

Non-Relational Databases on FPGAs: Survey, Design Decisions, Challenges

Jonas Dann · Daniel Ritter · Holger Frning

Abstract Non-relational database systems (NRDS), such as graph, document, key-value, and wide-column, have gained much attention in various trending (business) application domains like smart logistics, social network analysis, and medical applications, due to their data model variety and scalability. The broad data variety and sheer size of datasets pose unique challenges for the system design and runtime (incl. power consumption). While CPU performance scaling becomes increasingly more difficult, we argue that NRDS can benefit from adding field programmable gate arrays (FPGAs) as accelerators. However, FPGA-accelerated NRDS have not been systematically studied, yet.

To facilitate understanding of this emerging domain, we explore the fit of FPGA acceleration for NRDS with a focus on data model *variety*. We define the term NRDS class as a group of non-relational database systems supporting the same data model. This survey describes and categorizes the inherent differences and non-trivial trade-offs of relevant NRDS classes as well as their commonalities in the context of common design decisions when building such a system with FPGAs. For example, we found in the literature that for key-value stores the FPGA should be placed into the system as a smart network interface card (SmartNIC) to benefit from direct access of the FPGA to the network. However, more complex data models and processing of other classes (e. g., graph and document) commonly require more elaborate near-data or socket accelerator placements where the FPGA respectively has the only

or shared access to main memory. Across the different classes, FPGAs can be used as communication layer or for acceleration of operators and data access. We close with open research and engineering challenges to outline the future of FPGA-accelerated NRDS.

Keywords FPGA · hardware acceleration · non-relational databases · graph · document · key-value

1 Introduction

Recent business and socio-technical trends like smart applications (i. e., leveraging advanced analysis techniques), the internet of things, as well as business and social networks require applications to more efficiently deal with increasingly larger amounts of data in various, non-relational data models close to the underlying domain (cf. [13, 37]). For example, the predominant data exchange format in distributed business applications is JSON [72], requiring the processing and storage of nested object (i. e., key-value) and array data. Furthermore, working with JSON documents gains popularity in schema-less contexts that require flexible data models [37]. Similarly, applications like social network analysis require processing and storage capabilities of graph data that has a focus on entities and their relationships. At the same time the amount of data starkly increases into gigabytes of JSON documents (e. g., [72]) and big graph datasets (e. g., up to one trillion edges [29]). However, traditional relational database systems fall short on requirements of model *variety* (i. e., expressiveness, flexibility) and *efficiency* (e. g., data volume, scalability) of these applications [4, 25, 26, 111].

To address the requirements of variety and efficiency, new classes of systems emerged, namely non-relational database systems (NRDS), that provide a wide variety

J. Dann · D. Ritter
SAP SE, Germany
E-mail: {firstname.lastname}@sap.com

H. Frning
Heidelberg University, Germany
E-mail: holger.froening@ziti.uni-heidelberg.de

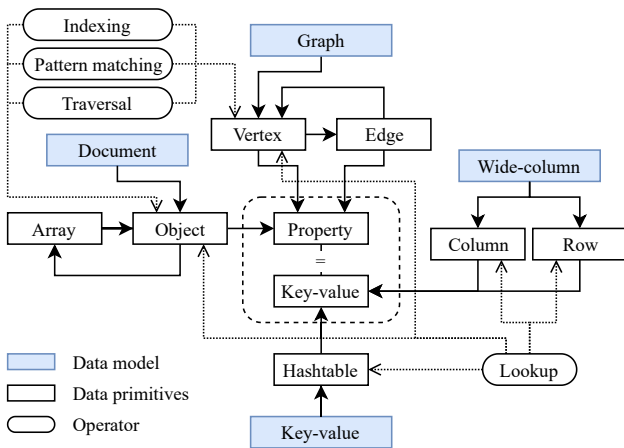


Fig. 1: NRDS data models and their relation

of database models (e. g., graphs, documents) for flexible on-the-fly and application-specific data modeling and efficient processing (e. g., by relaxing traditional relational database system constraints). Similar to recent non-relational data processing surveys [13, 37], we consider the following data models in this work: graph, document, key-value, and wide-column stores.

Variety On first sight, NRDS do not have many commonalities besides coarse-grained system design principles (e. g., scalability) and their application specificity. For example, complex graph query languages are conceptually different from the simple APIs of key-value stores (e. g., put, get [38]). Although acknowledging their differences, we argue that distinct non-relational data models share an underlying primitive in key-value pairs and types of operators, like lookup, traversal (e. g., BFS on graphs, path traversal on document hierarchies), pattern matching (e. g., Cypher [48] on graphs, XPath on documents), and indexing (e. g., PageRank on graphs, full text document indices). Figure 1 depicts those differences and commonalities regarding the data models (blue boxes). Each model has a differentiating structure layer of specific primitives (e. g., vertices and edges of a graph) connecting properties or key-value pairs. Key-value directly stores key-value pairs in a hash table (e. g., [64]), while wide-column uses a column name and a row name as key to the key-value pair (e. g., [71]). Graphs store vertices connected by edges and associate properties with both of those primitives [37], while documents are built from objects and arrays that again refer to properties [37] which are equivalent to key-values. Furthermore, all data models support simple lookup operations and the more expressive data models (document and graph) support structural traversals, pattern matching, and indexing.

Efficiency The NRDS-specific data models and operations are also relevant when considering efficient processing. For example, graph processing has highly irregular memory-access patterns and little or no temporal and spatial locality. This is challenging for current NRDS, mostly built on general-purpose CPUs (cf. Sect. 2), due to the CPU’s fixed deep cache hierarchy and coarse memory access granularity based on cache line sizes. Similar to the application specificity in database systems, there has been a fundamental change in computer architecture and technology trends based on ever rising performance requirements of applications (domain-specific architectures [58]). Since the end of Dennard scaling [40], frequency and therewith performance scaling of CPUs is reaching limits. Power consumption of computing systems is now a hard constraint leading to energy efficiency of operations such as computations and memory accesses dominating overall performance. While GPUs as one option offer massive parallelism, they exhibit significantly reduced performance when the internal cores do not execute the same instruction (i. e., warp divergence), e. g., common for graphs with varying degrees [107].

Thus, new computer architectures such as field programmable gate arrays (FPGAs) are explored for future performance scaling in relational database systems and NRDS (e. g., [64, 95, 109]). FPGAs are reconfigurable computing platforms that can implement custom massively parallel processor architectures. There are no structural restrictions on parallelism, all data on an FPGA is bit-addressable, and data types are not restricted to multiples of bytes. Additionally, FPGAs can be placed anywhere in a system which in particular is important for concepts like near-data processing or processing on the wire which allow for substantial reductions in the amount of data being moved. As a result, FPGAs are being widely used in various applications, including data centers [98], machine learning [122], and also data management which is the focus of this survey. For database systems specifically they might even outperform the currently more widely available GPUs [105]. The focus of this survey will thus be FPGA acceleration of NRDS which is not well studied yet.

1.1 Research gap and related surveys

While FPGA-accelerated NRDS may meet the requirements of emerging applications, there is only limited academic work in terms of survey research. Figure 2 depicts recent surveys related to FPGA-accelerated NRDS. We consider surveys on *relational* and *non-relational* data in the context of three kinds of acceleration: *no*

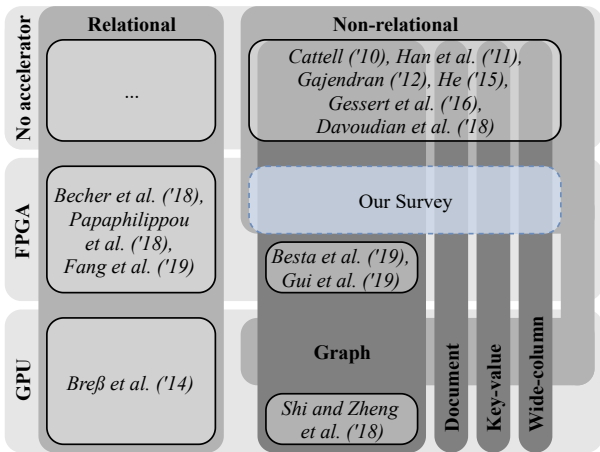


Fig. 2: Related surveys and our contribution

accelerator, FPGA-, and GPU-accelerated. The contribution to non-relational data processing is further specified by the NRDS classes (i. e., graph, document, key-value, wide-column).

Leaving the vast amount of work on non-accelerated relational database systems out of scope, there are several related surveys on acceleration of such systems. In [10], Becker et al. pose the challenge of using the on-the-fly reconfigurability of FPGAs in relational databases. This results in open questions on exploitation of heterogeneous hardware, query partitioning, and dynamic hardware reconfiguration. Papaphilippou et al. categorize the literature into frameworks (e. g., Centaur [95] or DoppioDB [109]) and specialized accelerators for common operators in a relational database, and also highlight upcoming cache coherent connectors for FPGAs (e. g., OpenCAPI). Most recently, Fang et al. [46] state main memory access, programmability, and GPUs as the three biggest factors holding back FPGAs in in-memory relational database systems in the past. Regarding GPU-accelerated relational database systems, [24] raises challenges like the I/O bottleneck and query planning but does not offer convincing solutions.

Non-accelerated NRDS are well-covered by surveys [26, 37, 50, 51, 57, 67] which are discussed in Sect. 2. Additionally, there are two surveys covering graph processing as a high-performance computing (HPC) workload [14, 55]. While this shows the feasibility of graph processing on FPGAs, these surveys only focus on the functional aspects and HPC applications. For our survey this literature has to be reinterpreted for the database perspective. GPU acceleration (e. g., [107] for graph) is out of scope of this work.

In summary, the related surveys support the relevance and timeliness of the topic but there is no survey on FPGA-accelerated NRDS. The objective of this

survey is to fill this gap in the intersection of FPGA acceleration and NRDS. Their inherent differences and commonalities regarding FPGA-accelerated NRDS result in overarching design decisions and patterns that will be instrumental to current and new systems.

1.2 Research method and contributions

To fill the FPGA-accelerated NRDS gap in the survey literature we rely on the design science methodology [132] as a method to collect, summarize, and evaluate FPGA-accelerated NRDS solutions. Our fundamental theory and motivation is: Non-relational database systems do not leverage modern FPGA hardware, from which we derive hypothesis (H1), i. e., *existing non-relational database systems do not realize the potential of FPGA acceleration*. This hypothesis is tested based on an introductory system review in Sect. 2 which aims at analyzing existing NRDS regarding system aspects (e. g., operators or scalability) We define system aspects as abstract requirements for a system to be regarded as an NRDS. Further, we believe that (H1) is grounded on missing fundamental research on FPGA-accelerated NRDS (cf. Tab. 2), which leads us to hypothesis (H2), i. e., *there are significant gaps in current research on non-relational FPGA acceleration*. We test this hypothesis on another observation artifact, i. e., a comprehensive literature analysis in Sect. 3. The literature analysis describes which system aspects have been previously studied and which solutions are provided along the system aspects. In hypothesis (H3), we argue that *there are patterns guiding the design of an FPGA-accelerated non-relational database system*. To address the detected gaps and missing FPGA-accelerated NRDS we propose new practical pattern categories on how to build FPGA accelerated NRDS or accelerate existing ones. To summarize, this survey makes the following contributions:

- We identify *system aspects* resulting from challenges in accelerator design for FPGAs and a review of existing NRDS (Sect. 2).
- We provide short solution summaries (*implementations*), a table classifying references by system aspect, and a list of *gaps* resulting from a comprehensive literature search (Sect. 3).
- We extract patterns from the literature with regard to *tasks, placements, non-trivial accelerator design decisions*, and accelerated architecture *justification* as guidelines to system architects (Sect. 4).
- From the gaps in the literature we derive *open research challenges* in the field of FPGA-accelerated NRDS (Sect. 5).

Section 6 concludes the survey.

1.3 Key insights – what will the reader learn?

In the course of this work we gain the following five key insights (highlighted as framed theorems in Sect. 4):

1. There are three accelerator task categories (operator, data access, and communication layer acceleration) that FPGAs are currently well suited for.
2. There are four fundamental patterns of FPGA placement (offload, SmartNIC, near-data, and socket).
3. The accelerator task in combination with the characteristics of operators of the NRDS class are sufficient to define the FPGA placement.
4. There are three operator switching strategies and six memory access optimization patterns guiding the development of accelerators for NRDS.
5. A portable, relevant benchmark suite that covers all necessary artifacts is missing for robust justification of accelerator usage decisions.

Additionally, the reader will learn about the considered NRDS classes and how they make use of accelerators (e.g., scalable key-value solutions were found, but no complete database solution for others).

2 Background

This section introduces fundamental concepts of FPGA hardware (Sect. 2.1) and NRDS classes (Sect. 2.2), required to understand the remainder of this work. The NRDS classes are discussed in the context of well-known, commercial NRDS which we review to give an answer to hypothesis (H1) i.e., *existing non-relational database systems do not realize the potential of FPGA acceleration*. Additionally, we collect important FPGA and NRDS system aspects of all NRDS classes in a taxonomy (Fig. 4) to guide the subsequent literature analysis.

2.1 Field Programmable Gate Arrays

FPGAs are a processor architecture platform that map to custom architecture designs, meaning a set of logic gates and their connection (circuit design). This is similar to application specific integrated circuits (ASICs) that are custom-made to represent a circuit design, but FPGAs are reprogrammable such that the design can be changed by a programmer at configuration time. This reconfigurability comes at a price: from an efficiency point of view, it is always preferable to implement a design as a custom ASIC. However, economic reasons as well as the need to adapt to changes in application behavior often prevent using ASICs. In recent

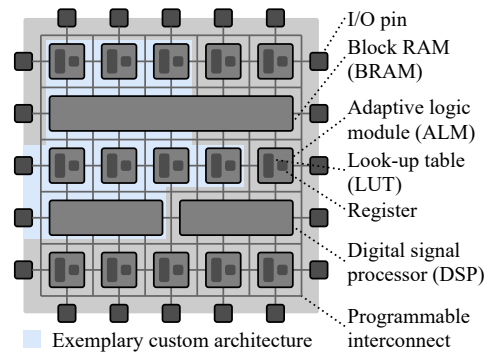


Fig. 3: FPGA on-chip resources

years, FPGAs emerged as accelerators for data processing (e.g., [118]). They provide unique opportunities to implement functionality, like custom single instruction multiple data (SIMD) units or processing and data structure hybrids, like systolic arrays.

Figure 3 shows an abstraction of an FPGA chip hiding vendor specifics. On an FPGA chip, a custom architecture is mapped to a grid of resources (e.g., look-up tables (LUTs) or registers) connected with a programmable interconnection network. Each custom architecture uses a certain amount of resources upon configuration (cf. Fig. 3), such that even multiple custom architectures might be deployed on the same chip. Additionally, FPGAs support partial dynamic reconfiguration where only parts of the FPGA can be reprogrammed. Data enters and leaves the FPGA chip through I/O pins. The logic itself is implemented as adaptive logic modules (ALMs) containing a LUT with multiple inbound and one outbound bit port and a register for optional storage of the outbound bit. Each individual LUT is programmed upon configuration. The FPGA further contains on-chip block RAM (BRAM) in the form of SRAM memory components for fast storage of data structures. BRAM in total is about as large as the caches on a CPU combined but finely configurable to the mapped architecture. Lastly, the FPGA contains hardened digital signal processors (DSPs) that allow fast arithmetics on floating-point numbers.

Similar to the CPU being placed on a mainboard, FPGA chips are soldered onto a board connecting components like network ports, DRAM chips, and PCIe connectors to the FPGA. Concerning DRAM, FPGAs support DDR3 and DDR4 as well as stacked memory enabling high bandwidth on-board data processing, like Hybrid Memory Cube (HMC). As an accelerator, boards containing an FPGA are most often placed in existing hardware systems with a CPU – the host system. Traditionally, accelerators are connected with the CPU over PCIe but recently there are also cache co-

herent interfaces (e.g., UPI, CXL, OpenCAPI). This allows deeper integration of the FPGA into the CPU’s data management and sets FPGAs apart from other accelerators like GPUs. As an example, the Contutto system explores deeper integration of FPGAs as programmable memory buffers into the datapath between CPU and memory [113].

Subsequently, the most relevant system aspects of FPGA accelerator design are discussed which are based on differences to designing CPU applications and also found to be challenging design decisions in related surveys [10, 14, 46], namely (i) software vs. hardware design as *design paradigm* (e.g., [46]), (ii) tight *CPU-FPGA collaboration* between existing CPU-based systems and FPGA accelerators (e.g., [46]), (iii) custom *memory access* controllers (e.g., [14]), and (iv) comparability through a *performance model* (e.g., [10]).

Design paradigm In principle, all LUTs ($> 1,000,000$ in modern FPGAs) can operate in parallel. This opens up a vast design space (when using hardware description languages like VHDL) compared to the well-formed design space of instruction-based processors. While there have been efforts to simplify development with higher-level languages (e.g., OpenCL), they do not seem to satisfy performance requirements for complex applications yet [139]. One key question of NRDS acceleration is thus, how to use constraints inherent to the NRDS class to reduce the accelerator design space without performance degradation.

CPU-FPGA collaboration Adding an FPGA to a data processing system is justified with sufficient improvement in overall performance or cost saving. However, most systems will still require a CPU, introducing data movement overhead between the two processors. Thus, effective accelerator integration requires not only high utilization of the FPGA but also little data movement. CPU-FPGA collaboration has to solve problems of task orchestration and data management.

Memory access As one unique feature when compared to CPUs, FPGAs do not access their connected on-board memory through a deep cache hierarchy that assumes temporal and spatial locality in memory accesses. Thus, FPGAs can implement unique caching strategies and placement of critical data on-chip.

Performance model The circuit-based programs of FPGAs have very different performance implications than instruction-based programs of CPUs. For big circuit designs it is not easy to comprehend the amount of parallelism and make performance predictions. Thus, performance models are about using the properties of the

Class	System	Refs	FPGA
Graph	Neo4j	[26, 37, 51, 57]	(♣)
	OrientDB	[37] (expert)	♣
Document	CouchDB	[26, 37, 50, 51, 57, 67]	♣
	MongoDB	[26, 37, 50, 51, 57, 67]	♣
Key-value	Redis	[26, 37, 51, 57, 67]	♣
	DynamoDB	[26, 37, 50, 51, 57]	♣
Wide-column	Cassandra	[26, 37, 50, 51, 57, 67]	(♣)
	HBase	[26, 37, 50, 51, 57, 67]	♣

♣: productive FPGA usage, (♣): explored FPGAs but no productive usage, ♣: no FPGA acceleration mentioned

Table 1: Commercial NRDS found in surveys

NRDS class to aid understanding of design decisions and comparing different approaches.

2.2 Non-relational database systems (NRDS)

NRDS are defined according to their data models (as established in Sect. 1). We found five surveys on non-accelerated NRDS (cf. Fig. 2) [26, 37, 50, 51, 57, 67]. According to these surveys, the predominant NRDS classes are graph, document, key-value, and wide-column. Other NRDS classes that we do not further consider in this work are less popular ones like spatial, object-oriented or timeseries database systems.

To review well-known commercial NRDS, we selected the systems of each class that were at least mentioned three times in the NRDS surveys and combined similar systems (e.g., Riak KV is similar to Amazon DynamoDB and Google Bigtable to Apache HBase). Additionally, we added OrientDB (marked as “expert” in Tab. 1) as the only NRDS that successfully combines graph, document, and object-oriented database concepts despite it only being named in [37]. Table 1 shows the selection of NRDS by class. The review will not only help to briefly recall the systems’ design choices and challenges but also discuss similarities. We considered public documentation and publications if available. Subsequently, the concepts of NRDS classes are introduced for the systems in Tab. 1.

2.2.1 Graph

Graph databases store entities (vertices V) and their relations (edges E). Naively, a graph is stored as adjacency matrix M of size $|V| \times |V|$ where a 1 at position $M_{v,w}$ represents an edge between vertex v and w . However, currently one of the most popular representations in the literature is compressed sparse row (CSR), as a compression technique for such sparse matrices. Moreover, graph partitioning is often applied in the literature, e.g., by simple horizontal (intervals of

source vertices) or vertical (intervals of destination vertices) partitioning or interval-shard [152, 142] where the graph is partitioned into p equally sized intervals I_0, I_1, \dots, I_{p-1} of vertices and shard $S_{i,j}$ is the set of all edges between vertices in I_i and I_j . We look at the graph database systems Neo4j and OrientDB (cf. Tab. 1). Neither of those NRDS deploys FPGA acceleration in production. However, Neo4j experimented with the Flash-accelerator CAPI SNAP [1], where an FPGA is inserted into the datapath between CPU and Flash storage to accelerate data accesses.

Neo4j¹ stores data as a property graph, where vertices and edges have properties (cf. Fig. 1). It supports Cypher [48] as a graph query language with a cost-based query optimizer, traversal patterns like breadth-first search (BFS), and algorithms to solve common graph problems, like shortest paths and centrality (e. g., PageRank). More details on graph algorithms can be found in [45]. For solving graph problems, Neo4j transforms the property graph into an in-memory projection that is optimized for traversal. Additionally, Neo4j supports different indexes. To provide horizontal scalability, Neo4j can run as a cluster of core nodes that replicate all changes between themselves. For additional read scaling, read replicas can be added to the cluster by registering at a core node and Neo4j drivers provide routing and load balancing. Changes to graphs follow causal consistency, where replication to a majority of core nodes has to finish to confirm a transaction. Optionally, ACID transaction guarantees can be enforced. For multi-tenant usage (multiple users working on same system), Neo4j provides role-based access control (RBAC) and intra-cluster encryption, and encrypted backups provide further security.

Another well-known graph database is OrientDB² that allows for queries and traversals on a property graph with native support for documents. OrientDB provides SQL-like language support with a graph extension and operators like BFS. For scalability, data can be replicated over a cluster and availability is guaranteed by multi-master replication (similar to Neo4j) and auto discovery of nodes. When multiple users are working on the database (OrientDB also supports RBAC), an optimistic multi-version concurrency control (MVCC) is used and ACID transaction guarantees can be applied. OrientDB supports encryption on disk.

2.2.2 Document

Document systems store formatted text documents in a document hierarchy. The two most popular document formats are Extensible Markup Language (XML) [22] and JavaScript Object Notation (JSON) [23]. All documents in a document store adhere to such a fixed formatting and are parsed upon ingestion from a string representation for quick processing afterwards. For XML stores, there is the XML Path Language (XPath) [31] as a query language. Apache CouchDB and MongoDB are the most-referenced document NRDS (cf. Tab. 1). We did not find FPGA acceleration options for any of these two document stores.

Apache CouchDB³ is a JSON document store operating as a cluster of masters with bidirectional, asynchronous replication. The API allows create, read, update, and delete (CRUD) operators on documents and more advanced view models for filtering and aggregation. Queries can be accelerated with B-tree indices. Writes to the database are isolated with MVCC and ACID transaction guarantees can be enforced. Like the graph databases, Apache CouchDB supports RBAC and provides security features, like update validation.

MongoDB⁴ is another JSON document store with similar to Apache CouchDB query possibilities but also adds capabilities for fulltext search and spatial queries. Availability is provided by a master-slave cluster setup, where nodes are placed into a replication set of a master that handles all writes. Writes are atomic, and recently multi-document transactions were added to the system. If the master fails, a new one is elected among the replicas. Data can also be sharded over multiple replication sets for scalability. Similar to the other systems, MongoDB provides RBAC for authentication of users. For security, communication between users and the database can be SSL-encrypted.

2.2.3 Key-value

Key-value stores operate on pairs of a key used for quick lookup and a value of arbitrary data. They may be used as an underlying persistence layer of another database or standalone depending on the use case. We subsequently discuss Redis and Amazon DynamoDB (cf. Tab. 1). While Amazon DynamoDB has no accelerator options, there is a recent extension for Redis by Algo-Logic [2] with the FPGAs directly attached to the

¹Neo4j documentation, visited 7/2020: <https://neo4j.com/docs/>

²OrientDB manual, visited 7/2020: <https://orientdb.com/docs/last/index.html>

³Apache CouchDB documentation, visited 7/2020: <https://docs.couchdb.org/en/stable>

⁴The MongoDB manual, visited 7/2020: <https://docs.mongodb.com/manual>

network for increased throughput, lower latency, and less energy consumption.

Redis⁵ may be used as an in-memory or a persistent key-value store, with a typical operator set with `set` and `get` operators. Key-value pairs are inserted into a hash table, where the key hash is used as the index in the table for fast lookup of values, but complex queries are not supported. Similar to the graph and document stores, Redis uses a master-slave setup with sharding for scalability. Additionally, it supports request routing on a proxy node and optional waiting for replication for consistency. Redis has security features like SSL encryption and action auditing.

Amazon DynamoDB⁶ [38] features a distribution scheme without a master. The key hashes denote a circular space and each node is randomly placed in this space upon entrance to the cluster. Each node is responsible for the keys preceding it in counter-clockwise order in this space and answers all requests to that partition. Data is replicated in clockwise order to a fixed number of nodes and upon failure, the node after a failed one is responsible for the failed nodes partition of the circular space. For load balancing, Amazon DynamoDB places many more virtual nodes on the circular space than there are nodes in the cluster. Multiple virtual nodes are then handled by each physical node. Amazon DynamoDB supports eventual consistency, MVCC, and is integrated with the RBAC of AWS accounts. For security, it features encryption and user authentication.

2.2.4 Wide-column

Wide-column databases also store data as pairs of key and value. However, different to key-value databases, the key has two predefined parts: a row name and a column name. Wide-column databases present their data as tables to the user but unlike relational databases the tables are not materialized, unstructured, and every row may have arbitrary columns. The most-referenced wide-column systems are Apache Cassandra and Apache HBase (cf. Tab. 1). For Apache Cassandra, FPGAs may be used to accelerate the data accesses where the FPGA denotes a data proxy [3].

Apache Cassandra⁷ [71] uses a multi-dimensional map, indexed with a key (everything with the same key constitutes a row) and columns that are grouped into column families. It supports `insert`, `get`, and `delete`

operations and additionally has a Cassandra query language (CQL) comparable to SQL. Similar to Amazon DynamoDB, the key hashes are used as a circular space to partition the data to the nodes in the cluster but load balancing is done by moving nodes in the circular space when imbalance is detected. Consistency is guaranteed with a replication quorum and RBAC is supported as well as SSL-encryption.

Apache HBase⁸ is a wide-column store based on Apache HDFS. The supported operator set as well as the binary representation is similar to Apache Cassandra, but Apache HBase distributes updates in a master-slave fashion. Similar to all other systems we looked at, Apache HBase supports RBAC.

2.3 Summary – NRDS on reprogrammable hardware

Although the architectures of the systems reviewed exhibit different feature sets, they cover a shared set of system aspects. For example, Neo4j provides an execution engine with Cypher query capabilities while Apache Cassandra provides them with CQL. Scalability in Amazon DynamoDB and Apache Cassandra is achieved with partitioning data with a circular hash space while e.g., MongoDB and Redis distribute work in a master-slave fashion. While the query languages and scalability schemes might differ in their concrete implementation, queries and scalability are two integral system aspects of any NRDS.

Figure 4 shows a generalized shared architecture as a component view of the reviewed NRDS capturing the system aspects we found in the system review. The architecture is based on a set of nodes in a networked environment, forming a cluster. Multiple different users or tenants send requests to the system, which distributes work with the request router and load balancer in the cluster manager. Changes to the underlying data are kept consistent by the consistency manager, which performs concurrency control and atomic writing. Optionally the transaction manager provides transactions (e.g., ACID) touching multiple data elements. The execution engine finally processes queries – made up of operators – on the data stored in the system. Queries performance may be improved with indexes. Persistent storage is optional, but data is partitioned over the nodes. Lastly, the replication manager aids discovery of new nodes in the system and fault handling. While the components in general might be similar to scale-out relational database systems, NRDS

⁵Redis documentation, visited 7/2020: <https://docs.redislabs.com/latest/index.html>

⁶Amazon DynamoDB developer guide, visited 7/2020: <https://docs.aws.amazon.com/amazondynamodb/latest/developerguide/Introduction.html>

⁷Apache Cassandra Documentation, visited 7/2020: <https://cassandra.apache.org/doc/latest/>

⁸Apache HBase reference guide, 7/2020: <https://hbase.apache.org/book.html>

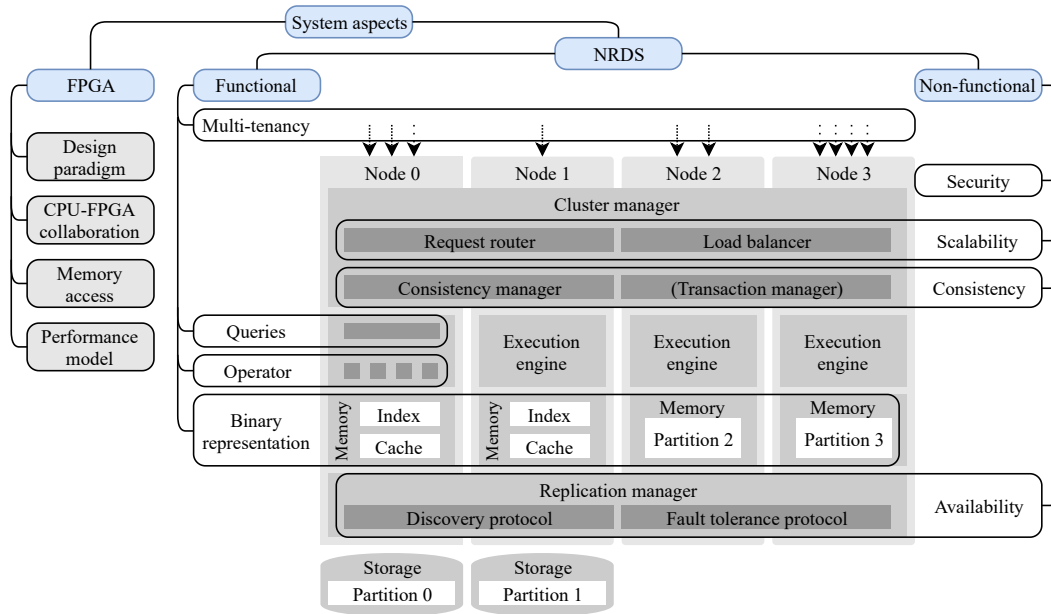


Fig. 4: Taxonomy of system aspects along common NRDS architecture

lay different emphasis on components because of different application requirements (cf. Sect. 1).

As an overlay in Fig. 4, we show the FPGA system aspects from Sect. 2.1 and functional and non-functional system aspects defining NRDS. A complete system would have to satisfy all (or most of the) combined system aspects. FPGAs are not well established enough in NRDS to place them in the shared architecture, such that their system aspects are added separately. The aspect taxonomy will be used in the literature analysis (Sect. 3) to classify literature. In summary, there are some first experiments with FPGAs in NRDS [1, 2, 3] that show the potential and feasibility, but the majority of systems has not considered FPGAs yet. Thus, we confirm hypothesis (H1) i. e., *existing non-relational database systems do not realize the potential of FPGA acceleration*.

3 Literature Review

In this section we conduct a literature analysis in order to answer hypothesis (H2), i. e., *there are significant gaps in current research on non-relational FPGA acceleration*, as set out in Sect. 1.2. The hypothesis raises three questions to be investigated in the literature analysis: (a) What are the most relevant of the identified NRDS classes? (b) Are there any system aspects that are not yet covered by literature? (c) Do existing approaches provide solutions to these topics?

The literature analysis is based on the guidelines described in [70]. The primary selection of references is

conducted in the domain of each individual NRDS class and with a focus on research articles (no patents and citations). This results in 351 hits before the following selection criteria are applied: (i) relation to computer science, FPGA, reconfigurable hardware, and acceleration, (ii) focus on data processing (excluding for example robotics, image processing, or graph-based FPGA design), (iii) availability of the document, (iv) written in English, (v) published (excluding Master and PhD theses). Overall, this results in 89 selected articles relevant to this survey.

Notably, we did not find dedicated literature for FPGA-acceleration of wide-column which, however, can be seen as very similar to key-value (cf. [37]).

3.1 Processing of selected literature – NRDS classes and trends

All selected articles were categorized by NRDS class giving an answer to question (a) What are the most relevant of the identified NRDS classes?

Figure 5 depicts the distribution of NRDS classes over time. It can be seen that hardware-accelerated graph processing already plays an early role from 1999–2002. After three years with only one article for document, graph sparks in the years 2006 and 2007. From 2009–2013 document moved into the focus, with graph joining in again from 2010 onward. Recently (i. e., 2013–2019), the picture seems to change, turning more or less exclusively to graph and key-value.

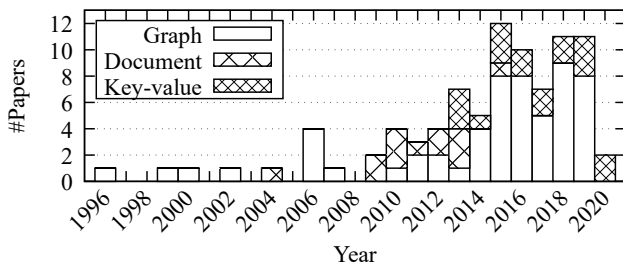


Fig. 5: Topics in literature over time

From Fig. 5 it can be concluded that most of the non-relational hardware work occurred in literature after 2006 with a dominance of graph. Apparently document lost significance, and other classes like key-value gained more attention. From the dominance of graph and the very different key-value type we conclude that a more fine-grained analysis of the mentioned database classes is meaningful.

3.2 Literature summaries

This section summarizes the approaches identified in the literature search, thus addressing questions (b) Are there any system aspects that are not yet covered by literature? and (c) Do existing approaches provide solutions to these topics? We organize the summaries chronologically (cf. Fig. 5) and by NRDS class. Additionally, we structure the literature by their strongest contribution(s) to the system aspects (cf. Fig. 4). Table 2 summarizes the found solutions and subclassification we identified for each NRDS class and system aspect pairing which we briefly introduce subsequently.

3.2.1 Graph

According to the timeline (cf. Fig. 5), the first non-relational accelerator solutions were provided for graph representing mostly HPC solutions not specifically tuned towards NRDS.

Operator In the literature, we identified solutions for five different operators which we discuss in the following order: shortest path, breadth-first search, maximum matchings, page rank, and centrality.

Shortest path Bondhugula et al. present a tiled Floyd-Warshall implementation solving the all-pairs shortest-paths problem using a pipeline of B processing elements (PE) which can each process l elements of the adjacency matrix [21]. They build a performance model with these two parameters where B is constrained by FPGA resources and l is constrained by I/O bandwidth.

Jagadeesh et al. use a parallel, synchronous Bellman-Ford algorithm where each PE on the FPGA represents one node in the graph (severely limiting graph size) [65]. Additionally, they model the internal computation time of their shortest path computation unit in cycles for performance prediction. A parallel implementation of Bellman-Ford that considers conflicting vertex updates is presented by Zhou et al. [146]. The edges are processed in parallel and conflicts are resolved by caching updates on-chip until they are applied in main memory. To optimize for sequential reads, the edges are sorted by destination vertex.

In [114, 115], Takei et al. combine Dijkstra with a SIMD distance comparison unit. A restriction of the approach is that all nodes have to fit into on-chip memory. Lei et al. propose an eager Dijkstra algorithm based on their own memory overflow extension of a priority queue and three memory channels (for overflow queue, graph, and output data) [75]. This resolves the graph size restriction of the earlier approach by Takei et al.

In [87], the authors propose a solution for all-pairs shortest-path by defining a partitioning that allows for processing the graph with a bidirectional systolic array with an optimal number $|V|$ of PEs. Betkaoui et al. solve all-pairs shortest-paths with parallel BFS kernels [16, 17]. Multiple parallel PEs issue many non-blocking memory requests to take advantage of the multi-channel memory system of their particular FPGA setup.

Breadth-first search (BFS) Wang et al. propose a solution for a parallel BFS with message passing on a fully-connected network between interval-partitioned soft cores [126]. The traversal levels are synchronized with a floating barrier. The number of random memory accesses is reduced by keeping the visited status in on-chip memory. Additionally, their approach allows switching traversal patterns (bottom-up and top-down) per level. TorusBFS [76] proposed by Lei et al. implements torus network-based message-passing between PEs in an interval-partitioned graph. The PEs are similar to those in [16], and auxiliary data structures are stored in BRAM.

In [123], BFS iterations are substituted by matrix-vector operations on a Boolean semi-ring (i.e. multiply and add are substituted with logical **and** and **or**). The data is horizontally partitioned but the matrix and vectors are never materialized. Dr. BFS [47] uses vertical partitioning to fit metadata of large graphs into on-chip memory. The tasks of computation and data access are separated into two different modules having a data access module for every memory bank. The computation modules use a pipelined combination network for sequential burst operations.

	Graph	Document	Key-value
NRDS – functional			
Operator	SP [16, 17, 21, 65, 75, 87, 114, 115, 146]; BFS [47, 76, 123, 126]; MM [12]; PR [147]; Cent. [53]	XML Parser [35, 59, 108]; Filter [27]; XPath [42, 88, 89]; Twig [90]	Insert [78]; Hash [80]
Binary representation	HP [148]; VP [28]; IS [33]; CSR [7]; HWM [60, 86]; Other [110, 116, 127, 128, 135]		Hash table [18, 61, 121, 140]; CAM [81]
Queries	Instructable processor [68]	Skel. autom. [119]; Other [83]	n/a
Multi-tenancy			Token bucket [62]
NRDS – non-functional			
Scalability	Part. [6, 9, 34]		Repl. [64]; Part. [100, 136]
Availability			Repl. [64]
Consistency	Locking [84]		MVCC [101]
Security			
FPGA			
Design paradigm	VC [44, 93, 96, 130, 141]; EC [149, 150]; Hybrid [138]; BSP [8, 39]; Other [9, 36, 94]		n/a
CPU-FPGA collaboration	Socket [20, 123, 129, 151]; Near-data [74]	Socket [124]; Near-data [120]; PCIe [125]	Near-data [73, 134]; DMA [77, 99]
Memory access	Req. merging [69]; Caching [106, 144]; Data placem. [145]; Other [15, 92, 137]		Bloom filter [30]; Flash [19, 136]
Performance model	Parallelism [21, 106]; Bottleneck [33, 34, 43, 144]		Bottleneck [100]

SP: Shortest path, BFS: Breadth-first search, MM: Maximum matchings, PR: PageRank, Cent.: Centrality, HP: Horizontal partitioning; VP: Vertical partitioning, IS: Interval-shard, CSR: Compressed sparse row, HWM: Hardware mapping, Part.: Partitioning, VC: Vertex-centric, EC: Edge-centric, BSP: Bulk-synchronous parallel, Req. merging.: Request merging, Data placem.: Data placement; Skel. autom.: Skeleton automaton; CAM: Content addressable memory, Repl.: Replication, MVCC: Multi-version concurrency control, DMA: Direct memory access

Table 2: Contributions by domain and system aspect

Maximum matchings Besta et al. propose a solution for maximum matchings in a graph which describes the maximum size set of edges that do not share a vertex [12]. Their substream-centric approach divides the incoming stream of edges by their weight into substreams, processes them in parallel, and merges the results.

PageRank Zhou et al. propose an implementation of PageRank (a measure for importance of vertices in a graph) split up into a scatter and gather phase using a source vertex interval (i. e., horizontal) graph partitioning with vertex, edge, and update sets for each partition [147]. The edge set of each partition is sorted by destination vertex to reduce the number of random writes.

Centrality In [53], the authors propose a stochastic matrix function estimator written in OpenCL and apply it to the subgraph centrality problem. Subgraph centrality (like PR) is another measure for importance of vertices in a graph.

Binary representation We found approaches that leverage different partitioning schemes (horizontal, vertical, interval-shard), an optimization of CSR for BFS, and other binary representations like hardware mapping.

Horizontal partitioning In [148] a horizontal partitioning scheme is proposed with an improved data layout enabling more sequential write accesses. This is achieved by sorting the edges in each partition by destination vertex which also allows update merging.

Vertical partitioning Chen et al. provide a solution that features a layout improvement of the partitioned data by storing edges inbound to the current vertex set and not storing an update list (i.e. directly streamed) [28]. The data is shuffled to graph PEs.

Interval-shard FPGP [33] is an edge-centric graph-processing framework based on the interval-shard graph partitioning scheme. Shards $S_{i,j}$ and their corresponding inbound I_i and outbound I_j vertex intervals are stored one after another in on-chip memory and processed by multiple PEs.

Compressed sparse row (CSR) The CyGraph architecture by Attia et al. proposes new optimizations for utilizing the memory bandwidth in parallel BFS by a custom CSR representation [7]. The visitation status of vertices is encoded in the row index which is replaced by the level after visitation (in-place BFS result).

Hardware mapping In [60], the complete graph with editable vertices and edges is represented as logic on the FPGA. An extension by Mencer et al. [86] draws parallels to content-addressable memory (CAM) and extends the idea of mapping the adjacency matrix to hardware by extending it to multi-context graph processing.

Other Skylarov et al. represent the graph as a matrix and use matrix operations to calculate graph colorings [110]. Wang et al. propose an edge-centric graph streaming model on an FPGA for partitioned graphs

to deal with load balancing issues of skewed graphs [127, 128]. The kp partitions are assigned to k PEs in a pre-processing step such that every PE processes approximately the same number of edges. The graph is compressed and streamed which results in a high effective bandwidth. When the application permits graph sampling, the data volume and irregular memory access can be tamed with a data structure derived from CSR [116]. The novel data structure allows storing multiple graphs and removal of vertices and edges with a second pointer array. Xu et al. propose a service-oriented accelerator that does asynchronous processing of BFS [135]. This is very different to other approaches in that batches of so called Batch Row Vectors are streamed into the accelerator and worked on. A Batch Row Vector contains for each edge in the batch the start and destination vertex and a property for the source and destination vertex.

Queries Queries in the context of accelerator design are about flexible chaining of operators. We did not find any accelerator designs explicitly handling query workloads. However, there is an instructable processor design that could be extended for query processing.

Instructable processor GraphSoC is a custom graph processor built from a 2D array of softcore processors (with send, receive, accumulate, and update instructions) connected by a packet-switched network [68].

Scalability There are no complete solutions for scalability either. However, there are three articles on graph partitioning to leverage multi-FPGA setups.

Partitioning Babb et al. presented an early work on scalability by compiling graph problems to whole arrays of FPGAs [9]. Virtual wires are used in between the FPGAs resulting in a multi-FPGA computing fabric. Attita et al. propose a vertex-centric graph processing framework based on the gather-apply-scatter (GAS) principle [6]. With partitioning and message passing, the workload can be distributed to multiple FPGAs. ForeGraph [34] is an edge-centric graph processing approach as a multi-FPGA extension to [33]. For p FPGAs, the graph is partitioned into p intervals resulting in p^2 shards. Each FPGA stores one interval and its outgoing p shards, additionally partitions its subgraph into sub-intervals and sub-shards, and subsequently does the same processing as [33] on the sub-shards. Updates to foreign intervals are propagated to the corresponding FPGA over the network.

Consistency The conflict management for highly concurrent systems (transactions on graphs) can be prob-

lematic. We only found one article proposing a locking scheme to achieve isolation.

Locking Ma et al. define a multi-threaded graph processing engine on FPGA using a global, transactional shared memory which allows fine-granular locking with an address signature table [84].

Design paradigm We found multiple graph accelerator design paradigms: vertex-, edge-centric, hybrid, and BSP, among others. Vertex- and edge-centric describe orthogonal approaches of traversing graphs: while the vertex-centric approach works on outgoing edges of active vertices, the edge-centric approach loads all edges from memory and discards those not needed currently.

Vertex-centric GraphGen [93] compiles graph algorithms from a domain specific language to RTL. The algorithms are built from user-defined instructions. RTL implementations of these instructions have to be provided together with an update function and a description on how those instructions play together. Weisz et al. provide programmability by using GraphGen and CoRAM in combination [130]. CoRAM allows to port the accelerator architecture to both Intel and Xilinx FPGAs. GraVF [44] shows how to compile Migen definitions to hardware. Ozda et al. provide a solution based on a configurable architecture template for vertex-centric graph algorithms [96]. In [141], conflicting updates on one vertex are resolved by accumulating updates in one cycle, parallelizing conflicting vertex updates, and removing sequential application of atomic protection.

Edge-centric HitGraph [149, 150] (based on [147, 148]) compiles edge-centric algorithms to RTL. The processing logic is split up into multiple PEs processing the graph by alternating scatter and gather phases. The vertices are partitioned and buffered in BRAM and partitions are skipped if they do not contain active vertices.

Hybrid Chengbo proposes a hybrid approach where there is a vertex-centric and an edge-centric module on the FPGA programmed in OpenCL [138]. A dispatcher switches the modules depending on the workload. The graph is vertically partitioned.

Bulk-synchronous parallel (BSP) DeLorimier et al. propose an accelerator, mapping the graph to sparse matrix operations executed in a BSP fashion [39]. The accelerator is split up into compute leaves with BRAMs attached. Data between the compute leaves is communicated over a network on chip. In [8], a design methodology based on the BSP model is proposed. Common architectural features are represented as templates which are specified with user-defined functions for GAS. All data flow is handled by the template.

Other In [9], problems on directed graphs are reformulated as closed semiring problems and compiled onto multiple FPGAs. For a graph instance the edges are mapped to summations and vertices are mapped to minimum operators. Dandali et al. address long synthesis times by a skeleton compilation with precompiled blocks that are adapted to the problem instance [36]. One PE is used for a vertex and has to fit onto the FPGA. GraphOps [94] is a general graph dataflow library that provides graph-specific building blocks for the generation of FPGA designs. It includes a locality-optimized property array.

CPU-FPGA collaboration We found the following contributions running on CPU-FPGA heterogeneous platforms with different task assignments.

Socket Bondhugula et al. propose a shortest-path hardware kernel communicating with the CPU via a shared memory region [20]. Another part of the solution is a graph data layout for the CPU to reduce cache misses. Umuroglu et al. propose an approach that distributes BFS iterations between CPU and FPGA: the iterations with few active vertices are performed on the CPU while the other iterations are performed on the FPGA [123]. Similarly, Zhou et al. address the respective drawbacks of vertex- and edge-centric graph processing by switching between the paradigms during execution [151]. The graph is horizontally partitioned (cf. [148]) and the paradigm is chosen individually for each partition in each iteration based on its active vertex ratio. Partitions with few active vertices are processed by the CPU in a vertex-centric paradigm while partitions with many active vertices are processed by the FPGA in an edge-centric paradigm. Wang et al. propose a general graph processing approach with a novel worklist (priority queue) based graph computation and software scheduler (reorders vertices to be processed) [129]. The FPGA inserts work items (vertices) into a pre-scheduled queue and the CPU reorders them into a scheduled queue.

Near-data ExtraV [74] proposes graph virtualization. The CPU accesses the data through a cache coherent FPGA attached to an SSD storing the graph. Transparently to the CPU, the FPGA applies compression and multi-versioning to writes and decompression and filtering to reads.

Memory access For memory access, we found solutions for request merging, custom caching, and custom data placement among others.

Request merging A CAM-based approach for the BFS memory access problem is proposed in [69]. The

solution is an architecture-aware software graph clustering algorithm that reduces bandwidth requirements for random requests to visited flags. The clustering is applied as an offline pre-processing step. A memory unit merges multiple requests (cache for storing recently checked flags implemented with CAM).

Caching Zhang et al. propose a map-reduce based BFS approach on HMC memory [144]. With a performance model they find the bottleneck in scanning the bitmap of the current frontier. Thus, a second caching layer for the bitmap is introduced where one bit in the caching layer represents the aggregate of many bits in the RAM bitmap. FabGraph [106] is an extension of ForeGraph [34] by a two-level vertex caching (level L1 attached to pipelines and shared L2; L1 only communicates with L2). The vertices of the current graph partition are stored in L2 and replaced in Hilbert order such that vertices can be used for multiple graph portions.

Data placement Zhang et al. address the problem of redundant memory accesses caused by high-degree vertices for graph traversals [145]. They correlate vertex degree with data access frequency and propose degree-aware data placement and degree-aware adjacency list compression (Exp-Golomb variable length coding) combined with hybrid traversal approach on HMC memory.

Other Betkaoui et al. address stalling pipelines caused by memory access latency through a memory crossbar to share off-chip memory [15]. The solution issues many parallel memory requests and decouples memory access and execution units. Ni et al. accelerate BFS by applying a horizontal partitioning allowing to distribute the graph and its associated metadata over multiple memory channels [92]. In this way, multiple PEs can traverse the graph in parallel utilizing a high memory bandwidth. Upon level synchronization, active vertices are exchanged between the PEs. A memory access improvement for vertex-centric graph processing is given by Yan et al. [137]. Random memory accesses are sequenced and graph pruning is applied to prevent ongoing traversal from the leaves. The solutions further includes an online pre-processing step during the apply phase when bandwidth is under-utilized.

Performance model Performance models are about understanding effects of design decisions on a conceptual level. We found the following solutions on modeling parallelism and bottlenecks.

Parallelism Bondhugula et al. define a performance model for parallelism in two orthogonal parameters B number of PEs (B loops at once), and l denoting the number of operators in PE (l elements at once) [21]. Then constraints for different FPGA resources are mod-

elled for the two parameters. Shao et al. model the runtime of their design as the sum of the time of the vertex transmission and the time of the edge streaming [106]. Both are modeled dependent on the internal parallelism parameters of their design. The model is used to determine the size of L1 and L2 caches in the design.

Bottleneck FPGP [33] finds the optimal number of PEs by modeling the runtime as the maximum of the time spent on interval loading, edge loading, and edge processing since those can be overlapped. ForeGraph [34] extends this model by also including time to load intervals from other boards in their multi-FPGA setup and compares theoretical performance against other systems. Zhang et al. specify a performance model that is a slight variation of the network model [144]. The memory system is represented by the packet size, packet overhead, bandwidths, and internal latency of memory. An upper-bound performance model for vertex-centric graph processing on multiple FPGA systems is proposed by Engelhardt et al. [43]. The solution includes an architecture generator for multiple FPGAs with an application kernel and fitting dataset.

3.2.2 Document

Another early NRDS class is document (cf. Fig. 5) with a focus on XML in the FPGA-accelerated literature. While JSON is the predominant document format in commercial document stores (cf. Sect. 2.2.2), we did not find any solutions in the literature. In the following, we will discuss solutions we found to the system aspects for the document NRDS class.

Operator In the literature, we found operators for XML parsing, document filtering, XPath evaluation, Twig, and XML projection which we subsequently discuss.

XML parser Parsing of XML includes tasks like well-formed checking, schema validation, and tree construction. In [35], Dai et al. propose a solution that leverages recurring idioms in XML processing (one-to-one string match, one-to-many string membership test, one-to-many string search), and a speculative pipeline structure for tree construction skewed for common case high-throughput and edge case pipeline stalls. The FPGA is placed close to the network on a SmartNIC. An alternative approach by Sidhu implements tree automata as a pair of a lexical and a tree automaton where the states of the lexical automaton form the transitions of the tree automaton [108]. Huang et al. propose a sliding window XML parsing accelerator [59]. They assume that the XML is valid and based on that can process multiple non-delimiter characters in one cycle. Delimiter characters are still processed one by one.

Filter Chalamalasetti et al. [27] provide an implementation of document filtering with scoring against topic profiles. The work mainly improves the bloom filter design of [124]. The new design leverages multiple banks and reduces contention.

XPath Mitra et al. propose a solution for XPath-based filtering of XML documents by mapping the XPath queries to regular expressions [88]. These expressions are clustered by common profile prefixes and mapped to FPGA state machines (one per XPath). A global stack is used for the inherent parent-child relationships. In [42], publish-subscribe systems are extended to become an XML broker using an XPath processor with CAM. The article also provides a hardware-based XML parser. Moussalli et al. [89] address the challenge of recursive XML filtering. For that, each XPath is mapped into a stack whose width matches the XPath depths in bits and the height corresponds to the depth of the document. Open tags are handled as push events and close tags as pop events.

Twig The same authors propose an FPGA-based solution for twig matching on XML documents based on [89] combined with dynamic programming [90].

Queries For document stores, two different solutions for microsecond reprogrammable state machines as integral parts of string matching were proposed. This does not allow query processing in itself but chaining of operators which can be used for query processing.

Skeleton automata Teubner et al. propose the idea of skeleton automata as a fixed finite state automaton structure with parameterized transitions, allowing dynamic workload change in microseconds, instead of long (partial) reprogramming of the FPGA.

Other The ZuXA system [83] implements a general programmable state machine with hash index on a rule table and a clustering scheme for very large automata that is applied as an XML acceleration engine.

CPU-FPGA collaboration We found three different collaboration schemes: socket, near-data, and PCIe.

Socket Vanderbauwhede et al. address the challenge of power consumption of information filtering on streams of documents with a multi-FPGA setup [124]. A CPU places the document stream into main memory from where it is fetched by the FPGAs. An on-chip BRAM Bloom filter is used to quickly discard irrelevant documents before the documents are matched against profiles stored in a hash table in on-board memory.

Near-data The XLynx system [120] by Teubner et al. provides a solution for hybrid CPU-FPGA XQuery processing with dynamic XML projection. The FPGA im-

plements the same XML projector as in [119] placed into the data path between the document server and XQuery engine such that data is filtered before being queried, reducing load on the XQuery engine.

PCIe In [125], previous work on document filtering [27, 124] is advanced by embedding it into a hybrid CPU-FPGA system. The CPU handles parsing the network document stream passing it to the FPGA as 64bit words using separator words between documents. Additionally, the words are dictionary encoded.

3.2.3 Key-value

The most recent, NRDS-related research on FPGAs concerns key-value (cf. Fig. 5), arguably the conceptually simplest class of NRDS. We subsequently present the system aspect solutions for key-value stores.

Operator Solutions are found for the externally visible insert and integral internal hash operator.

Insert Liang et al. address the unpredictable insert performance of cuckoo hashing. It avoids hash collisions by computing multiple hashes per key and reinserts one key-value pair when all hash positions are already occupied [78]. The to-be-reinserted pair possibly triggers another pair to be reinserted into the hash table stalling naive pipelines. The proposed solution splits up the pipeline into an insert and reinsert pipeline.

Hash Fast key-based value access requires reliable hash-functions like Murmurhash2⁹ missing on FPGAs. Liu et al. contribute an implementation of Murmurhash2 on FPGA with different kernels for different key sizes that are applied through dynamic reconfiguration [80].

Binary representation The predominant binary representations are based on hash tables but there is also a CAM implementation.

Hash table Istvan et al. leverage the FPGA's scalability and energy efficiency for a hash table implementation sustaining 10Gbps by implementing a sophisticated pipeline with concurrency control and separate key-hash and value store [18, 61]. Collision handling is done with buckets by chaining fixed length pre-allocated memory regions for a tradeoff between probability of collisions and memory bandwidth. Tong et al. present a hash table with one operation per clock cycle throughput [121]. They propose two hash table access schemes. The first one provides multiple slots for each hash value to reduce collisions where each operation

scans all slots. For bandwidth limited deployments, instead of scanning all slots, a second hash function decides the slot to work on similar to cuckoo hashing. The tradeoff is again between probability of collisions and memory bandwidth. FASTHash provides even higher throughput by processing p queries in each cycle [140]. This is achieved by using p parallel, data-independent PEs with eventual consistency of updates. The hash table is split up into p partitions each owned by one of the PEs. Each partition is replicated to all PEs but only the PE that owns it inserts new values into it.

Content-addressable memory (CAM) Lockwood et al. propose a CAM-based, network-attached solution [81]. They define an own message format for optimized hardware parsing of requests. The key-value data is either stored in BRAM or DRAM while only BRAM guarantees low latency access and the value addresses are looked up with an emulated CAM.

Multi-tenancy Multi-tenancy is about performance and data isolation between multiple users working concurrently on a system.

Token bucket Istvan et al. achieve this by defining traffic shapers with token buckets (from networking) and introducing tenant-specific registers for temporary query data [62].

Scalability Data center level scalability is required by many applications using key-value stores. There are solutions on data replication for increased read throughput and data partitioning to aid work distribution.

Replication The Caribou system addresses scalability and fault-tolerance with data replication [64]. The FPGAs are deployed as a distributed storage layer in the storage nodes and allow for operator push-down (full scan and value predication) for near-data processing. For scalability, key-value pairs are replicated between the nodes such that read requests can be served by any node in the storage layer.

Partitioning FULL-KV [100] is a network-attached accelerator for CPU-FPGA hybrid systems, extending [99] to two nodes. The key-value store is partitioned to the nodes and requests are routed by a proxy. BlueCache [136] acts as a caching layer of distributed network-attached FPGAs. For operation with BlueCache, all application servers are equipped with a PCIe-attached FPGA that is connected to the other FPGAs via network and Flash for data storage. Each CPU collects key-value store requests and passes them to its accelerator card in batches. After parsing, requests are routed to the FPGAs containing the data and answered there.

⁹Murmurhash, visited 7/2020: <https://github.com/appleby/smhasher>

Availability Availability is about fault tolerance meaning responsiveness even when single nodes fail.

Replication Besides its contribution to scalability, the Caribou system [64] provides availability through replication. Writes are replicated from a master node to all other nodes in the storage layer. Data is still available when one node fails.

Consistency For the consistency system aspect, there is one solution on isolation with MVCC, thus covering only one facet of consistency.

Multi-version concurrency control (MVCC) Ren et al. define and implement MVCC support on top of [99] by storing a version-value B-tree for every key [101]. They also define atomic operations like compare and swap, compare and get, and predecessor version.

CPU-FPGA collaboration We found near-data solutions with the FPGA between CPU and data and solutions based on direct memory access to expand value storage.

Near-data In [73], the FPGA is used as an in-line accelerator for Memcached acceleration processing 96% of the requests. The CPU is only used as a fallback. Based on [99], Xie et al. design an FPGA distributed memory proxy with four pipelined data paths and a pipelined consistent hashing processor [134]. The orchestrator for the key-value store reaches up to 100Gbps.

Direct memory access (DMA) Li et al. propose a network attached key-value store based on SmartNICs that accesses the main memory over PCIe DMA [77]. The problems of the CPU cache hierarchy and FPGA low storage capacity are addressed by Qiu et al. [99]. Similar to [77], the hash table is in on-board memory and the data in main memory (accessed with PCIe DMA). The solution features novel memory allocation and fragmentation schemes.

Memory access For key-value stores, we found solutions for Bloom filters and Flash storage.

Bloom filter Cho et al. propose a key-value store with cuckoo hashing and decoupled hash table and values [30]. A Bloom filter is used to control hash table read access and thus reduce the amount of memory requests.

Flash Flash storage was proposed as a viable storage medium for values by Blott et al. [19]. Values are much larger in size than keys and are only very selectively accessed after their address has already been found over the hash table. Thus, storing values in Flash storage allows much larger amounts of data to be stored. Blott et al. scale out a Memcached server to 40 Terabytes of data this way [19]. Similarly, BlueCache stores the hash table as an index cache in RAM and stores the corresponding values on the attached Flash storage [136].

Performance model We found one solution modeling the performance bottleneck of the accelerator design.

Bottleneck Qiu et al. examine the theoretical performance of their system based on the network performance as the bottleneck of the system [100]. They split up their analysis by put and get operator.

3.3 Synthesis and discussion of system aspects

The system review in Sect. 2 resulted in a taxonomy of relevant FPGA and NRDS system aspects (cf. Fig. 4) that guided the literature analysis that addresses hypothesis (H2), i. e., *there are significant gaps in current research on non-relational FPGA acceleration* as set out in Sect. 1.2.

To address (H2), we first answer question (a) What are the most relevant of the identified NRDS classes? by studying potential contributions chronologically. We found that while document still plays a role in recent research efforts, current emphasis lies on accelerating graph and key-value. This also supports the necessity of this work to consider FPGA-accelerated NRDS and their design decisions.

Secondly, we strive to answer questions (b) Are there any system aspects that are not yet covered by literature? and (c) Do existing approaches provide solutions to these topics? through the selection and summarizing of research literature. Table 2 sets the relevant FPGA and NRDS system aspects into context to the NRDS classes. The focus on system aspects comes from our objective to guide practitioners and research towards FPGA-accelerated NRDS.

We found a large body of work on general graph operators (e. g., BFS) and binary representation in the HPC domain. While HPC does not focus on database processing, the work denotes a starting point for the graph database research. Notably, when it comes to database-specific system aspects like queries (i. e., GraphSoC [68]), scalability (i. e., most notably ForeGraph for multi-FPGA graph traversal [34]), and consistency (i. e., simple locking [84]) only few solutions are provided. Due to the focus on FPGAs, there is an equally large body of work on design paradigms (e. g., HitGraph [148, 149, 150]), CPU-FPGA collaboration (e. g., most notably ExtraV [74]), memory access, and performance models. No solutions were found for availability, multi-tenancy, and security.

In the document class, there are only few solutions for functional system aspects. Most notable is the XLynx system [119, 120]. Non-functional system aspects are not covered. FPGA system aspect solutions were

only found for CPU-FPGA collaboration [120, 124, 125]. All other categories remain without a solution.

For key-value stores, important functional topics like operators and binary representation are covered (e.g., hash tables [61, 121, 140] and CAM [81]). Queries and sophisticated design paradigms are not applicable to key-value since there are only CRUD operators. The non-functional system aspects scalability, availability, and multi-tenancy are mainly studied in the Caribou system [62, 64]. Further solutions are provided for consistency (i.e., isolation with MVCC [101]) as well as CPU-FPGA collaboration, memory access, and one performance model. We found no solutions for security.

In summary, we only found FPGA-accelerated NRDS for key-value databases, while in the other NRDS classes, solutions for functional system aspects are available. For NRDS, these solutions would have to be adapted to databases to be usable. There were only few non-functional solutions provided.

4 FPGA-accelerated NRDS

This section gives several perspectives on NRDS class commonalities, addressing hypothesis (H3), i.e., *there are patterns guiding the design of an FPGA-accelerated non-relational database system*. We revisit existing practical FPGA-accelerated NRDS solutions [1, 2, 3] (cf. Sect. 2) and the current state in research (cf. Sect. 3) and discuss overarching patterns that we found along four questions on building an FPGA-accelerated NRDS starting from the idea or need of accelerating a system, over its design and implementation, and finally the evaluation of the accelerator’s impact:

- Q1 Which problems are FPGAs able to solve? (accelerator task: Sect. 4.1)
- Q2 Where to put the FPGA(s)? (accelerator placement: Sect. 4.2)
- Q3 How to implement NRDS operators? (accelerator design: Sect. 4.3)
- Q4 How to measure improvements of the NRDS? (accelerator justification: Sect. 4.4)

Subsequently, we give answers to questions Q1–4 in the form of a practitioners guide.

4.1 Accelerator task

When dealing with a concrete system implementation, either an accelerator is added to an existing system or becomes relevant to the design of a new one. However, in both cases one has to decide whether an FPGA is suitable. As a first step, we answer question Q1 (Which

problems are FPGAs able to solve?) by summarizing the problems typically solved by FPGAs in the different NRDS classes and differentiating task categories that are well suited to FPGAs. Notice, however, that in this survey we cannot represent all possible motivations and tasks that are met by practitioners in practice.

Graph For graph processing, FPGAs are mainly used to offload operators because of inefficiencies in the cache hierarchy and coarse-grained memory access of CPUs resulting from the irregular memory access patterns of graph workloads (cf. [7, 34, 149]). An FPGA has more control over its memory access and thus helps to alleviate the problem of irregular memory accesses through custom memory controller design and full control of data placement in on-chip memory (cf. Sect. 4.3).

Document FPGAs for document processing are mainly used as bandwidth amplifiers for the CPU. Much of document processing is parsing and filtering with large data movement costs putting heavy load on the CPU even though a lot of the data is discarded and pollutes the cache hierarchy. Thus, FPGAs can be used as a flexible stream processing accelerator in the data path to the memory (e.g., [27]), disk, or network (e.g., [35]). Sometimes, however, the CPU is completely bypassed and the FPGA is used as a standalone accelerator in the network (e.g., [119]).

Key-value For key-value stores, the literature mainly shows two schemes for motivation of FPGA usage. Either a single function instrumental to key-value stores is accelerated (e.g., insertion [78] or hashing [80]) because the CPU does not meet the latency requirements, or building a full system is motivated by large roundtrip latencies of CPUs from the network through the operating system network stack and back (e.g., [64, 77, 100]).

Summary – task categories The two biggest problems that FPGAs solve for the different NRDS classes are: (1) data movement from peripherals (e.g., network or disk) to the CPU and (2) memory access inefficiencies by the fixed CPU architecture. Figure 6 shows a simplified version of the shared system architecture from Fig. 4. The tasks that FPGAs might take over in an NRDS to address these problems fall into three categories: operator acceleration, data access acceleration, and communication layer.

Operator acceleration focuses on improving performance for one or multiple operators of the NRDS class. A lot of the literature focuses on this task category with implementations of specific operators.

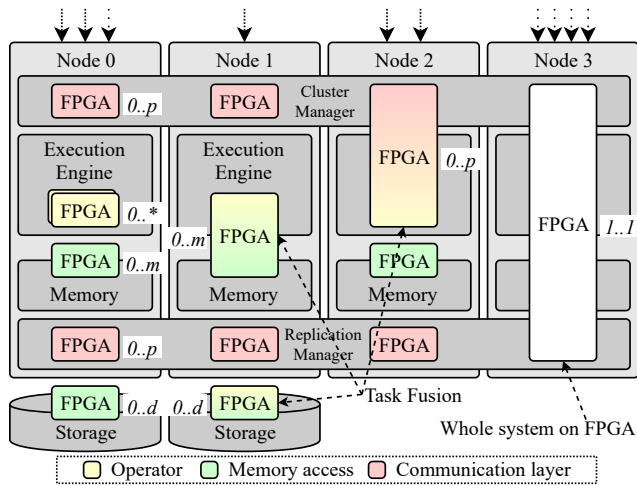


Fig. 6: Potential FPGA tasks in NRDS

Particularly interesting for NRDS are accelerator implementations in the context of hybrid CPU-FPGA systems (e.g., [20, 65, 124, 125, 129]). In data-centric applications, the FPGA can take pressure off the CPU with data access acceleration. One already existing example is graph virtualization, where non application-specific access patterns are accelerated on an FPGA near storage [74].

Placing FPGAs in the communication layer is another promising option. One example is a proxy layer for key-value stores where the request router (cf. Fig. 4) is placed in an FPGA outside the other nodes which routes traffic to the correct CPU nodes [134].

On node 1 and 2 in Fig. 6 we show possible combined acceleration of tasks that we call *task fusion*. Task fusion is possible when the resources on the FPGA fit both tasks and should be done if more than one task benefits from FPGA acceleration. The tasks combined on one FPGA may even accelerate the overall system more than if they would be accelerated separately because data movement costs are reduced.

In the relational database literature we found examples of operator push down where operator acceleration and data access are fused into one (e.g., [49, 133]). This worked especially well for filter operators pushed to the FPGA that reduced the amount of data communicated to the CPU. This kind of accelerators use their close proximity to memory to reduce the data load on the CPU. The newly presented Enzian system [5] also enables this with a cache-coherent attachment of the FPGA and opens up questions of near-data processing where the FPGA is used for on-the-fly conversion of binary representation. One prominent example of fusing operators with the communication layer at data center scale is the Microsoft Catapult project [98].

There, servers push document classification workloads to a communication layer of multiple FPGAs connected to each other. Another example from the key-value literature is [99]. They use the FPGA as an entry point from the network and do pre-processing of queries in the FPGA while – due to size – storing the actual data values in the memory of the CPU.

In extreme cases, one to all nodes in the NRDS cluster can be mapped to FPGAs (node 3) (e.g., [64]). This works well for key-value systems (and might work for wide-column systems) since the overall system is simple enough, and also works for providing standalone services on FPGAs (e.g., [120]). In this case the FPGA has to be network-attached which saves a lot of overhead by not going through the CPU cache hierarchy and operating system network stack.

To aid the decision of how many FPGAs to put into a system we added cardinalities to each possible FPGA task (Fig. 6). In the confines that are set by the system hardware, independent FPGAs could be added for different tasks. One FPGA can be added for each network port (p is number of network ports), memory subsystem, and disk (m is number of memory subsystems; d is number of disks) in the system. For operators, there is no such restriction. There can be as many FPGAs as fit into the hardware system. Insight 1 answers question Q1 (Which problems are FPGAs able to solve?):

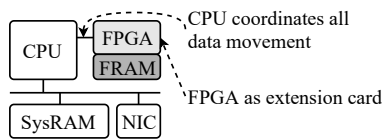
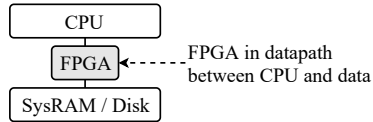
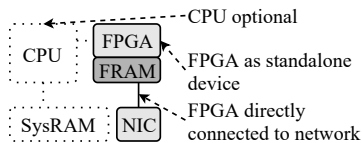
Insight 1 *There are three accelerator task categories (operator, data access, and communication layer acceleration) that FPGAs are currently well suited for.*

4.2 Accelerator placement

After deciding for which task the FPGA accelerator is used in the system, in this section, we discuss FPGA placement patterns in the context of one cluster node. In the literature, we discovered four FPGA placement patterns which we discuss before we show how to chose a placement based on the task (cf. Sect. 4.1) and properties of the workload. This answers question Q2 (Where to put the FPGA(s)?).

Placement patterns We differentiate between the main memory of the overall system (SysRAM), attached to the CPU, and RAM directly attached to the FPGA on the board (FRAM).

The *offload* accelerator placement (Fig. 7) is defined in being attached only to the CPU (e.g., over PCIe) and an on-board FRAM much smaller than SysRAM. Input data is directly written to FRAM by the CPU and execution is triggered by the CPU. The FPGA

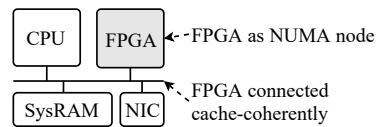
Fig. 7: Placement pattern: *Offload*Fig. 8: Placement pattern: *Near-data*Fig. 9: Placement pattern: *SmartNIC*

works on the input data in FRAM, and the results are transferred by the CPU to the SysRAM on notification of the CPU by the FPGA. This placement only allows master-slave setups, and thus can introduce a lot of overhead for acceleration because the FPGA cannot move data in the system on its own, and CPU cycles are wasted on data movement. The FPGA-accelerated in-memory database survey [46] shows this placement as *IO-attached accelerator*.

The second placement we found in the literature is the *near-data* placement (Fig. 8). It is defined by the way the FPGA is inserted into the data path between the CPU and the SysRAM or disk. In this way, the FPGA provides an interface to the CPU to interact with the underlying resource. In a more restricted way, [46] define this as a *bandwidth amplifier* which decompresses data, however this placement can accelerate more workloads than just decompression, e. g., filtering and binary representation conversion. In [91] there is a similar placement where the FPGA is placed in the data path between the disk and the CPU.

Figure 9 shows the *SmartNIC* placement where the FPGA is directly attached to the network interface controller (NIC). This placement option may even completely eliminate the CPU in the system if there are no tasks besides what is implemented on the FPGA. The SmartNIC placement optimizes for low latency of the overall system by saving multiple round trips through the operating system kernel on the CPU.

In emerging systems (e. g., [5]), the FPGA may be placed as a *socket* (Fig. 10) at least conceptually with a cache coherent access to SysRAM. The three previ-

Fig. 10: Placement pattern: *Socket*

ously discussed placements can be represented with the FPGA being a socket with little overhead. However, the socket placement also enables new work distribution strategies where the CPU does not have to coordinate execution of the FPGA and data movement. In [46] this placement option is called *co-processor*.

Placement decision FPGA placements discussed in the literature can be reduced to the four fundamental patterns from Sect. 4.2. Based on these, Fig. 11 shows a decision tree guiding the practitioner towards choosing an accelerator placement pattern depending on task and workload properties. Additionally, we added CPUs and GPUs as alternative accelerator options.

For tasks that incur large data movement overheads from either the memory, disk, or network we have introduced shortcuts (shown as blue arrows) to the near-data and SmartNIC placements respectively.

For workloads that do not exhibit massive parallelization opportunities we do not see much potential in applying an accelerator. Thus, this leads to adding more CPUs to the system or alternatively adding more nodes to the cluster. For compute-bound problems with structured parallelism, meaning large numbers of homogeneous threads running in parallel, we would chose GPUs over FPGAs because they are specifically made to handle these workloads [79]. Similarly, for tasks with heavy reliance on unstructured floating-point operations, we would most of the time advise against using an FPGA as an accelerator because the DSP floating point units on the FPGA will quickly become the bottleneck.

For compute-bound problem instances that are not better suited to GPUs or multi-CPU, an *offload* accelerator approach is chosen. The data movement is a big source of overhead in this placement model such that it only works for compute-bound problems where data movement costs, on the slow link between CPU and FPGA, are negligible compared to the duration of the computation. The offload accelerator approach is implemented by most of the graph literature (e. g., [34, 47, 149]). While we think that the offload pattern could be a viable option for database systems, we focus on more significant improvements through acceleration.

In the category of memory-bound problems, we differentiate first between workloads with simple operators (e. g., lookup) and second workloads with complex op-

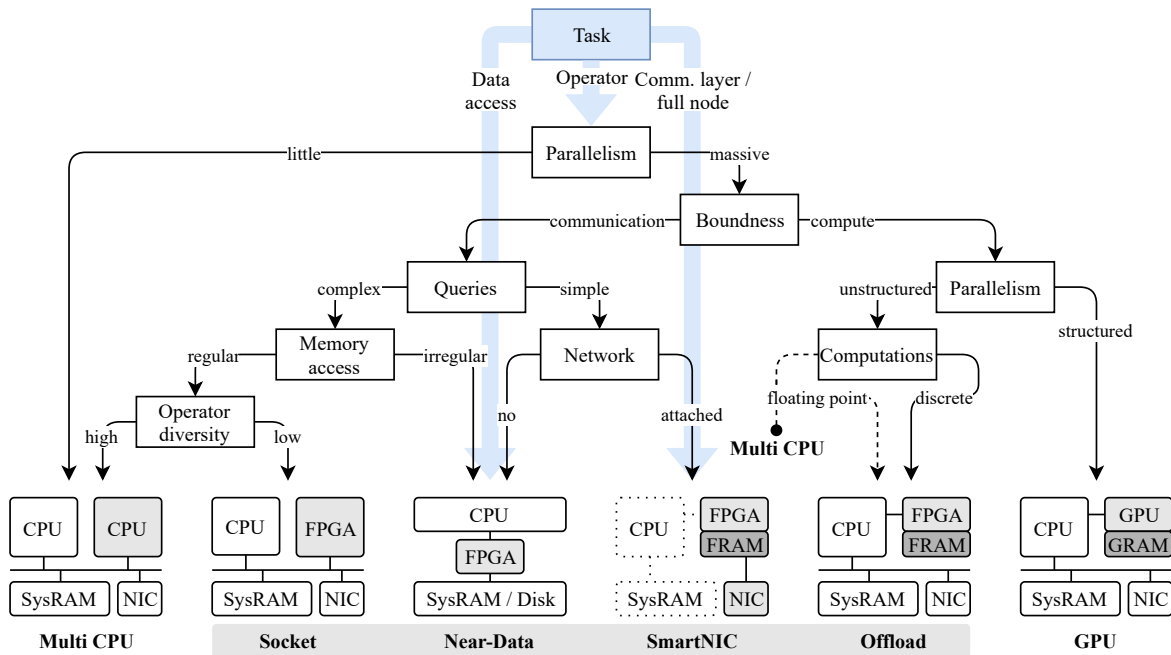


Fig. 11: Accelerator placement decision tree

erators (e.g., graph traversal). Simple queries are defined as a combination of few simple operators with few predicate expressions. This category includes key-value and wide-column stores and can include document and graph stores in certain scenarios (e.g., data provider for graph neural networks). If the database system is network-attached, we choose a *SmartNIC*. This placement was also found to be efficient in related data processing domains like data-intensive messaging [103, 104]. If the database system is only part of a larger architecture and not network-attached, we chose the *near-data* approach. This placement option is also chosen when there are complex queries with irregular memory accesses (e.g., in graph traversal). We think that the *socket* placement will benefit FPGA accelerated systems especially in the database context.

As shown in Fig. 11, adding an FPGA to the system is not always the best strategy to improving performance of a system. Moreover, the traditional approach of placing an accelerator in a system as an offload accelerator is not the best option for database systems.

Table 3 shows all systems found in the system review and literature analysis and how the placement decisions look like. We exclude the systems describing frameworks without a description of how they are deployed (i.e., GraphGen [93], GraVF [44], and GraphOps [94]). For tasks that have a quick path to the placement patterns, the decision from the decision tree is always correct. This includes the commercial solutions from the system review, e.g., CAPI SNAP [1] which uses near-

data deployment of FPGAs for graph workloads. For pure operator acceleration we propose near-data placement for graph and socket placement for document workloads. This is only implemented by CyGraph indicating a lack of consideration for the whole system and therefore data movement in the offload placement pattern. The offload placement pattern can be seen as the initial step an inexperienced practitioner might take, since it trades off performance against easy system integration.

Summary – accelerator placement In this section we first showed how FPGAs can be attached to the other hardware components in a system. Especially compared to a CPU, FPGAs can be placed close to the data, whether in memory, disk, or network.

Insight 2 *There are four fundamental patterns of FPGA placement (offload, SmartNIC, near-data, and socket).*

Thereafter, we established a decision tree guiding the practitioner from the tasks (cf. Sect. 4.1) and the characteristics of the operators towards choosing a placement pattern. We validated this decision tree with the systems found in the system review and literature analysis by comparing our and their placement decision.

Insight 3 *The accelerator task in combination with the characteristics of operators of the NRDS class are sufficient to decide the FPGA placement.*

NRDS class	System identifier	Task	Queries	Mem. access	Op. div.	NRDS	Decision	Correct
Graph	CAPI SNAP [1]	D	complex	irregular	n/a		Near-data	
	HAGAR [86]	O	complex	irregular	n/a		Near-data	(Offload)
	GraphStep [39]	CO	complex	irregular	n/a		SmartNIC	
	CyGraph [7]	O	complex	irregular	n/a		Near-data	
	TorusBFS [76]	CO	complex	irregular	n/a		SmartNIC	
	FPGP [33]	O	complex	irregular	n/a		Near-data	(Offload)
	ForeGraph [34]	CO	complex	irregular	n/a		SmartNIC	
	ExtraV [74]	D	complex	irregular	n/a		Near-data	
	Dr. BFS [47]	O	complex	irregular	n/a		Near-data	(Offload)
	FabGraph [106]	O	complex	irregular	n/a		Near-data	(Offload)
	HitGraph [149]	O	complex	irregular	n/a		Near-data	(Offload)
Document	ZuXA [83]	O	complex	regular	low		Socket	(Offload)
	XLynx [120]	CO	complex	regular	low		SmartNIC	
Key-value	Algo-Logic [2]	C	simple	n/a	n/a		SmartNIC	
	BlueCache [136]	D	simple	n/a	n/a		Near-data	
	Caribou [64]	F	simple	n/a	n/a		SmartNIC	
	KV-Direct [77]	F	simple	n/a	n/a		SmartNIC	
	FULL-KV [100]	F	simple	n/a	n/a		SmartNIC	
Wide-column	rENIAC [3]	C	simple	n/a	n/a		SmartNIC	

: commercially available; task options (D: data access, O: operator, C: communication layer, F: full system), combinations are permitted; n/a: “not applicable”; Mem. access: Memory access, Op. div.: Operator diversity; : yes, : no

Table 3: Accelerator decision validation for systems in the literature (all systems are communication bound)

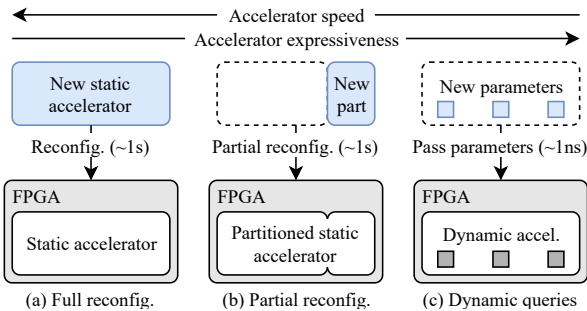


Fig. 12: Strategies for operator switching

4.3 Accelerator design

Now that we know what the FPGA should do and where it is placed in the system we need to answer question Q3 (How to implement NRDS operators?). Although implementation details largely depend on specific algorithms and data structures of the NRDS classes, we found patterns for two critical accelerator design considerations. In the following we introduce operator switching strategies and memory access optimization patterns common for all NRDS classes on FPGAs.

Operator switching strategies In the literature analysis we saw little about accelerators that are able to switch to different operators without compiling a new accelerator each time. However, it is required of a NRDS accelerator to process multiple different operators in parallel and in very quick succession on multiple different datasets in memory, since it is not sufficient for an accelerator to improve the performance of one infrequently used operator to offset added cost and complexity. This is easy to achieve on instruction-based architectures,

like CPUs, since their programs are easily switched out but difficult to achieve on FPGAs since they cannot switch their architecture without significant overhead. This is shown in Fig. 12(a) as *full reconfiguration* where an operator switch takes seconds [97].

To alleviate the overhead of full reconfiguration, the relational database community pursued *partial reconfiguration* (Fig. 12(b)) where only parts of the accelerator architecture are switched (e. g., [41, 85, 95, 131, 153]). While this works for coarse-grained functionality switching during runtime, the part that is reconfigured still is unavailable for seconds.

Thus, we advocate for a more elegant and expressive solution in what we call *dynamic queries* (Fig. 12(c)). The dynamic queries switching strategy is based on a dynamic, parameterized or instructable, domain-specific accelerator (e. g., [63, 68, 83, 112, 119, 120]) that allows to process multiple different operators in parallel and with only nanosecond switching delay on multiple datasets by just passing new parameters or instructions. The added accelerator expressiveness might come at the loss of some accelerator speed [117] but saves magnitudes in operator switching delay. The difficulty of designing such an accelerator lies in finding abstractions of high generality without introducing much overhead that slows down performance. The examples we found above focus on a small set of domain-specific primitives that are combined into a dynamic accelerator.

Although we did not find references for it in the literature, the operator switching strategies can be applied in combination. Depending on the task it might be beneficial to partially or even fully reconfigure the FPGA as long as it is done at a frequency low enough (e. g., if lasting workload changes are detected).

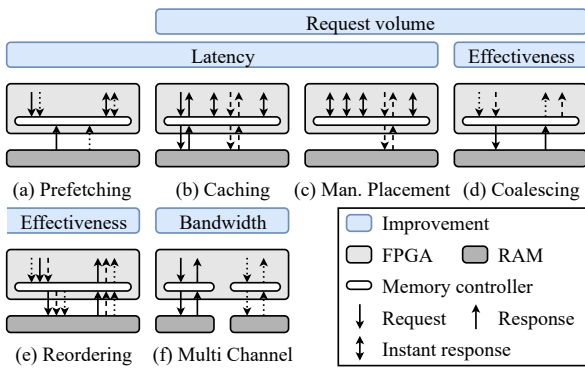


Fig. 13: Memory access optimization patterns

Memory access optimization patterns For NRDS, memory accesses are one of the most instrumental challenges to good performance. We found six memory access optimization patterns (Fig. 13) in the literature that are applicable to all NRDS classes. Each pattern is implemented in the memory controller (endpoint to memory on the FPGA) of the accelerator design and is independent of the accelerators algorithms and data structures. The memory controllers performance can be improved along four axes: by reducing *latency* of accesses, reducing the number of accesses (*request volume*), enhancing the amount of effective data per access (*effectiveness*), and increasing raw memory *bandwidth*. In the following we introduce these six combinable patterns.

Prefetching (cf. Fig. 13(a)) Prefetching is a technique to hide memory access latency that starves processing of input data. Therefore, memory requests are issued before the data is processed to overlap computation and loading of data if the access locations are known beforehand. By the time the data is processed, it already resides on chip to be consumed. One example is partitioning the data and overlapping partition loading with partition processing [106, 147]. Another is issuing large amounts of non-blocking memory requests such that there is always data to process [16].

Caching (cf. Fig. 13(b)) Caching (or automatic data placement) reduces high-bandwidth utilization by storing accessed values in on-chip memory (cache). If the cache is full, there is an automated policy (e. g., least recently used) replacing values in the cache with new ones [66]. If a value in the cache is requested, it is instantly served from on-chip memory, but this only works well for workloads with strong temporal locality. One example is multi-level caching in [106]. In [129], Wang et al. combine caching with the before mentioned reordering to increase the spatial locality of memory accesses.

Manual data placement (cf. Fig. 13(c)) FPGAs provide the practitioner with full control over what data

resides in quickly accessible on-chip memory not only custom caching techniques can be employed but critical data can be permanently placed on the FPGA. This may be done for frequently accessed data structures critical to the performance of the accelerator (e. g., [16, 78, 144, 145]). Another example is storing a highly efficient data structure like a bloom filter that allows filtering of memory accesses for presence of values in a dataset (e. g., [11, 30, 124]).

Coalescing (cf. Fig. 13(d)) Coalescing means merging multiple memory requests of single data items into one memory access (e. g., [69]). Since modern DRAM operates on rows of memory that are in the kilobyte range, accessing single data items is wasteful. If the workload exhibits strong spatial and temporal locality, coalescing can reduce the number of memory accesses per single data item and thus decrease the request volume by simultaneously increasing the effectiveness of each memory access. In [47] this is done by having many more compute units than data access units that issue many accesses enabling data access units to coalesce some of the accesses. Another example is combining write requests before they are written to memory and thus reducing the overall number of writes [146, 148].

Reordering (cf. Fig. 13(e)) Usable memory bandwidth suffers from irregular accesses. Reordering of memory requests can improve upon this if the workload exhibits spatial and temporal access locality. This can be done online (e. g., [137]), at the cost of increased latency, or offline (e. g., [147, 148]), if there is a correlation between data and memory access order.

Multi Channel (cf. Fig. 13(f)) If multiple memory channels are available, the memory bandwidth can be increased by distributing memory accesses over those channels. This is a meta-pattern that can be combined with any of the aforementioned patterns. One example of using multiple channels is placing different data structures on different channels [75, 92].

Summary – accelerator design While there was no focus on operator switching in the HPC-motivated literature, there has been some work on the topic especially in the document class. We see it as a crucial consideration in FPGA-accelerated NRDS. Furthermore, the memory access optimization patterns are especially important to NRDS since memory access acceleration is one of the big motivations to use FPGAs. Thus, we conclude with the following insight:

Insight 4 *There are three operator switching strategies and six memory access optimization patterns guiding the development of accelerators for NRDS.*

Domain	Name	Dataset	Workload	#Refs
Graph	graph500	👍	👍	1
	live-journal	👍	👎	7
	rmat	👍	👎	4
	pokec	👍	👎	4
	orkut	👍	👎	3
	twitter	👍	👎	3
	wiki-talk	👍	👎	2
	indochina	👍	👎	2
	flickr	👍	👎	2
	roadnet-ca	👍	👎	2
Doc.	xmark	👍	👍	4
	trec aquaint	👍	👎	3
	toxgene	👍	👎	3
	yfilter	👎	👍	3

👍: yes, 👎: no

Table 4: Benchmark statistics

4.4 Justification

In an FPGA-accelerated system, the FPGA’s improvement on performance must not be evaluated on the accelerated part of the workload in isolation but the whole system. An FPGA introduces a new component into the hardware system that entails costs having to be justified by the performance improvement. Costs occur in the form of cost of ownership (which might benefit from FPGA energy efficiency) but also cost of programming and operating a whole new hardware architecture. While the decision if the performance improvement outweighs the cost lies in the judgement of the practitioner, in this section we will introduce measuring performance improvement with benchmarking and performance models for FPGA-accelerated NRDS. There-with we answer question Q4 (How to measure improvements of the NRDS system?).

Benchmarks As a guideline to good benchmarking we follow the four key criteria for domain-specific benchmarks from [54]. Benchmarks should be easy to understand (*simple*) and scale from small to powerful systems in the present and towards the future (*scalable*). However, the two criteria most critical to FPGA-accelerated NRDS are that benchmarks are *portable* and *relevant* which we discuss in the following.

Benchmarking only works well when either comparing different systems with the same program or different programs with the same underlying system. Either the program or the system as variables have to be fixed for comparability. This especially poses a big problem for benchmarking on FPGAs since different FPGAs have very different specifications, and implementations are often tuned to one specific board. We did not find any solutions for this problem in literature.

Regarding relevance, we found different datasets and workloads listed in Tab. 4. Excluding “indochina”, all

non-synthetic graph datasets can also be found in the Stanford Network Analysis Project (SNAP) graph collection¹⁰. Considering the overall number of articles we found by NRDS class, it does not seem as if any workload or dataset is established especially when compared to e.g., TPC in the database literature. For key-value systems, benchmarking is currently done in literature as a sequence of a subset of the three basic API functions (i.e., get, put, delete) without a standardized dataset. However, a well-formed benchmark establishes not only datasets and workloads but all of the following artifacts: (1) Workloads (e.g., for graph: BFS, shortest path, weakly connected components) (2) Datasets (e.g., for graph: twitter, rmat, live-journal) (3) Domain-specific performance measures (e.g., for graph: traversed edges per second (TEPS)) (4) Benchmark (reference) implementation details. Thus, current performance measurements not only lack relevance and artifacts to be regarded as benchmarks.

A solution could be provided by existing comprehensive benchmarks that are not yet used in literature [102]. The YCSB program suite [32] offers capabilities for benchmarking key-value and document stores. For graph stores there are the LDBC Graphalytics Benchmark, LDBC Social Network Benchmark, and GAP Benchmark Suite. The GAP Benchmark Suite and the LDBC Graphalytics Benchmark cover kernels common in graph analytics (e.g., BFS or PageRank) while the LDBC Social Network Benchmark covers querying workloads. Recently, a cross NRDS class benchmark was proposed [143]. These benchmarks, designed for CPU-based systems initially, could also be used to benchmark new designs on FPGAs with modifications for FPGA-specific problems to make them portable.

Performance models As established in Sect. 2.1, FPGAs exhibit unstructured parallelism exacerbating comprehension of performance and algorithm complexity on an abstract performance model level. In the literature we found different ways (cf. Tab. 2) to model performance of the proposed solutions breaking down to modeling pipeline and data parallelism in the face of constrained resources (logic resources and memory bandwidth). This means, the system architecture is broken down into components with known performance (i.e., pipeline steps or replicated PEs) where performance is measured in throughput of data which scales linearly with pipeline steps and data parallelism. Sometimes this is embedded into a roofline model where the performance is first capped by the amount of parallelism and later by the available memory bandwidth.

¹⁰SNAP dataset collection, visited 7/2020: <https://snap.stanford.edu/data/index.html>

Summary – justification Justification in the form of workloads and datasets performing well on FPGAs is provided by the literature but the performance measurement landscape is scattered and sometimes specifically tuned to the contribution of the article. The biggest challenge, however, is the lacking portability of current benchmarks because measurements are performed on vastly different FPGA setups without accounting for e. g., different memory bandwidths. It is sometimes unclear if performance improvements stem from better design or just better hardware. Thus, we conclude with the following insight:

Insight 5 *A portable, relevant benchmark suite that covers all necessary artifacts is missing for robust justification of accelerator usage decisions.*

4.5 Discussion – insights

Over the course of this section, we gained five insights into building an FPGA-accelerated NRDS answering questions Q1–4. From the motivations of the different NRDS classes in Sect. 4.1 and resulting tasks we saw that two of the biggest challenges of current CPU-based systems are data movement and memory access. With the patterns we found for placement and memory access optimization, we guide the practitioner towards addressing these challenges with FPGAs regardless of NRDS class augmented with common operator switching strategies. However, performance largely depends on specific algorithms and data structures designed on a use case basis. In this regard, we were not able to uncover even more inter-class structure on such a high abstraction level. Nonetheless, the insights we gained help the practitioner to apply FPGAs to existing or newly designed NRDS, regardless of data model, to take pressure off the CPU or eliminate it from the system architecture completely. These common patterns might even foster building FPGA accelerators supporting multiple data models at once.

5 Open research challenges

The literature analysis did not only summarize many interesting solutions but showed several gaps lacking solutions (cf. Tab. 2). In this section we discuss the important open challenges and research questions we think should be pursued in the near future.

Regarding the differences between the NRDS classes the open challenges vary. We found that key-value systems are exceptional in that there already exist complete NRDS (research and commercial). However, key-value systems are the simplest NRDS class, and their

system design leaves many questions unanswered that come up in other classes. Wide-column stores are not represented in the literature, but solutions from the key-value class should be applicable (cf. [37]).

The literature on document and graph stores as well as existing accelerator prototypes indicate feasibility. However, the found solutions in the literature are rather HPC-specific and cannot directly be applied to NRDS. Thus, there are in general many challenges to be solved towards a complete NRDS.

Besides these rather broad considerations, we identified several open research challenges in the course of this survey that we discuss subsequently.

Non-functional NRDS aspects As a broad trend in the literature, the coverage of non-functional NRDS aspects are an open challenge. While consistency protocols may be transferred from the non-accelerated NRDS literature, it is broadly unclear how to provide production-grade scalability, availability, multi-tenancy, and security with an FPGA-accelerated NRDS. FPGAs as a relatively new processor architecture for data processing are not integrated as deeply into current systems and do, in contrast to CPUs, not have widely used operating systems providing basic functionality. Possibly some solutions can be transferred from FPGA-accelerated relational database systems.

Flexibility – dynamic queries The static implementations presented in the current literature need to become more flexible. There are elegant solutions (e. g., skeleton automata [119]) to be found for intractable or parametrizable accelerators that can process more than one rigid operator. This follows the proposal of domain-specific architectures in [58] and will also largely improve the projected ratio of the workload processed by the accelerator leading to better overall performance of the system.

Comparability – standard benchmarks As discussed in Sect. 4.4, there are no commonly used benchmarks in the NRDS classes on FPGAs yet. A standardized benchmark will be instrumental in gaining more credibility in performance claims, comparability, and justification of FPGAs as an NRDS accelerator. Moreover, better performance measures help uncover performance impediments in other domains like DRAM (bank parallelism utilization analyzed in [52]).

Collaborative memory usage Data movement overhead dominates accelerated systems and narrows their potential for improvements. The emergence of cache coherent attachments of FPGAs to the system main mem-

ory might alleviate this. FPGA-directed data movement could take pressure off the CPU and also make more fine-grained acceleration possible. However, we did not find any solution in the literature for the collaborative usage of the system main memory and smart movement of data.

Near-data memory usage for graph systems Especially for document stores, we found deployments of the FPGA in the datapath between CPU and memory or disk (i. e., near-memory) to reduce the volume of data being moved to the CPU. Since graph workloads are highly memory-bound, we see big potential for a tighter integration of FPGA and memory to hide inefficiencies of irregular memory access patterns.

Cross data model processing A relatively new trend in NRDS are cross data model systems, i. e., systems that allow storing and accessing data in multiple data models simultaneously [82]. One example is OrientDB that supports polymorphic queries over graph and document data in one unified system. FPGAs could, e. g., be used near-data to transform and change the binary representation on-the-fly as data is loaded to the CPU.

Cross-accelerator architecture As GPUs and FPGAs get more popular and ever more present in the data center, there might be more performance gains in heterogeneous systems using both accelerator types at the same time. There are first works on those systems in other research areas (e. g., [56]), but none in the NRDS literature. This kind of acceleration might be especially beneficial to NRDS supporting multiple data models, where different workloads are particularly well-suited to different processor architectures.

6 Summary

FPGAs are an instrumental tool in achieving performance gains in data-centric systems in the near future. In this survey we open up the field of FPGA-accelerated non-relational database systems (NRDS) by studying and answering three hypotheses formulated in Sect. 1.2.

To start with, for hypothesis (H1), i. e., *existing non-relational database systems do not realize the potential of FPGA acceleration*, we conducted a system review of commercial NRDS. We found three experimental extensions to existing systems using FPGAs as accelerators [1, 2, 3], showing the NRDS-acceleration feasibility but no mainstream adoptions of FPGAs in NRDS, yet. Hence, we confirm the potential of FPGAs as accelerators for NRDS, but also conclude that this potential is not yet realized in commercial systems.

To give an answer to hypothesis (H2), i. e., *there are significant gaps in current research on non-relational FPGA acceleration*, we derive a system aspect taxonomy that guides an extensive literature analysis by categorizing the large body of research, for which we provide an insightful overview of existing contributions.

Taking the results of the literature analysis as a knowledge base, we derived common patterns, confirming hypothesis (H3), i. e., *there are patterns guiding the design of an FPGA-accelerated non-relational database system*. Therefore, we compile a list of four relevant questions that practitioners, like system architects, have to ask themselves when designing and constructing an FPGA-accelerated NRDS. This survey gives answers to these questions by easy-to-apply patterns for FPGA task definition, FPGA placement, accelerator design considerations, and benchmarking.

In summary, we provide a comprehensive introduction of FPGA acceleration and CPU offload potential for NRDS and present it in a form suitable to everybody interested in the field. However, we especially regard this survey as a guide for system architects in their decision making and a reference for researchers to guide and conduct new research.

Acknowledgements

We thank Norman May and Wolfgang Lehner for various valuable discussions in the context of this article.

References

1. Accelerating Neo4j with CAPI SNAP from IBM Power Systems. <https://neo4j.com/blog/neo4j-capi-snap-ibm-power-systems/>. Visited 7/2020
2. Achieving a half-billion IOPs in a 1U Redis server with FPGA acceleration. http://algorlogic.com/sites/default/files/2020_05_12-Half-Billion%20IOPs%20in%20a%201U%20REDIS%20server%20with%20FPGA%20Acceleration_Algo-Logic_JWLockwood-1M.pdf. Visited 7/2020
3. rENIACs FPGA-based distributed data engine boosts Cassandra database performance as much as 10x when deployed as data proxy. <https://blogs.intel.com/psg/reniacs-fpga-based-distributed-data-engine-boosts-cassandra-database-performance-as-much-as-10x-when-deployed-as-data-proxy/>. Visited 7/2020
4. Abadi, D.: Consistency tradeoffs in modern distributed database system design: CAP is only part of the story. *IEEE Computer* **45**(2), 37–42 (2012)
5. Alonso, G., Roscoe, T., Cock, D., Ewaida, M., Kara, K., Korolija, D., Sidler, D., Wang, Z.: Tackling hardware/software co-design from a database perspective. In: CIDR, Online Proceedings

6. Attia, O.G., Grieve, A., Townsend, K.R., Jones, P.H., Zambreno, J.: Accelerating all-pairs shortest path using a message-passing reconfigurable architecture. In: *ReConFig*, pp. 1–6 (2015)
7. Attia, O.G., Johnson, T., Townsend, K., Jones, P.H., Zambreno, J.: Cygraph: A reconfigurable architecture for parallel breadth-first search. In: *IPDPS*, pp. 228–235 (2014)
8. Ayupov, A., Yesil, S., Ozdal, M.M., Kim, T., Burns, S.M., Ozturk, O.: A template-based design methodology for graph-parallel hardware accelerators. *IEEE Trans. on CAD of Integrated Circuits and Systems* **37**(2), 420–430 (2018)
9. Babb, J., Frank, M.I., Agarwal, A.: Solving graph problems with dynamic computation structures. In: *Other Conferences* (1996)
10. Becher, A., G., L.B., Broneske, D., Drewes, T., Gurusurthy, B., Meyer-Wegener, K., Pionteck, T., Saake, G., Teich, J., Wildermann, S.: Integration of FPGAs in database management systems: challenges and opportunities. *Datenbank-Spektrum* **18**(3), 145–156 (2018)
11. Becher, A., Ziener, D., Meyer-Wegener, K., Teich, J.: A co-design approach for accelerated SQL query processing via FPGA-based data filtering. In: *FPT*, pp. 192–195 (2015)
12. Besta, M., Fischer, M., Ben-Nun, T., de Fine Licht, J., Hoefler, T.: Substream-centric maximum matchings on FPGA. In: *FPGA*, pp. 152–161 (2019)
13. Besta, M., Peter, E., Gerstenberger, R., Fischer, M., Podstawski, M., Barthels, C., Alonso, G., Hoefler, T.: Demystifying graph databases: analysis and taxonomy of data organization, system designs, and graph queries. *CoRR abs/1910.09017* (2019)
14. Besta, M., Stanojevic, D., de Fine Licht, J., Ben-Nun, T., Hoefler, T.: Graph processing on FPGAs: taxonomy, survey, challenges. *CoRR abs/1903.06697* (2019)
15. Betkaoui, B., Thomas, D.B., Luk, W., Przulj, N.: A framework for FPGA acceleration of large graph problems: Graphlet counting case study. In: *FPT*, pp. 1–8 (2011)
16. Betkaoui, B., Wang, Y., Thomas, D.B., Luk, W.: Parallel FPGA-based all pairs shortest paths for sparse networks: a human brain connectome case study. In: *FPL*, pp. 99–104 (2012)
17. Betkaoui, B., Wang, Y., Thomas, D.B., Luk, W.: A reconfigurable computing approach for efficient and scalable parallel graph exploration. In: *ASAP*, pp. 8–15 (2012)
18. Blott, M., Karras, K., Liu, L., Vissers, K.A., Bär, J., István, Z.: Achieving 10Gbps line-rate key-value stores with FPGAs. In: *HotCloud* (2013)
19. Blott, M., Liu, L., Karras, K., Vissers, K.A.: Scaling out to a single-node 80Gbps memcached server with 40 terabytes of memory. In: *HotStorage* (2015)
20. Bondhugula, U., Devulapalli, A., Dinan, J., Fernando, J., Wyckoff, P., Stahlberg, E., Sadayappan, P.: Hardware/software integration for FPGA-based all-pairs shortest-paths. In: *FCCM*, pp. 152–164 (2006)
21. Bondhugula, U., Devulapalli, A., Fernando, J., Wyckoff, P., Sadayappan, P.: Parallel FPGA-based all-pairs shortest-paths in a directed graph. In: *IPDPS* (2006)
22. Bray, T., Paoli, J., Sperberg-McQueen, C.M., Maler, E., Yergeau, F., et al.: Extensible markup language (XML) 1.0 (2000)
23. Bray, T., et al.: The javascript object notation (JSON) data interchange format (2014)
24. Breß, S., Heimes, M., Siegmund, N., Bellatreche, L., Saake, G.: GPU-accelerated database systems: Survey and open challenges. *Trans. Large-Scale Data- and Knowledge-Centered Systems* **15**, 1–35 (2014)
25. Brewer, E.: Cap twelve years later: How the “rules” have changed. *IEEE Computer* **45**(2), 23–29 (2012)
26. Cattell, R.: Scalable SQL and nosql data stores. *SIGMOD Rec.* **39**(4), 12–27 (2010)
27. Chalamalasetti, S.R., Margala, M., Vanderbauwhede, W., Wright, M., Ranganathan, P.: Evaluating FPGA-acceleration for real-time unstructured search. In: *IEEE ISPASS*, pp. 200–209 (2012)
28. Chen, X., Bajaj, R., Chen, Y., He, J., He, B., Wong, W., Chen, D.: On-the-fly parallel data shuffling for graph processing on OpenCL-based FPGAs. In: *FPL*, pp. 67–73 (2019)
29. Ching, A., Edunov, S., Kabiljo, M., Logothetis, D., Muthukrishnan, S.: One trillion edges: graph processing at facebook-scale. *PVLDB* **8**(12), 1804–1815 (2015)
30. Cho, J.M., Choi, K.: An FPGA implementation of high-throughput key-value store using bloom filter. In: *International Symposium on VLSI Design, Automation and Test*, pp. 1–4 (2014)
31. Clark, J., DeRose, S., et al.: XML path language (XPath) (1999)
32. Cooper, B.F., Silberstein, A., Tam, E., Ramakrishnan, R., Sears, R.: Benchmarking cloud serving systems with YCSB. In: *ACM Symposium on Cloud Computing*, pp. 143–154 (2010)
33. Dai, G., Chi, Y., Wang, Y., Yang, H.: FPGP: graph processing framework on FPGA A case study of breadth-first search. In: *FPGA*, pp. 105–110 (2016)
34. Dai, G., Huang, T., Chi, Y., Xu, N., Wang, Y., Yang, H.: ForeGraph: Exploring large-scale graph processing on multi-FPGA architecture. In: *FPGA*, pp. 217–226 (2017)
35. Dai, Z., Ni, N., Zhu, J.: A 1 cycle-per-byte XML parsing accelerator. In: *FPGA*, pp. 199–208 (2010)
36. Dandalis, A., Mei, A., Prasanna, V.K.: Domain specific mapping for solving graph problems on reconfigurable devices. In: *IPPS*, pp. 652–660 (1999)
37. Davoudian, A., Chen, L., Liu, M.: A survey on NoSQL stores. *ACM Comput. Surv.* **51**(2), 40:1–40:43 (2018)
38. DeCandia, G., Hastorun, D., Jampani, M., Kakulapati, G., Lakshman, A., Pilchin, A., Sivasubramanian, S., Vosshall, P., Vogels, W.: Dynamo: Amazon’s highly available key-value store. In: *SOSP*, pp. 205–220 (2007)
39. DeLorimier, M., Kapre, N., Mehta, N., Rizzo, D., Eslick, I., Rubin, R., Uribe, T.E., Jr., T.F.K., DeHon, A.: GraphStep: a system architecture for sparse-graph algorithms. In: *FCCM*, pp. 143–151 (2006)
40. Dennard, R.H., Gaensslen, F.H., Rideout, V.L., Bassous, E., LeBlanc, A.R.: Design of ion-implanted MOSFET’s with very small physical dimensions. *IEEE Journal of Solid-State Circuits* **9**(5), 256–268 (1974)
41. Dennl, C., Ziener, D., Teich, J.: Acceleration of SQL restrictions and aggregations through FPGA-based dynamic partial reconfiguration. In: *FCCM*, pp. 25–28 (2013)
42. El-Hassan, F., Ionescu, D.: A hardware architecture of an XML/XPath broker for content-based publish/subscribe systems. In: *ReConFig*, pp. 138–143 (2010)
43. Engelhardt, N., Hung, C.D., So, H.K.: Performance-driven system generation for distributed vertex-centric graph processing on multi-FPGA systems. In: *FPL*, pp. 215–218 (2018)

44. Engelhardt, N., So, H.K.: GraVF: a vertex-centric distributed graph processing framework on FPGAs. In: FPL, pp. 1–4 (2016)
45. Even, S.: Graph algorithms. Cambridge University Press (2011)
46. Fang, J., Mulder, Y.T.B., Hidders, J., Lee, J., Hofstee, H.P.: In-memory database acceleration on FPGAs: a survey. VLDB J. **29**(1), 33–59 (2020)
47. Finnerty, E., Sherer, Z., Liu, H., Luo, Y.: Dr. BFS: data centric breadth-first search on FPGAs. In: Annual Design Automation Conference, p. 208 (2019)
48. Francis, N., Green, A., Guagliardo, P., Libkin, L., Linddaaker, T., Marsault, V., Plantikow, S., Rydberg, M., Selmer, P., Taylor, A.: Cypher: An evolving query language for property graphs. In: SIGMOD, pp. 1433–1445 (2018)
49. Francisco, P., et al.: The Netezza data appliance architecture: A platform for high performance data warehousing and analytics (2011)
50. Gajendran, S.K.: A survey on NoSQL databases. University of Illinois (2012)
51. Gessert, F., Wingerath, W., Friedrich, S., Ritter, N.: NoSQL database systems: a survey and decision guidance. Comput. Sci. Res. Dev. **32**(3-4), 353–365 (2017)
52. Ghose, S., Li, T., Hajinazar, N., Cali, D.S., Mutlu, O.: Understanding the interactions of workloads and DRAM types: A comprehensive experimental study. CoRR [abs/1902.07609](https://arxiv.org/abs/1902.07609) (2019)
53. Giefers, H., Staar, P.W.J., Polig, R.: Energy-efficient stochastic matrix function estimator for graph analytics on FPGA. In: FPL, pp. 1–9 (2016)
54. Gray, J.: The benchmark handbook for database and transaction systems. Morgan Kaufmann (1993)
55. Gui, C., Zheng, L., He, B., Liu, C., Chen, X., Liao, X., Jin, H.: A survey on graph processing accelerators: challenges and opportunities. J. Comput. Sci. Technol. **34**(2), 339–371 (2019)
56. Hampel, V., Pionteck, T., Maehle, E.: An approach for performance estimation of hybrid systems with FPGAs and GPUs as coprocessors. In: Architecture of Computing Systems, pp. 160–171 (2012)
57. He, C.: Survey on NoSQL database technology. J. of Applied Sci. and Eng. Innovation **2**(2) (2015)
58. Hennessy, J.L., Patterson, D.A.: A new golden age for computer architecture. Commun. ACM **62**(2), 48–60 (2019)
59. Huang, L., Jiang, J., Wang, C., Wang, Y., Pei, Y.: A slide-window-based hardware XML parsing accelerator. In: CCF National Conference on Computer Engineering and Technology, pp. 108–117. Springer (2014)
60. Huelsbergen, L.: A representation for dynamic graphs in reconfigurable hardware and its application to fundamental graph algorithms. In: FPGA, pp. 105–115 (2000)
61. István, Z., Alonso, G., Blott, M., Vissers, K.A.: A flexible hash table design for 10Gbps key-value stores on FPGAs. In: FPL, pp. 1–8 (2013)
62. István, Z., Alonso, G., Singla, A.: Providing multi-tenant services with FPGAs: Case study on a key-value store. In: FPL, pp. 119–124 (2018)
63. István, Z., Sidler, D., Alonso, G.: Runtime parameterizable regular expression operators for databases. In: FCCM, pp. 204–211 (2016)
64. István, Z., Sidler, D., Alonso, G.: Caribou: Intelligent distributed storage. PVLDB **10**(11), 1202–1213 (2017)
65. Jagadeesh, G.R., Srikanthan, T., Lim, C.M.: Field programmable gate array-based acceleration of shortest-path computation. IET Computers & Digital Techniques **5**(4), 231–237 (2011)
66. Jain, A., Lin, C.: Cache replacement policies. Synthesis Lectures on Computer Architecture **14**(1), 1–87 (2019)
67. Jing Han, Haihong E, Guan Le, Jian Du: Survey on NoSQL database. In: International Conference on Pervasive Computing and Applications, pp. 363–366 (2011)
68. Kapre, N.: Custom FPGA-based soft-processors for sparse graph acceleration. In: ASAP, pp. 9–16 (2015)
69. Khoram, S., Zhang, J., Strange, M., Li, J.: Accelerating graph analytics by co-optimizing storage and access on an FPGA-HMC platform. In: FPGA, pp. 239–248 (2018)
70. Kitchenham, B.: Procedures for performing systematic reviews. Keele, UK, Keele University **33**(2004), 1–26 (2004)
71. Lakshman, A., Malik, P.: Cassandra: a decentralized structured storage system. Operating Systems Review **44**(2), 35–40 (2010)
72. Langdale, G., Lemire, D.: Parsing gigabytes of JSON per second. VLDB J. **28**(6), 941–960 (2019)
73. Lavasani, M., Angepat, H., Chiou, D.: An FPGA-based in-line accelerator for Memcached. IEEE Comput. Archit. Lett. **13**(2), 57–60 (2014)
74. Lee, J., Kim, H., Yoo, S., Choi, K., Hofstee, P., Nam, G., Nutter, M., Jamsek, D.A.: ExtraV: Boosting graph processing near storage with a coherent accelerator. PVLDB **10**(12), 1706–1717 (2017)
75. Lei, G., Dou, Y., Li, R., Xia, F.: An FPGA implementation for solving the large single-source-shortest-path problem. IEEE Trans. on Circuits and Systems **63-II**(5), 473–477 (2016)
76. Lei, G., Rong-chun, L., Guo, S.: TorusBFS: a novel message-passing parallel breadth-first search architecture on FPGAs. In: Engineering Science and Technology: An International Journal, vol. 5 (2015)
77. Li, B., Ruan, Z., Xiao, W., Lu, Y., Xiong, Y., Putnam, A., Chen, E., Zhang, L.: KV-Direct: High-performance in-memory key-value store with programmable NIC. In: SOSP, pp. 137–152 (2017)
78. Liang, W., Yin, W., Kang, P., Wang, L.: Memory efficient and high performance key-value store on FPGA using Cuckoo hashing. In: FPL, pp. 1–4 (2016)
79. Lindholm, E., Nickolls, J., Oberman, S.F., Montrym, J.: NVIDIA Tesla: A unified graphics and computing architecture. IEEE Micro **28**(2), 39–55 (2008)
80. Liu, K., Zhao, M., Ju, L., Jia, Z.: Dynamically reconfigurable architecture for high-throughput hash function in key-value store. In: IEEE HPCC, pp. 1964–1970 (2019)
81. Lockwood, J.W., Monga, M.: Implementing ultra low latency data center services with programmable logic. In: HOTI, pp. 68–77 (2015)
82. Lu, J., Holubová, I.: Multi-model databases: A new journey to handle the variety of data. ACM Comput. Surv. **52**(3), 55:1–55:38 (2019)
83. Lunteren, J.V., Engbersen, T., Bostian, J., Carey, B., Larsson, C.: XML accelerator engine. In: International Workshop on High Performance XML Processing. ACM (2004)
84. Ma, X., Zhang, D., Chiou, D.: FPGA-accelerated transactional execution of graph workloads. In: FPGA, pp. 227–236 (2017)
85. Manev, K., Vaishnav, A., Kritikakis, C., Koch, D.: Scalable filtering modules for database acceleration on FPGAs. In: HEART, pp. 4:1–4:6 (2019)
86. Mencer, O., Huang, Z., Huelsbergen, L.: HAGAR: efficient multi-context graph processors. In: FPL, pp. 915–924 (2002)

87. Milovanovic, E.I., Milovanovic, I.Z., Bekakos, M.P., Tselepis, I.N.: Computing all-pairs shortest paths on a linear systolic array and hardware realization on a reprogrammable FPGA platform. *The Journal of Supercomputing* **40**(1), 49–66 (2007)
88. Mitra, A., Vieira, M.R., Bakalov, P., Najjar, W.A., Tsotras, V.J.: Boosting XML filtering with a scalable fpga-based architecture. *CoRR* **abs/0909.1781** (2009)
89. Moussalli, R., Salloum, M., Najjar, W.A., Tsotras, V.J.: Accelerating XML query matching through custom stack generation on FPGAs. In: *HiPEAC*, pp. 141–155 (2010)
90. Moussalli, R., Salloum, M., Najjar, W.A., Tsotras, V.J.: Massively parallel XML twig filtering using dynamic programming on FPGAs. In: *ICDE*, pp. 948–959 (2011)
91. Müller, R., Teubner, J.: FPGA: what's in it for a database? In: *SIGMOD*, pp. 999–1004 (2009)
92. Ni, S., Dou, Y., Zou, D., Li, R., Wang, Q.: Parallel graph traversal for FPGA. *IEICE Electronic Express* **11**(7), 20130987 (2014)
93. Nurvitadhi, E., Weisz, G., Wang, Y., Hurkat, S., Nguyen, M., Hoe, J.C., Martínez, J.F., Guestrin, C.: GraphGen: An FPGA framework for vertex-centric graph computation. In: *FCCM*, pp. 25–28 (2014)
94. Oguntebi, T., Olukotun, K.: GraphOps: a dataflow library for graph analytics acceleration. In: *FPGA*, pp. 111–117 (2016)
95. Owaida, M., Sidler, D., Kara, K., Alonso, G.: Centaur: a framework for hybrid CPU-FPGA databases. In: *FCCM*, pp. 211–218 (2017)
96. Ozdal, M.M., Yesil, S., Kim, T., Ayupov, A., Greth, J., Burns, S.M., Özturk, Ö.: Energy efficient architecture for graph analytics accelerators. In: *ISCA*, pp. 166–177 (2016)
97. Papadimitriou, K., Dollas, A., Hauck, S.: Performance of partial reconfiguration in FPGA systems: A survey and a cost model. *ACM Trans. Reconfigurable Technol. Syst.* **4**(4), 36:1–36:24 (2011)
98. Putnam, A., et al.: A reconfigurable fabric for accelerating large-scale datacenter services. In: *ISCA*, pp. 13–24 (2014)
99. Qiu, Y., Lv, H., Xie, J., Yin, W., Wang, L.: Ultra-low-latency and flexible in-memory key-value store system design on CPU-FPGA. In: *FPT*, pp. 142–149 (2018)
100. Qiu, Y., Xie, J., Lv, H., Yin, W., Luk, W., Wang, L., Yu, B., Chen, H., Ge, X., Liao, Z., Shi, X.: FULL-KV: Flexible and ultra-low-latency in-memory key-value store system design on CPU-FPGA. *IEEE Transactions on Parallel and Distributed Systems* **31**(8), 1828–1444 (2020)
101. Ren, Y., Liao, Z., Shi, X., Xie, J., Qiu, Y., Lv, H., Yin, W., Wang, L., Yu, B., Chen, H., He, X.: A low-latency multi-version key-value store using b-tree on an FPGA-CPU platform. In: *FPL*, pp. 321–325 (2019)
102. Reniers, V., Landuyt, D.V., Rafique, A., Joosen, W.: On the state of nosql benchmarks. In: *International Conference on Performance Engineering*, pp. 107–112 (2017)
103. Ritter, D.: Hardware accelerated application integration: challenges and opportunities. In: *Active Workshop at ACM Middleware*, p. 15 (2017)
104. Ritter, D., Dann, J., May, N., Rinderle-Ma, S.: Hardware accelerated application integration processing: Industry paper. In: *DEBS*, pp. 215–226 (2017)
105. Roozmeh, M., Lavagno, L.: Implementation of a performance optimized database join operation on FPGA-GPU platforms using OpenCL. In: *IEEE NORCHIP*, pp. 1–6 (2017)
106. Shao, Z., Li, R., Hu, D., Liao, X., Jin, H.: Improving performance of graph processing on FPGA-DRAM platform by two-level vertex caching. In: *FPGA* (2019)
107. Shi, X., Zheng, Z., Zhou, Y., Jin, H., He, L., Liu, B., Hua, Q.: Graph processing on GPUs: a survey. *ACM Comput. Surv.* **50**(6), 81:1–81:35 (2018)
108. Sidhu, R.: High throughput, tree automata based XML processing using FPGAs. In: *FPT*, pp. 74–81 (2013)
109. Sidler, D., István, Z., Owaida, M., Kara, K., Alonso, G.: DoppioDB: a hardware accelerated database. In: *SIGMOD*, pp. 1659–1662 (2017)
110. Sklyarov, V., Skliarova, I., Pimentel, B.: Modeling and FPGA-based implementation of graph coloring algorithms. In: *ICARA*, pp. 443–448 (2006)
111. Stonebraker, M.: SQL databases v. nosql databases. *Commun. ACM* **53**(4), 10–11 (2010)
112. Sukhwani, B., Min, H., Thoennes, M., Dube, P., Iyer, B., Brezzo, B., Dillenberger, D., Asaad, S.W.: Database analytics acceleration using FPGAs. In: *PACT*, pp. 411–420 (2012)
113. Sukhwani, B., Roewer, T., Haymes, C.L., Kim, K., McPadden, A.J., Dreps, D.M., Sanner, D., Lunteren, J.V., Asaad, S.W.: Contutto: a novel FPGA-based prototyping platform enabling innovation in the memory subsystem of a server class processor. In: *MICRO*, pp. 15–26 (2017)
114. Takei, Y., Hariyama, M., Kameyama, M.: An SIMD architecture for shortest-path search and its FPGA implementation. In: *PDPTA* (2014)
115. Takei, Y., Hariyama, M., Kameyama, M.: Evaluation of an FPGA-based shortest-path-search accelerator. In: *PDPTA*, p. 613 (2015)
116. Tariq, M.U., Saeed, F.: Parallel sampling-pipeline for indefinite stream of heterogeneous graphs using OpenCL for FPGAs. In: *Big Data*, pp. 4752–4761 (2018)
117. Teubner, J.: Fpgas for data processing: Current state. *Inf. Technol.* **59**(3), 125 (2017)
118. Teubner, J., Woods, L.: *Data Processing on FPGAs. Synthesis Lectures on Data Management.* Morgan & Claypool Publishers (2013)
119. Teubner, J., Woods, L., Nie, C.: Skeleton automata for FPGAs: reconfiguring without reconstructing. In: *SIGMOD*, pp. 229–240 (2012)
120. Teubner, J., Woods, L., Nie, C.: *XLynx* - an FPGA-based XML filter for hybrid XQuery processing. *ACM Trans. Database Syst.* **38**(4), 23:1–23:39 (2013)
121. Tong, D., Zhou, S., Prasanna, V.K.: High-throughput online hash table on FPGA. In: *IPDPS*, pp. 105–112 (2015)
122. Umuroglu, Y., Fraser, N.J., Gambardella, G., Blott, M., Leong, P.H.W., Jahre, M., Vissers, K.A.: FINN: a framework for fast, scalable binarized neural network inference. In: *FPGA*, pp. 65–74 (2017)
123. Umuroglu, Y., Morrison, D., Jahre, M.: Hybrid breadth-first search on a single-chip FPGA-CPU heterogeneous platform. In: *FPL*, pp. 1–8 (2015)
124. Vanderbauwhede, W., Azzopardi, L., Moadeli, M.: FPGA-accelerated information retrieval: high-efficiency document filtering. In: *FPL*, pp. 417–422 (2009)
125. Vanderbauwhede, W., Frolov, A., Chalamalasetti, S.R., Margala, M.: A hybrid CPU-FPGA system for high throughput (10Gb/s) streaming document classification. *SIGARCH Computer Architecture News* **41**(5), 53–58 (2013)
126. Wang, Q., Jiang, W., Xia, Y., Prasanna, V.K.: A message-passing multi-softcore architecture on FPGA for breadth-first search. In: *FPT*, pp. 70–77 (2010)

127. Wang, X., Huang, L., Zhu, Y., Zhou, Y., Peng, H., Xiong, H.: Addressing memory wall problem of graph computation in reconfigurable system. In: HPCCC, pp. 302–307 (2015)
128. Wang, X., Zhu, Y., Huang, L.: A comprehensive reconfigurable computing approach to memory wall problem of large graph computation. *J. Syst. Archit.* **70**, 59–69 (2016)
129. Wang, Y., Hoe, J.C., Nurvitadhi, E.: Processor assisted worklist scheduling for FPGA accelerated graph processing on a shared-memory platform. In: FCCM, pp. 136–144 (2019)
130. Weisz, G., Nurvitadhi, E., Hoe, J.: GraphGen for CoRAM: Graph computation on FPGAs. In: Workshop on the Intersections of Computer Architecture and Reconfigurable Logic (2013)
131. Werner, S., Heinrich, D., Pionteck, T., Groppe, S.: Semi-static operator graphs for accelerated query execution on FPGAs. *Microprocessors and Microsystems - Embedded Hardware Design* **53**, 178–189 (2017)
132. Wieringa, R.J.: *Design Science Methodology for Information Systems and Software Engineering*. Springer (2014)
133. Woods, L., István, Z., Alonso, G.: Ibex - an intelligent storage engine with support for advanced SQL off-loading. *PVLDB* **7**(11), 963–974 (2014)
134. Xie, J., Qiu, Y., Yin, W., Wang, L.: High-throughput and low-latency distributed management proxy for key-value store over 100Gbps Ethernet on FPGA. In: FPT, pp. 224–230 (2019)
135. Xu, C., Wang, C., Gong, L., Jin, L., Li, X., Zhou, X.: Domino: graph processing services on energy-efficient hardware accelerator. In: ICWS, pp. 274–281 (2018)
136. Xu, S., Lee, S., Jun, S.W., Liu, M., Hicks, J., Arvind: BlueCache: a scalable distributed Flash-based key-value store. *PVLDB* **10**(4), 301–312 (2016)
137. Yan, M., Hu, X., Li, S., Akgun, I., Li, H., Ma, X., Deng, L., Ye, X., Zhang, Z., Fan, D., Xie, Y.: Balancing memory accesses for energy-efficient graph analytics accelerators. In: ISLPED, pp. 1–6 (2019)
138. Yang, C.: An efficient dispatcher for large scale graph processing on OpenCL-based FPGAs. *CoRR abs/1806.11509* (2018)
139. Yang, C., Sheng, J., Patel, R., Sanaullah, A., Sachdeva, V., Herbordt, M.C.: OpenCL for HPC with FPGAs: Case study in molecular electrostatics. In: HPEC, pp. 1–8 (2017)
140. Yang, Y., Kuppannagari, S.R., Srivastava, A., Kannan, R., K, P.V.: FASTHash: FPGA-based high throughput parallel hash table (2020)
141. Yao, P.: An efficient graph accelerator with parallel data conflict management. *CoRR abs/1806.00751* (2018)
142. Young, J.S., Romera, J., Hauck, M., Fröning, H.: Optimizing communication for a 2d-partitioned scalable BFS. In: HPEC, pp. 1–7 (2016)
143. Zhang, C., Lu, J., Xu, P., Chen, Y.: Unibench: A benchmark for multi-model database management systems. In: TPC Technology Conference, pp. 7–23 (2018)
144. Zhang, J., Khoram, S., Li, J.: Boosting the performance of FPGA-based graph processor using Hybrid Memory Cube: a case for breadth first search. In: FPGA, pp. 207–216 (2017)
145. Zhang, J., Li, J.: Degree-aware hybrid graph traversal on FPGA-HMC platform. In: FPGA, pp. 229–238 (2018)
146. Zhou, S., Chelmiss, C., Prasanna, V.K.: Accelerating large-scale single-source shortest path on FPGA. In: IPDPS, pp. 129–136 (2015)
147. Zhou, S., Chelmiss, C., Prasanna, V.K.: Optimizing memory performance for FPGA implementation of PageRank. In: ReConfig, pp. 1–6 (2015)
148. Zhou, S., Chelmiss, C., Prasanna, V.K.: High-throughput and energy-efficient graph processing on FPGA. In: FCCM, pp. 103–110 (2016)
149. Zhou, S., Kannan, R., Prasanna, V.K., Seetharaman, G., Wu, Q.: HitGraph: high-throughput graph processing framework on FPGA. *IEEE Trans. Parallel Distrib. Syst.* **30**(10), 2249–2264 (2019)
150. Zhou, S., Kannan, R., Zeng, H., Prasanna, V.K.: An FPGA framework for edge-centric graph processing. In: Proceedings of the 15th ACM International Conference on Computing Frontiers, pp. 69–77 (2018)
151. Zhou, S., Prasanna, V.K.: Accelerating graph analytics on CPU-FPGA heterogeneous platform. In: 29th International Symposium on Computer Architecture and High Performance Computing, pp. 137–144 (2017)
152. Zhu, X., Han, W., Chen, W.: Gridgraph: Large-scale graph processing on a single machine using 2-level hierarchical partitioning. In: USENIX ATC, pp. 375–386 (2015)
153. Ziener, D., Bauer, F., Becher, A., Dennl, C., Meyer-Wegener, K., Schürfeld, U., Teich, J., Vogt, J., Weber, H.: FPGA-based dynamically reconfigurable SQL query processing. *TRETS* **9**(4), 25:1–25:24 (2016)