

Privacy Preserving Data Access in Edge Computing

Ujjwal Goel

Advisor: **Dr. Ferdous Ahmed Barbhuiya**

Department of Computer Science and Engineering
Indian Institute of Information Technology, Guwahati

November 2022

Declaration

I hereby declare that except where specific reference is made to the work of others, the contents of this dissertation are original and have not been submitted in whole or in part for consideration for any other degree or qualification in this, or any other university. This dissertation is my own work and contains nothing which is the outcome of work done in collaboration with others, except as specified in the text and Acknowledgements.

Ujjwal Goel
November 2022

Acknowledgements

I would like to acknowledge my supervisor, Dr. Ferdous Ahmed Barbhuiya for his constant guidance and support throughout the project.

I would also like to thank Dr. Rakesh Matam, Meghali Nandi, Mehbub Alam and Aditya Sharma for their help and guidance.

Abstract

In this project, we are working on the problem of secure data access in the edge computing environment. Specifically, we are interested in the methods for resource discovery in edge while preserving the data privacy. We study how Distributed Hash Tables (DHTs) and Learned Index Structures have been used by researchers to solve similar problems in other domains. In our work, we are considering using the XOR based routing algorithm used in Kademlia, a kind of Distributed Hash Table. Though Kademlia is generally used for the storage and retrieval of data, we are considering using only the routing algorithm used in Kademlia. For retrieval itself, we are interested in Learned Index Structures (LIS), which have been gaining popularity lately.

Table of contents

List of figures	xi
1 Introduction	1
1.1 Edge computing	1
1.1.1 Why use edge computing in place of cloud computing?	1
1.1.2 Challenges with edge computing	1
1.1.3 Edge computing for vehicles	2
1.2 Problem statement	3
2 Literature review	5
2.1 Distributed Hash Tables	6
2.1.1 How does Kademlia work?	6
2.2 Learned indexes	8
3 Proposed Methodology	9
3.1 Using Distributed Hash Tables for Edge Computing	9
3.2 Using Learned Indices for Edge Computing	10
3.3 Conclusion	10
References	11

List of figures

- 1.1 Vehicular edge network architecture 2
- 2.1 Kademlia 7
- 2.2 Recursive Model Index 8

Chapter 1

Introduction

1.1 Edge computing

Edge computing is a newer computation paradigm as compared to cloud computing. The difference between them is that edge computing brings computation closer to the edge of the network. As more of the data is being generated near the edge of the network, edge computing enables the data to be processed near the data source. [17]

1.1.1 Why use edge computing in place of cloud computing?

The most important reason for using edge computing is that it reduces the latency as the data has to travel a shorter distance. [17] In addition to this, with the emergence of Internet of Things, the data generated at the edge is expected to grow rapidly. For example, it is expected that the autonomous vehicles will generate data in the order of 1Tb per second. In such a case, it may not be possible for the cloud to keep pace with the data generated at the edge. Edge computing can help in this situation as it can process the data at the edge itself.

1.1.2 Challenges with edge computing

One of the major challenges with edge computing is that the edge devices are resource constrained. They have limited processing power, memory and storage. Often, this problem is solved using offloading. A node that is low on resources may enlist the other nodes for help. For example, if a node is low on storage space, then the data generated on this node or the data this node was supported to store may be offloaded to other nodes in the network. Similarly, data processing can be offloaded to other nodes.

The other challenge is that, no node in the edge network has complete knowledge of the edge network. This is in contrast with cloud computing where a central node is responsible for all the processing and does not need to coordinator with others. As such, it can become difficult to do any computation on the edge as the nodes may not be having all the information required for the computation.

1.1.3 Edge computing for vehicles

Researchers have been working on using edge computing for vehicles for a while now. In autonomous vehicles, the decision must be taken quickly as there is very low tolerance for latency. Further, the data that a vehicle generate is rarely needed outside the geographical vicinity in which the vehicle is present. Using edge computing makes the common use-case faster. As such, edge computing can be used to improve the performance of the autonomous vehicles.

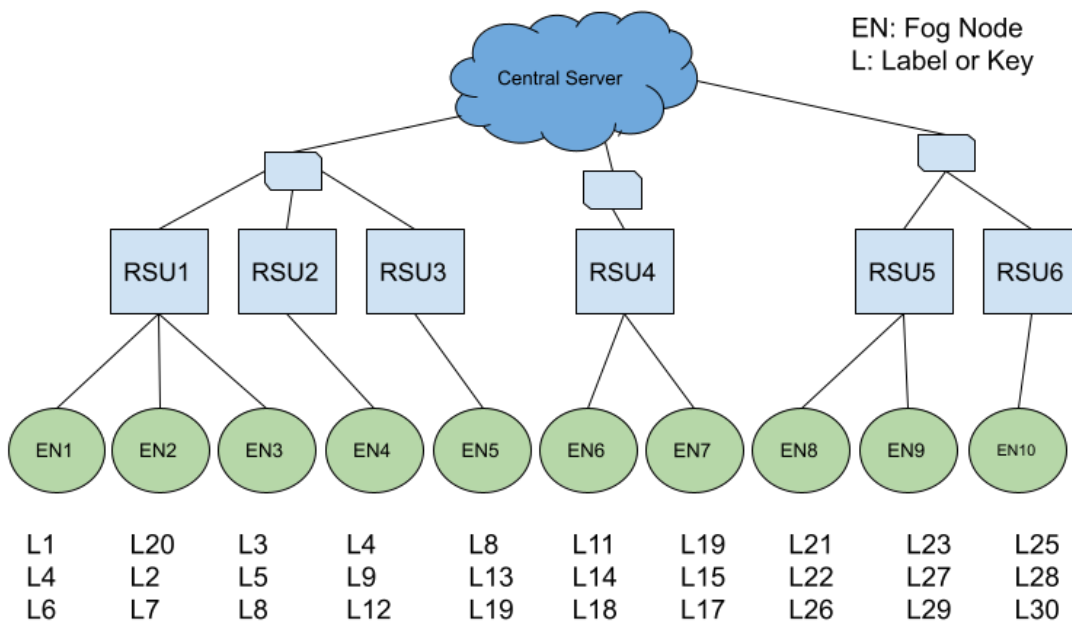


Fig. 1.1 Vehicular edge network architecture

Vehicular edge networks usually consists of a central cloud to which a number of gateway devices called roadside units (RSUs) are connected. The RSUs are further connected to the vehicles. Vehicular edge networks are highly dynamic and vehicles change their position quite often. Because of this, hand-offs may be needed very frequently.

1.2 Problem statement

Because of the issues described above, we need to have some kind of method to enable nodes to find out where the required data may be present in the edge network.

We are trying to come up with a method for answering resource queries in edge networks while keeping in mind the constraints on resources that come with edge computing.

Requirements for the solution:

- Data privacy is maintained and the data itself should not be used for indexing or querying
- The solution should be scalable
- The algorithm should preferably be online – it should be able to index or answer queries for the data that is being generated

Further, it is not possible to simply boardcast the query to all the nodes in the network, as it would very inefficient. Doing so will lead to privacy issues, as the queries would be known to every node in the network. Also, caching the data on other nodes should be avoided because of the privacy issues.

Use cases

The problem of resource or data discovery is by no mean unique to edge computing or IoT. Even in the computer architecture, Flat Cache-Only Memory Architecture (COMA) faces the similar problem. In Flat COMA, memory lines can freely migrate. Whenever there is a cache miss, the memory line must be located in the network and migrated to where it is needed. [6]

But in our work, we are trying to solve this problem in the context of edge computing, where the nodes may be static or dynamic. Weather crowd-sensing would need only static nodes, but in vehicular edge networks, the nodes are dynamic.

Chapter 2

Literature review

Numerous surveys have been done on edge computing as well as data access in IoT or distributed systems. These include surveys by Zineddine Kouahla et al. [20] and Yasmin Fathy et al. [18].

Yasmin Fathy et al. [18] prescribes the following taxonomy for queries in IoT:

- **Data queries:** Here, we are only interested in the data, regardless of its source. They may be indexed by:
 - Spatial features: The data may be indexed using spatial attributes, e.g. geo-coordinates. Space filling curves may be used to reduce dimensionality.
 - Thematic features: data is indexed based on the keywords or its attributes.
 - Temporal features: data is indexed based on the time at which it was generated.
- **Resource queries:** Rather than finding the data stored on the node, the nodes are only interested in finding a physical resource on the network.
- **Higher-level abstractions:** the queries are based on higher-level abstractions like events, activities and patterns which require data from several nodes to be combined and processed.

Zineddine Kouahla et al. [20] also classifies indexing techniques as below:

- Indexing structures in a Multidimensional Space
 - Hashing-based techniques
 - Tree-based techniques
 - Bitmap-based techniques

- Indexing structures in a Metric Space
 - Tree-based techniques

In this work, the focus is on resource queries. The advantage of focusing on resource queries over data queries is that the data queries can be trivially answered once the resource queries are answered. Further, it allows the nodes to implement some kind of authentication mechanism to ensure that the data is not publically available.

In 2012, Federica Paganelli and David Parlanti [5] proposed a scheme that used DHTs for service discovery in Internet of Things. Their method allowed multi-attribute queries as well as range queries. The method consists of three layers. They first used a space-filling curve to linearize and map the multidimensional domain into a one-dimensional one. Then a Prefix-Hash Table is used on top of a generic DHT with standard interface for get and put. At the end, the DHT is an implementation based on the Kademlia algorithm. Federica Paganelli and David Parlanti [5], however, did not apply DHTs to edge networks.

Researchers have also attempted to use Gaussian Mixture Models (GMM) for searching and indexing in edge computing.

2.1 Distributed Hash Tables

The Distributed Hash Tables (DHTs) are similar to the conventional hash tables, except that the key-value pairs are stored across multiple nodes in the network. Distributed hash tables are highly scalable and fault-tolerant. Like traditional hash-tables, DHTs also support “get” and “put” operations. They are commonly used in the peer-to-peer networks. There exist multiple variants of DHTs including Chord [1], Tapestry [19] and Kademlia [14].

In DHTs, the data is placed deterministically on the nodes in the network. DHTs are self-organizing and self-healing.

Of these DHTs, Kademlia is the most widely used. It is used in BitTorrent [4] as well as in Ethereum. Every key-value pair is stored with multiple nodes in the network. This is done to ensure that the network is robust to node failures. Kademlia is not affected if a node leaves or joins the network. Each node joining the Kademlia randomly generates a 160-bit identifier. [14]

2.1.1 How does Kademlia work?

Whenever a key-value pair is to be stored in the network, a 160 bit hash of the key is calculated. Then, nodes with the identifier closest to the hash are selected for storing the

key-value pair. For measuring the distance and finding the closest nodes to a given key, Kademia uses XOR metric. Even though it is non-euclidean, the XOR metric satisfies all four properties of a distance metric. [3]

Similarly, when finding the value for a given key, the nodes with the identifier closest to the hash of the key are selected. An attempt is then made to retrieve the corresponding values from these nodes.

The authors also propose an XOR based routing algorithm for finding the node with the identifier closest to a given key. With the algorithm proposed, in each iteration, we get one bit closer to the target node. Thus, if the network size is n , then the algorithm takes $O(\log(n))$ iterations to find the target node.

Nodes in Kademia can be visualized as leaves in a binary tree, where the position of each of the nodes is given by its identifier. In the algorithm, each node is required to maintain a routing table. Each of the node is expected to store information about k nodes from each of the maximal tree which it is not a part of in its routing table. The information includes the identifier, IP address and the port number. The requirement of storing information about k nodes from each is to bring redundancy and to ensure that the network is robust to node failures.

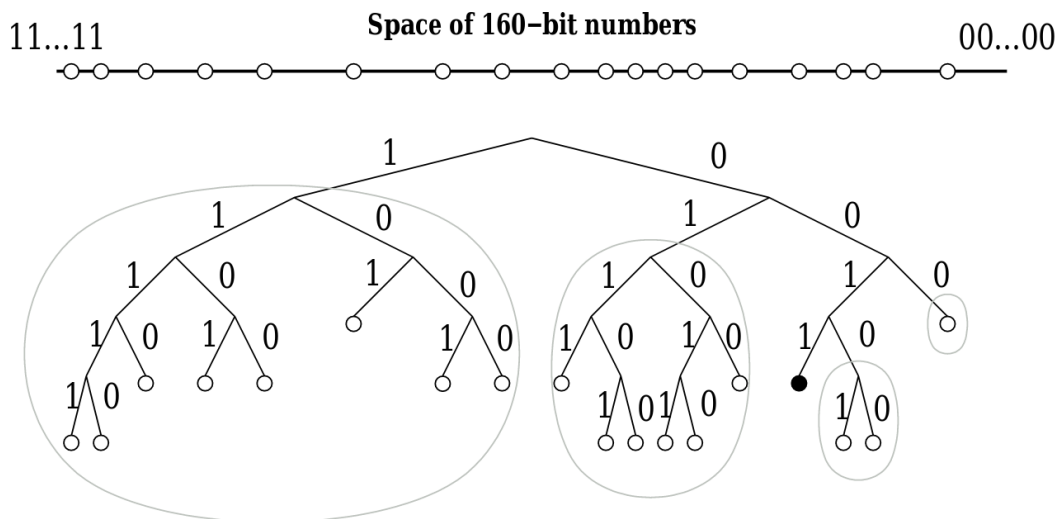


Fig. 2.1 The node with identifier 0011 stores information about nodes in four subtrees, which are encircled. [14]

2.2 Learned indexes

Kraska et al. [10] proposed that machine learning models may be used along with traditional data structures to improve performance. For example, in a sorted array, instead of using binary search or an index, a machine learning model may be used to predict the appropriate index of the element. This can be applied to many structures like B-trees, bloom filters, etc.

Several learned indexes have been proposed – Learned Spatial Index [13] for answering spatial range queries, Learned Secondary Index [9] for unsorted data, RadixSpline [8], an index that can be built in a single pass and Practical Learned Index (PLEX) [15] which was introduced to make the learned index more practical by reducing the number of hyperparameters to just one.

Learned indices have recently been gaining popularity in the database community recently. The idea of using machine learning models has been so influential that researchers are applying the same approach to even the sorting algorithms. [11]

Kraska et al. [10] also proposed a Recursive Model Index which is hierarchical in nature. They are similar to B-trees except, a machine learning model is used in each stage. In each stage, the model predicts another model until the leaf node is reached where the data is stored. It can also be viewed as each node picking a better expert to answer the query.

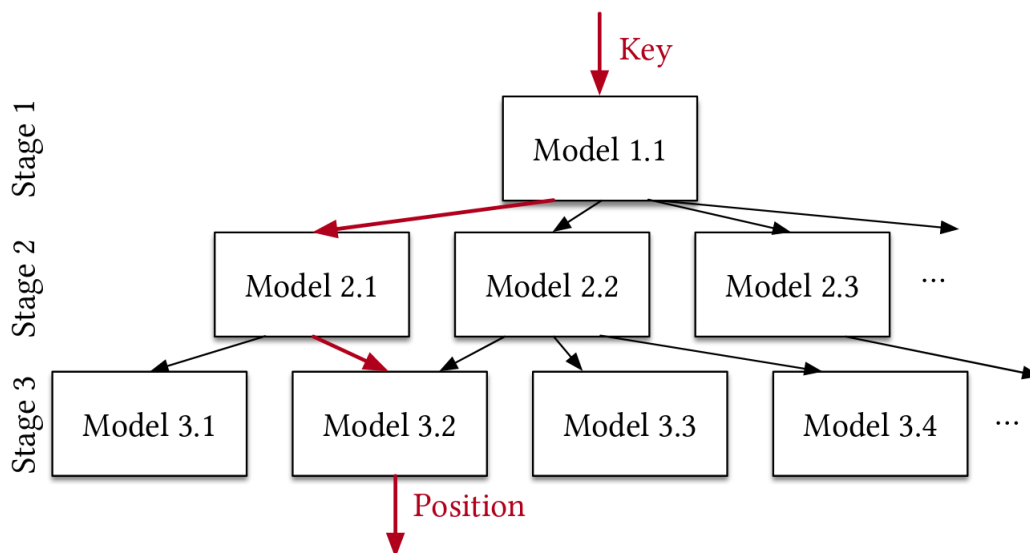


Fig. 2.2 The Recursive Model Index. [10]

Chapter 3

Proposed Methodology

3.1 Using Distributed Hash Tables for Edge Computing

We tried to extend DHTs and learned indices to edge networks. We intend to use Kademlia not for content discovery, but only the XOR based routing algorithm proposed by Petar Maymounkov and David Mazières [14] for routing, once the target node is determined using the learned index.

Junjie Xie et al. [7] observed that DHTs may use a longer path than the shortest possible path between two nodes. This happens because the nodes which may be logically close to each as determined by XOR metric on node identifier may, in fact, be far apart physically. One possible way of solving this problem is to modify the node identifier used by Kademlia so that it is based on the location of the node. Meng Wei and Guoqing Dong [12] did this by using the IP address of the nodes as the prefix in the node identifier. While this should make the routing in Kademlia more efficient, this method suffers from the drawback that the IP address of the node may change.

The another way in which node identifiers can be modified is to encode the services offered by a node into its identifier. For example, if a node is offering weather sensing service, its identifier may be prefixed by “10” or if it is offering road-traffic information, its identifier may be prefixed by “01”. Similarly, if both the services are offered by the node, its identifier may start with “11”. Because of how routing works in Kademlia, it would be trivial to search for nodes offering a particular service if this scheme is used. While we could not find any work that has done this, we believe that this too will have limitations.

We used an open-source Python library to test if Kademlia would be appropriate for this use case. [2]

3.2 Using Learned Indices for Edge Computing

While Kraska et al. [10] proposed that the model should be trained only once, we believe it would be better if the model is updated periodically so that it can serve the queries better.

To test the performance of learned indexes, we simulated an edge network using Python. Inside each edge node, we put five key value pairs. The keys were randomly generated strings and the values were randomly generated integers. First, keys were encoded using ordinal encoding. Then a one-vs-rest multi-class classifier was used to predict the location in which the key-value pairs might be stored. Though the results were not conclusive, using a real world dataset by Taxi & Limousine Commission [16] should give conclusive results. This is because random data carries the most amount of entropy, it is the most difficult to compress, and real world data usually has a lot of redundancy and patterns which the model can learn.

3.3 Conclusion

We believe learned indexes can be a good fit for this problem. Along with learned indexes, the routing algorithm proposed in Kademia can be used to routing after modifying it so that it favors the shortest path between the nodes based on the scheme proposed by Meng Wei and Guoqing Dong [12].

References

- [1] Chord: A scalable peer-to-peer lookup service for internet applications: ACM SIGCOMM Computer Communication Review: Vol 31, No 4. URL <https://dl.acm.org/doi/abs/10.1145/964723.383071>.
- [2] Kademia/index.rst at master · bmuller/kademia. URL <https://github.com/bmuller/kademia>.
- [3] Metric space. *Wikipedia*, August 2022. URL https://en.wikipedia.org/w/index.php?title=Metric_space&oldid=1107612588.
- [4] Andrew Loewenstern and Arvid Norberg. DHT Protocol, January 2008. URL https://www.bittorrent.org/beps/bep_0005.html.
- [5] Federica Paganelli and David Parlanti. A DHT-Based Discovery Service for the Internet of Things. *Journal of Computer Networks and Communications*, 2012: e107041, October 2012. ISSN 2090-7141. doi: 10.1155/2012/107041. URL <https://www.hindawi.com/journals/jcnc/2012/107041/>.
- [6] Josep Torrellas. Cache-Only Memory Architecture. page 4.
- [7] Junjie Xie, Chen Qian, Deke Guo, Minmei Wang, Shouqian Shi, and Honghui Chen. Efficient Indexing Mechanism for Unstructured Data Sharing Systems in Edge Computing. In *IEEE INFOCOM 2019 - IEEE Conference on Computer Communications*, pages 820–828, April 2019. doi: 10.1109/INFOCOM.2019.8737617.
- [8] Andreas Kipf, Ryan Marcus, Alexander van Renen, Mihail Stoian, Alfons Kemper, Tim Kraska, and Thomas Neumann. RadixSpline: A Single-Pass Learned Index, May 2020. URL <http://arxiv.org/abs/2004.14541>.
- [9] Andreas Kipf, Dominik Horn, Pascal Pfeil, Ryan Marcus, and Tim Kraska. LSI: A Learned Secondary Index Structure, May 2022. URL <http://arxiv.org/abs/2205.05769>.
- [10] Tim Kraska, Alex Beutel, Ed H. Chi, Jeffrey Dean, and Neoklis Polyzotis. The Case for Learned Index Structures. In *Proceedings of the 2018 International Conference on Management of Data*, pages 489–504, Houston TX USA, May 2018. ACM. ISBN 978-1-4503-4703-7. doi: 10.1145/3183713.3196909. URL <https://dl.acm.org/doi/10.1145/3183713.3196909>.
- [11] Ani Kristo, Kapil Vaidya, Ugur Çetintemel, Sanchit Misra, and Tim Kraska. The Case for a Learned Sorting Algorithm. In *Proceedings of the 2020 ACM SIGMOD International Conference on Management of Data*, SIGMOD '20, pages 1001–1016,

- New York, NY, USA, May 2020. Association for Computing Machinery. ISBN 978-1-4503-6735-6. doi: 10.1145/3318464.3389752. URL <https://doi.org/10.1145/3318464.3389752>.
- [12] Meng Wei and Guoqing Dong. Improvement of Kademia Based on Physical Location. In *2013 10th Web Information System and Application Conference*, pages 119–122, November 2013. doi: 10.1109/WISA.2013.31.
- [13] Varun Pandey, Alexander van Renen, Andreas Kipf, Ibrahim Sabek, Jialin Ding, and Alfons Kemper. The Case for Learned Spatial Indexes, August 2020. URL <http://arxiv.org/abs/2008.10349>.
- [14] Petar Maymounkov and David Mazières. Kademia: A Peer-to-Peer Information System Based on the XOR Metric. In Gerhard Goos, Juris Hartmanis, Jan van Leeuwen, Peter Druschel, Frans Kaashoek, and Antony Rowstron, editors, *Peer-to-Peer Systems*, volume 2429, pages 53–65. Springer Berlin Heidelberg, Berlin, Heidelberg, 2002. ISBN 978-3-540-44179-3 978-3-540-45748-0. doi: 10.1007/3-540-45748-8_5. URL http://link.springer.com/10.1007/3-540-45748-8_5.
- [15] Mihail Stoian, Andreas Kipf, Ryan Marcus, and Tim Kraska. Towards Practical Learned Indexing, November 2021. URL <http://arxiv.org/abs/2108.05117>.
- [16] Taxi & Limousine Commission. New York Taxi Rides open dataset. URL <https://www.nyc.gov/site/tlc/about/tlc-trip-record-data.page>.
- [17] Weisong Shi, Jie Cao, Quan Zhang, Youhuizi Li, and Lanyu Xu. Edge Computing: Vision and Challenges. *IEEE Internet of Things Journal*, 3(5):637–646, October 2016. ISSN 2327-4662. doi: 10.1109/JIOT.2016.2579198.
- [18] Yasmin Fathy, Payam Barnaghi, and Rahim Tafazolli. Large-Scale Indexing, Discovery, and Ranking for the Internet of Things (IoT). *ACM Computing Surveys*, 51(2):29:1–29:53, March 2018. ISSN 0360-0300. doi: 10.1145/3154525. URL <https://doi.org/10.1145/3154525>.
- [19] B.Y. Zhao, Ling Huang, J. Stribling, S.C. Rhea, A.D. Joseph, and J.D. Kubiatowicz. Tapestry: A resilient global-scale overlay for service deployment. *IEEE Journal on Selected Areas in Communications*, 22(1):41–53, January 2004. ISSN 1558-0008. doi: 10.1109/JSAC.2003.818784.
- [20] Zineddine Kouahla, Ala-Eddine Benrazek, Mohamed Amine Ferrag, Brahim Farou, Hamid Seridi, Muhammet Kurulay, Adeel Anjum, and Alia Asheralieva. A Survey on Big IoT Data Indexing: Potential Solutions, Recent Advancements, and Open Issues. *Future Internet*, 14(1):19, January 2022. ISSN 1999-5903. doi: 10.3390/fi14010019. URL <https://www.mdpi.com/1999-5903/14/1/19>.