

# Substream-Centric Maximum Matchings on FPGA

MACIEJ BESTA, Department of Computer Science, ETH Zurich, Switzerland  
 MARC FISCHER, Department of Computer Science, ETH Zurich, Switzerland  
 TAL BEN-NUN, Department of Computer Science, ETH Zurich, Switzerland  
 DIMITRI STANOJEVIC, Department of Computer Science, ETH Zurich, Switzerland  
 JOHANNES DE FINE LICHT, Department of Computer Science, ETH Zurich, Switzerland  
 TORSTEN HOEFLER, Department of Computer Science, ETH Zurich, Switzerland

Developing high-performance and energy-efficient algorithms for maximum matchings is becoming increasingly important in social network analysis, computational sciences, scheduling, and others. In this work, we propose the first maximum matching algorithm designed for FPGAs; it is energy-efficient and has provable guarantees on accuracy, performance, and storage utilization. To achieve this, we forego popular graph processing paradigms, such as vertex-centric programming, that often entail large communication costs. Instead, we propose a *substream-centric* approach, in which the input stream of data is divided into substreams processed independently to enable more parallelism while lowering communication costs. We base our work on the *theory of streaming graph algorithms* and analyze 14 models and 28 algorithms. We use this analysis to provide theoretical underpinning that matches the physical constraints of FPGA platforms. Our algorithm delivers high performance (more than 4× speedup over tuned parallel CPU variants), low memory, high accuracy, and effective usage of FPGA resources. The substream-centric approach could easily be extended to other algorithms to offer low-power and high-performance graph processing on FPGAs.

CCS Concepts: • **Computer systems organization** → **Reconfigurable computing**; • **Hardware** → **Reconfigurable logic and FPGAs**; **Reconfigurable logic applications**; • **Mathematics of computing** → **Matchings and factors**; **Graph algorithms**; *Approximation algorithms*; • **Theory of computation** → *Parallel algorithms*; Streaming models; • **Computing methodologies** → Vector / streaming algorithms.

## ACM Reference Format:

Maciej Besta, Marc Fischer, Tal Ben-Nun, Dimitri Stanojevic, Johannes De Fine Licht, and Torsten Hoefler. 2020. Substream-Centric Maximum Matchings on FPGA. *ACM Trans. Reconfig. Technol. Syst.* 1, 1, Article 1 (January 2020), 32 pages. <https://doi.org/10.1145/3377871>

## 1 INTRODUCTION

Analyzing large graphs has become an important task. Example applications include investigating the structure of Internet links, analyzing relationships in social media, or capturing the behavior of proteins [2, 95]. There are various challenges related to the efficient processing of such graphs. One of the most prominent ones is the size of the graph datasets, reaching trillions of edges [39]. Another one is the fact that processing such graphs can be very power-hungry [5].

---

Authors' addresses: Maciej Besta, Department of Computer Science, ETH Zurich, Switzerland; Marc Fischer, Department of Computer Science, ETH Zurich, Switzerland; Tal Ben-Nun, Department of Computer Science, ETH Zurich, Switzerland; Dimitri Stanojevic, Department of Computer Science, ETH Zurich, Switzerland; Johannes De Fine Licht, Department of Computer Science, ETH Zurich, Switzerland; Torsten Hoefler, Department of Computer Science, ETH Zurich, Switzerland.

---

Publication rights licensed to ACM. ACM acknowledges that this contribution was authored or co-authored by an employee, contractor or affiliate of a national government. As such, the Government retains a nonexclusive, royalty-free right to publish or reproduce this article, or to allow others to do so, for Government purposes only.

© 2020 Copyright held by the owner/author(s). Publication rights licensed to ACM.

1936-7406/2020/1-ART1 \$15.00  
<https://doi.org/10.1145/3377871>

Deriving and approximating *maximum matchings* (MM) [35] are important graph problems. A matching in a graph is a set of edges that have no common vertices. Maximum matchings are used in computational sciences, image processing, VLSI design, or scheduling [35, 130]. For example, a matching of the carbon skeleton of an aromatic compound can be used to show the locations of double bonds in the chemical structure [130]. As deriving the exact MM is usually computationally expensive, significant focus has been placed on developing fast approximate solutions [46].

To enable high-performance graph processing, various schemes were proposed, such as vertex-centric approaches [58], streaming [115], and others [124]. They are easily deployable in combination with the existing processing infrastructure such as Spark [139]. However, they were shown to be often inefficient [102] and they are not explicitly optimized for power-efficiency.

To enable power-efficient graph processing, several graph algorithms and paradigms for FPGAs were proposed [32, 47, 48, 57, 78, 105, 107, 135, 141, 145, 147, 149]. Yet, *none targets maximum matchings*. In addition, the established paradigms for designing graph algorithms that were ported to FPGAs, for example the vertex-centric paradigm, *are not straightforwardly applicable to the MM problem* [116].

In this work, we propose *the first design and implementation of approximating maximum matchings on FPGAs*. Our design is power-efficient *and* high-performance. For this, we forego the established vertex-centric paradigm that may result in complex MM codes [116]. Instead, basing on *streaming theory* [60], we propose a *substream-centric* FPGA design for deriving MM. In this approach, we ① divide the incoming stream of edges into *substreams*, ② process each substream independently, and ③ merge these results to form the final algorithm outcome.

For highest power-efficiency, phases ①–② run on the FPGA; both phases work in the streaming fashion and offer much parallelism, and we identify the FPGA as the best environment for these phases. Conversely, the final gathering phase, that usually takes  $< 1\%$  of the processing time as well as consumed power and exhibits little parallelism, is conducted on the CPU for more performance.

To provide formal underpinning of our design and thus enable guarantees of correctness, memory usage, or performance, we base our work on the family of *streaming models* that were developed to tackle large graph sizes. A special case is the *semi-streaming model* [60], created specifically for graph processing. It assumes that the input is a sequence of edges (pairs of vertices), which can be accessed only sequentially in one direction, as a stream. The main memory (can be randomly accessed) is assumed to be of size  $O(n \text{ polylog}(n))$ <sup>1</sup> ( $n$  is the number of vertices in the graph). Usually, only one pass over the input stream is allowed, but some algorithms assume a small (usually constant or logarithmic) number of passes. We investigate *a total of 14 streaming models* and *a total of 28 MM algorithms* created in these models, and use the insights from this investigation to develop our MM FPGA algorithm, ensuring both empirical speedups and provable guarantees on runtime, used memory, and correctness.

Towards these goals, we contribute:

- the first design and implementation of the maximum matching algorithm on FPGAs,
- an in-depth analysis of the potential of using streaming theory (14 models and 28 algorithms) for accelerating graph processing on FPGAs,
- a substream-centric paradigm that combines the advantages of semi-streaming theory and FPGA capabilities,
- detailed performance analysis demonstrating significant speedups over state-of-the-art baselines on both CPUs and FPGAs.

This paper is an extended version of our work published in the FPGA'19 proceedings [19].

<sup>1</sup> $O(\text{polylog}(n)) = O(\log^c(n))$  for some constant  $c \in \mathbb{N}$

## 2 BACKGROUND AND NOTATION

We first present the necessary concepts.

### 2.1 Graph-Related Concepts

**2.1.1 Graph Model.** We model an undirected graph  $G$  as a tuple  $(V, E)$ ;  $V = \{v_1, \dots, v_n\}$  is a set of vertices and  $E \subseteq V \times V$  is a set of edges;  $|V| = n$  and  $|E| = m$ . Vertex labels are  $\{1, 2, \dots, n\}$ . If  $G$  is weighted, it is modeled by a tuple  $(V, E, w)$ ;  $w(e)$  or  $w(u, v)$  denote the weight of an edge  $e = (u, v) \in E$ . The maximum and minimum edge weight in  $G$  are denoted with  $w_{max}$  and  $w_{min}$ .  $G$ 's adjacency matrix is denoted by  $A$ .

**2.1.2 Compressed Sparse Row (CSR).** In the well-known CSR format,  $A$  is represented with three arrays: *val*, *col*, and *row*. *val* contains all  $A$ 's non-zeros (that correspond to  $G$ 's edges) in the row major order. *col* contains the column index for each corresponding value in *val*. Finally, *row* contains starting indices in *val* (and *col*) of the beginning of each row in  $A$ .  $D$  denotes the diameter of  $G$ . CSR is widely adopted for its simplicity and low memory footprint for sparse matrices.

**2.1.3 Graph Matching.** A *matching*  $M \subseteq E$  in a graph  $G$  is a set of edges that share no vertices.  $M$  is called *maximal* if it is no longer a matching once any edge not in  $M$  is added to it.  $M$  is *maximum* if there is no matching with more edges in it. Maximum matchings (MM) in unweighted graphs are called *maximum cardinality* matchings (MCM). Maximum matchings in weighted graphs are called *maximum weighted* matchings (MWM). Example matchings are illustrated in Figure 1.

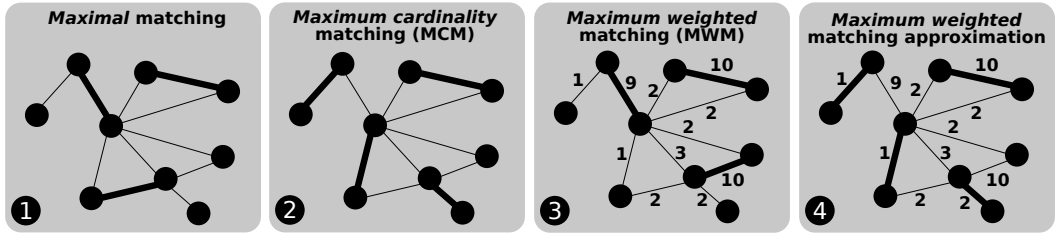


Fig. 1. **Example matchings.** Edges in matchings are represented by bold lines, edge weights are represented with numbers.

**2.1.4 Maximum Weighted Matching.** Given a weighted graph  $G = (V, E, w)$ , a maximum weighted matching is a matching  $M^*$ , such that its weight  $w(M^*) = \sum_{e \in M^*} w(e)$  is maximized. An algorithm provides an  $\alpha$ -approximation of  $M^*$ , if – for any derived weighted matching  $M$  – it holds that  $w(M^*)/w(M) \leq \alpha$  (therefore  $\alpha \geq 1$ ). Note that the approximation ratio of the MCM is defined inversely compared to the MWM: We say that an algorithm, that returns a matching  $M_C$ , provides an  $\alpha$ -approximation to the maximum cardinality  $M_C^*$  if it holds that  $|M_C|/|M_C^*| \geq \alpha$  (therefore  $\alpha \leq 1$ ). We do this to conform to the general approximation notation of maximum cardinality matching [76, 81, 86] and maximum weighted matching [46, 62, 66].

### 2.2 Architecture-Related Concepts

**2.2.1 FPGAs.** FPGAs aim to combine the advantages of Application Specific Integrated Circuits (ASICs) and CPUs: they offer ASIC's high performance and low power usage, and they can be reconfigured to compute arbitrary functions, similarly to CPUs. Usually, the FPGA clock frequency is  $\approx 10$ –600MHz, dependent on the algorithm and the FPGA platform. This is an order of magnitude less compared to high-end CPUs (up to 4.7GHz [73]) and below GPUs (up to 1.5GHz [106]). Yet,

due to the custom design deployed directly in hardware, multiple advantages such as low power consumption arise [50].

**2.2.2 FPGA Components and Fundamental Building Blocks.** Xilinx uses vanilla look-up tables (LUTs) while Intel employs Adaptive Logic Modules as fundamental building blocks [121, 131]. Their micro-architecture is different but fundamental functionality is similar. Xilinx combines multiple LUTs and associated registers into CLBs [121], while Intel combines ALMs into LABs [131]. Next, Block Random Access Memory (BRAM) allows to store small amounts of data (up to 20 kbits per BRAM [71]) and provides fast data access, acting similarly to a CPU cache. Today, thousands of BRAM units are distributed over a single FPGA.

**2.2.3 FPGA+CPU.** Hybrid computation systems consist of a host CPU and an attached FPGA. First ❶, an FPGA can be added to the system as an accelerator; the host main memory is separated from the FPGA private DRAM memory and data must be transferred over PCIe. Often, the FPGA is configured as a PCIe endpoint with a direct memory access (DMA) controller, allowing to move data between the host and the FPGA without the need of CPU resources. PCIe is high-bandwidth oriented, but exhibits high overhead and latency for small packets [41]. This drawback is overcome by storing often accessed data in the private DRAM using the memory controller, or storing the data on chip in the FPGA’s BRAM. Second ❷, the CPU and the FPGA can be directly linked by an interconnect, such as Intel’s QuickPath Interconnect (QPI), providing a coherent view to a single shared main memory. Examples of these systems include Intel HARP [108] and the Xilinx Extensible Processing Platform [117]. The direct main memory access allows to share data without the need to copy it to the FPGA. To prevent direct physical main memory accesses, HARP provides a translation layer, allowing the FPGA to operate on virtual addresses. It is implemented in both hardware as a System Protocol Layer (SPL) and in software, for example as a part of the Centaur framework [109]. Moreover, a cache is available to reduce access time. According to Choi et al. [41], systems with direct interconnect exhibit lower latency and higher throughput than PCIe connected FPGAs. *In our substream-centric FPGA design for deriving MM, we use a hybrid CPU+FPGA system to take advantage of both the CPU and the FPGA in the context of graph processing.*

### 3 FROM SEMI-STREAMING TO FPGAS

We first summarize the analysis into the theory of streaming models and algorithms. We conducted the analysis to provide formal underpinning of our work and thus **ensure provable properties, for example correctness, approximation, or performance**. Towards this goal, *we analyzed 14 different models of streaming* (simple streaming [68], semi-streaming [60], insert-only [60], dynamic [8], vertex-arrival [44], adjacency-list [101], cash-register [103], Turnstile [103], sliding window [49], annotated streaming [37], StreamSort [3], W-Stream [52], online [80], and MapReduce [51]) and *28 different MM algorithms*. Moreover, to understand whether streaming itself is the best option for developing MWM on FPGAs, we extensively analyzed existing graph processing works on FPGAs [47, 48, 57, 84, 89, 90, 96, 107, 110, 137, 138, 142, 143, 147–149]. The outcome of this analysis is as follows: the best candidates for adoption in the FPGA setting are **semi-streaming** graph algorithms that **expose parallelism by decomposing the incoming stream of edges for independent processing**, for example the MWM algorithm by Crouch and Stubbs [46].

#### 3.1 Why Streaming (aka Edge-Centric)?

Before analyzing different streaming models and algorithms, we first investigate whether streaming itself is the best paradigm for developing the MWM algorithm for FPGAs. In the streaming paradigm, also known as the *edge-centric* paradigm, the “first class citizen” is an edge. A sequence of edges is streamed in and out from the memory to the FPGA, and a specified operation is performed on each

edge and possibly some associated edge properties. This way of accessing the graph has major advantages because of its sequential memory access pattern, which improves spatial locality [30].

However, there also exist other paradigms for processing graphs, most importantly the established vertex-centric paradigm [97], where the “first class citizen” is a vertex. Here, one programs an algorithm by developing a (usually small) routine that is executed *for each vertex in the graph concurrently*. In this routine, one usually has access to the neighbors of a given vertex. Such an approach can lead to many random memory accesses as neighboring vertices may be stored in different regions of the memory. Still, it is often used because many important algorithms such as BFS or PageRank can easily be implemented in this model [30].

To identify the best graph processing paradigm for implementing the MWM algorithm in FPGA, we first analyze the existing FPGA graph processing implementations, focusing on the used paradigm. Table 1 illustrates the most relevant designs. We group separately generic graph processing frameworks and specific algorithm implementations. Each group is sorted chronologically. Selected columns in this table constitute criteria used to categorize the surveyed FPGA works (the full results of this analysis are in a separate extended survey [30]).

The first such criterion is **generality**, i.e., whether a given FPGA scheme is focused on a particular graph problem or whether it constitutes a generic framework that facilitates implementing different graph algorithms. Another criterion is a used **graph programming paradigm**. We also distinguish between works that target a single FPGA and ones that scale to **multiple FPGAs**. Finally, we consider the used **programming language** and the **storage location** of the *whole* processed graph datasets. In the latter, “DRAM” indicates that the input dataset is located in DRAM and it is streamed in and out of the FPGA during processing (i.e., only a part of the input dataset is stored in BRAM at a time). Contrarily, “BRAM” indicates that the whole dataset is assumed to be located in BRAM. To investigate the scalability of the analyzed solutions, we provide sizes  $(n, m)$  of the largest processed graphs.

The analysis indicates that the *streaming (edge-centric) paradigm and its variants have so far been the most successful in processing large graphs*. The only vertex-centric design that processed a graph with  $m > 1\text{B}$  required multiple FPGAs [34]. Contrarily, two recent edge-centric designs based on single FPGAs were able to conduct computations on such graphs [89, 148].

Moreover, although the vertex-centric paradigm facilitates developing simple algorithms such as BFS or PageRank, it is often difficult to use for solving more complex graph problems. For example, as Salihoglu and Widom state [116], “(...) *implementing graph algorithms efficiently on Pregel-like systems (...) can be surprisingly difficult and require careful optimizations*.” For example, when describing graph problems as fundamental as deriving Strongly Connected Components (SCCs), Bi-Connected Components (BCCs), and Minimum Spanning Forest (MSF), Salihoglu and Widom [116] observe that “(...) *implementing even the basic versions of SCC and MSF is quite challenging, taking more than 700 lines of code*.” while Yan et al. [136] state that “*It is challenging to design Pregel algorithms for problems such as BCCs and SCCs*.” The problem is not only related to Pregel. Yan et al. [136] make similar observations about the established GraphLab [94] and PowerGraph [65] vertex-centric frameworks, stating that they do “*not support algorithms in which a vertex needs to communicate with a non-neighbor, such as the [Shiloach-Vishkin] algorithm [for Connected Components], the list ranking algorithm, and the [Case Checking] algorithm*.” They make similar observations for BCCs and SCCs. Thus, when developing vertex-centric graph algorithms, one resorts to algorithms that fit the vertex-centric paradigm well. An example is Label Propagation for Connected Components [136]. Yet, this algorithm takes  $O(D)$  time in the PRAM model [67] while the Shiloach-Vishkin algorithm [122], hard to express in the vertex-centric paradigm [136], uses only  $O(\log n)$  time in PRAM. Similar observations are made for other graph problems [83]. This

indicates that *it is difficult to design efficient vertex-centric formulations of graph algorithms that require accessing more than neighbors of each vertex.*

**Consequently, in our work, we use edge streaming for developing the MWM algorithm for FPGAs.** This is because (1) it has been shown to scale to large graphs and (2) it straightforwardly enables pipelining of edges, thus facilitating the utilization of FPGA hardware resources. Finally (3), the existing rich theory of streaming graph processing for complex graph problems such as matchings, random walks, and others [3, 8, 37, 44, 49, 51, 52, 60, 60, 80, 101, 103, 103] indicates that it is easier to develop fast MWM schemes with edge streaming than with the vertex-centric paradigm.

### 3.2 Why Semi-Streaming?

The *semi-streaming model* [60] was created specifically for graph processing. However, there are numerous other streaming models that are also used for developing graph algorithms, namely simple streaming [68], insert-only [60], dynamic [8], vertex-arrival [44], adjacency-list [101], cash-register [103], Turnstile [103], sliding window [49], annotated streaming [37], StreamSort [3], W-Stream [52], online [80], and MapReduce [51]. A detailed comparison of these models and analysis of their relationships is outside the scope of this paper<sup>2</sup>. Here, we briefly justify why we selected semi-streaming as the basis for our MWM algorithm for FPGAs. First, semi-streaming enables a generic streaming setting in which one can arbitrarily process the incoming stream of edges. This will enable our substream-centric approach where the incoming edges are divided into independent substreams. Second, there exists a rich body of algorithms and techniques developed for the semi-streaming setting. Finally, most importantly, in semi-streaming one assumes that processing the incoming stream of edges can utilize at most  $O(n \text{ polylog}(n))$  random memory. **Thus, algorithms under this model may address the limited FPGA BRAM capacity better than algorithms in models with weaker memory-related constraints.**

### 3.3 Which Semi-Streaming MM Algorithm?

Table 2 compares the considered semi-streaming and related MM algorithms. We identify those with properties suggesting an effective and versatile FPGA design: low space consumption, one pass, and applicability to general graphs. Finally, virtually all designed algorithms are approximate. Yet, as we show later (§ 5), in practice they deliver near-accurate results.

We conjecture that the majority of the considered MM algorithms deliver limited performance on FPGA because *their design is strictly sequential*: every edge in the incoming stream can only be processed after processing the previous edge in the stream is completed. However, we identify some algorithms that introduce a certain amount of parallelism. Here, we focus on the algorithm by Crouch and Stubbs [46], used as a basis for our FPGA design (last row of Table 2). We first outline this algorithm and then justify our selection. We also discuss other considered MWM algorithms.

**3.3.1 Algorithm Intuition.** The MWM algorithm by Crouch and Stubbs [46] delivers a  $(4 + \epsilon)$ -approximation of MWM. It consists of two parts. In Part 1, one selects  $L$  subsets of the incoming (streamed) edges and computes a *maximum cardinality* matching for each such subset. In Part 2, the derived maximum matchings are combined into the final *maximum weighted matching*. The approach is visualized in Figure 2.

**3.3.2 Algorithm Details.** The algorithm of Crouch and Stubbs [46] provides a  $(4 + \epsilon)$ -approximation to the MWM problem assuming an arbitrarily ordered stream of incoming edges with possible graph

<sup>2</sup>This comparison will be provided in a separate survey on streaming graph processing; this survey will be released upon the publication of this work.

Reference (scheme name)	Venue	Generic Design <sup>1</sup>	Considered Problems <sup>2</sup>	Programming Paradigm <sup>4</sup>	Used Language	Multi FPGAs <sup>4</sup>	Input Location <sup>5</sup>	$n^\dagger$	$m^\dagger$
Kapre [79] (GraphStep)	FCCM'06	👍	spreading activation [93]	BSP (similar to vertex-centric)	unsp.	👍	BRAM	220k	550k
Weisz [135] (GraphGen)	FCCM'14	👍	TRW-S, CNN [127]	Vertex-Centric	unsp.	👎	DRAM	110k	221k
Kapre [78] (GraphSoC)	ASAP'15	👍	SpMV	Vertex-Centric	C++ (HLS)	👍	BRAM	17k	126k
Dai [47] (FPGP)	FPGA'16	👍	BFS	Edge-Centric*	unsp.	👍	DRAM	41.6M	1.4B
Oguntebi [107] (GraphOps)	FPGA'16	👍	BFS, SpMV, PR, Vertex Cover	Edge-Centric*	MaxJ (HLS)	👎	BRAM	16M	128M
Zhou [147] (GraVF)	FCCM'16	👍	SSSP, WCC, MST	Edge-Centric	unsp.	👎	DRAM	4.7M	65.8M
Engelhardt [57] (ForeGraph)	FPL'16	👍	BFS, PR, SSSP, CC	Vertex-Centric	Migen (HLS)	👎	BRAM	128k	512k
Dai [48] (ForeGraph)	FPGA'17	👍	PR, BFS, WCC	Edge-Centric*	unsp.	👍	DRAM	41.6M	1.4B
Zhou [149] (ExtraV)	SBAC-PAD'17	👍	BFS, SSSP	Hybrid (Vertex +Edge-Centric)	unsp.	👎	DRAM	10M	160M
Lee [89] (ExtraV)	FPGA'17	👍	BFS, PR, CC	Edge-Centric*	C++ (HLS)	👎	DRAM	124M	1.8B
Zhou [148] report (2018)	CF'18	👍	SpMV, PR	Edge-Centric	unsp.	👎	DRAM	41.6M	1.4B
Yang [137] report (2018)	report (2018)	👍	BFS, PR, WCC	Edge-Centric*	OpenCL	👎	DRAM	4.85M	69M
Yao [138]	report (2018)	👍	BFS, PR, WCC	Vertex-Centric	unsp.	👎	BRAM	4.85M	69M
Betkaoui [32]	FTP'11	👎	GC	Vertex-Centric	Verilog	👍	DRAM	300k	3M
Betkaoui [33]	FPL'12	👎	APSP	Vertex-Centric	Verilog	👍	DRAM	38k	72M
Betkaoui [34]	ASAP'12	👎	BFS	Vertex-Centric	Verilog	👍	DRAM	16.8M	1.1B
Attia [13] (CyGraph)	IPDPS'14	👎	BFS	Vertex-Centric	VHDL	👍	DRAM	8.4M	536M
Zhou [146]	ReConFig'15	👎	PR	Edge-Centric	unsp.	👎	DRAM	2.4M	5M
Besta [19]	FPGA'19	👎	MWM	Substream-Centric	Verilog	👎	DRAM	4.8M	117M

Table 1. Summary of the features of selected works sorted by publication date. <sup>1</sup>**Generic Design**: this criterion indicates whether a given scheme provides a graph processing framework that supports more than one graph algorithm (👍) or whether it focuses on concrete graph algorithm(s) (👎). <sup>2</sup>**Considered Problems**: this column lists graph problems (or algorithms) that are explicitly considered in a given work; they are all described in detail in an extended survey [30] (BFS: Breadth-First Search [43], SSSP: Single-Source Shortest Paths [43], APSP: All-Pairs Shortest Paths [43], PR: PageRank [111], CC: Connected Components [43], WCC: Weakly Connected Components [43], MST: Minimum Spanning Tree [43], SpMV: Sparse Matrix and Dense Vector product, TC: Triangle Counting [118], BC: Betweenness Centrality [104], GC: Graphlet Counting [118], TRW-S: Tree-Reweighted Message Passing [128], CNN: Convolutional Neural Networks [17, 144]). <sup>3</sup>**Used Programming Paradigm**: this column specifies programming paradigms and models used in each work; \*The star indicates that a given scheme uses a paradigm similar to the edge-centric streaming paradigm, for example sharding as used in GraphChi [88], where edges are first preprocessed and divided into shards, with shards being streamed in and out of the main memory. <sup>4</sup>**Multi FPGAs**: this criterion indicates whether a given scheme scales to multiple FPGAs (👍) or not (👎). <sup>5</sup>**Input Location**: this column indicates the location of the *whole* input graph dataset. “DRAM” indicates that it is located in DRAM and it is streamed in and out of the FPGA during processing (i.e., only a part of the input dataset is stored in BRAM at a time). Contrarily, “BRAM” indicates that the whole dataset is assumed to be located in BRAM.  $n^\dagger, m^\dagger$ : these two columns contain the numbers of vertices and edges used in the largest graphs considered in respective works. In any of the columns, “unsp.” indicates that a given value is not specified.

Reference	Approx.	Space	#Passes	Wgh <sup>1</sup>	Gen <sup>2</sup>	Par <sup>3</sup>
[60]	1/2	$O(n)$	1	🚫	👍	👍
[86, Theorem 6]	1/2 + 0.0071	$O(n \text{ polylog}(n))$	2	🚫	👍	👍
[86, Theorem 2]	1/2 + 0.003*	$O(n \text{ polylog}(n))$	1	🚫	👍	👍
[77, Theorem 1.1]	$O(\text{polylog}(n))$	$O(\text{polylog}(n))$	1	🚫	👍	👍
[60, Theorem 1]	2/3 - $\epsilon$	$O(n \log n)$	$O(\log(1/\epsilon)/\epsilon)$	🚫	🚫	👍
[6, Theorem 19]	1 - $\epsilon$	$O(n \text{ polylog}(n)/\epsilon^2)$	$O(\log \log(1/\epsilon)/\epsilon^2)$	🚫	🚫	👍
[86, Theorem 5]	1/2 + 0.019	$O(n \text{ polylog}(n))$	2	🚫	🚫	👍
[86, Theorem 1]	1/2 + 0.005*	$O(n \log n)$	1	🚫	🚫	👍
[86, Theorem 4]	1/2 + 0.0071*	$O(n \text{ polylog}(n))$	2	🚫	🚫	👍
[81]	1 - 1/e	$O(n \text{ polylog}(n))$	1	🚫	🚫	👍
[64, Theorem 20]	1 - 1/e	$O(n)$	1	🚫	🚫	👍
[76, Theorem 2]	$1 - \frac{e^{-k} k^{k-1}}{(k-1)!}$	$O(n)$	$k$	🚫	🚫	👍
[40]	1	$\tilde{O}(k^2)$	1	🚫	👍	👍
[40]	1/ $\epsilon$	$\tilde{O}(n^2/\epsilon^3)$	1	🚫	👍	👍
[12, Theorem 1]	1/ $n^\epsilon$	$\tilde{O}(n^{2-3\epsilon} + n^{1-\epsilon})$	1	🚫	🚫	👍
[60, Theorem 2]	6	$O(n \log n)$	1	👍	👍	🤔
[98, Theorem 3]	2 + $\epsilon$	$O(n \text{ polylog}(n))$	$O(1)$	👍	👍	🤔
[98, Theorem 3]	5.82	$O(n \text{ polylog}(n))$	1	👍	👍	🤔
[140]	5.58	$O(n \text{ polylog}(n))$	1	👍	👍	🤔
[59]	4.911 + $\epsilon$	$O(n \text{ polylog}(n))$	1	👍	👍	🤔
[66]	3.5 + $\epsilon$	$O(n \text{ polylog}(n))$	1	👍	👍	🤔
[113]	2 + $\epsilon$	$O(n \log^2 n)$	1	👍	👍	🤔
[62]	2 + $\epsilon$	$O(n \log n)$	1	👍	👍	🤔
[60, Section 3.2]	2 + $\epsilon$	$O(n \log n)$	$O(\log_{1+\epsilon/3} n)$	👍	👍	🤔
[6, Theorem 28]	$\frac{1}{1-\epsilon}$	$O(n \log(n)/\epsilon^4)$	$O(\epsilon^{-4} \log n)$	👍	👍	👍
[6, Theorem 22]	$\frac{1}{\frac{2}{3}(1-\epsilon)}$	$O(n(\frac{\epsilon \log n - \log \epsilon}{\epsilon^2}))$	$O(\epsilon^{-2} \log(\epsilon^{-1}))$	👍	👍	👍
[6, Theorem 22]	$\frac{1}{1-\epsilon}$	$O(n(\frac{\epsilon \log n - \log \epsilon}{\epsilon^2}))$	$O(\epsilon^{-2} \log(\epsilon^{-1}))$	👍	🚫	👍
[46]	4 + $\epsilon$	$O(n \text{ polylog}(n))$	1	👍	👍	👍

Table 2. (§ 3) **Comparison of algorithms for maximum matching.** \*Approximation in expectation, <sup>1</sup>Wgh: accepted weighted graphs, <sup>2</sup>Gen: accepted general (non-bipartite) graphs, <sup>3</sup>Par: Potential for parallelization;  $k$  is the size of a given maximum matching. 👍: A given feature is offered. 🚫: A given feature is not offered. In the context of parallelization: 👍: a given algorithm is based on a method that is easily parallelizable (e.g., sampling), 🚫: a given algorithm uses a method that may be complex to parallelize (e.g., augmenting paths), 🤔: it is unclear how to parallelize a given algorithm (e.g., it is based on a greedy approach).



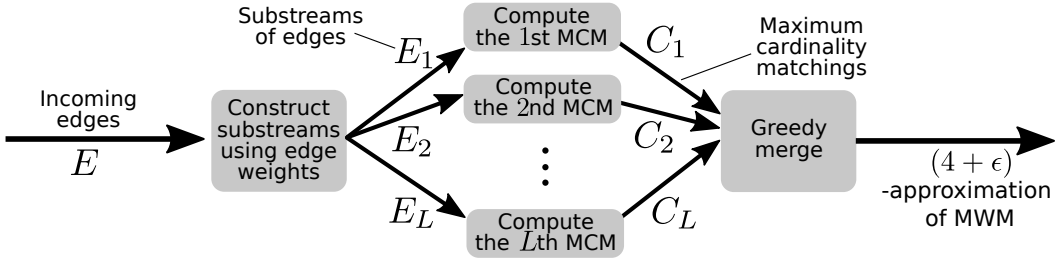


Fig. 2. **The design of the MWM algorithm** of Crouch and Stubbs [46].

updates (edge insertions). The basic idea is to reduce the MWM problem to  $L \equiv O(\text{polylog}(n))$  instances of the MCM problem. Given the input stream of incoming edges  $E$ ,  $O\left(\frac{1}{\epsilon} \log n\right)$  many substreams are generated. Each substream  $E_i$  is created by filtering the edges according to their weight. Specifically, we have  $E_i = \{e \in E \mid w(e) \geq (1 + \epsilon)^i\}$ . Since an edge that belongs to substream  $i + 1$  also belongs to substream  $i$ , it holds that  $E_{i+1} \subseteq E_i$ . Next, for each substream, an MCM  $C_i$  is constructed. The final  $(4 + \epsilon)$ -approximation to MWM is greedily constructed by considering the edges of every  $C_i$ , in the descending order of  $i$ .

We select this algorithm as the basis of our substream-centric FPGA design as it (1) can be straightforwardly parallelized, (2) ensures only  $O(n \text{ polylog} n)$  memory footprint as it belongs to the semi-streaming model, (3) targets general weighted graphs, (4) its structure matches well the design of a hybrid FPGA+CPU system: while substreams can be processed in parallel on the FPGA, the greedy sequential merging of substreams into the final MWM can be done on the CPU, and (5) it needs only one pass over the streamed dataset of size  $O(m + n)$ , limiting expensive data transfers between the FPGA and DRAM.

### 3.4 How To Adapt Semi-Streaming to FPGAs?

In § 4, we describe the FPGA adaptation, design, and implementation of the selected semi-streaming MM algorithm. We stream the edges as stored in the CSR representation. Our substream-centric design implements a staged pipeline with throughput of up to one edge per cycle.

### 3.5 Other Considered Matching Algorithms

We could not find other algorithms that would clearly satisfy all the above five criteria (stated in § 3.3) simultaneously. However, we do not conclude that other algorithms are unsuitable for an efficient FPGA implementation, and leave developing such designs as future work. Here, for completeness, we describe other streaming matching algorithms considered in our analysis. Details of these algorithms (models, techniques, random vs. deterministic design) can be found in Table 3.

**3.5.1 Generic Techniques for Deriving Matchings.** First, we shortly describe generic techniques that we identified while analyzing matching algorithms. We identified these techniques to investigate which algorithms are easily parallelizable (when a given technique is easily parallelizable, this implies that algorithms relying on this technique may also be easily parallelizable).

**Sampling and Unbiased Estimators** Sampling [54] in general is used to estimate some quantity, for example the number of triangles in a graph. One first samples edges that are then used to generate an unbiased estimator of the graph property in question. After that, the challenging task is to show that the space constraint of  $O(n \text{ polylog}(n))$  is not exceeded, and that the estimator succeeds to give an estimation within a small error range with some probability. Often, the error

range and success probability can be controlled and are assumed to be constants. This method is abbreviated with a simple *General sampling*.

**Sketching and  $L_0$  Sampling** Graph sketching is a technique that reduces dimensionality of given data while preserving certain properties. It is commonly used in dynamic graph streams because in such streams a sampled edge might be removed afterwards, and thus a simple sampling scheme as described above cannot be used. Specifically, sketching based on  $L_0$  sampling proved to be useful in different cases [7, 8] when the dynamic streaming model with inserts and deletes is assumed.

**Simulation of Local Distributed Algorithms** An approach used in several cases is based on porting a local distributed algorithm solving a given graph problem to the streaming model to simulate its behaviour and solve the same problem in the streaming setting [77].

**Augmenting Paths** An  $M$ -augmenting path [112] is a path that starts and ends with vertices that are not adjacent to edges that belong to  $M$ , and all other vertices in this path are adjacent to such edges. Moreover, every second edge in the path belongs to  $M$ . It is easy to see that by removing all the edges of such a path from  $M$ , and by adding those in this path that were not in  $M$ , one increases the size of  $M$  by 1. This technique is used to improve the size of maximal matchings.

**Greedy Approach** A traditional greedy approach [42] is also used in streaming settings.

**Linear Programming (LP)** Some problems can be reduced to a linear program [112] which also provides a solution to the initial problem.

**Local Ratio** The local ratio technique [15] is an approach for solving combinatorial optimization problems and is related to linear programming.

*3.5.2 Maximum Cardinality Matching.* We start with algorithms for MCM algorithms. A simple *maximal* matching can be obtained by iterating over edges in some arbitrary order and greedily adding an edge only if both its endpoints are not used yet; this scheme requires  $O(n)$  space [60]. Since every maximal matching is a  $1/2$ -approximation to a maximum matching, this scheme leads to a  $1/2$ -approximation for the maximum cardinality matching. For a long time, this approach was used to derive the best approximation of a maximum matching, using only one pass.

Konrad et al. [86] present a variety of algorithms, taking either one or two passes over the input stream of edges. The general idea is to simulate a three pass algorithm. The original three pass algorithm relies on the refinement of a maximum matching in a bipartite graph using  $M$ -augmenting paths, as already used by Feigenbaum et al. [60, Theorem 1].

Kapralov et al. [77] simulate a local distributed algorithm in the semi-streaming model. The local algorithm is able to distinguish graphs with an  $\Omega(n)$  size matching from the graphs having no  $n/\text{polylog}(n)$  size matching. This approach is transformed into a one-pass semi-streaming algorithm, requiring only  $O(\text{polylog}(n))$  space.  $O(\log(n))$  many instances of this algorithm are executed in parallel, resulting in an  $O(\text{polylog}(n))$ -approximation.

For bipartite graphs,  $(1 - 1/e)$ -approximations are possible in both the online [81] and vertex model [64]. Both assume that both vertex classes belonging to sets  $U, W$  of the bipartite graph  $G = (U, W, E)$  have the same size, and that one set of the vertices is known beforehand (say  $U$ ). The other set (say  $W$ ) is then streamed in and at the same time the edges are revealed to the other set. However, the two approaches differ by the fact that one is online, so must make a decision as soon as the edges arrive, and the other can defer the decision to a later point in time. Additionally, the algorithm of Goel et al. [64] is deterministic. A refinement of Kapralov's scheme [76] allows multiple passes on the input and also achieves a  $(1 - 1/e)$ -approximation for one pass. Differently, the  $|U| = |W|$  constraint is not mentioned.

In the dynamic graph stream model, Chitnis et al. [40] present an exact approximation using  $\tilde{O}(k^2)$  space (where  $k$  is the size of the matching  $|S^*|$ ), requiring only one pass. The approach is

Reference	Model	Determinism	Technique	Matching type
[60]	Insertion-Only	Deterministic	-	Cardinality
[86, Theorem 6]	Insertion-Only	Deterministic	Augmenting Paths	Cardinality
[86, Theorem 2]	Insertion-Only	Deterministic	Augmenting Paths	Cardinality
[77, Theorem 1.1]	Insertion-Only	Deterministic	Local Algorithm	Cardinality
[60, Theorem 1]	Insertion-Only	Deterministic	Augmenting Paths	Cardinality
[6, Theorem 19]	Insertion-Only	Deterministic	LP	Cardinality
[86, Theorem 5]	Insertion-Only	Deterministic	Augmenting Paths	Cardinality
[86, Theorem 1]	Insertion-Only	Deterministic	Augmenting Paths	Cardinality
[86, Theorem 4]	Insertion-Only	Randomized	Augmenting Paths	Cardinality
[81]	Online	Randomized	-	Cardinality
[64, Theorem 20]	Vertex-Arrival	Deterministic	-	Cardinality
[76, Theorem 2]	Vertex-Arrival	Deterministic	-	Cardinality
[40]	Dynamic Graph Stream	Randomized	General Sampling	Cardinality
[40]	Dynamic Graph Stream	Randomized	General Sampling	Cardinality
[12, Theorem 1]	Dynamic Graph Stream	Randomized	$L_0$ Sampling	Cardinality
[60, Theorem 2]	Insertion-Only	Deterministic	Greedy	Weighted
[98, Theorem 3]	Insertion-Only	Deterministic	Greedy	Weighted
[98, Theorem 3]	Insertion-Only	Deterministic	Greedy	Weighted
[140]	Insertion-Only	Deterministic	-	Weighted
[59]	Insertion-Only	Deterministic	-	Weighted
[66]	Insertion-Only	Deterministic	-	Weighted
[113]	Insertion-Only	Deterministic	Local Ratio	Weighted
[62]	Insertion-Only	Deterministic	Local Ratio	Weighted
[60, Section 3.2]	Insertion-Only	Deterministic	-	Weighted
[6, Theorem 28]	Insertion-Only	Deterministic	LP	Weighted
[6, Theorem 22]	Insertion-Only	Deterministic	LP	Weighted
[6, Theorem 22]	Insertion-Only	Deterministic	LP	Weighted
[46]	Insertion-Only	Deterministic	-	Weighted

Table 3. (§ 3.5) **Comparison of algorithms for maximum matching.** **Model:** model used to construct a given algorithm, **Determinism:** whether a given algorithm is deterministic or randomized, **Technique:** used general technique (see § 3.5.1);

refined for  $(1/\alpha)$ -approximative matchings, using  $\tilde{O}(n^2/\alpha^3)$  memory. Both algorithms rely on a sampling primitive, which runs in parallel and is also applicable in the MapReduce [51] setting.

A one-pass algorithm in the dynamic graph stream model is presented by Assadi et al. [12, Theorem 1.1]. The algorithm uses a bipartite graph and two approximations of the matching as the input. Note that one can run the algorithm for  $O(\log(n))$  many estimates of the matching, to determine the correct approximation value. The algorithm relies on  $L_0$ -sampling to process the input and succeeds with probability of at least 0.15. Despite the fact that the algorithm runs for

bipartite graphs only, by choosing a random bipartition of the vertices it is possible to run the algorithm for arbitrary graphs, reducing the approximation by a factor of at most two.

Works on maximum matchings in low arboricity graphs also exist [45, 99, 100].

**3.5.3 Maximum Weighted Matching.** Feigenbaum et al. [60] presented the first 6-approximation in 2005. In the same year, the bound was improved to 5.82 [98], which also allows a  $(2 + \epsilon)$ -approximation using a constant number of passes, assuming  $\epsilon$  is small. Note that both of these one-pass algorithms decide at edge arrival, if the edge is kept or not by comparing the weight of the incoming edge to some value that depends on the matching computed so far. Using so called shadow-edges, which may be reinserted into the matching later on [140], the approximation value can be reduced to 5.58. Epstein et al. [59] partition the input edges into  $O(\log(n))$  edge sets. For each set, a separate maximal cardinality matching is computed. Finally, a greedy algorithm is applied to merge the cardinality matchings. This randomized method allows a  $(4.911 + \epsilon)$ -approximation, and can be derandomized by running all possible outcomes of the randomized algorithm in parallel. Note that there is a constant number of parallel executions for a fixed  $\epsilon$ . A similar approach is used to lower the one-pass approximation to  $4 + \epsilon$  [46]: The algorithm reduces the maximum weight matching problem to a polylog number of copies of the maximum cardinality matching problem. At the end, a greedy merge step is applied to get the final result. It is also proven that this specific approach cannot provide a better approximation than  $3.5 + \epsilon$ . This lower bound of  $3.5 + \epsilon$  was achieved two years later [66]. Recently, the  $2 + \epsilon$  approximation ratio was achieved [113] using the local ratio technique [14, 16]. Ghaffari [62] improved the algorithm and reduced the space required from  $O(n \log^2(n))$  to  $O(n \log(n))$ . The proof is done differently using a blaming-charging argument.

Different multi-pass algorithms exist: Feigenbaum et al. [60] noticed that multiple passes allow to emulate an already existing algorithm [132] solving the maximum weighted matching problem. This approach uses only  $O(n \log(n))$  space resulting in an  $(2 + \epsilon)$ -approximation with  $O(\log(n))$  passes. Ahn et al. [6] rely on linear programming: given a graph  $G = (V, E, w)$ , a suitable linear program is defined, which needs to be solved. The Multiplicative Weights Update Meta-Method [11] is used to solve the linear program. Different approaches are presented to lower the amount of space and passes on the input stream.

## 4 MAXIMUM MATCHING ON FPGA

We now describe the design and implementation of the substream-centric MM for FPGAs.

### 4.1 Overview of the Algorithm

We start with a high-level overview of the MWM algorithm. A pseudo code is shown in Listing 1. For each edge, we iterate in the descending order of  $i$  over the  $L$  substreams, identifying them by their respective weights (Line 11). The  $i$ -th substream weight is given by  $(1 + \epsilon)^i$ . For each maximum matching  $C_i$ , we use a bit matrix  $MB$  to track if a vertex has an incident edge to ensure that  $C_i$  remains a matching (i.e., that no two vertices share an edge). Bits included in  $MB$  are called *matching bits*. Bits in  $MB$  associated with a vertex  $u$ , the source vertex of a processed edge, ( $u$ -matching bits) determine if  $u$  has an incident edge included in some matching; they are included in column  $mb_u$  of matrix  $MB$ . Matching bits associated with vertex  $v$ , the destination vertex of a processed edge, ( $v$ -matching bits) track the incident edges of  $v$ ; they are included in column  $mb_v$  of matrix  $MB$ . Since there are  $L$  matchings and  $n$  vertices, the bit matrix  $MB$  is a matrix of size  $L \times n$ . Furthermore, every matching stores its edges in a list. If an edge is added, a flag is set to true to prevent that the edge is added to multiple lists (Line 16). This reduces the runtime of the post-processing part, in which we iterate in the descending order over the  $L$  lists of edges to generate the  $(4 + \epsilon)$ -approximation to the maximum weighted matching.

```

1 //Input:  $\epsilon$ ,  $E$ ,  $L$ . Output:  $T$  (a  $(4+\epsilon)$ -approximation of MWM). I/O
2
3 //PART 1 (Stream processing): compute  $L$  maximum matchings
4 C: List of Lists; //  $L$  lists to store edges in  $L$  substreams
5 MB: Matrix; //The matching bits matrix of size  $L \times n$ 
6 substream_weights: List; //The list of substream weights;
7 //substream_weights[i] =  $(1+\epsilon)^i$ .
8 has_added: bool; //Controlling adding an edge to only one MCM
9 foreach(WeightedEdge e : E) {
10   has_added = false;
11   for(i = L-1; i >= 0; i--) {
12     if(e.weight >= substream_weights[i]) {
13       if(!MB[e.u][i] && !MB[e.v][i]) {
14         MB[e.u][i] = 1; MB[e.v][i] = 1;
15         if(!has_added) { //Add e only once to the matchings
16           C[i].add(e); has_added = true;
17         } } } } }
18
19 //PART 2 (Post processing): combine  $L$  matchings into a MWM
20 T: List; //A list with the edges of the final MWM
21 tbits: List; //An array containing the matching bits of T
22 for(i = L-1; i >= 0; i--) {
23   foreach(WeightedEdge e : C[i]) {
24     if(!tbits[e.u] && !tbits[e.v]) {
25       tbits[e.u] = 1; tbits[e.v] = 1;
26       T.add(e);
27     } } }
28 return T; CPU

```

Listing 1. (§ 4.1) The high-level overview of the substream-centric Maximum Weighted Matching algorithm, based on the scheme by Crouch and Stubbs [46]

**4.1.1 Time & Space Complexity.** The space complexity is  $O(nL)$  to track the matching bits, and  $O(\min(m, n/2)L \log(n))$  to store the edges of  $L$  maximum matchings. The time complexity is  $O(mL)$  for substream processing on the FPGA and  $O(nL)$  for substream merging on the CPU, resulting in the total complexity of  $O(mL + nL)$ .

**4.1.2 Reducing Data Transfer with Matching Bits Storage.** We assume that the input is streamed according to the CSR order corresponding to the input adjacency matrix. If we process a matrix row, we load the edges from DRAM to the FPGA. Further, we can store the matching bits  $mb_u$  of vertex  $u$  in BRAM on the FPGA, since they are reused multiple times (temporal locality). The matching bits of  $v$  are streamed in from DRAM. Since the matching bits for  $v$  are not used afterwards for the same matrix row, we write them back to DRAM. Using this approach, we can process the whole graph row by row and need to store only the  $u$ -matching bits in BRAM.

## 4.2 Blocking Design for More Performance

**4.2.1 Problem of Data Dependency.** We cannot start processing the next row of the adjacency matrix until the last matching bits of the previous row have been written to DRAM, because we might require accessing the same  $v$ -matching bits again (read after write dependency). In such a design, the waiting time required after each row could grow, decreasing performance.

**4.2.2 Solution with Blocking Rows.** We alleviate the data dependency by applying *blocking*. We merge  $K$  adjacent rows to become one stream; we call the merged stream of  $K$  rows an *epoch*, and

```

1 //Input and Output: as in Listing 1. I/O
2
3 //PART 1 (Stream processing): compute  $L$  maximum matchings
4 for(Epoch  $k = 1$ ;  $k \leq \lceil n/K \rceil$ ;  $k++$ ) {
5   Load  $u$ -matching bits from DRAM into double-buffered BRAM
6   Merge the  $K$  rows of edges (loaded from DRAM into one stream  $S$ )
7   with a merging network (Figure 4), apply lexicographic order
8   //Process each edge
9   foreach(WeightedEdge  $e : S$ ) {
10    Matching bits requester loads matching bits ( $e.v$ ) from DRAM
11    //Apply the 8 stage pipeline for each edge
12    Stage 1: extract  $v$ -matching bits from a data chunk,
13    determine BRAM address
14    Stage 2: load the matching bits for  $e.u$  from BRAM
15    Stage 3: wait for one cycle due to the latency of the BRAM
16    Stage 4: store the arriving BRAM data in a register, select
17    the correct matching bits, compute  $e_l[i] = e.w = (1+\epsilon)^i$ 
18    Stage 5: compute the matchings
19    Stage 6: write  $u$ -matching bits to BRAM, write
20     $v$ -matching bits to double-buffered BRAM if required
21    Stage 7: determine the least significant bit in  $te$ ,
22    store them in variable  $i$ 
23    Stage 8: write the edge to DRAM at  $C[i]$ 
24    (if part of a matching), write  $v$ -matching bits to DRAM
25  }
26  Wait till all writes to DRAM are committed} FPGA
27
28 //PART 2 (Post processing): As in Listing 1. CPU

```

Listing 2. (§ 4.1–§ 4.4) The pseudocode of the substream-centric MWM algorithm, enhanced with the blocking optimization and a lexicographic ordering.

denote the  $k$ -th epoch (starting counting from 1) as  $k$ . There are  $\lceil n/K \rceil$  epochs in total. To enable merging the rows, we define a *lexicographic ordering* over all edges.

**4.2.3 Lexicographic Ordering.** Let a tuple  $(u, v, w, k)$  denote an edge with vertices  $u, v$ , weight  $w$ , and associated epoch  $k = \lfloor (u-1)/K \rfloor + 1$ . Then, the lexicographic ordering is given by:  $(u_a, v_a, w_a, k_a) < (u_b, v_b, w_b, k_b)$  iff  $k_a < k_b \vee (k_a = k_b \wedge v_a < v_b) \vee (k_a = k_b \wedge v_a = v_b \wedge u_a < u_b)$ ; the edge weight is ignored. An example is in Figure 3 (top). The lexicographic ordering is implemented by a simple merging network.

**4.2.4 Advantages of Blocking.** At the end of each epoch,  $v$ -matching bits are written to DRAM. This reduces the number of such transfers from  $n$  to  $n/K$ . Moreover, if edges in different rows share the same  $v$ -matching bits, only one load from DRAM is required. Finally,  $u$ -matching bits can be kept in BRAM, since they are reused multiple times.

**4.2.5 Further Optimizations.** To achieve a performance of *up to one processed edge per cycle*, we pipeline the processed edges, we distribute the  $u$ -matching bits over multiple BRAMs (to facilitate reading data from different addresses), and we double buffer  $u$ -matching bits to reduce latencies.

### 4.3 Input and Output Format

The input to the FPGA algorithm is a custom variant of the Compressed Sparse Row (CSR) format. An example is given in Figure 3 (bottom). The format has two parts: The `pointer_data` and the `graph_data`. First, the `pointer_data` stores information about the start and end of each row of the

adjacency matrix. An entry contains three parts: the ID of the *data chunk* with information about where the first edge is stored, the data chunk offset denoting the offset of the first edge from the start of the data chunk, and the number of associated edges (a data chunk refers to data of a given size at an aligned memory address). Each entry uses 32 bits, making an entry of the *pointer\_data* 96 bits. We fit five entries (480 bits) in a data chunk. Second, the *graph\_data* is a stream of edges. One entry consists of the column index and the edge weight. The row identifier is given by the corresponding entry in the *pointer\_data*. One *graph\_data* entry requires 64 bits, allowing to store eight edges in a data chunk.

Our custom data layout has different advantages over the usual CSR format. First, a single entry of the *pointer\_data* already gives all required information about the start and length of the row of the adjacency matrix. This entails some redundancy compared to the traditional CSR, but only requires one load from DRAM to resolve a given edge. Further, CSR splits the column indices and values. We merge them together in one stream, reducing the number of random accesses.

The output of the FPGA consists of  $L$  substreams of edges. The  $i$ -th stream contains edges of the maximum matching  $C_i$ . We use 128 bits for each edge: 32 bits each for the vertex IDs, the edge weight, and the assigned index  $i$  of the maximum matching (which could be omitted). A single data chunk therefore contains four output edges.

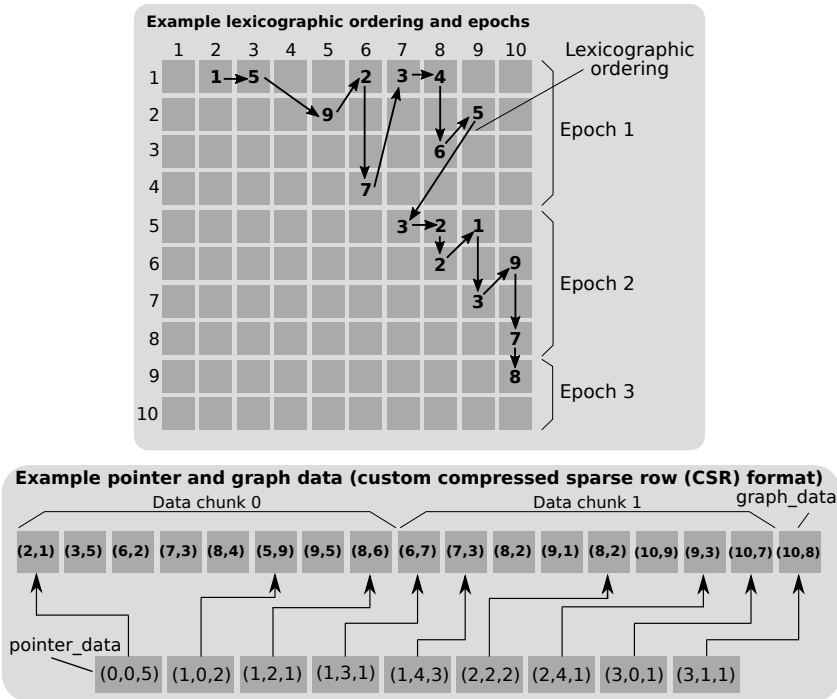


Fig. 3. An example input adjacency matrix, its annotated lexicographic ordering illustrated by arrows ( $K = 4$ ), and and its custom compressed sparse row (CSR) format. The entries of the adjacency matrix denote the weight of an edge.

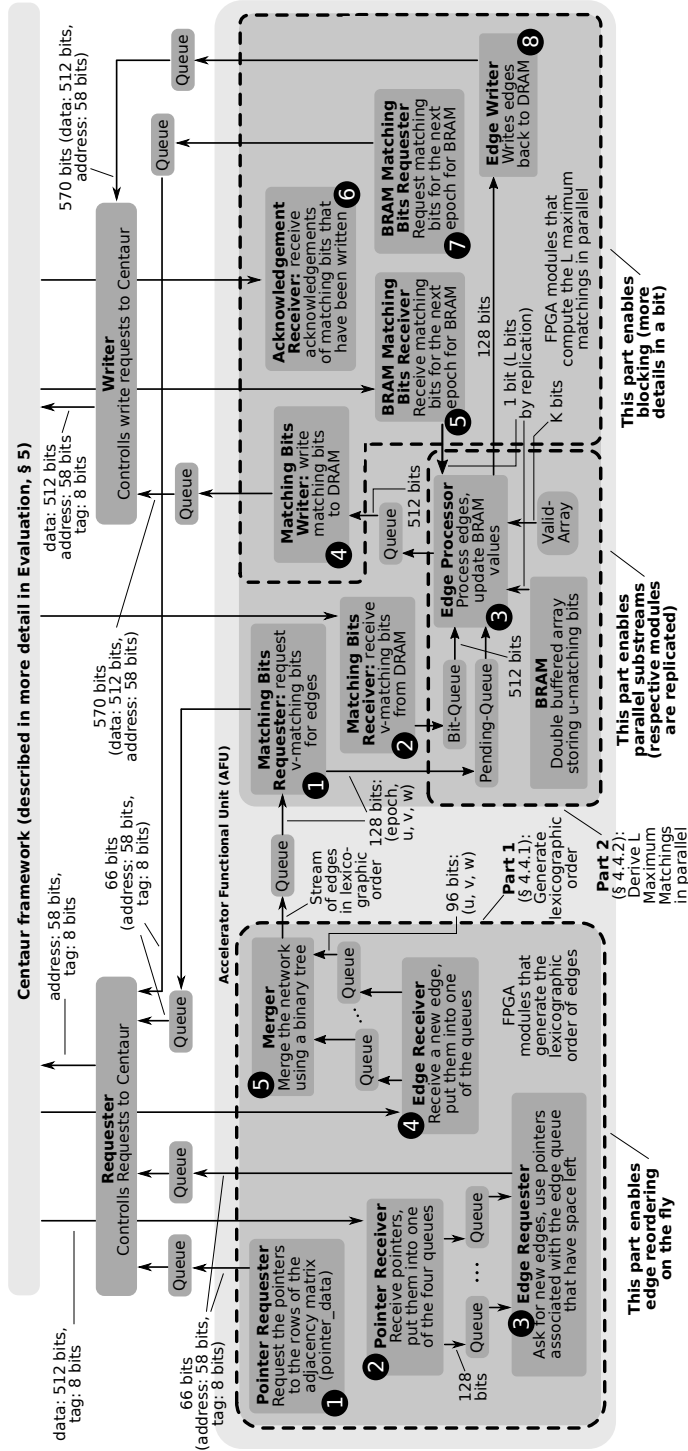


Fig. 4. (§ 4.4) The interaction of the FPGA modules to approximate MWM. For clarity, the State Controller is omitted. The wires of incoming data from Centaur consists of 512 bits for data and 8 bits for tag. All modules are connected using AXI interfaces. All valid bits are omitted. The merger network is in Figure 5.



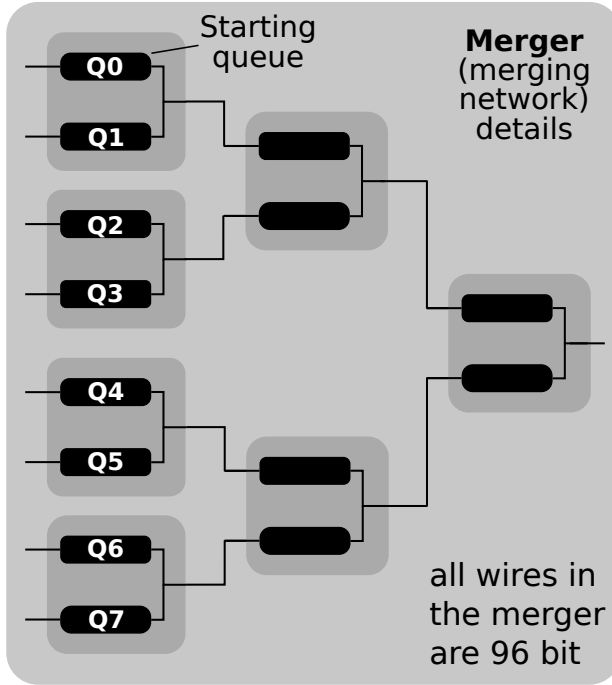


Fig. 5. (§ 4.4) The **merger network** from Figure 4 for  $K = 8$ .

#### 4.4 Details of Processing Substreams on FPGA

We explain the interaction of the FPGA modules dedicated to generating the lexicographic ordering (Part 1) and computing the maximum matchings (Part 2); see Figure 4, Figure 5, and Listing 2.

**4.4.1 Generating Lexicographic Ordering.** As input to the FPGA, we get the address pointing to the start of `pointer_data`, the number of vertices  $n$ , the number of edges  $m$ , a pointer  $p_{out}$  where we write the output to, and an offset value  $o$  to distinguish the  $L$  output streams (start of output stream  $i$  is at  $p_{out} + i \cdot o$ ). The **pointer requester** is responsible for requesting the data chunks holding the `pointer_data`. The requested pointers arrive at the **pointer receiver**. Given a data chunk, the pointer receiver unwraps the five pointers, and passes them to the **edge requester**. The `pointer_data` from the pointer receiver is passed to four different queues  $Q_0, Q_1, Q_2, Q_3$ , where every queue gets a subset of the pointers dependent on  $K$ . Assume for simplicity that the vertex IDs start at 0 and  $(K \bmod 4) = 0$ . Then, to be precise, given a pointer  $p(u)$  pointing to row  $u$ , we assign  $p(u)$  to  $Q_i$  if  $(u \bmod K) \geq K/4 \cdot i \wedge (u \bmod K) < K/4 \cdot (i + 1)$ . For example, with  $K = 16$ ,  $Q_0$  stores  $p(u)$  with  $u = 0, 1, 2, 3, 16, 17, 18, 19, 31, \dots$ , and  $Q_1$  stores  $p(u)$  with  $u = 4, 5, 6, 7, 20, \dots$ . The `pointer_data` is loaded from the queues into a BRAM array  $BP$  of size  $K$ , where every entry holds two pointers ( $2K$  pointers are therefore stored in total). If an entry  $i$  of  $BP$  has pointers  $p(u')$  and  $p(u'')$ , it holds that  $i = (u' \bmod K) = (u'' \bmod K)$  and  $p(u'')$  requests edges for an epoch after  $p(u')$ . Therefore, only the first pointer in an entry is valid to use and we have random access to  $K$  valid pointers in total. To describe the mechanism that determines the selection of the next pointer to request new edges, we first inspect further processing steps.

The **edge receiver** gets data chunks containing `graph_data` from the framework and unwraps them (we use the Centaur framework [109] to access main memory independently of the CPU).

Information regarding the offset and number of edges which are valid for a data chunk request is also passed from the edge requester to the edge receiver. Next, an edge  $e = (u, v, w)$  is passed from the edge receiver to the **merger**. There, the edge is inserted in a *starting queue* (with ID  $(u \bmod K)$ ). The merger merges the  $K$  streams in lexicographic ordering. It consists of a series of merging elements, where each element has two input queues and an output port. The element compares edges in its queues and outputs the edges according to the lexicographic ordering. The merging elements form a binary tree, such that for a given  $K$ , there are  $K/2$  starting elements with  $K$  starting queues in total.

The edge requester can observe the size of the starting queues of the merger. It operates in two modes to determine a pointer to new edges. In mode 1, one selects a pointer  $p(u)$  from queue  $Q_i$  as the next candidate if the corresponding starting (merger) queue  $(u \bmod K)$  does not overflow, and store the pointer in BRAM  $BP$  at position  $(u \bmod K)$ . If mode 1 fails (for example, if there is no empty space at the appropriate position in  $BP$ ), then mode 2 selects the pointer according to the merger starting queue which has the least amount of edges. Note that the edge requester also takes the requests which are in flight into account to predict the future size of the starting queue. This approach ensures that the merger queues do not overflow and their load is balanced.

For a row  $u$  which has no edges, a special information is passed from the edge requester to the edge receiver. It then inserts an artificial edge in the merger. This allows to overcome problems, where a merging element waits for new input, but does not receive any, since the adjacency matrix row is empty. The merging network filters these edges at the output port (they are not passed on).

**4.4.2 Deriving  $L$  Maximum Matchings.** The stream in lexicographic ordering is passed to the **matching bits requester**. This module requests the  $v$ -matching bits from DRAM. It can only operate when the bits of the epoch before have been acknowledged. Also, it only processes edges belonging to the current epoch which is defined by the **state controller**. The requested data is received in the **matching bits receiver**. It passes the full data chunk to the **edge processor**. Using the matching bits and the ordered stream of edges, the edge processor computes the  $L$  maximum matchings in parallel in an 8-stage pipeline (Listing 2, Lines 10–24). In Stage 1,  $v$ -matching bits for a given edge are extracted from a data chunk. Further, the address of the  $u$ -matching bits in BRAM is computed. Since the more up-to-date  $v$ -matching bits might also be stored in BRAM, this address is also determined. In Stage 2, read requests to fetch the matching bits from BRAM are issued. Stage 3 only waits one clock cycle for BRAM to return the data. In Stage 4, the BRAM data arrives and is stored in a register. The stage also decides if  $v$ -matching bits are taken from the data chunk or from BRAM. Further, the stage computes the matching value  $te$  indicating if an edge  $e = (u, v, w)$  belongs to substream  $E_i$ ;  $te[i] = w \geq (1 + \varepsilon)^i$  for  $i \in \{0, \dots, L - 1\}$ . In Stage 5, the actual matching is computed. As the BRAM data from Stage 4 may already be obsolete, the computed values are also stored in registers for instant access in the next cycle. The result is passed to Stage 6, in which the updated  $u$ -matching bits (and if required also the  $v$ -matching bits) are written back to BRAM. In Stage 7, the maximum matching with the highest index, to which the edge is assigned, is determined. Finally, Stage 8 passes the edge to the **edge writer** to write it back to DRAM (if the edge is used in a matching) and also passes the updated  $v$ -matching bits to the **matching bits writer** for writing back to DRAM.

The BRAM storing the  $u$ -matching bits is double buffered. While the first BRAM buffer is used in the edge processor, the matching bits for the next epoch are loaded from DRAM to the second BRAM buffer. Since an epoch can alter the  $u$ -matching bits required for the next epoch, we write the according updates also in the double buffered BRAM if required. To prevent that stale data from DRAM overwrites the more up-to-date data, we use a register (the valid-array) as flag. After an epoch, the access is redirected to the BRAM containing the loaded data. The **BRAM matching**

**bits requester** requests the according data from DRAM, and the **BRAM matching bits receiver** unwraps the data chunks. It passes the data to the edge processor. There, Stage 6 checks for data from the BRAM matching bits receiver and updates the according entry in the BRAM.

The **acknowledgement receiver** tracks the count of write acknowledgements from the framework and determines if all  $v$ -matching bits are committed to DRAM when an epoch ends. When all edges from the epoch are processed, the state controller indicates the start of the next epoch.

#### 4.5 Substream Merging on the CPU

After the  $L$  MCMs are written to DRAM, the CPU inspects them in the decreasing order to compute the final maximum matching  $(4 + \varepsilon)$ -approximation. This part is a simple greedy scheme that exposes little parallelism, thus we execute it on the CPU. It takes  $O(Ln)$  time and  $O(Ln)$  work.

#### 4.6 Summary of Optimizations

In Figure 4, we also use dashed rectangles to illustrate which modules are responsible for the most important optimizations: edge reordering on the fly, parallel substreams (pipelining), and blocking. Modules responsible for pipelining are appropriately replicated.

#### 4.7 Interactions with DRAM

We use the Centaur framework [109] as the interface to the Accelerator Functional Unit (AFU), the custom FPGA implementation, allowing to access main memory independently of the CPU. Centaur consists of a software and a hardware part. The software part allows to start and stop hardware functions, to allocate and deallocate the shared memory, and pass input parameters to the FPGA. The hardware part is responsible for bootstrapping the FPGA, setting up the QPI endpoint, and handling reads and writes to the main memory. In our design, we use dedicated arbiter modules for all read and write requests to Centaur: the requester and the writer. The requester has four queues. The pointer requester, the edge requester, the matching bits requester, and the BRAM matching bits requester can all write the DRAM address (from which a data chunk of 512 bits should be read) to these four requester queues. The requester uses a fixed priority order to send the requests to the Centaur framework. Centaur provides an 8 bit tag to identify the requests. For simplicity, fixed tags are used for each of the four modules emitting requests. The modules listening for incoming data (pointer receiver, edge receiver, matching bits receiver, and BRAM matching bits receiver) process the data chunk of size 512 bits only when the tag matches the expected value. The design relies on the FIFO behavior of Centaur, such that requests with the same tag are not reordered. The writer orchestrates the modules (the matching bits writer and the edge writer) which issue writes to DRAM. Similarly to the requester, the writer has a queue for each writing module (two in total) and it uses a fixed priority order. The modules issuing the requests use fixed tags. Note that data is written in chunks of 512 bits. The acknowledgement receiver monitors the Centaur interface for writes that have been written. This information is passed to the state controller (not shown in Figure 4) to orchestrate the modules. Since Centaur allows only to access data in chunks of 512 bits, the addresses passed to the framework have 58 bits.

## 5 EVALUATION

We now illustrate the advantages of our hybrid (CPU+FPGA) MWM design and inspect resource and energy consumption. *For every benchmark, each tested algorithm was synthesized, routed, and executed on the hybrid FPGA platform specified below.*

## 5.1 Setup, Methodology, Baselines

**5.1.1 Compared Algorithms.** Since to our best knowledge *no MWM algorithms for FPGAs are available*, we compare our design to three state-of-the-art CPU implementations. In total, we evaluate three CPU and two CPU+FPGA algorithms; see Table 4. First ❶, we implement a sequential CPU-only version of the substream-centric MWM, based on the scheme by Crouch and Stubbs [46], as presented in Listing 1 (CS-SEQ). Second ❷, we parallelize the algorithm with OpenMP’s `parallel-for` statement to compute different maximum matchings in parallel (CS-PAR). Third ❸, we implement the algorithm by Ghaffari [62] (G-SEQ) that provides a  $(2 + \epsilon)$ -approximation to MWM with time complexity of  $O(m)$  and space complexity of  $O(n \log(n))$  bits. Thus, this algorithm is optimal in the asymptotic time and space complexity. We compare these three algorithms to our optimized FPGA+CPU implementation, SC-OPT ❹. Finally, we also tested SC-SIMPLE ❺, a variant of our implementation that uses no blocking. SC-SIMPLE delivers more performance than the comparison baselines but it is consistently outperformed by SC-OPT, we thus usually exclude it for clarity of presentation. However, we use it in power consumption experiments to illustrate how much additional power is used by the design optimizations in SC-OPT. *To our best knowledge, we report the first performance data for deriving maximum matchings on the FPGA.*

**5.1.2 Implementation Details and Reproducibility.** We implement our algorithms in Verilog on a hybrid CPU+FPGA system using the Centaur framework [109]. The modules outlined in Figure 4 are connected using AXI interfaces. To facilitate reproducibility and interpretability [69], we make the whole code publicly available<sup>3</sup>.

Algorithm	Platform	Time complexity
Crouch et al. [46] Sequential (CS-SEQ)	CPU	$O(mL + nL)$
Crouch et al. [46] Parallel (CS-PAR)	CPU	$O(mL/T + nL)$
Ghaffari [62] Sequential (G-SEQ)	CPU	$O(m)$
Substream-Centric, no blocking (SC-SIMPLE)	Hybrid	$O(m + nL^2)$
Substream-Centric, with blocking (SC-OPT)	Hybrid	$O(m + n/K + nL)$

Table 4. (§ 5) Overview of **the evaluated MWM algorithm implementations.**

**5.1.3 Setup.** We use Intel HARP 2 [108], a hybrid CPU+FPGA system. It is a dual socket platform where one socket is occupied by an Intel Broadwell Xeon E5-2680 v4 CPU [72] with 14 cores (28 threads) with up to 3.3 GHz clock frequency. Each core has 32 KByte L1 cache and there is 35 MByte L3 cache in total. An Arria-10 FPGA (10AX115U3F45E2SGE3) is in the other socket. The used FPGA has speed grade 2 [74]. It provides 55 Mbit in 2,713 BRAM units and 427,200 ALMs. The FPGA is connected to the CPU by one QPI and two PCIe links. The system runs Ubuntu 16.04.3 LTS with kernel 4.4.0-96 as the operating system. All host code is compiled with gcc 5.4.0 and the `-O3` compile flag.

**5.1.4 Datasets.** The input graphs are shown in Table 5. We use both synthetic (Kronecker) power-law graphs of size up to  $n = 2^{21}$ ,  $m = 48n$  from the 10th DIMACS challenge [1] and real world KONECT [85] and SNAP [91] graphs. For unweighted graphs, we assigned weights uniformly at random with a fixed seed. The value range is given by  $[1, (1 + \epsilon)^{L-1} + 1]$ .

<sup>3</sup>[https://spl.inf.ethz.ch/Parallel\\_Programming/Matchings-FPGA](https://spl.inf.ethz.ch/Parallel_Programming/Matchings-FPGA)

Graph	Type	Reference	$m$	$n$
Kronecker	Synthetic power-law	DIMACS 10 [1]	$\approx 48n$	$2^k; k = 16, \dots, 21$
Gowalla	Social network	KONECT [85]	950,327	196,591
Flickr	Social network	KONECT [85]	33,140,017	2,302,925
LiveJournal1	Social network	SNAP [91]	68,993,773	4,847,571
Orkut	Social network	KONECT [85]	117,184,899	3,072,441
Stanford	Hyperlink graph	KONECT [85]	2,312,497	281,903
Berkeley	Hyperlink graph	KONECT [85]	7,600,595	685,230
arXiv hep-th	Citation graph	KONECT [85]	352,807	27,770

Table 5. Selected used graph datasets.  $Kx$  denotes a Kronecker graph with  $2^x$  vertices.

**5.1.5 Measurements.** The runtime is measured by `clock_gettime` with parameter `CLOCK_MONOTONIC_RAW`, allowing the nanosecond resolution. The runtime of the FPGA implementations is determined by the Centaur framework. We execute each benchmark ten times to gather statistics and we use box plot entries to visualize data distributions.

## 5.2 Scaling Size of Synthetic Graphs

We first evaluate the impact from varying graph sizes (synthetic power-law Kronecker graphs), for the fixed amount of parallelism (the *weak scaling* experiment). The results are illustrated in Figure 6. The throughput for CS-SEQ and CS-PAR stays approximately constant below  $\approx 12M$  edges/s. G-SEQ decreases in performance as the graph size increases. We conjecture that this is due to the increasing size of the hash map used to track pointers. This increases the time for inserts and deletes, and might also require re-allocations to increase the space. The performance for SC-OPT increases from  $\approx 135M$  to  $\approx 140M$  edges/s. This is because the initial (constant) overhead (due to reading from DRAM) becomes less significant with larger graphs. **We conclude that the substream-centric SC-OPT beats comparison targets for all considered sizes of power-law Kronecker graphs.**

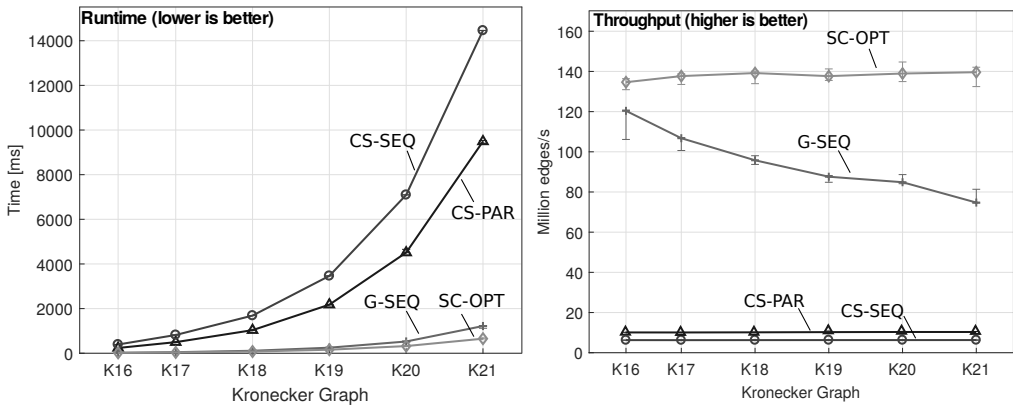


Fig. 6. (§ 5.2) **Influence of graph size  $n$  on performance** (synthetic power-law graphs).  $K = 32, L = 64, T = 4, \varepsilon = 0.1$ .

### 5.3 Processing Different Real-World Graphs

We next analyze the performance of the considered designs for different real-world graphs; the results are illustrated in Figure 7. CS-SEQ and CS-PAR achieve sustained  $\approx 3\text{M}$  edges/s and  $\approx 10\text{M}$  edges/s, respectively. The performance of SC-OPT is  $\approx 45\text{M}$  edges/s for small graphs due to the initial overhead of reading data from DRAM. Compared to the experiment with Kronecker graphs, the performance of both SC-OPT and G-SEQ is lower for all graphs except Orkut. The reason is the average vertex degree: it equals  $\approx 48$  in Kronecker graphs compared to  $\approx 14$  in Flickr and LiveJournal1. If the ratio is high, G-SEQ can drop many edges without further processing in an early phase. This reduces expensive updates to the hash map and lists. For SC-OPT, the waiting time (of data dependencies) lowers the performance. Still, **stream-centric SC-OPT ensures highest performance for all considered real-world graphs.**

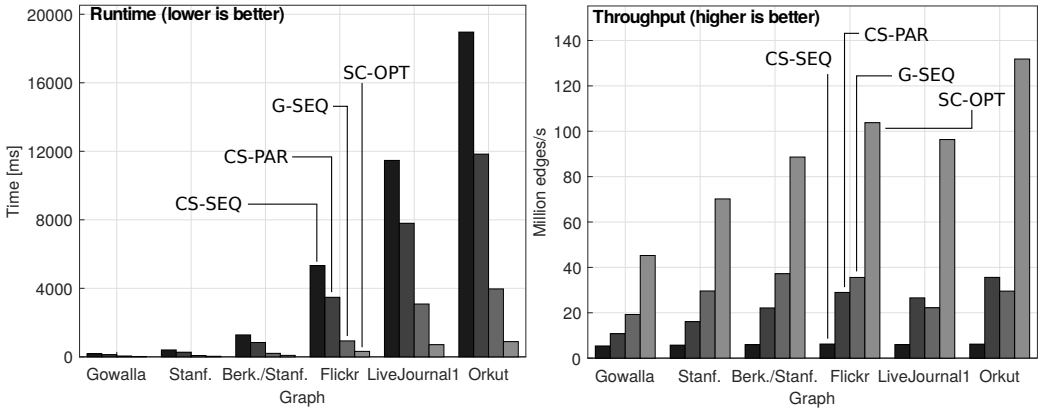


Fig. 7. (§ 5.3) Influence of graph dataset  $G$  on performance (real-world graphs).  $K = 32, L = 64, T = 4, \epsilon = 0.1$ .

### 5.4 Scaling Number of Threads $T$

In the CPU versions, one can compute in parallel different maximum matchings in SC-PAR using  $T$  threads. In the following, we run a *strong scaling* experiment (fixed graph size, varying  $T$ ) for a power-law Kronecker graph. Figure 8 illustrates the results. Since G-SEQ and CS-SEQ are not multi-threaded, they do not scale with  $T$ . The parallelized CS-PAR reaches up to  $\approx 40\text{M}$  edges/s, a  $\approx 6\times$  improvement over the sequential version, and an  $\approx 14\times$  improvement over the parallel version with one thread. Therefore, the algorithm is still  $\approx 3\times$  slower than SC-OPT which achieves up to  $\approx 140\text{M}$  edges/s on the K20 Kronecker graph. Scaling is limited since the parallel version takes  $L$  passes over the stream, whereas the other CPU algorithms process the input in one pass. The bandwidth usage of the parallel version with  $T = 64$  threads is  $\approx 32$  GB/s ( $\approx 44\text{M}$  edges and 64 passes in one second), assuming no data sharing. Note that we only parallelize the stream-processing part which computes the  $L = 64$  maximum matchings. However, as our analysis shows that the post-processing part takes  $<1\%$  of the computation time of the maximum matching, parallelization of post-processing would provide hardly any benefit. We conjecture that the scaling of SC-PAR stops due to bandwidth limitations and the limited computational resources of 14 cores. Finally, **SC-OPT is the fastest regardless of  $T$  used by other schemes.**

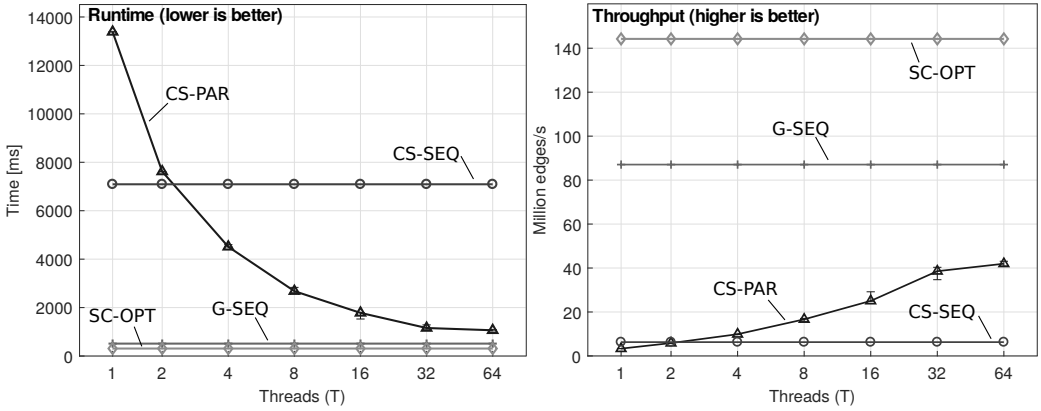


Fig. 8. (§ 5.4) Influence of the number of threads  $T$  on performance. The input graph is Kronecker with  $n = 2^{20}$ ,  $K = 32$ ,  $L = 64$ ,  $\epsilon = 0.1$ .

### 5.5 Approximation Analysis

We briefly analyze how well in practice SC-OPT approximates the exact MWM. The results are in Figure 9 (SC-OPT, SC-SIMPLE, CS-SEQ, and CS-PAR produce the same results). The accuracy is negligibly ( $\approx 3\%$ ) lower than that of G-SEQ for a fixed  $\epsilon$  and varying  $n$  (Kronecker graphs). The higher  $\epsilon$  becomes, the more SC-OPT has advantage over G-SEQ. As higher  $\epsilon$  entails less circuit complexity (fewer substreams are processed independently, assuming a fixed  $w_{max}$  [46]), we conclude that the substream-centric MWM SC-OPT scheme provides better approximation than G-SEQ when physical resources become more constrained.

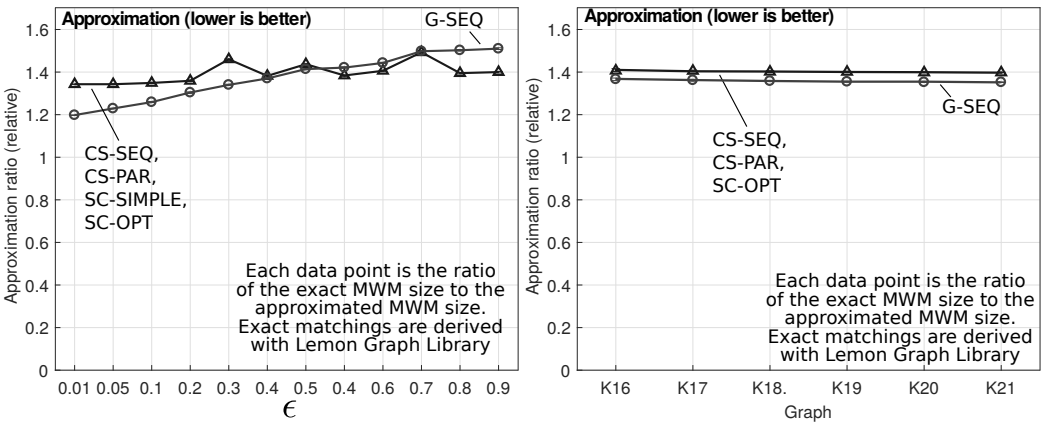


Fig. 9. (§ 5.5) Approximation analysis. The input graph is Kronecker with  $n = 2^{19}$  (left).  $L = 128$  and  $\epsilon = 0.1$  (right). Explanation of the approximation ratio is provided in Section 2.

### 5.6 Influence of Blocking Parameter $K$

We also analyze the performance impact from  $K$ , a parameter that determines how many rows in the streamed-in adjacency matrix are merged together using a lexicographic ordering. Figure 10 illustrates the results. On one hand, the CPU schemes cannot take significant advantage when  $K$

increases, showing that no cache locality is exploited. On the other hand, FPGA-based SC-OPT accelerates from  $\approx 125\text{M}$  to  $\approx 175\text{M}$  edges/s. This is up to  $2\times$  faster than the work-optimal G-SEQ and up to  $55\times$  faster than CS-SEQ. This is expected as the amount of stalling is reduced by a factor of  $n/K$ . Moreover, increasing  $K$  allows to share more matching bits between edges. The performance impact is reduced when  $K$  reaches 256. We conjecture this is because of the random access to the matching bits, approaching the peak random bandwidth. Furthermore, G-SEQ outperforms all other CPU implementations with up to  $\approx 90\text{M}$  edges/s. Compared to CS-SEQ ( $\approx 3.15\text{M}$  edges/s) and CS-PAR ( $\approx 5.6\text{M}$  edges/s), this is  $>15\times$ . Finally, parallelization comes with high overhead, such that the four threads in CS-PAR achieve less than  $2\times$  speedup compared to CS-SEQ. **We conclude that our blocking scheme enables SC-OPT to achieve even higher speedups.**

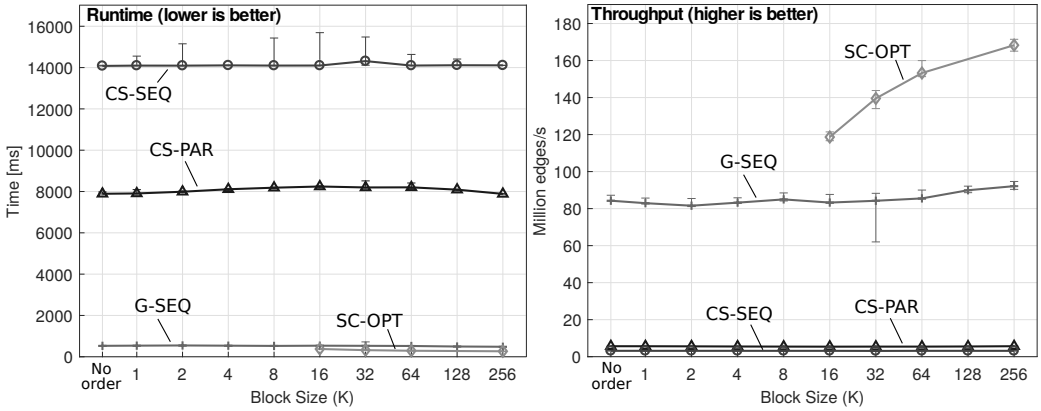


Fig. 10. (§ 5.6) **Influence of epoch size  $K$  on the performance.** The input graph is Kronecker with  $n = 2^{20}$ .  $L = 128$ ,  $T = 4$ , and  $\varepsilon = 0.1$ .

## 5.7 Influence of Maximum Matching Count $L$

Finally, we analyze the impact of  $L$  on performance.  $L$  is the number of substreams and thus maximum matchings computed independently. CS-SEQ and CS-PAR achieve high performance with up to  $\approx 400\text{M}$  edges/s for  $L = 1$ . The performance drops linearly with  $L$  ( $X$ -axis has a logarithmic scale) to  $\approx 800\text{k}$  edges/s for CS-SEQ and  $\approx 1.3\text{M}$  edges/s for CS-PAR. G-SEQ also drops in performance as  $L$  increases due to  $\varepsilon$  and  $w_{max}$ . Since  $L$  increases, we also increase the range of the weight ( $L$  influences the approximation by  $\varepsilon = \sqrt[L]{w_{max}} - 1$ ). Thus, for  $L = 1$  the maximum edge weight is given by  $w_{max} = 1$ , allowing G-SEQ to drop many edges in an early phase. The drop of performance between  $L = 32$  and  $L = 64$  are due to a change in  $\varepsilon$ , requiring G-SEQ to store more data. Similarly, we change  $\varepsilon$  between  $L = 128$  and  $L = 256$ . **SC-OPT keeps its performance at  $\approx 140\text{M}$  edges/s ( $\approx 330\text{ms}$ ) and outperforms other schemes.**

## 5.8 FPGA Resource Utilization

Table 6 shows the usage of FPGA resources. As maximum matchings are computed on the FPGA in one clock cycle, the number of computed matchings  $L$  influences the amount of used logic. Moreover, for SC-OPT,  $K$  and  $L$  determine the FPGA layout. Specifically,  $K$  influences the BRAM usage, since every element in the merging network requires two queues which are each mapped to one BRAM unit. We also consider the amount  $B$  [bits] of BRAM allocated to storing the matching bits. SC-OPT requires only 21% of Arria-10's BRAM and 32% out of all ALMs for a design that



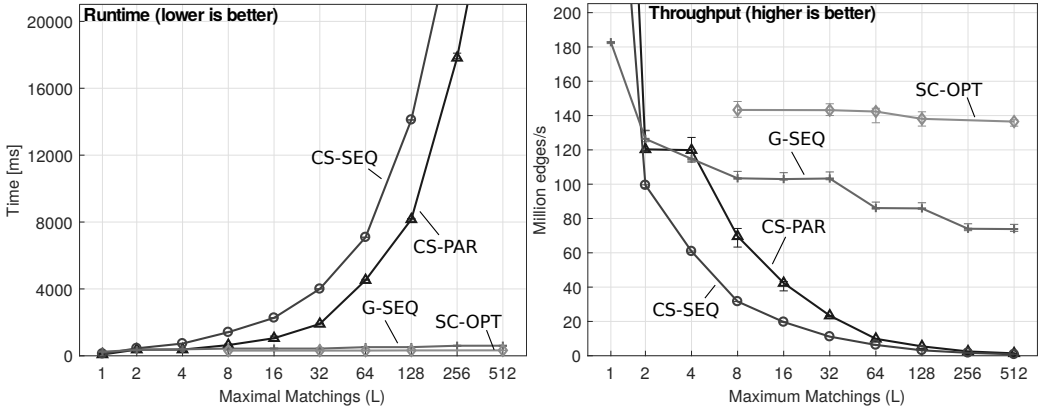


Fig. 11. (§ 5.7) **Influence of  $L$  on performance.** The input graph is Kronecker ( $n = 2^{20}$ ,  $K = 32$ ,  $T = 4$ ). As  $L$  changes,  $\epsilon$  changes as follows: for  $1 \leq L \leq 32$ , we select  $\epsilon = 0.6$ , for  $64 \leq L \leq 128$  we select  $\epsilon = 0.1$ , and for  $256 \leq L \leq 512$  we select  $\epsilon = 0.03$ ;  $w_{max}$  is given by  $w_{max} = (1 + \epsilon)^L$ . We restricted the range of  $L$  for SC-OPT due to the significant runtime required to generate different bitstreams for evaluation.

outperforms other targets by at least  $\approx 2\times$  (Figures 10–11); these speedups can be increased even further by maximizing circuitry utilization. Finally, we include SC-SIMPLE in the analyses to illustrate the impact of additional optimizations in SC-OPT. CS-SIMPLE uses less ALMs (approx. 21%) than CS-OPT (usage varies, up to 82%).

FPGA Algorithm	Parameters	Used BRAM	Used ALMs
SC-SIMPLE	$\log B = 12, L = 8$	5.6 MBit (10%)	89,388 (21%)
SC-SIMPLE	$\log B = 18, L = 6$	21 MBit (38%)	88,920 (21%)
SC-OPT	$K = 32, L = 512$	11.5 MBit (21%)	151,998 (32%)
SC-OPT	$K = 256, L = 128$	24.8 MBit (45%)	350,556 (82%)

Table 6. (§ 5.8) **FPGA resource usage** for different parameters.

## 5.9 Energy Consumption

We estimate the energy consumption of SC-SIMPLE and SC-OPT using the Altera PowerPlay Power Analyzer Tool; see Table 7 (we use 200MHz and include static power). Furthermore, the host CPU (Broadwell Xeon E5-2680 v4) has TDP of 120 Watt [72] when all cores are in use (We use TDP as the baseline for the CPU because the utilized server is physically located elsewhere and we are unable to directly measure the used power). The TDP is an upper bound for CS-PAR at  $T = 64$ . **FPGA designs reduce consumed energy by at least  $\approx 88\%$  compared to the CPU.**

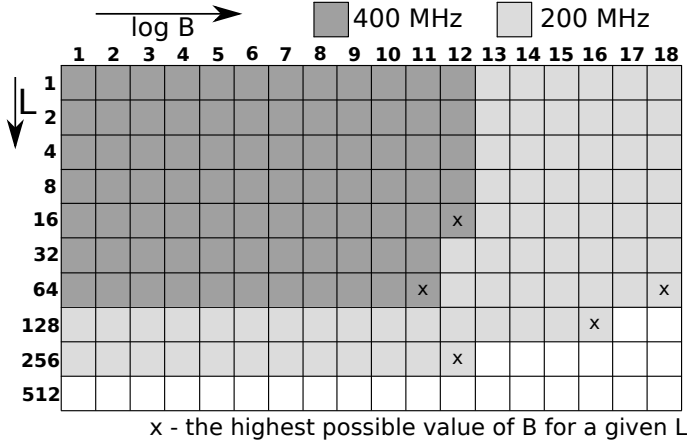
## 5.10 Design Space Exploration

We now briefly analyze the interaction between the performance of our FPGA design and the limitations due to the clock frequency. The resource usage, determined by  $L$  and  $B$ , influences the frequency upper bound due to wiring and logic complexity. We applied a grid search to derive feasible frequencies for SC-SIMPLE; see Figure 12 (we exclude SC-OPT as our analysis shown that the design is too complex to run at 400MHz and we were only able to use it with 200MHz). Dark

Algorithm	Parameters	Energy Consumption [W]
SC-SIMPLE	$\log B = 18, L = 6$	14.714
SC-SIMPLE	$\log B = 12, L = 8$	14.598
SC-OPT	$K = 32, L = 512$	14.789
SC-OPT	$K = 256, L = 128$	14.789
SC-OPT	$K = 32, L = 64$	14.657
CS-PAR	$T = 64$	120

Table 7. (§ 5.9) **Estimated energy consumption for different parameters.**

grey indicates 400MHz, light grey indicates 200MHz (we use only the two frequencies as only those two were supported by the Centaur framework at the time of evaluation). Two factors have shown to limit the performance. First, while computing the matching, we use an addition with a variable that uses  $L$  bits. Thus, the addition complexity grows linearly with  $L$ . More importantly, the BRAM signal propagation limits the frequency. For example, for SC-SIMPLE and  $\log B = 13$ , the place and route report shows that the reset signal to set all BRAM units to zero becomes the critical path. As alleviating these two issues would make our final design even more complex, we leave it for future work. Specifically, we are now working on a general FPGA substream centric engine that will feature pipelined reset BRAM signals.

Fig. 12. (§ 5.10) **Design space exploration: the used (available) frequencies.**

### 5.11 Optimality Analysis

We also discuss how far the obtained results are from the maximum achievable performance numbers; we focus on the most optimized SC-OPT. SC-OPT can process up to  $\approx 175$ M edges/s. This is close to the optimum due to different reasons: Firstly, the implementation can process up to 1 edge per cycle (200M edges/s). Thus, the achieved performance is optimal within only  $\approx 12\%$ . Second, assuming that edges are read aligned from memory, it allows to read 8 edges per read request. Further, if every edge requires its own data chunk with matching bits, it needs 1 request per edge. Overall, this results in  $1 + 1/8 = 1.125$  read requests per edge. Under this assumption, the performance is limited to 178M edges/s. SC-OPT performs close to this bound, which is possible because *the matching bits can be shared between edges*.

## 6 BEYOND SUBSTREAM-CENTRIC MM

We now briefly discuss how to apply our substream-centric FPGA design to other streaming graph algorithms. First, we identify some MM schemes that also divide the streamed dataset into substreams and can straightforwardly be adopted to the hybrid CPU+FPGA system. The **MWM algorithm by Grigorescu et al. [66]** reduces the MWM problem to  $O(\epsilon^{-1} \log(n))$  instances of maximum matchings, which could be processed on the FPGA analogously to our design; its merging phase could also be executed on the CPU. All our optimizations, such as blocking, are applicable in this case. Moreover, the **MWM algorithm by Feigenbaum et al. [60, Algorithm 4]** does not divide the stream of edges into substreams, but its design would potentially allow for applying our blocking scheme. A key part of this algorithm is maintaining a certain value  $q_e$  associated with each edge  $e$ . Given an edge  $e = (u, v, w)$ ,  $q_e$  depends on values  $q_u$  and  $q_v$  associated with vertices  $u$  and  $v$ . We can apply the blocking pattern by storing  $q_u$  for  $u$  in BRAM, and streaming in  $q_v$  for  $v$ . Next, the **MWM algorithm by Ghaffari [62]** provides a  $(2 + \epsilon)$ -approximation. The algorithm compares the weights of incoming edges to values  $\varphi$ , indexed by  $u$  or  $v$ . Therefore, it can be computed on the FPGA using the blocking pattern by storing the values  $\varphi_u$  in BRAM, and streaming  $\varphi_v$  from DRAM, similarly to matching bits in our design. Further, as the algorithm requires postprocessing to derive the final result, it could be also delegated to the CPU.

We also identify algorithms unrelated to matching that could be enhanced with our design. The random **triangle counting algorithm by Buriol et al. [36]** is also a suitable candidate for the presented blocking pattern. The algorithm requires three passes. In pass 1, the number of paths of length two in the input graph is computed. In pass 2, a random path of length two is selected. In pass 3, the stream is searched for a certain edge, dependent on the randomly selected path. To reduce variance, passes 2–3 are run in parallel using a pre-determined number of random variables (up to a million). This also implies that in pass 3 every edge in the stream must be checked against a million edges. To reduce the workload, a hash map is used. The map is filled with edges which are expected to occur. We propose the following approach to exploit the blocking pattern: the CPU fills a hash map for each epoch with edges expected to arrive. The map is passed to the FPGA. The edges for this epoch are streamed in and compared to the pre-filled hash map. If the epoch changes, the next hash map is passed over.

## 7 RELATED WORK

Our work touches on various areas. We now discuss related works, briefly summarizing the ones covered in previous sections (streaming models in § 3 and streaming maximum matching algorithms in § 3.3, Table 2, and § 6).

### 7.1 Graph Processing on FPGAs

The FPGA community has recently gained interest in processing graphs [18–20, 23, 25, 27–31, 63, 125] and other forms of general irregular computations [21, 22, 24, 53, 61, 82, 119, 120, 129]. First, some established CPU-related schemes were ported to the FPGA setting, for example vertex-centric [57, 58], GAS [145], edge-centric [149], BSP [78], and MapReduce [141]. There are also efforts independent of the above, such as FPGP [47], ForeGraph [48], and others [32, 78, 105, 107, 135, 147]. These works target popular graph algorithms such as BFS or PageRank. Multiplication of matrices and vectors [26, 87] has also been addressed in the context of FPGAs [55, 56, 92, 126, 134, 151]; these efforts could be used for energy-efficient and high-performance graph analytics on FPGAs due to the possibility of expressing graph algorithms in the language of linear algebra [82]. *Our work differs from these designs as we focus on the problem of finding graph matchings. For more detailed analysis of related work, please refer to Table 1.*

## 7.2 Graph Matchings and FPGAs

The only work related to matchings and FPGAs that we are aware of merely uses matchings to enhance FPGA segmentation design [38], *which is unrelated to deriving matchings and graph processing in general.*

## 7.3 Streaming Models and Algorithms

We investigate the rich theory of streaming models [3, 8, 37, 44, 49, 51, 52, 60, 68, 80, 101, 103] and identify the semi-streaming model [60] as the best candidate for using together with FPGAs to deliver algorithms with provable properties that match FPGA characteristics such as limited memory. We then investigate semi-streaming algorithms for maximum matchings [6, 12, 40, 46, 59, 60, 62, 64, 66, 76, 77, 81, 86, 98, 113, 140] and identify the scheme by Crouch and Stubbs [46] that we use as the basis for our substream-centric design *that ensures low-power, high-performance, and high-accuracy general maximum weighted matchings on FPGAs.*

## 7.4 Hybrid FPGA+CPU Platforms

Finally, our work is related to the study of hybrid CPU+FPGA platforms [4, 9, 10, 70, 75, 109, 114, 123, 133, 141, 150, 150]. We illustrate a case study with maximum matchings and show that *hybrid platforms can outperform state-of-the-art parallel CPU designs in both performance and power consumption.* Other works on graph processing on hybrid FPGA-CPU systems include a hybrid scheme for BFS [133].

## 8 CONCLUSION

An important problem in today's graph processing is developing high-performance and energy-efficient algorithms for approximating maximum matchings. Towards this goal, we propose the first maximum matching algorithm for FPGAs. Our algorithm is *substream-centric*: the input stream is divided into substreams that are processed independently on the FPGA and merged into the final outcome on the CPU. This exposes parallelism while keeping communication costs low: only  $O(m)$  data must be streamed from DRAM to the FPGA. Our algorithm is energy-efficient (88% less consumed energy over a tuned CPU variant) and provably accurate, fast (speedups of  $>4\times$  over parallel CPU baselines), and memory-efficient ( $O(n \log^c n)$  required storage).

The underlying FPGA design uses several novel optimizations, such as merging rows of the graph adjacency matrix and ordering resulting blocks lexicographically. This enables low utilization of FPGA resources (only 21% of Arria-10's BRAM and 32% out of all ALMs) while outperforming CPU baselines by at least  $\approx 2\times$ . Both the FPGA implementation and the substream-centric approach could be extended to other graph problems.

Finally, to the best of our knowledge, the proposed design is the first to combine the theory of streaming with the FPGA setting. Our insights coming from the analysis of 14 streaming models and 28 streaming matching algorithms can be used to develop more efficient FPGA designs.

**Acknowledgements** We thank Mohsen Ghaffari for inspiring discussions that helped us better understand graph streaming theory. We also thank David Sidler for his help with the FPGA infrastructure. Funded by the European Research Council (ERC) under the European Union's Horizon 2020 programme grant No. 678880. TBN is supported by the ETH Zurich Postdoctoral Fellowship and Marie Curie Actions for People COFUND program.

## REFERENCES

- [1] 10th DIMACS Challenge. Kronecker Generator Graphs, 2011.
- [2] C. Aggarwal and K. Subbian. Evolutionary network analysis: A survey. *CSUR*, 2014.
- [3] G. Aggarwal, M. Datar, S. Rajagopalan, and M. Ruhl. On the streaming model augmented with a sorting primitive. In *Foundations of Computer Science, 2004. Proceedings. 45th Annual IEEE Symposium on*, pages 540–549. IEEE, 2004.

- [4] J. Agron, W. Peck, E. Anderson, D. Andrews, E. Komp, R. Sass, F. Baijot, and J. Stevens. Run-time services for hybrid CPU/FPGA systems on chip. In *2006 27th IEEE International Real-Time Systems Symposium (RTSS'06)*, pages 3–12. IEEE, 2006.
- [5] J. Ahn, S. Hong, S. Yoo, O. Mutlu, and K. Choi. A scalable processing-in-memory accelerator for parallel graph processing. *ACM SIGARCH Computer Architecture News*, 43(3):105–117, 2016.
- [6] K. J. Ahn and S. Guha. Linear programming in the semi-streaming model with application to the maximum matching problem. In *ICALP*, 2011.
- [7] K. J. Ahn, S. Guha, and A. McGregor. Analyzing graph structure via linear measurements. In *Proceedings of the twenty-third annual ACM-SIAM symposium on Discrete Algorithms*, pages 459–467. Society for Industrial and Applied Mathematics, 2012.
- [8] K. J. Ahn, S. Guha, and A. McGregor. Graph sketches: sparsification, spanners, and subgraphs. In *PODS*, 2012.
- [9] D. Andrews, D. Niehaus, and P. Ashenden. Programming models for hybrid CPU/FPGA chips. *Computer*, 37(1):118–120, 2004.
- [10] D. Andrews, D. Niehaus, R. Jidin, M. Finley, W. Peck, M. Frisbie, J. Ortiz, E. Komp, and P. Ashenden. Programming models for hybrid fpga-cpu computational components: a missing link. *IEEE micro*, 24(4):42–53, 2004.
- [11] S. Arora, E. Hazan, and S. Kale. The multiplicative weights update method: a meta-algorithm and applications. *Theory of Computing*, 8(1):121–164, 2012.
- [12] S. Assadi, S. Khanna, Y. Li, and G. Yaroslavtsev. Maximum matchings in dynamic graph streams and the simultaneous communication model. In *Proceedings of the Twenty-Seventh Annual ACM-SIAM Symposium on Discrete Algorithms*, pages 1345–1364. Society for Industrial and Applied Mathematics, 2016.
- [13] O. G. Attia, T. Johnson, K. Townsend, P. Jones, and J. Zambreno. Cygraph: A reconfigurable architecture for parallel breadth-first search. In *Parallel & Distributed Processing Symposium Workshops (IPDPSW), 2014 IEEE International*, pages 228–235. IEEE, 2014.
- [14] A. Bar-Noy, R. Bar-Yehuda, A. Freund, J. Naor, and B. Schieber. A unified approach to approximating resource allocation and scheduling. *Journal of the ACM (JACM)*, 48(5):1069–1090, 2001.
- [15] R. Bar-Yehuda, K. Bendel, A. Freund, and D. Rawitz. Local ratio: A unified framework for approximation algorithms. in memoriam: Shimon even 1935-2004. *ACM Computing Surveys (CSUR)*, 36(4):422–463, 2004.
- [16] R. Bar-Yehuda and S. Even. A local-ratio theorem for approximating the weighted vertex cover problem. *North-Holland Mathematics Studies*, 109:27–45, 1985.
- [17] T. Ben-Nun, M. Besta, S. Huber, A. N. Ziogas, D. Peter, and T. Hoefler. A modular benchmarking infrastructure for high-performance and reproducible deep learning. *arXiv preprint arXiv:1901.10183*, 2019.
- [18] M. Besta et al. Slim graph: Practical lossy graph compression for approximate graph processing, storage, and analytics. 2019.
- [19] M. Besta, M. Fischer, T. Ben-Nun, J. D. F. Licht, and T. Hoefler. Substream-Centric Maximum Matchings on FPGA. Feb. 2019. In *Proceedings of the 27th ACM/SIGDA International Symposium on Field-Programmable Gate Arrays*.
- [20] M. Besta, M. Fischer, V. Kalavri, M. Kapralov, and T. Hoefler. Practice of streaming and dynamic graphs: Concepts, models, systems, and parallelism, 2019.
- [21] M. Besta, S. M. Hassan, S. Yalamanchili, R. Ausavarungnirun, O. Mutlu, and T. Hoefler. Slim noc: A low-diameter on-chip network topology for high energy efficiency and scalability. In *Proceedings of the Twenty-Third International Conference on Architectural Support for Programming Languages and Operating Systems*, pages 43–55. ACM, 2018.
- [22] M. Besta and T. Hoefler. Fault tolerance for remote memory access programming models. In *ACM HPDC*, pages 37–48, 2014.
- [23] M. Besta and T. Hoefler. Accelerating irregular computations with hardware transactional memory and active messages. In *Proceedings of the 24th International Symposium on High-Performance Parallel and Distributed Computing*, pages 161–172. ACM, 2015.
- [24] M. Besta and T. Hoefler. Active access: A mechanism for high-performance distributed data-centric computations. In *ACM ICS*, 2015.
- [25] M. Besta and T. Hoefler. Survey and taxonomy of lossless graph compression and space-efficient graph representations. *arXiv preprint arXiv:1806.01799*, 2018.
- [26] M. Besta, R. Kanakagiri, H. Mustafa, M. Karasikov, G. Rättsch, T. Hoefler, and E. Solomonik. Communication-efficient jaccard similarity for high-performance distributed genome comparisons. *arXiv preprint arXiv:1911.04200*, 2019.
- [27] M. Besta, F. Marending, E. Solomonik, and T. Hoefler. Slimsell: A vectorizable graph representation for breadth-first search. In *Proc. IEEE IPDPS*, volume 17, 2017.
- [28] M. Besta, E. Peter, R. Gerstenberger, M. Fischer, M. Podstawski, C. Barthels, G. Alonso, and T. Hoefler. Demystifying graph databases: Analysis and taxonomy of data organization, system designs, and graph queries. *arXiv preprint arXiv:1910.09017*, 2019.
- [29] M. Besta, M. Podstawski, L. Groner, E. Solomonik, and T. Hoefler. To push or to pull: On reducing communication and synchronization in graph computations. In *Proceedings of the 26th International Symposium on High-Performance Parallel and Distributed Computing*, pages 93–104. ACM, 2017.
- [30] M. Besta, D. Stanojevic, J. D. F. Licht, T. Ben-Nun, and T. Hoefler. Graph processing on FPGAs: Taxonomy, survey, challenges. *arXiv preprint arXiv:1903.06697*, 2019.
- [31] M. Besta, D. Stanojevic, T. Zivic, J. Singh, M. Hoerold, and T. Hoefler. Log (graph): a near-optimal high-performance graph representation. In *Proceedings of the 27th International Conference on Parallel Architectures and Compilation Techniques*, page 7. ACM, 2018.
- [32] B. Betkaoui, D. B. Thomas, W. Luk, and N. Potherszujl. A framework for FPGA acceleration of large graph problems: Graphlet counting case study. In *FPT*, 2011.
- [33] B. Betkaoui, Y. Wang, D. B. Thomas, and W. Luk. Parallel FPGA-based all pairs shortest paths for sparse networks: A human brain connectome case study. In *22nd International Conference on Field Programmable Logic and Applications (FPL)*, pages 99–104, Aug 2012.
- [34] B. Betkaoui, Y. Wang, D. B. Thomas, and W. Luk. A reconfigurable computing approach for efficient and scalable parallel graph exploration. In *Application-Specific Systems, Architectures and Processors (ASAP), 2012 IEEE 23rd International Conference on*, pages 8–15. IEEE, 2012.
- [35] J. A. Bondy, U. S. R. Murty, et al. *Graph theory with applications*. 1976.
- [36] L. S. Burkol, G. Frahling, S. Leonardi, A. Marchetti-Spaccamela, and C. Sohler. Counting triangles in data streams. In *PODS*, 2006.
- [37] A. Chakrabarti, G. Cormode, and A. McGregor. Annotations in data streams. In *ICALP*, 2009.
- [38] Y.-W. Chang, J.-M. Lin, and D. Wong. Graph matching-based algorithms for FPGA segmentation design. In *ICCAD*, 1998.
- [39] A. Ching, S. Edunov, M. Kabiljo, D. Logothetis, and S. Muthukrishnan. One trillion edges: Graph processing at Facebook-scale. *VLDB*, 2015.
- [40] R. Chitnis, G. Cormode, H. Esfandiari, M. Hajiaghayi, A. McGregor, M. Monemizadeh, and S. Vorotnikova. Kernelization via sampling with applications to finding matchings and related problems in dynamic graph streams. In *Proceedings of the Twenty-Seventh Annual ACM-SIAM Symposium on Discrete Algorithms*, pages 1326–1344. Society for Industrial and Applied Mathematics, 2016.

- [41] Y.-k. Choi, J. Cong, Z. Fang, Y. Hao, G. Reinman, and P. Wei. A quantitative analysis on microarchitectures of modern CPU-FPGA platforms. In *DAC*, 2016.
- [42] T. H. Cormen. *Introduction to algorithms*. MIT press, 2009.
- [43] T. H. Cormen, C. Stein, R. L. Rivest, and C. E. Leiserson. *Introduction to Algorithms*. McGraw-Hill Higher Education, 2nd edition, 2001.
- [44] G. Cormode, J. Dark, and C. Konrad. Independent sets in vertex-arrival streams. *arXiv:1807.08331*, 2018.
- [45] G. Cormode, H. Jowhari, M. Monemizadeh, and S. Muthukrishnan. The sparse awakens: Streaming algorithms for matching size estimation in sparse graphs. *arXiv preprint arXiv:1608.03118*, 2016.
- [46] M. Crouch and D. M. Stubbs. Improved streaming algorithms for weighted matching, via unweighted matching. In *LIPICs-Leibniz Inf.*, 2014.
- [47] G. Dai, Y. Chi, Y. Wang, and H. Yang. FPGP: Graph Processing Framework on FPGA. In *FPGA*, 2016.
- [48] G. Dai, T. Huang, Y. Chi, N. Xu, Y. Wang, and H. Yang. ForeGraph: Exploring large-scale graph processing on multi-FPGA architecture. In *FPGA*, 2017.
- [49] M. Datar, A. Gionis, S. Indyk, and R. Motwani. Maintaining stream statistics over sliding windows. *SIAM J. on Comp.*, 2002.
- [50] J. de Fine Licht, S. Meierhans, and T. Hoefler. Transformations of high-level synthesis codes for high-performance computing. *arXiv:1805.08288*, 2018.
- [51] J. Dean and S. Ghemawat. Mapreduce: simplified data processing on large clusters. *Communications of the ACM*, 51(1):107–113, 2008.
- [52] C. Demetrescu, I. Finocchi, and A. Ribichini. Trading off space for passes in graph streaming problems. *TALG*, 2009.
- [53] S. Di Girolamo, K. Taranov, A. Kurth, M. Schaffner, T. Schneider, J. Beránek, M. Besta, L. Benini, D. Roweth, and T. Hoefler. Network-accelerated non-contiguous memory transfers. *arXiv preprint arXiv:1908.08590*, 2019.
- [54] W. J. Dixon and F. J. Massey Jr. *Introduction to statistical analysis*. 1957.
- [55] R. Dorrance, F. Ren, and D. Marković. A scalable sparse matrix-vector multiplication kernel for energy-efficient sparse-blas on fpgas. In *Proceedings of the 2014 ACM/SIGDA international symposium on Field-programmable gate arrays*, pages 161–170. ACM, 2014.
- [56] H. ElGindy and Y.-L. Shue. On sparse matrix-vector multiplication with fpga-based system. In *Proceedings. 10th Annual IEEE Symposium on Field-Programmable Custom Computing Machines*, pages 273–274. IEEE, 2002.
- [57] N. Engelhardt and H. K.-H. So. Grafv: A vertex-centric distributed graph processing framework on FPGAs. In *FPL*, 2016.
- [58] N. Engelhardt and H. K.-H. So. Vertex-centric Graph Processing on FPGA. In *FCCM*, 2016.
- [59] L. Epstein, A. Levin, J. Mestre, and D. Segev. Improved approximation guarantees for weighted matching in the semi-streaming model. *SIAM Journal on Discrete Mathematics*, 25(3):1251–1265, 2011.
- [60] J. Feigenbaum, S. Kannan, A. McGregor, S. Suri, and J. Zhang. On graph problems in a semi-streaming model. *Theoretical CS*, 2005.
- [61] R. Gerstenberger, M. Besta, and T. Hoefler. Enabling highly-scalable remote memory access programming with mpi-3 one sided. *Scientific Programming*, 22(2):75–91, 2014.
- [62] M. Ghaffari. Space-optimal semi-streaming for  $(2 + \epsilon)$ -approximate matching. *arXiv:1701.03730*, 2017.
- [63] L. Gianinazzi, P. Kalvoda, A. De Palma, M. Besta, and T. Hoefler. Communication-avoiding parallel minimum cuts and connected components. In *Proceedings of the 23rd ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming*, pages 219–232. ACM, 2018.
- [64] A. Goel, M. Kapralov, and S. Khanna. On the communication and streaming complexity of maximum bipartite matching. In *Proceedings of the twenty-third annual ACM-SIAM symposium on Discrete Algorithms*, pages 468–485. SIAM, 2012.
- [65] J. E. Gonzalez, Y. Low, H. Gu, D. Bickson, and C. Guestrin. Powergraph: distributed graph-parallel computation on natural graphs. In *OSDI*, volume 12, page 2, 2012.
- [66] E. Grigorescu, M. Monemizadeh, and S. Zhou. Streaming weighted matchings: Optimal meets greedy. *arXiv:1608.01487*, 2016.
- [67] T. J. Harris. A survey of PRAM simulation techniques. *ACM Computing Surveys (CSUR)*, 26(2):187–206, 1994.
- [68] M. R. Henzinger, P. Raghavan, and S. Rajagopalan. Computing on data streams. *External Mem. Alg.*, 1998.
- [69] T. Hoefler and R. Belli. Scientific benchmarking of parallel computing systems: twelve ways to tell the masses when reporting performance results. In *Proceedings of the international conference for high performance computing, networking, storage and analysis*, page 73. ACM, 2015.
- [70] R. Inta, D. J. Bowman, and S. M. Scott. The Chimera: an off-the-shelf CPU/GPGPU/FPGA hybrid computing platform. *International Journal of Reconfigurable Computing*, 2012:2, 2012.
- [71] Intel. Intel Core i7-8700K Processor, 2017.
- [72] Intel. Intel Xeon Processor E5-2680 v4, 2017.
- [73] Intel. Stratix 10 GX/SX Device Overview, 2017.
- [74] Intel Arria. Intel Arria 10 Device Overview, 2017.
- [75] R. Jidin. Extending the Thread Programming Model Across Hybrid FPGA/CPU Architectures. *Information Technology and Telecommunications Center (ITTC)*, University of Kansas, 2005.
- [76] M. Kapralov. Better bounds for matchings in the streaming model. In *Proceedings of the twenty-fourth annual ACM-SIAM symposium on Discrete algorithms*, pages 1679–1697. SIAM, 2013.
- [77] M. Kapralov, S. Khanna, and M. Sudan. Approximating matching size from random streams. In *Proceedings of the twenty-fifth annual ACM-SIAM symposium on Discrete algorithms*, pages 734–751. SIAM, 2014.
- [78] N. Kapre. Custom FPGA-based soft-processors for sparse graph acceleration. In *ASAP*, 2015.
- [79] N. Kapre, N. Mehta, D. Rizzo, I. Eslick, R. Rubin, T. E. Uribe, F. Thomas Jr, A. DeHon, et al. Graphstep: A system architecture for sparse-graph algorithms. In *Field-Programmable Custom Computing Machines, 2006. FCCM'06. 14th Annual IEEE Symposium on*, pages 143–151. IEEE, 2006.
- [80] C. Karande, A. Mehta, and P. Tripathi. Online bipartite matching with unknown distributions. In *STOC*, 2011.
- [81] R. M. Karp, U. V. Vazirani, and V. V. Vazirani. An optimal algorithm for on-line bipartite matching. In *Proceedings of the twenty-second annual ACM symposium on Theory of computing*, pages 352–358. ACM, 1990.
- [82] J. Kepner, P. Aaltonen, D. Bader, A. Buluç, F. Franchetti, J. Gilbert, D. Hutchison, M. Kumar, A. Lumsdaine, H. Meyerhenke, et al. Mathematical foundations of the graphblas. *arXiv preprint arXiv:1606.05790*, 2016.
- [83] A. Khan. Vertex-centric graph processing: The good, the bad, and the ugly. *arXiv preprint arXiv:1612.07404*, 2016.
- [84] S. Khoram, J. Zhang, M. Strange, and J. Li. Accelerating graph analytics by co-optimizing storage and access on an FPGA-hmc platform. In *Proceedings of the 2018 ACM/SIGDA International Symposium on Field-Programmable Gate Arrays*, pages 239–248. ACM, 2018.

- [85] KONECT. Konec network dataset, 2017.
- [86] C. Konrad, F. Magniez, and C. Mathieu. Maximum matching in semi-streaming with few passes. *Approximation, Randomization, and Combinatorial Optimization. Algorithms and Techniques*, pages 231–242, 2012.
- [87] G. Kwasniewski, M. Kabić, M. Besta, J. VandeVondele, R. Solcà, and T. Hoefler. Red-blue pebbling revisited: near optimal parallel matrix-matrix multiplication. In *ACM/IEEE Supercomputing*, page 24. ACM, 2019.
- [88] A. Kyrola, G. E. Blelloch, and C. Guestrin. Graphchi: Large-scale graph computation on just a pc. USENIX, 2012.
- [89] J. Lee, H. Kim, S. Yoo, K. Choi, H. P. Hofstee, G.-J. Nam, M. R. Nutter, and D. Jamsek. Extrav: boosting graph processing near storage with a coherent accelerator. *Proceedings of the VLDB Endowment*, 10(12):1706–1717, 2017.
- [90] G. Lei, Y. Dou, R. Li, and F. Xia. An FPGA implementation for solving the large single-source-shortest-path problem. *IEEE Transactions on Circuits and Systems II: Express Briefs*, 63(5):473–477, 2016.
- [91] J. Leskovec and A. Krevl. SNAP Datasets: Stanford large network dataset collection, 2014.
- [92] J. d. F. Licht, G. Kwasniewski, and T. Hoefler. Flexible communication avoiding matrix multiplication on fpga with high-level synthesis. *arXiv preprint arXiv:1912.06526*, 2019.
- [93] H. Liu and P. Singh. Conceptnet—a practical commonsense reasoning tool-kit. *BT technology journal*, 22(4):211–226, 2004.
- [94] Y. Low, J. Gonzalez, A. Kyrola, D. Bickson, C. Guestrin, and J. M. Hellerstein. Graphlab: A new framework for parallel machine learning. *preprint arXiv:1006.4990*, 2010.
- [95] A. Lumsdaine, D. Gregor, B. Hendrickson, and J. Berry. Challenges in Parallel Graph Processing. *Par. Proc. Let.*, 2007.
- [96] X. Ma, D. Zhang, and D. Chiou. FPGA-accelerated transactional execution of graph workloads. In *Proceedings of the 2017 ACM/SIGDA International Symposium on Field-Programmable Gate Arrays*, pages 227–236. ACM, 2017.
- [97] G. Malewicz, M. H. Austern, A. J. Bik, J. C. Dehnert, I. Horn, N. Leiser, and G. Czajkowski. Pregel: a system for large-scale graph processing. In *Proceedings of the 2010 ACM SIGMOD International Conference on Management of data*, pages 135–146. ACM, 2010.
- [98] A. McGregor. Finding graph matchings in data streams. In *APPROX-RANDOM*, volume 3624, pages 170–181. Springer, 2005.
- [99] A. McGregor and S. Vorotnikova. Planar matching in streams revisited. In *LIPICs-Leibniz International Proceedings in Informatics*, volume 60. Schloss Dagstuhl-Leibniz-Zentrum fuer Informatik, 2016.
- [100] A. McGregor and S. Vorotnikova. A simple, space-efficient, streaming algorithm for matchings in low arboricity graphs. In *OASIS-Open Access Series in Informatics*, volume 61. Schloss Dagstuhl-Leibniz-Zentrum fuer Informatik, 2018.
- [101] A. McGregor, S. Vorotnikova, and H. T. Vu. Better algorithms for counting triangles in data streams. In *PODS*, 2016.
- [102] F. McSherry, M. Isard, and D. G. Murray. Scalability! but at what COST? In *HotOS*, 2015.
- [103] S. Muthukrishnan. Data streams: Algorithms and applications. *Foundations and Trends® in Theoretical Computer Science*, 1(2):117–236, 2005.
- [104] M. E. Newman. A measure of betweenness centrality based on random walks. *Social networks*, 27(1):39–54, 2005.
- [105] E. Nurvitadhi, G. Weisz, Y. Wang, S. Hurkat, M. Nguyen, J. C. Hoe, J. F. Martinez, and C. Guestrin. Graphgen: An FPGA framework for vertex-centric graph computation. In *Field-Programmable Custom Computing Machines (FCCM), 2014 IEEE 22nd Annual International Symposium on*, pages 25–28. IEEE, 2014.
- [106] Nvidia. GEFORCE GTX 1080 Ti, 2017.
- [107] T. Oguntebi and K. Olukotun. Graphops: A dataflow library for graph analytics acceleration. In *FPGA*, 2016.
- [108] N. Oliver, R. R. Sharma, S. Chang, B. Chitlur, E. Garcia, J. Grecco, A. Grier, N. Ijith, Y. Liu, P. Marolia, et al. A reconfigurable computing system based on a cache-coherent fabric. In *ReConFig*, 2011.
- [109] M. Owaida, D. Sidler, K. Kara, and G. Alonso. Centaur: A framework for hybrid CPU-FPGA databases. In *FCCM*, 2017.
- [110] M. M. Ozdal, S. Yesil, T. Kim, A. Ayupov, J. Greth, S. Burns, and O. Ozturk. Energy efficient architecture for graph analytics accelerators. In *Computer Architecture (ISCA), 2016 ACM/IEEE 43rd Annual International Symposium on*, pages 166–177. IEEE, 2016.
- [111] L. Page, S. Brin, R. Motwani, and T. Winograd. The pagerank citation ranking: Bringing order to the web. Technical report, Stanford InfoLab, 1999.
- [112] C. H. Papadimitriou and K. Steiglitz. *Combinatorial optimization: algorithms and complexity*. Courier Corporation, 1998.
- [113] A. Paz and G. Schwartzman. A  $(2+\epsilon)$ -approximation for maximum weight matching in the semi-streaming model. In *Proceedings of the Twenty-Eighth Annual ACM-SIAM Symposium on Discrete Algorithms*, pages 2153–2161. SIAM, 2017.
- [114] A. Putnam, D. Bennett, E. Dellinger, J. Mason, P. Sundararajan, and S. Eggers. CHiMPS: A C-level compilation flow for hybrid CPU-FPGA architectures. In *2008 International Conference on Field Programmable Logic and Applications*, pages 173–178. IEEE, 2008.
- [115] A. Roy, I. Mihailovic, and W. Zwaenepoel. X-stream: Edge-centric graph processing using streaming partitions. In *Proceedings of the Twenty-Fourth ACM Symposium on Operating Systems Principles*, pages 472–488. ACM, 2013.
- [116] S. Salihoglu and J. Widom. Optimizing graph algorithms on Pregel-like systems. *VLDB*, 2014.
- [117] M. Santarini. Zynq-7000 EPP sets stage for new era of innovations. *Xcell*, 2011.
- [118] T. Schank. Algorithmic aspects of triangle-based network analysis. 2007.
- [119] P. Schmid, M. Besta, and T. Hoefler. High-performance distributed rma locks. In *Proceedings of the 25th ACM International Symposium on High-Performance Parallel and Distributed Computing*, pages 19–30. ACM, 2016.
- [120] H. Schweizer, M. Besta, and T. Hoefler. Evaluating the cost of atomic operations on modern architectures. In *2015 International Conference on Parallel Architecture and Compilation (PACT)*, pages 445–456. IEEE, 2015.
- [121] L. Shang, A. S. Kaviani, and K. Bathala. Dynamic power consumption in Virtex™-II FPGA family. In *FPGA*, 2002.
- [122] Y. Shiloach and U. Vishkin. An  $O(\log n)$  parallel connectivity algorithm. Technical report, Computer Science Department, Technion, 1980.
- [123] D. Sidler, Z. István, M. Owaida, and G. Alonso. Accelerating pattern matching queries in hybrid CPU-FPGA architectures. In *Proceedings of the 2017 ACM International Conference on Management of Data*, pages 403–415. ACM, 2017.
- [124] Y. Simmhan, A. Kumbhare, C. Wickramaarachchi, S. Nagarkar, S. Ravi, C. Raghavendra, and V. Prasanna. Goffish: A sub-graph centric framework for large-scale graph analytics. In *EuroPar*, 2014.
- [125] E. Solomonik, M. Besta, F. Vella, and T. Hoefler. Scaling betweenness centrality using communication-efficient sparse matrix multiplication. In *Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis*, page 47. ACM, 2017.
- [126] J. Sun, G. Peterson, and O. Storaasli. Sparse matrix-vector multiplication design on fpgas. In *15th Annual IEEE Symposium on Field-Programmable Custom Computing Machines (FCCM 2007)*, pages 349–352. IEEE, 2007.
- [127] J. Sun, N.-N. Zheng, and H.-Y. Shum. Stereo matching using belief propagation. *IEEE Transactions on pattern analysis and machine intelligence*, 25(7):787–800, 2003.

- [128] R. Szeliski, R. Zabih, D. Scharstein, O. Veksler, V. Kolmogorov, A. Agarwala, M. Tappen, and C. Rother. A comparative study of energy minimization methods for markov random fields with smoothness-based priors. *IEEE transactions on pattern analysis and machine intelligence*, 30(6):1068–1080, 2008.
- [129] A. Tate, A. Kamil, A. Dubey, A. Größlinger, B. Chamberlain, B. Goglin, C. Edwards, C. J. Newburn, D. Padua, D. Unat, et al. Programming abstractions for data locality. PADAL Workshop 2014, April 28–29, Swiss National Supercomputing Center . . . , 2014.
- [130] N. Trinajstić, D. J. Klein, and M. Randić. On some solved and unsolved problems of chemical graph theory. *International Journal of Quantum Chemistry*, 1986.
- [131] J. Tyhach, M. Hutton, S. Atsatt, A. Rahman, B. Vest, D. Lewis, M. Langhammer, S. Shumarayev, T. Hoang, A. Chan, et al. Arria™ 10 device architecture. In *CICC*, 2015.
- [132] R. Uehara and Z.-Z. Chen. Parallel approximation algorithms for maximum weighted matching in general graphs. *Information Processing Letters*, 76(1-2):13–17, 2000.
- [133] Y. Umuroglu, D. Morrison, and M. Jahre. Hybrid breadth-first search on a single-chip FPGA-CPU heterogeneous platform. In *FPL*, 2015.
- [134] X. Wang and S. G. Ziavras. Performance-energy tradeoffs for matrix multiplication on fpga-based mixed-mode chip multiprocessors. In *8th International Symposium on Quality Electronic Design (ISQED'07)*, pages 386–391. IEEE, 2007.
- [135] G. Weisz, E. Nurvitadhi, and J. Hoe. Graphgen for coram: Graph computation on FPGAs. In *CARL*, 2013.
- [136] D. Yan, J. Cheng, K. Xing, Y. Lu, W. Ng, and Y. Bu. Pregel algorithms for graph connectivity problems with performance guarantees. *Proceedings of the VLDB Endowment*, 7(14):1821–1832, 2014.
- [137] C. Yang. An efficient dispatcher for large scale graphprocessing on opencl-based FPGAs. *arXiv preprint arXiv:1806.11509*, 2018.
- [138] P. Yao. An efficient graph accelerator with parallel data conflict management. *arXiv preprint arXiv:1806.00751*, 2018.
- [139] M. Zaharia, R. S. Xin, P. Wendell, T. Das, M. Armbrust, A. Dave, X. Meng, J. Rosen, S. Venkataraman, M. J. Franklin, et al. Apache Spark: a unified engine for big data processing. *CACM*, 2016.
- [140] M. Zelke. Weighted matching in the semi-streaming model. *Algorithmica*, 62(1-2):1–20, 2012.
- [141] J. Zhang, S. Khoram, and J. Li. Boosting the performance of FPGA-based graph processor using hybrid memory cube: A case for breadth first search. In *FPGA*, 2017.
- [142] J. Zhang, S. Khoram, and J. Li. Boosting the performance of FPGA-based graph processor using hybrid memory cube: A case for breadth first search. In *Proceedings of the 2017 ACM/SIGDA International Symposium on Field-Programmable Gate Arrays, FPGA '17*, pages 207–216, New York, NY, USA, 2017. ACM.
- [143] J. Zhang and J. Li. Degree-aware hybrid graph traversal on FPGA-hmc platform. In *Proceedings of the 2018 ACM/SIGDA International Symposium on Field-Programmable Gate Arrays*, pages 229–238. ACM, 2018.
- [144] S. Zhang, Z. Du, L. Zhang, H. Lan, S. Liu, L. Li, Q. Guo, T. Chen, and Y. Chen. Cambricon-x: An accelerator for sparse neural networks. In *2016 49th Annual IEEE/ACM International Symposium on Microarchitecture (MICRO)*, pages 1–12, Oct 2016.
- [145] J. Zhou, S. Liu, Q. Guo, X. Zhou, T. Zhi, D. Liu, C. Wang, X. Zhou, Y. Chen, and T. Chen. Tunao: A high-performance and energy-efficient reconfigurable accelerator for graph processing. In *CCGRID*, 2017.
- [146] S. Zhou, C. Chelmiss, and V. K. Prasanna. Optimizing memory performance for FPGA implementation of pagerank. In *ReConFig*, pages 1–6, 2015.
- [147] S. Zhou, C. Chelmiss, and V. K. Prasanna. High-throughput and energy-efficient graph processing on FPGA. In *Field-Programmable Custom Computing Machines (FCCM), 2016 IEEE 24th Annual International Symposium on*, pages 103–110. IEEE, 2016.
- [148] S. Zhou, R. Kannan, H. Zeng, and V. K. Prasanna. An FPGA framework for edge-centric graph processing. In *Proceedings of the 15th ACM International Conference on Computing Frontiers*, pages 69–77. ACM, 2018.
- [149] S. Zhou and V. K. Prasanna. Accelerating Graph Analytics on CPU-FPGA Heterogeneous Platform. *SBAC-PAD*, 2017.
- [150] J. Zhu, I. Sander, and A. Jantsch. Buffer minimization of real-time streaming applications scheduling on hybrid CPU/FPGA architectures. In *Proceedings of the Conference on Design, Automation and Test in Europe*, pages 1506–1511. European Design and Automation Association, 2009.
- [151] L. Zhuo and V. K. Prasanna. Sparse matrix-vector multiplication on fpgas. In *Proceedings of the 2005 ACM/SIGDA 13th international symposium on Field-programmable gate arrays*, pages 63–74. ACM, 2005.