wearable and Implantable Body Sensor Networks*, vol. 5, London, UK, 2004.
- [13] A. Page, O. Kocabas, S. Ames, M. Venkatasubramanian, and T. Soyata, "Cloud-based secure health monitoring: optimizing fully-homomorphic encryption for streaming algorithms," in *Proceedings of the IEEE Globecom Workshops (GC Wkshps '14)*, pp. 48–52, IEEE, Austin, Tex, USA, December 2014.
- [14] J. Bange, M. Gryzwa, K. Hoyme, D. C. Johnson, J. LaLonde, and W. Mass, "Medical data transport over wireless life critical network," *US Patent 7,978,062*, 2011.
- [15] N. Kara and O. A. Dragoi, "Reasoning with contextual data in telehealth applications," in *Proceedings of the 3rd IEEE International Conference on Wireless and Mobile Computing, Networking and Communications (WiMoB '07)*, p. 69, IEEE, October 2007.
- [16] G. Li, J. Liu, X. Li, L. Lin, and R. Wei, "A multiple biomedical signals synchronous acquisition circuit based on over-sampling and shaped signal for the application of the ubiquitous healthcare," *Circuits, Systems, and Signal Processing*, vol. 33, no. 10, pp. 3003–3017, 2014.
- [17] A. Bar-Or, J. Healey, L. Kontothanassis, and J. M. van Tong, "BioStream: a system architecture for real-time processing of physiological signals," in *Proceedings of the 26th Annual International Conference of the IEEE Engineering in Medicine and Biology Society (EMBC' 04)*, vol. 2, pp. 3101–3104, September 2004.
- [18] W. Raghupathi and V. Raghupathi, "Big data analytics in healthcare: promise and potential," *Health Information Science and Systems*, vol. 2, article 3, 2014.
- [19] S. Ahmad, T. Ramsay, L. Huebsch et al., "Continuous multiparameter heart rate variability analysis heralds onset of sepsis in adults," *PLoS ONE*, vol. 4, no. 8, Article ID e6642, 2009.
- [20] A. L. Goldberger, L. A. Amaral, L. Glass et al., "Physiobank, physiokit, and physionet components of a new research resource for complex physiologic signals," *Circulation*, vol. 101, no. 23, pp. e215–e220, 2000.
- [21] E. J. Siachalou, I. K. Kitsas, K. J. Panoulas et al., "ICASP: an intensive-care acquisition and signal

processing integrated framework,” *Journal of Medical Systems*, vol. 29, no. 6, pp. 633–646, 2005.

[22] M. Saeed, C. Lieu, G. Raber, and R. G. Mark, “Mimic ii: amassive temporal icu patient database to support research in intelligent patient monitoring,” in *Proceedings of the Computers in Cardiology*, pp. 641–644, IEEE, September 2002.

[23] A. Burykin, T. Peck, and T. G. Buchman, “Using ‘off-the-shelf’ tools for terabyte-scale waveform recording in intensive care: computer system design, database description and lessons learned,” *Computer Methods and Programs in Biomedicine*, vol. 103, no. 3, pp. 151–160, 2011.

[24] G. Adrian, G. E. Francisco, M. Marcela, A. Baum, L. Daniel, and G. B. de Quiros Fern ´ an, “Mongodb: an open source alternative for HL7-CDA clinical documents management,” in *Proceedings of the Open Source International Conference (CISL ’13)*, Buenos Aires, Argentina, 2013.

[25] K. Kaur and R. Rani, “Managing data in healthcare information systems: many models, one solution,” *Computer*, vol. 48, no. 3, pp. 52–59, 2015.