



# Parallelising Control Flow in Dynamic-scheduling High-level Synthesis

JIANYI CHENG, Imperial College London, United Kingdom

LANA JOSIPOVIĆ, ETH Zürich, Switzerland

JOHN WICKERSON and GEORGE A. CONSTANTINIDES, Imperial College London, United Kingdom

Recently, there is a trend to use high-level synthesis (HLS) tools to generate dynamically scheduled hardware. The generated hardware is made up of components connected using handshake signals. These handshake signals schedule the components at runtime when inputs become available. Such approaches promise superior performance on “irregular” source programs, such as those whose control flow depends on input data. This is at the cost of additional area. Current dynamic scheduling techniques are well able to exploit parallelism among instructions *within* each basic block (BB) of the source program, but parallelism *between* BBs is under-explored, due to the complexity in runtime control flows and memory dependencies. Existing tools allow some of the operations of different BBs to overlap, but to simplify the analysis required at compile time they require the BBs to *start* in strict program order, thus limiting the achievable parallelism and overall performance.

We formulate a general dependency model suitable for comparing the ability of different dynamic scheduling approaches to extract maximal parallelism at runtime. Using this model, we explore a variety of mechanisms for runtime scheduling, incorporating and generalising existing approaches. In particular, we precisely identify the restrictions in existing scheduling implementation and define possible optimisation solutions. We identify two particularly promising examples where the compile-time overhead is small and the area overhead is minimal and yet we are able to significantly speed up execution time: (1) parallelising consecutive independent loops; and (2) parallelising independent inner-loop instances in a nested loop as individual threads. Using benchmark sets from related works, we compare our proposed toolflow against a state-of-the-art dynamic-scheduling HLS tool called Dynamatic. Our results show that, on average, our toolflow yields a 4× speedup from (1) and a 2.9× speedup from (2), with a negligible area overhead. This increases to a 14.3× average speedup when combining (1) and (2).

CCS Concepts: • **Hardware** → **High-level and register-transfer level synthesis; Modeling and parameter extraction;**

Additional Key Words and Phrases: FPGA, high-level synthesis, dynamic scheduling, static analysis

## ACM Reference format:

Jianyi Cheng, Lana Josipović, John Wickerson, and George A. Constantinides. 2023. Parallelising Control Flow in Dynamic-scheduling High-level Synthesis. *ACM Trans. Reconfig. Technol. Syst.* 16, 4, Article 55 (September 2023), 32 pages.

<https://doi.org/10.1145/3599973>

This work is supported by the EPSRC (EP/P010040/1, EP/R006865/1).

Authors' addresses: J. Cheng, J. Wickerson, and G. A. Constantinides, Imperial College London, London, United Kingdom; emails: [jianyi.cheng17@imperial.ac.uk](mailto:jianyi.cheng17@imperial.ac.uk), [j.wickerson@imperial.ac.uk](mailto:j.wickerson@imperial.ac.uk), [g.constantinides@imperial.ac.uk](mailto:g.constantinides@imperial.ac.uk); L. Josipović, ETH Zürich, Zürich, Switzerland; email: [ljospovic@ethz.ch](mailto:ljospovic@ethz.ch).

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner/author(s).

© 2023 Copyright held by the owner/author(s).

1936-7406/2023/09-ART55 \$15.00

<https://doi.org/10.1145/3599973>

## 1 INTRODUCTION

FPGAs are now widely used as a reconfigurable device for custom high-performance computing, such as in datacentres including Microsoft Project Catapult [43] and Amazon EC2 F1 instances [1]. However, users need to understand low-level hardware details for directly programming on FPGAs. To lift this restriction for software engineers, **high-level synthesis (HLS)** tools automatically translate a software program in a high-level software language, such as C, into low-level hardware descriptions. This could also significantly reduce the design effort compared to manual **register transfer level (RTL)** implementation. Today, various HLS tools have been developed in both academia, including LegUp from the University of Toronto [7], Bambu from the Politecnico di Milano [9] and Dynamatic from EPFL [33], and industry, including Intel HLS compiler [31], Xilinx Vivado HLS [51], Cadence Stratus HLS [46], and Siemens Catapult HLS [10].

A central step of the HLS process is scheduling, which maps each operation in the input program to a clock cycle. This mapping can be decided either at compile time (statically) or at runtime (dynamically). There has been recent interest in dynamic scheduling, because it enables the hardware to adapt its behaviour at runtime to particular input values, memory access patterns, and control flow decisions. Therefore, it potentially achieves better performance compared to the conservative schedule produced by static analysis.

Dynamic-scheduling HLS tools, such as Dynamatic [33], transform a sequential program into a circuit made up of components that are connected by handshaking signals. Each component can start as soon as all of its inputs are ready. Although these tools aim to allow out-of-order execution as much as possible, they must take care to respect dependencies in the source program. There are two kinds of dependencies: memory dependencies (i.e., dependency via a memory location) and data dependencies (i.e., dependency via a program variable). There are also two scopes of dependency: between instructions in the same **basic block (BB)** and between instructions in different BBs. This leads to four cases to consider:

- (1) **Intra-BB data dependencies:** these can be respected by placing handshaking connections between the corresponding hardware operations in the circuit.
- (2) **Intra-BB memory dependencies:** these can be kept in the original program order using elastic components named *load-store queues (LSQs)* [32]. An LSQ is a hardware component that schedules memory operations at runtime.
- (3) **Inter-BB data dependencies:** these can be respected using handshaking connections, as in (1), and additionally by starting BBs in strict program order, so the inputs of each BB are accepted in program order [35].
- (4) **Inter-BB memory dependencies:** these can be respected by starting BBs in strict program order and using an LSQ, as in (2).

In all cases, existing dynamic-scheduling HLS tools well exploit parallelism for cases (1) and (2) above. This allows out-of-order execution *within* a BB but requires different BBs to start in order, even when some BBs are independent and could start in parallel. This, naturally, leads to missed opportunities for performance improvements.

The existing dependency model only analyses intra-BB dependencies, which exploits parallelism at the data flow level. The inter-BB dependencies are resolved by maintaining the original program order. The order is preserved by sequentially starting BB execution even though these BBs are executing in parallel, as shown in Figure 1(b). Analysis for inter-BB dependencies is limited. In this article, we propose a dependency model that formalises all four cases above. As a demonstration of applications, we use the dependency model to focus on further exploiting parallelism for cases (3) and (4). We find BBs that can be started out-of-order or simultaneously, as shown in Figures 1(c) and 1(d), and use static analysis (powered by the Microsoft Boogie verification engine [36]) to

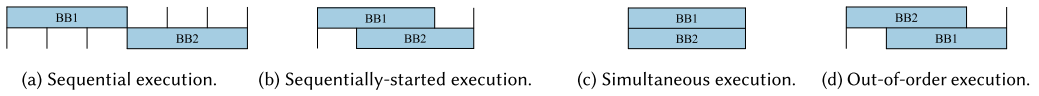


Fig. 1. Basic block schedules. Assume BB1 executes before BB2 in the original program order. Dynamic only supports (a) and (b). We show how to support (c) and (d) using analysis in our proposed model.

ensure that inter-BB dependencies are still respected. We then achieve a parallel BB schedule by two techniques: (1) parallelising independent consecutive loops (inter-block scheduling); and (2) parallelising independent consecutive inner-loop instances in a nested loop (C-slow pipelining). Our main contributions include:

- a general dependency model that formalises both data flow dependencies and control flow dependencies for dynamic-scheduling high-level synthesis;
- a technique that automatically identifies the absence of dependencies between consecutive loops using the Microsoft Boogie verifier for parallelism;
- a technique that automatically identifies the absence of dependencies between consecutive iterations of a loop using the Microsoft Boogie for C-slow pipelining; and
- results and analysis showing that our techniques, compared to original Dynamatic, achieve, on average,  $14.3\times$  speedup with 10% area overhead.

The rest of our article is organised as follows: Section 2 introduces existing works on dynamic-scheduling HLS, parallelising **control-flow graphs (CFGs)** for HLS, and C-slow pipelining. Section 3 explains our general dependency model for dynamic-scheduling HLS. Section 4 demonstrates inter-block scheduling as an application of the proposed dependency model. Section 5 demonstrates C-slow pipelining as another application of the proposed dependency model. Section 7 evaluates the effectiveness of these two techniques.

## 2 BACKGROUND

This section first reviews related work on existing HLS tools that use dynamic scheduling. We then compare existing works on static analysis of memory dependencies for HLS with our work. Finally, we review related works on parallelising CFGs for HLS and C-slow pipelining on FPGAs.

### 2.1 Dynamic-scheduling High-level Synthesis

Most HLS tools such as Xilinx Vivado HLS [51] and Dynamatic [33] translate an input program into an **intermediate representation (IR)** such as LLVM IR and then transform the IR into a **control data flow graph (CDFG)** for scheduling [22]. A CDFG is a two-level directed graph that contains a set of vertices connected by edges. The top level is a CFG, where each vertex represents a **basic block (BB)** in the IR, and each edge represents the control flow. At the lower level, each vertex is also a **data flow graph (DFG)**, where each sub-vertex inside the DFG represents an operation in the BB, and each sub-edge represents a data dependency. The CDFG is used as part of the dependency constraints in both static and dynamic scheduling [20, 33].

In dynamic-scheduling HLS, initial work was studied by Page and Luk [28], which maps oc-cam programs into hardware and has been extended to support a commercial language named Handel-C [11]. The idea of mapping a C program into a netlist of pre-defined hardware components has been studied in both asynchronous and synchronous worlds. In the asynchronous world, Venkataramani et al. [48] propose a toolflow that maps ANSI-C programs into asynchronous hardware designs. Li et al. [39] propose a dynamic-scheduling HLS tool named Fluid that supports the synthesis of complex CFGs into asynchronous hardware designs. In the synchronous world,

Sayuri and Nagisa [44] propose a method that synthesises single-level loops into dynamically scheduled circuits. Josipović et al. [33] propose an open-sourced HLS tool named “Dynamatic” that automatically translates a program into a dynamically pipelined hardware.

Dynamatic uses pre-defined components with handshake connections formalised by Carloni et al. [8]. Each edge in the CDFG of the input program is translated to a handshake connection between components. This allows a component to execute at the earliest time when all its inputs are valid. The memory dependency is controlled by **load-store queues (LSQs)**. An LSQ exploits out-of-order memory accesses by checking memory dependency in program order at runtime [32] and early executing those independent memory accesses.

Dynamatic parallelises DFGs within and across BBs for high performance, but the CFG still starts BBs sequentially. Sequentially starting BBs is required to respect inter-BB dependencies at runtime. An unverified BB schedule may cause an error. Our toolflow uses Boogie to formally prove that the transformed BB schedule cannot break any dependency, such that the synthesised hardware is still correct.

## 2.2 Parallelising Control Flows for HLS

Dependency analysis for parallelising a CFG of a sequential program has been well-studied in the software compiler world [29]. Traditional approaches exploit BB parallelism using polyhedral analysers such as Pluto [3] and Polly [25]. These tools automatically parallelise code that contains affine memory accesses [4, 23] and have been widely used in HLS to parallelise hardware kernels [40, 41, 49, 55]. However, polyhedral analysis is not applicable when analysing irregular memory patterns such as non-affine memory accesses, which are commonly seen in applications amenable for dynamic scheduling, such as tumour detection [52] and video rendering [47].

Recently, there are works that use formal verification to prove the absence of dependency to exploit hardware parallelism. Compared with affine or polyhedral analysis, formal verification can analyse non-affine memory accesses but takes a longer time. Zhou et al. [54] propose a **satisfiability-modulo theory (SMT)**-based [24] approach to verify the absence of memory contention in banked memory among parallel kernels. Cheng et al. propose a Boogie-based approach for simplifying memory arbitration for multi-threaded hardware [13]. Microsoft Boogie [36] is an automated program verifier on top of SMT solvers. The Boogie verifier does not run a Boogie program but generates a set of specifications for verification. It uses its own intermediate verification language to describe the behaviour of a program to be verified, which can be automatically decoded into SMT queries. An SMT solver under Boogie then reasons the program behaviour, including the values that its variables may take. Our work also uses Boogie but for parallelising BBs in dynamically scheduled hardware. Boogie has its own constructs and, here, we list the ones used in this article:

- (1) `if (*) {A} else {B}` is a non-deterministic choice. The program arbitrarily does A or B.
- (2) `havoc x` assigns an arbitrary values to a variable or an array `x`, used to capture all the possible values of `x`.
- (3) `assert c` proves the condition `c` for all the values that the variables in `c` may take.

In this article, we use Boogie to verify the absence of memory dependency between different iterations of the same loop or two consecutive loops. Our approach can generate Boogie programs from arbitrary programs based on the formulation by Reference [19].

Mapping a parallel BB schedule into hardware has also been widely studied. Initial work by Cabrera et al. [5] proposes an OpenMP extension to off-load computation to an FPGA. Leow et al. [38] propose a framework that maps OpenMP code in Handel-C [11] to VHDL programs. Choi et al. [18] propose a plugin that synthesises both OpenMP and Pthreads C programs into multi-threaded

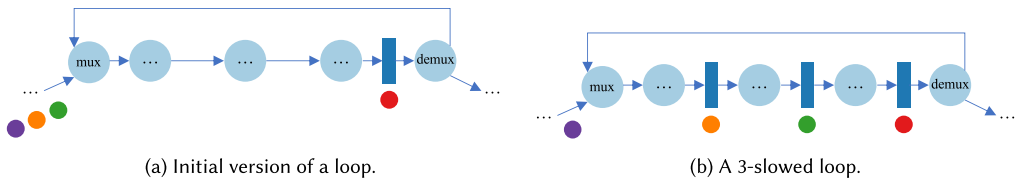


Fig. 2. An example of 3-slowing a loop. The 3-slowed loop has tripled latency, the same throughput, and one-third critical path compared to the initial version.

hardware, used in an open-sourced HLS tool named LegUp [7]. Gupta et al. propose an HLS tool named SPARK that parallelises control flow with speculation [26]. Existing commercialised HLS tools [10, 30, 46, 51] support multi-threaded hardware synthesis using manually annotated directives by users. These works either require user annotation or only use static scheduling, while our approach only uses automated dynamic scheduling.

Finally, there are works on simultaneously starting BB in dynamic-scheduling HLS. Cheng et al. [14] propose an HLS tool named DASS that allows each statically scheduled component to act as a static island in a dynamically scheduled circuit. Each island is still statically scheduled, while our toolflow only uses dynamic scheduling.

### 2.3 C-Slow Pipelining

C-slow pipelining is a technique that replaces each register in the circuit with  $C$  registers to construct  $C$  independent threads [42]. The circuit then operates as  $C$ -thread hardware while keeping one copy of resources. For instance, a stream of data enters a pipelined loop in Figure 2(a). We use **initiation interval (II)** for evaluating the hardware performance. An II of a loop is defined as the time difference in block cycles between the start of the same operation in two consecutive iterations. The loop computes with an II of 1, as illustrated by the presence of one register in the cycle. Assume the loop trip count is  $N$ , then the latency of the loop is approximately  $N$  cycles for a large  $N$ . The overall throughput of the hardware is  $1/N$ , and the critical path is the delay of the cycle. Assume that each set of data is independent of other sets. Figure 2(b) demonstrates a 3-slowed loop that is functionally equivalent to the one in Figure 2(a). There are three registers in the cycle, evenly distributed in the path. This increases the latency of the hardware to  $3N$  cycles. The loop can iterate with three sets of data in the cycle concurrently. Then, the overall throughput of the hardware is approximately  $3/(3N) = 1/N$ , and the critical path is nearly  $1/3$  of the one in Figure 2(a). A  $C$ -slow loop can have a better throughput or clock frequency to achieve approximately  $C$  times speedup.

C-slow pipelining was first proposed by Leiserson et al. for optimising the critical path of synchronous circuits [37]. Markovskiy and Patel [42] propose a C-slow based-approach to improve the throughput of a microprocessor. Weaver et al. [50] propose an automated tool that applies C-slow retiming on a class of applications for certain FPGA families. Our work brings the idea of C-slow pipelining into the dynamic HLS world. We analyse nested loops at the source level to determine  $C$  for each loop by checking the dependency between inputs to the loop and then apply hardware transformations to achieve C-slow pipelining.

## 3 DEPENDENCY MODEL FOR DYNAMIC-SCHEDULING HLS

In this section, we formalise the dependency model for dynamic-scheduling HLS and demonstrate the restriction of the state-of-the-art dynamic-scheduling HLS tool. The dependency formulation for static scheduling has been well-studied [6, 16, 53]. Here, we extend the dependency model in Reference [16] to support dynamic scheduling.

### 3.1 Scheduling Specifications

We first formulate the fundamental specifications of the dependency constraints for scheduling. For two runtime events  $x$  and  $y$ ,<sup>1</sup> we introduce the following terms:

- $d(x, y)$  denotes whether executions of  $x$  and  $y$  have dependencies,
- $x < y$  denotes whether  $x$  executes before  $y$  in strict program order,
- $t(x) \in \mathbb{N}$  denotes the start time of the execution  $x$  in clock cycles, and
- $l(x) \in \mathbb{N}$  denotes the latency of the execution  $x$  in clock cycles.

The general dependency constraint is listed as follows:

$$\forall x, y. d(x, y) \wedge x < y \Rightarrow t(x) + l(x) \leq t(y), \quad (1)$$

where  $x$  and  $y$  are the operations in the input program. The total latency of the execution is defined as  $t_T$ :

$$\forall x. t_T \geq t(x) + l(x). \quad (2)$$

The goal is to determine a schedule with a minimum  $t_T$  that must not break Constraint 1.

### 3.2 Dynamic Implementation

Now, we show how the state-of-the-art dynamic-scheduling HLS tool, Dynamatic [33], parallelises the instruction execution dynamically.

**3.2.1 Control Flow Graph.** We assume that the input program is sequential. As introduced in Section 2.1, the input program is initially lowered from C/C++ to a CFG, where each BB contains instructions in sequential order. A CFG illustrates the control flows of a program. Each vertex represents a BB, and each edge represents a control transition between two BBs. Each of these vertices corresponds to a subgraph at the lower level, known as **data flow graph (DFG)**. Each subgraph vertex represents an operation, and each edge between these vertices represents a data dependency between two operations. For a given program and its inputs, we define the following terms for the input source:

- $B = \{b_1, b_2, \dots\}$  denotes the set of all the BBs in the program, and
- $I_b = \{i_1, i_2, \dots\}$  denotes the set of all the instructions within a BB  $b$ .

The original execution order of the input program can be defined as follows:

- $E_B \subseteq B \times \mathbb{N}$  denotes the set of execution of BBs, where  $(b_h, k)$  denotes execution of BB  $h$  ( $b_h$ ) in its  $k$ th iteration,
- $< \subseteq E_B \times E_B$  denotes the original program order of BB execution, where  $(b_h, k) < (b_{h'}, k')$  denotes  $(b_h, k)$  executes before  $(b_{h'}, k')$  in strict program order,
- $E_I \subseteq \bigcup_{b \in B} I_b \times \mathbb{N}$  denotes the set of execution of instructions, and
- $<_b : I_b \times I_b$  denotes the original program order of instruction execution within a BB  $b$ .

The execution of BBs  $<$  can be dynamic, where a BB may have a different number of iterations than another BB. However, inside each BB, the execution of instructions  $<_b$  is static, where they always have the same number of iterations and execution order inside the BB, as shown in the following constraints:

$$\forall b, i, k. (b, k) \in E_B \wedge i \in I_b \Rightarrow (i, k) \in E_I, \quad (3)$$

$$\forall i, i', k. (i, k) \in E_I \wedge (i', k) \in E_I \wedge i < i' \Rightarrow (i, k) < (i', k). \quad (4)$$

<sup>1</sup>The runtime event could be the execution of an instruction, a basic block or a loop.

By combining  $<$  and  $<_b$  lexicographically, the original program order of the instruction execution can be obtained. This is used as a reference for correctness checks. A schedule being correct is defined as the execution result by this schedule is always the same as the result by sequential execution in the original program order.

**3.2.2 Dependency Constraints.** There are mainly four types of dependencies to solve as introduced in Section 1. First, the intra-BB data dependencies are represented as edges in DFGs inside BBs and can be directly mapped to *handshake signals* between hardware operations. An operation may have variable latency, depending on its inputs. Operations not connected through handshake signals are independent and can execute in parallel or out-of-order. Although an operation may have a variable latency and execute out-of-order, each data path propagates in-order data through the edges as formalised by Carloni et al. [8]. Such a preserved data order inside each BB preserves the intra-BB data dependencies.

Second, the intra-BB memory dependencies are dynamically scheduled using LSQs. Dynamic analyses  $<_b$  and statically encodes the sequential memory order of each BB into the LSQ. An LSQ allocates the memory operations in a BB into its queue at the start of the corresponding BB execution. It dynamically checks these memory operations following the order of  $<_b$  and executes a memory operation if it is independent of all its priorly executing operations. The LSQ ensures that the intra-BB memory dependencies are resolved by statically encoding  $<_b$  into the LSQ for dependency check.

We now need to resolve the inter-BB dependencies. Completely resolving inter-BB dependencies for concurrent execution at runtime is still an open question. The most dynamic approach so far is by Dynamic, which enables dynamic scheduling with only one restriction. The restriction requires BB executions must start sequentially in strict program order, even they can execute in parallel. Let  $t : E_I \cup E_B \rightarrow \mathbb{N}$  denote the start times of an instruction execution or a BB execution in clock cycles.

$$\forall i, b, k. i \in I_b \wedge (b, k) \in E_B \Rightarrow t(b, k) \leq t(i, k) \quad (5)$$

$$D : \forall e, e'. e \in E_B \wedge e' \in E_B \wedge e < e' \Rightarrow t(e) < t(e') \quad (6)$$

This preserves the original program order  $<$  during runtime hardware execution and provides a reference of correctness for dynamic scheduling when combined with  $<_b$ .

Third, with Constraint 6, the inter-BB data dependencies are preserved using muxes. Dynamic uses a mux to select data input to a BB at the start of the BB by its preceding BB execution. Since the BBs are restricted to start sequentially, there can only be at most one BB receiving at most one starting signal from its multiple preceding BBs. The input data of a BB is selected by the muxes based on the starting signal and accept the correct input data for the computation. The starting signal in the sequential BB execution order ensures that the data sent to each BB is also in strict program order. Such property ensures in-order data flow between BBs, which preserves inter-BB data dependencies.

Finally, the inter-BB memory dependencies are resolved by the LSQ by dynamically allocating memory operations between these BB executions. LSQ dynamically monitor the start signals of BB executions and allocate groups of memory operations in the same BB order, also known as  $<$ . It checks the memory operations in an order that combines  $<$  and  $<_b$  lexicographically, the same as the original program order. With the logic that resolves the intra-BB memory dependencies, the LSQ ensures that the out-of-order memory execution always has the same results as the execution in the allocated order, i.e., the program order. Such an approach of preserving  $<$  for allocation in the LSQ resolves the inter-BB memory dependencies.

With the implementation above,  $<_b$  is directly encoded to hardware at compile time, and  $<$  is recovered at runtime. All four kinds of dependencies can be resolved by checking dependencies

between operations in strict program order. Let  $d : (E_I \cup E_B)^2 \rightarrow \{0, 1\}$  denote whether two executions have dependencies, and let  $l : E_I \rightarrow \mathbb{N}$  denote the latencies of instruction execution in clock cycles. Dynamatic ensures the following dependency constraint always holds:

$$\forall e, e'. e \in E_I \wedge e' \in E_I \wedge e < e' \wedge d(e, e') \Rightarrow t(e) + l(e) \leq t(e'). \quad (7)$$

**3.2.3 Hardware Restriction.** The dependency constraints above ensure the correctness of the schedule with optimised performance. However, the hardware architecture could also bring restrictions on performance.

Here, we only consider the case of pipelining. HLS uses hardware pipelining to exploit parallelism among different iterations of a BB using a single hardware instance. A resource restriction is then added to this scheduling model.

$$\forall b, k, k'. (b, k) \in E_B \wedge (b, k') \in E_B \wedge k \neq k' \Rightarrow t(b, k) \neq t(b, k') \quad (8)$$

This means that two iterations of the same BB cannot start simultaneously in pipelining, which requires multiple hardware instances. In Dynamatic implementation, Constraint 8 is covered by Constraint 6. We keep it here, as Constraint 6 will be relaxed in the later sections. Also, Dynamatic generates a data flow hardware architecture where the input data order is always the same as the output data order for each BB. This could also restrict the execution of each individual operation:

$$\forall i, k, k', b. b \in B \wedge i \in I_b \wedge (i, k) \in E_I \wedge (i, k') \in E_I \wedge t(b, k) < t(b, k') \Rightarrow t(i, k) < t(i, k'). \quad (9)$$

**3.2.4 Summary.** In summary, Dynamatic automatically generates efficient dynamically scheduled hardware with minimal static analysis. The generated hardware must satisfy four constraints. First, Constraint 6 ensures the program order is preserved in the hardware design. Second, given Constraint 6, the hardware logic ensures Constraint 7 for dynamically resolving dependencies using the sequential program order. Third, the performance of pipelined architecture is restricted by the resource Constraint 8. Finally, the performance is also restricted by the Constraint 9 introduced from the data flow architecture.

Constraint 6, which restricts BBs to start sequentially, significantly reduces the need for static analysis and simplifies the dynamic scheduling problem. However, this is too conservative for dependency analysis. For instance, when all the BB executions do not have dependencies, Constraint 6 could cause sub-optimal performance. We now demonstrate how Constraint 6 can be relaxed in this model and what analysis can be applied to achieve a schedule with a better performance.

### 3.3 Possible Relaxations

We seek to relax Constraint 6 and only restrict the order between the BB executions that may have dependencies. If two instruction executions in two BB executions have dependencies, then these two BB executions have dependencies.

$$d((b, k), (b', k')) = (\exists i, i'. (i, k) \in E_I \wedge (i', k') \in E_I \Rightarrow d((i, k), (i', k'))) \quad (10)$$

Then, Constraint 6 can be relaxed to:

$$D' : \forall e, e'. e \in E_B \wedge e' \in E_B \wedge e < e' \wedge d(e, e') \Rightarrow t(e) < t(e'). \quad (11)$$

However, analysing the dependency between two BB execution at runtime is still challenging.

To enable static analysis for dependency analysis, we over-approximate analysing two particular BB executions to all the executions of two particular BBs. This means that each dependency is checked between statements in the source instead of runtime events. Let  $d' : B \times B \rightarrow \{0, 1\}$



denote whether two basic blocks may have dependency during their executions.

$$d'(b, b') = (\exists k, k'. d((b, k), (b', k'))) \quad (12)$$

We then relax Constraint 11 to the following:

$$D'' : \forall b, b', k, k'. (b, k) \in E_B \wedge (b', k') \in E_B \wedge (b, k) < (b', k') \wedge d'(b, b') \Rightarrow t(b, k) < t(b', k') \quad (13)$$

This ensures that only some BBs must start sequentially, and BBs that cannot have dependency during the whole execution can start in parallel or out-of-order. The dependency set of the program from our formulation lies between Constraint 6 and Constraint 11, where  $D' \subseteq D'' \subseteq D$ . With this new constraint, the dependencies are still respected with existing muxes and LSQs, since the dependent BB executions remain starting sequentially.

In the rest of the article, we demonstrate how to apply static analysis for relaxing Constraint 6 towards Constraint 13 and enable two hardware optimisation techniques for dynamic-scheduling HLS.

## 4 DYNAMIC INTER-BLOCK SCHEDULING

This section demonstrates an application of the proposed dependency model for achieving the simultaneous execution of two independent BBs by parallellising independent sequential loops. It formalises our prior conference paper [15] in the proposed model. Section 4.1 demonstrates a motivating example of parallellising independent sequential loops. Section 4.2 explains how to formulate the problem into the proposed model in Section 3 and shows how to use Microsoft Boogie to automatically determine the absence of dependency between these sequential loops. Section 4.3 illustrates efficient hardware transformation for parallellising sequential loops.

### 4.1 Motivating Example

Here, we illustrate a motivating example of parallellising two sequential loops in dynamically scheduled hardware. Figure 3(a) shows an example of two sequential loops, `loop_0` and `loop_1`. In each iteration of `loop_0`, an element at index  $f(i)$  of array `a` is loaded and processed by a function `op0`. The result is stored to an element at index  $i$  of array `b`. In each iteration of `loop_1`, an element at index  $h(j)$  of array `a` is loaded and processed by a function `op1`. The result is stored back to array `a` at index  $g(j)$ . For simplicity, let  $f(i) = 0$ ,  $g(j) = j*j+1$  and  $h(j) = j$ . Hence, there is no memory dependency between two loops, that is,  $\forall 0 \leq i < X. \forall 0 \leq j < Y. f(i) \neq g(j)$ .

Dynamatic [33], the state-of-the-art dynamic scheduled HLS tool, synthesises hardware that computes in a schedule shown in Figure 3(b). The green bars represent the pipeline schedule of `loop_0`, and the blue bars represent the pipeline schedule of `loop_1`. In `loop_1`, the interval between the starts of consecutive iterations, known as the *initiation interval* ( $\Pi$ ), is variable because of the dynamic inter-iteration dependency between loading from `a[h(j)]` and storing to `a[g(j)]`. For instance, if we suppose that  $g$  and  $h$  are defined such that  $g(0) = h(1)$ , then the first two iterations must be executed sequentially, and if we further suppose that  $g(1) \neq h(2)$ , then the second and third iterations are pipelined with an  $\Pi$  of 1.

However, `loop_1` is stalled until all the iterations in `loop_0` have started, even though it has no dependency on `loop_0`. The reason is that Dynamatic forces all the BBs to start sequentially to preserve any potential inter-BB dependency, such as the inter-iteration memory dependency in `loop_1`. For this example, each loop iteration is a single BB, and at most one loop iteration starts in each clock cycle.

An optimised schedule is shown in Figure 3(c). In the figure, both loops start from the first cycle and iterate in parallel, resulting in better performance. Existing approaches cannot achieve the optimised schedule: static scheduling can start `loop_0` and `loop_1` simultaneously such as using

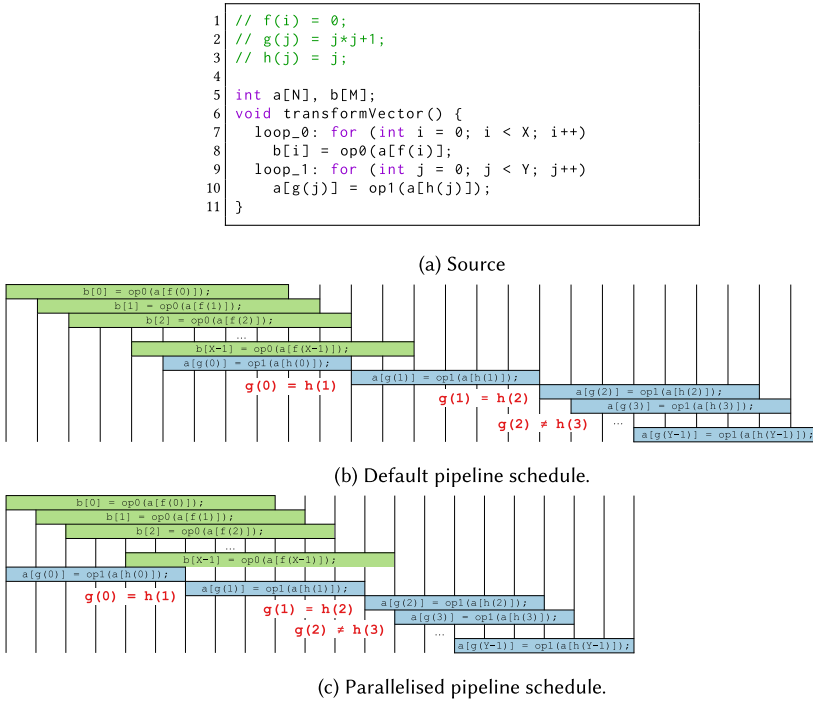


Fig. 3. Motivating example. Assume no dependence between two loops. The dynamically scheduled hardware from the original Dynamatic [33] a schedule in (b). Our work achieves an optimised schedule in (c).

multi-threading in LegUp HLS [7], but loop<sub>1</sub> is sequential, as the static scheduler assumes the worst case of dependency and timing; dynamic scheduling has a better throughput of loop<sub>1</sub>, but cannot start it simultaneously with loop<sub>0</sub>.

Besides, determining the absence of dependency between these two loops for complex  $f(i)$ ,  $g(j)$ , and  $h(j)$  is challenging. In this section, our toolflow (1) generates a Boogie program to formally prove that starting loop<sub>0</sub> and loop<sub>1</sub> simultaneously cannot break memory dependency and (2) parallelises these loops in dynamically scheduled hardware if they are proved independent. The Boogie program generated for this example is explained later (in Figure 5).

The transformation for the example in Figure 3(a) is demonstrated in Figures 4(a) and 4(b). Figure 4(a) shows the CFG generated by the original Dynamatic. The CFG consists of a set of pre-defined components, as listed in Table 1. As indicated by the red arrows, a control token enters the upper block and triggers all the operations in the first iteration of loop<sub>0</sub>. It circulates within the upper block for  $X$  cycles and then enters the lower block to start loop<sub>1</sub>. Figure 4(b) shows a parallelised CFG by our toolflow. Initially, a control token is forked into two tokens. These two tokens simultaneously trigger loop<sub>0</sub> and loop<sub>1</sub>. A join is used to synchronise the two tokens when they exit these loops. Both designs use the same hardware, yet, Figure 4(b) uses these resources in a more efficient way by allowing the two loops to be used in parallel, reducing the overall execution time. The rest of Section 4 explains the details of our approach.

## 4.2 Problem Formulation and Dependency Analysis

Here, we first show how to formalise the problem based on the model in Section 3. We then show how to extract sets of subgraphs from a sequential program, where subgraphs in the same set

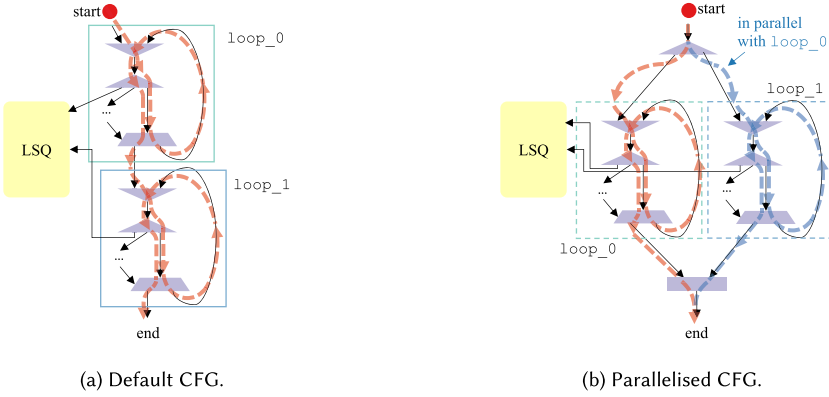

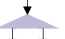
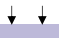



Fig. 4. Hardware transformation of the motivating example in Figure 3.

Table 1. Elastic Components for Dynamically Scheduled HLS

	<b>Merge:</b> takes the input data from an arbitrary predecessor and propagates it to its single successor.
	<b>Fork:</b> takes the input data from its single predecessor and replicates it to each of its multiple successors.
	<b>Join:</b> triggers its single successor only when the input data of its all predecessors is available.
	<b>Branch:</b> takes the data from its data predecessor and propagates it to one of its multiple successors based on the select value from its control predecessor.

may start in parallel. The absence of dependency between these parallelised subgraphs is formally verified using the generated Boogie program by our tool.

**4.2.1 Problem Formulation.** The search space for BBs that can start in parallel could be huge, and it scales exponentially with the code size. To increase scalability, we limit our scope to loops. Each loop forms a subgraph in the CFG for analysis. Parallelising BBs outside any loop adds significant search time but has negligible improvement in latency. We define the following terms:

- $G = \{g_1, g_2, \dots\}$  denotes a set of consecutive subgraphs in the CFG of the program,
- $E_G \subseteq G \times \mathbb{N}$  denotes the executions of subgraphs,
- $\prec_G \subseteq E_G \times E_G$  denotes the original program order of subgraph execution, and
- $B_g \subseteq B$  denotes the set of all the BBs in subgraph  $g$ .

The subgraph set  $G$  must satisfy the following constraints based on the original sequential program order. First, the BB sets of all the subgraphs in  $G$  must be disjoint.

$$\forall g, g'. g \in G \wedge g' \in G \wedge g \neq g' \Rightarrow B_g \cap B_{g'} = \emptyset \quad (14)$$

Second, no BB outside a subgraph executes during the execution of the subgraph in sequential program order. Let  $b_0(e)$  and  $b_n(e)$  denote the first BB execution and the last BB execution in a subgraph execution, where  $e \in E_G$ .

$$\nexists b, k, g, m. (g, m) \in E_G \wedge b \notin B_g \wedge (b, k) \in E_B \Rightarrow b_0(g, m) < (b, k) < b_n(g, m) \quad (15)$$

Finally, all the subgraphs in  $G$  must consecutively execute, i.e., directly connected to at least one subgraph in CFG. That means that they are sequentially executed in each iteration. Let  $c(g, g')$  denote whether the execution of a subgraph  $g'$  is consecutive after the execution of subgraph  $g$ .

$$c(g, g') = (\nexists m, g''. (g, m) \in E_G \wedge (g', m) \in E_G \wedge (g'', m) \in E_G \Rightarrow (g, m) < (g'', m) < (g', m)) \quad (16)$$

$$\nexists g, g', m, b, k. (g, m) \in E_G \wedge (g', m) \in E_G \wedge c(g, g') \wedge (b, k) \in E_B \Rightarrow b_n(g, m) < (b, k) < b_0(g', m) \quad (17)$$

Now, we start to map the execution of subgraphs into a hardware schedule. As explained in Constraint 6 in Section 3, Dynamatic forces BB to start sequentially. This leads to the following constraint regardless any dependency:

$$\forall g, g', m. c(g, g') \wedge (g, m) < (g', m) \Rightarrow t(b_n(g, m)) < t(b_0(g', m)). \quad (18)$$

If it is proven that the execution  $(g, m)$  cannot have any dependency with the execution  $(g', m)$  of its consecutively following subgraph, then there is no need to use muxes and LSQs to resolve the dependency between these two subgraphs. Then,  $(g', m)$  can start execution early, such as  $t(b_0(g, m)) = t(b_0(g', m))$ , which leads to a correct schedule with a better performance. Let  $d'(g, g')$  denote that two subgraphs  $g$  and  $g'$  may have a dependency, and let  $d'_c(g, g')$  denote that there exists a subgraph between the execution of  $g$  and  $g'$  in the same iteration including  $g$  that may have a dependency with subgraph  $g'$ .

$$d'(g, g') = (\exists b, b'. b \in B_g \wedge b' \in B_{g'} \Rightarrow d'(b, b')) \quad (19)$$

$$d'_c(g, g') = d'(g, g') \vee (\exists g'', m. (g, m) < (g'', m) < (g', m) \Rightarrow d'(g'', g')) \quad (20)$$

Constraint 6 is now relaxed to:

$$\forall g, g', m. (g, m) < (g', m) \wedge d'_c(g, g') \Rightarrow t(b_n(g, m)) < t(b_0(g', m)), \quad (21)$$

$$\forall b, b', k, k', g. b \in B_g \wedge b' \in B_{g'} \wedge g \in G \wedge (b, k) < (b', k') \Rightarrow t(b, k) < t(b', k'). \quad (22)$$

Constraint 21 restricts the starting time of a subgraph execution by its most recently executed subgraph that has dependencies. Inside each subgraph, BB execution remains to start sequentially, as shown in Constraint 22.

The optimised schedule still respects all the inter-BB dependencies, since it only modifies the start times of independent subgraph execution. The intra-BB dependencies remain unchanged, since the transformation is applied at only the subgraph level. Therefore, Constraint 7 still holds for the optimised schedule. Also, Constraint 8 and Constraint 9 still hold, as the hardware pipelining and dataflow architecture remain the same.

The following sections explain how to solve two main problems: (1) How to efficiently determine a large set of  $G$  and a highly parallelised schedule for  $G$ ? (2) How to map a parallelised schedule into efficient hardware?

**4.2.2 Subgraph Extraction.** Given an input program, our toolflow analyses sequential loops in each loop depth and constructs a number of sets of subgraphs. Each set contains several consecutive sequential loops at the same depth, where each loop forms a subgraph. For instance, the example in Figure 3 has a set of two subgraphs, corresponding to `loop_0` and `loop_1`. Our toolflow then checks the dependency among the subgraphs for each set. Dynamatic translates data dependency into handshake connections in hardware for correctness. Our toolflow does not change these connections, so the data dependency is still preserved. For memory dependencies, our toolflow generates a Boogie program to prove the absence of dependency among subgraphs.

```

1 procedure pickOneMemoryAccess() returns (valid: bool,
2   addr: Index, array: Array, subgraphID: Index,
3   type: MemoryType) {
4   loop_0: for (i = 0; i < X; i++) {
5     // b[i] = op0(a[f(i)]);
6     if (*) { valid := true; addr := f(i); array := a;
7       subgraphID := 0; type := LOAD; return; }
8     if (*) { valid := true; addr := i; array := b;
9       subgraphID := 0; type := STORE; return; } }
10  loop_1: for (j = 0; j < Y; j++) {
11    // a[g(j)] = op1(a[h(j)]);
12    if (*) { valid := true; addr := h(j); array := a;
13      subgraphID := 1; type := LOAD; return; }
14    if (*) { valid := true; addr := g(j); array := a;
15      subgraphID := 1; type := STORE; return; } }
16  valid := false; return; }

```

(a) Procedure that arbitrarily picks a memory access.

```

1 procedure main() {
2   // assume that all the arrays have arbitrary values
3   havoc a, b;
4   // valid: whether the returned memory access is valid
5   // addr: which address the memory access touches
6   // array: which array the memory access touches
7   // subgraphID: which subgraph the memory access is in
8   // type: the type of memory access, either load/store
9   call valid_0, addr_0, array_0, subgraphID_0, type_0 :=
10  pickOneMemoryAccess();
11  call valid_1, addr_1, array_1, subgraphID_1, type_1 :=
12  pickOneMemoryAccess();
13  assert !valid_0 || !valid_1 ||
14  subgraphID_0 == subgraphID_1 ||
15  array_0 != array_1 ||
16  (type_0 == LOAD && type_1 == LOAD) ||
17  addr_0 != addr_1;
18 }

```

(b) Main procedure that proves the absence of dependency.

Fig. 5. A Boogie program generated for the example in Figure 3. It tries to prove the absence of memory dependency between two sequential loops loop\_0 and loop\_1.

For this example, Boogie proves that the two loops do not conflict on any memory locations and can be safely reordered.

For example, Figure 5 shows the Boogie program that proves the absence of a dependency between loop\_0 and loop\_1 in Figure 3. The Boogie program consists of two procedures. First, the procedure in Figure 5(a) describes the behaviour of function transformVector and arbitrarily picks a memory access during the whole execution. The procedure returns a few parameters for analysis, as listed in lines 4–8 in Figure 5(b). The for loop structures are automatically translated using an open-sourced tool named EASY [13]. In the rest of the function, each memory operation is translated to a non-deterministic choice if(\*). It arbitrarily returns the parameters of a memory operation or continues the program. If all the memory operations are skipped, then the procedure returns an invalid state in line 16. The non-deterministic choices over-approximate the exact memory locations to a set of potential memory locations. For instance, any memory location accessed by the code in Figure 3(a) is reachable by the procedure in Figure 5(a). The assertions in Figure 5(b) must hold for any possible memory location returned by the procedure in Figure 5(a) to pass verification.

Figure 5(b) shows the main procedure. In line 3, the verifier assumes both arrays hold arbitrary values, making the verification input independent. Then, the verifier arbitrarily picks two memory accesses in lines 9–10. Each memory access can capture any memory access during the execution of transformVector. The assertion describes the dependency constraint to be proved that for any two valid memory accesses (line 11), if they are in different subgraphs (line 12), they must be independent. Lines 13–15 describe the independency, where they either touch different arrays or different indices, or they are both load operations. If the assertion always holds, then it is safe to parallelise loop\_0 and loop\_1.

Our toolflow generates  $\frac{k(k-1)}{2}$  assertions for  $k$  subgraphs, because that is the number of ways of picking two subgraphs from  $k$ . The subgraphs are rescheduled based on the verification results. If a subgraph is independent of any of its preceding subgraphs within a distance of  $n$ , then it can simultaneously start with its  $(m - n)$ th last subgraph. In this article, we analyse at loop level for parallelism instead of the BB level for better scalability, and the search space is already huge. In the cases of two or more consecutive subgraphs that are all mutually independent, it is straightforward to schedule them all in parallel. However, a sequence of subgraphs that are neither completely independent nor completely dependent may result in several possible solutions. For instance, the CFG in Figure 6(a1) contains three consecutive loops, BB1, BB2, and BB3. BB1 and BB2 can be parallelised, as can BB2 and BB3, but BB1 and BB3 cannot. We, therefore, have to choose between

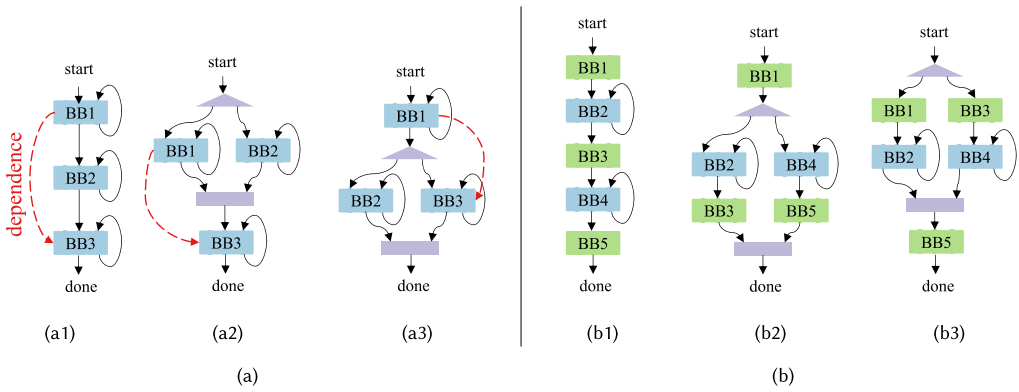


Fig. 6. A CFG may be parallelised differently, depending on (a) parallelising in top-to-bottom/bottom-to-top order, and (b) grouping BBs with the loop before/after. The dashed arrows represent memory dependency.

parallelising them as in Figure 6(a2) or in Figure 6(a3). Our current approach greedily parallelises BBs in top-to-bottom order, so it yields Figure 6(a2) by default, but this order can be overridden via a user option. It may be profitable in future work to consider Figure 6(a3) as an alternative if BB2 and BB3 have more closely matched latencies.

Second, the BBs between sequential loops can be included in a subgraph of either loop, resulting in several solutions. For instance, Figure 6(b1) can be parallelised to Figure 6(b2) or to Figure 6(b3). In Figure 6(b2), BB2 is grouped with its succeeding loop BB3 and so is BB4. In Figure 6(b3), the BBs are grouped with their preceding loops. This may result in different verification results, which affect whether the subgraphs can be parallelised. For instance, if BB3 depends on BB2, then Figure 6(b2) is memory-legal and Figure 6(b3) is invalid (our toolflow will keep the CFG as in Figure 6(b1)). This grouping can be controlled via a user option.

### 4.3 Hardware Transformation

We here explain how to construct dynamically scheduled hardware in which BBs can start simultaneously. First, we illustrate how to insert additional components to enable BB parallelism. Second, we show how to simplify the data flow to avoid unnecessary stalls for subgraphs.

**4.3.1 Components Insertion for Parallelism.** With given sets of subgraphs that start simultaneously, our toolflow inserts additional components into the dynamically scheduled hardware to enable parallelism. For each set, our toolflow finds the start of the first subgraph and the exit of the last subgraph in the program order. The trigger of the first subgraph is forked to trigger the other subgraphs in the set. The exit of the last subgraphs is joined with the exits of the other subgraphs and then triggers its succeeding BB. For the example in Figure 4(b), the start of the function is forked to trigger both loop\_0 and loop\_1. A join is used to synchronise the BB starting signals in loop\_0 and loop\_1. The join waits for all the BBs in both loops to start and then starts the succeeding BB of the loops.

The BB starting order is now out-of-order, but the computed data must be in-order. The transformation above ensures the order of data does not affect the correctness. Since we only target loops, only the muxes at the header of the loops are affected. Outside of the loops to be parallelised, the order remains unmodified. When each parallelised loop starts, a token enters the loop and circulates through the loop exactly in the program order. The parallelised loop outputs are synchronised by the join, thus, everything that happens later remains in order. Only the BB orders among these parallelised loops are out-of-order, which have been proven independent.

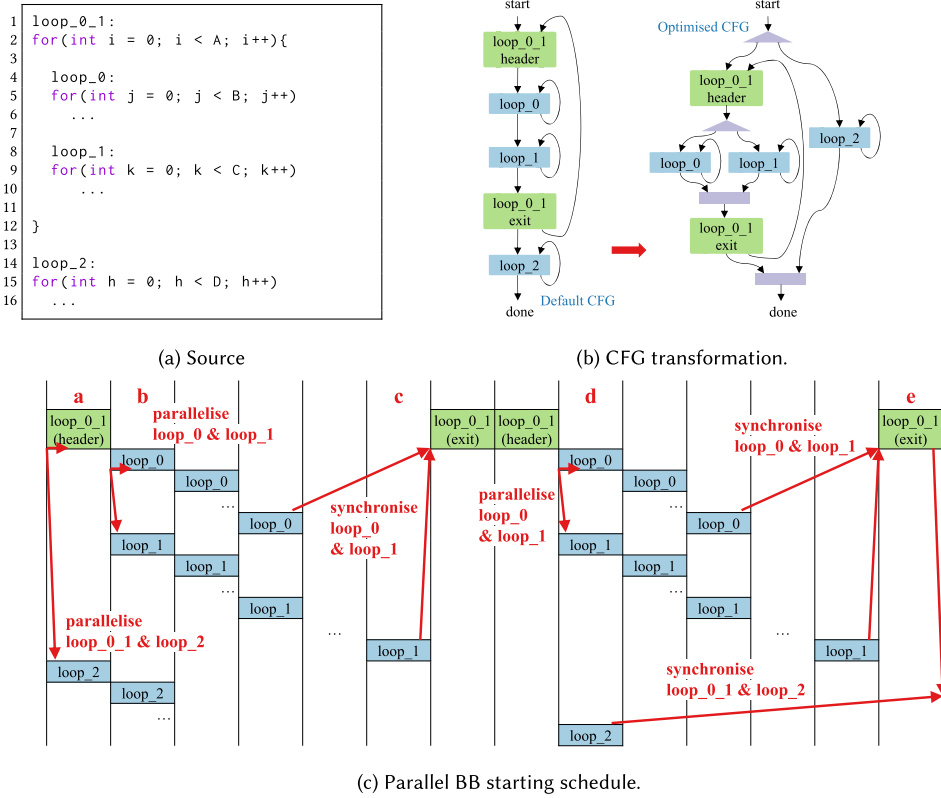


Fig. 7. An example of parallelising BB starting schedule by CFG transformation. There are two sets of sequential loops in different depths. Assuming all the loops are independent, each set of sequential loops starts simultaneously after the transformation in (b). (c) only shows the time when a BB starts, where a BB may take multiple cycles to execute.

An advantage of such transformation is that the execution of parallelised subgraphs and their succeeding BB are in parallel, although they still start in order. The memory dependencies between these subgraphs and the succeeding BB are still respected at runtime, as they start in order. This effect qualitatively corresponds to what standard dynamically scheduled hardware exhibits, yet, in that case, only on a single BB at a time. Compared to traditional static scheduling, which only starts the succeeding BB when all the subgraphs finish execution, our design can achieve better performance.

Figure 7 shows an example of parallelising nested parallel subgraphs. The code contains two sequential loops, loop\_0\_1 and loop\_2. Loop loop\_0\_1 is a nested loop that contains two sequential loops, loop\_0 and loop\_1. For simplicity, assume that there is no dependency between any two loops.

Our toolflow constructs two sets of subgraphs in two depths, allowing more parallelism in the CFGs. One set contains loop\_0\_1 and loop\_2, and the other set contains loop\_0 and loop\_1. The transformation of CFG is illustrated in Figure 7(b). loop\_0\_1 and loop\_2 are parallelised at the start the program, and loop\_0 and loop\_1 are further parallelised inside loop\_0\_1. The corresponding BB starting schedule is demonstrated in Figure 7(c), which only shows the time when each BB starts. A BB may have a long latency and execute in parallel with other BBs.

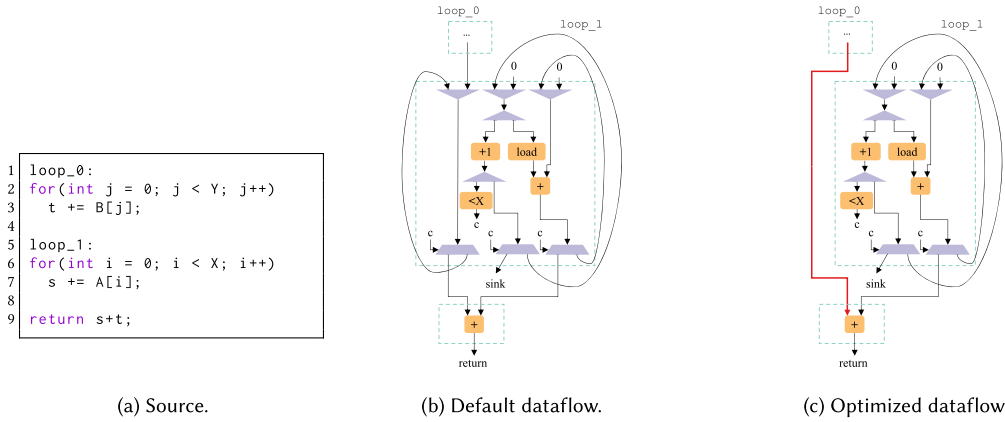


Fig. 8. An example of simplified data flow for live variables.  $t$  is a live variable in line 3, but not used in  $\text{loop}_1$ .  $t$  circulates in  $\text{loop}_1$  to preserve liveness but is seen as data dependencies, stalling  $\text{loop}_1$  before  $t$  is valid. Our toolflow identifies and removes these cycles, such that  $\text{loop}_1$  can start earlier.

**4.3.2 Forwarding Variables in Data Flow.** The second step is to simplify the data flow of live variables for parallelising sequential loops. Dynamic directly translates the CFGD of an input program into a hardware dataflow graph. In the data flow graph, each vertex represents a hardware operation, and each edge represents a data dependency between two operations.

The data flow of a loop uses cycles for each variable that has carried dependency. The data circulates in the cycle and updates its value in each iteration. However, such an approach also maintains all the live variables in these cycles while executing a loop, even when they are not used inside the loop. The edges of these cycles are seen as data dependencies in the hardware, where the edges for unused live variables could cause unnecessary pipeline stalls.

For example, the loops in Figure 8(a) can be parallelised.  $\text{loop}_0$  accumulates array B onto  $t$ , and  $\text{loop}_1$  accumulates array A onto  $s$ . The sum of  $s$  and  $t$  is returned. The dataflow graph of  $\text{loop}_1$  is shown in Figure 8(b). The loop iterator  $i$  and the variable  $s$  have carried dependency in  $\text{loop}_1$ . They are kept and updated in the middle and right cycles. The result of  $\text{loop}_0$ ,  $t$ , is still live and required by addition in line 9.  $t$  is kept in the left cycle, circulating with  $i$  and  $s$ .

$\text{loop}_1$  is stalled by the absence of  $t$  even when parallelised with  $\text{loop}_0$ , but  $t$  is not needed by  $\text{loop}_1$ . To remove these unnecessary cycles, our toolflow checks whether a live variable is used in the loop. If it is not, then our toolflow removes the corresponding cycle and directly forwards the variable to its next used BB. Figure 8(c) illustrates the transformed dataflow graph.  $t$  is now directly forwarded to the final adder, enabling two loops to start simultaneously.

**4.3.3 LSQ Handling.** The parallel BB schedule also affects the LSQs. First, the original Dynamic starts BB sequentially, whereas the LSQ expects sequential BB allocation. Our parallelised schedule allows multiple BBs to start simultaneously; therefore, we place a round-robin arbiter for the LSQ to serialise the allocations. The out-of-order allocation still preserves correctness, as the simultaneous BB requests have been statically proven independent by our approach in Section 4.2.2.

Second, the arbiter may cause deadlock if the LSQ depth is not sufficient to consume and reorder all memory accesses, (e.g., a later access may be stuck in an LSQ waiting for a token from an earlier access, but the earlier access cannot enter the LSQ if it is full, thus never supplying the token). This issue has been extensively explored in the context of shared resources in dataflow circuits [34];



similarly to what is suggested in this work, the appropriate LSQ size could be determined based on the number of overlapping loop iterations and their IIs. Although systematically determining the minimal allowed LSQ depth is out of the scope of this work, we here assume a conservative LSQ size that ensures that deadlock never occurs in the benchmarks we consider. We note that minimising the LSQ is orthogonal to our contribution and could only positively impact our results (by reducing circuit area and improving its critical path).

## 5 DYNAMIC C-SLOW PIPELINING

This section demonstrates another application of the proposed dependency model for achieving out-of-order execution of independent and consecutive iterations of the same BB by C-slow pipelining in nested loops. It formalises our prior conference paper [17] in the proposed model. Section 5.1 demonstrates a motivating example of C-slow pipelining the innermost loop of a loop nest. Section 5.2 explains how to formulate the problem into the proposed model in Section 3 and shows how to use Microsoft Boogie to automatically determine the absence of dependency between the outer-loop iterations of the loop nest for C-slow pipelining. Section 5.3 explains how to realise C-slow pipelining in hardware.

### 5.1 Motivating Example

In this section, we use a motivating example to demonstrate the problem of pipelining a nested loop. In Figure 9(a), a loop nest updates the elements in an array  $a$ . The outer loop  $\text{loop}_0$  loads an element at address  $f(i)$ . The inner loop  $\text{loop}_1$ , bounded by  $N-i$ , computes  $s$  with a row in an array  $b$ , shown as function  $g$ . The result is then stored back to array  $a$  at address  $h(i)$  at the end of each outer-loop iteration.

For simplicity, assume there is no inter-iteration dependency in the outer loop  $\text{loop}_0$ . Assume the latency of function  $g$  is three cycles. An inter-iteration dependency of the inner loop  $\text{loop}_1$  on  $s$  causes a minimum II of 3. The pipeline schedule of the hardware from vanilla Dynamatic is shown in Figure 9(b). The first iteration of  $\text{loop}_0$  is optimally pipelined with an II of 3 shown as the green bars. However, the second iteration, shown as blue bars, can only start after the last iteration of  $\text{loop}_1$  in the first iteration of  $\text{loop}_0$  starts. Although the loop contains three registers, the control flow is not parallelised because of Constraint 6.

The schedule shown in Figure 9(c) is also correct and achieves better performance. Since the II of  $\text{loop}_0$  is 3, the two empty slots between every two consecutive iterations allow the next two iterations of  $\text{loop}_0$  to start earlier. The second iteration of  $\text{loop}_0$  now starts one cycle after the start of the first iteration of  $\text{loop}_0$ , followed by the third iteration shown as orange bars. After the last iteration of  $\text{loop}_1$  starts, the current inner-loop instance leaves new empty pipeline slots spare. This triggers the start of the fourth iteration, shown as yellow bars, filling into the new empty slot.

The reason that Dynamatic cannot achieve the schedule in Figure 9(c) is that the control flow in the latter schedule is out-of-order, which breaks Constraint 6. The LSQ cannot retain the original program order of memory accesses and cannot verify the correctness of memory order from the out-of-order control flow, which may lead to wrong results. In this section, we use static analysis to prove such control flow will still maintain a legal memory access order for a given program, so the LSQ still works correctly for this new program order.

This is an example for which traditional techniques such as loop interchange and loop unrolling do not help because they only work under stringent constraints. In this example, loop interchanging cannot be applied, because the bound of the inner loop depends on its outer loop. Also, loop unrolling does not change the control flow and cannot improve performance. In this section, we propose a general approach that works for arbitrary nested loops.

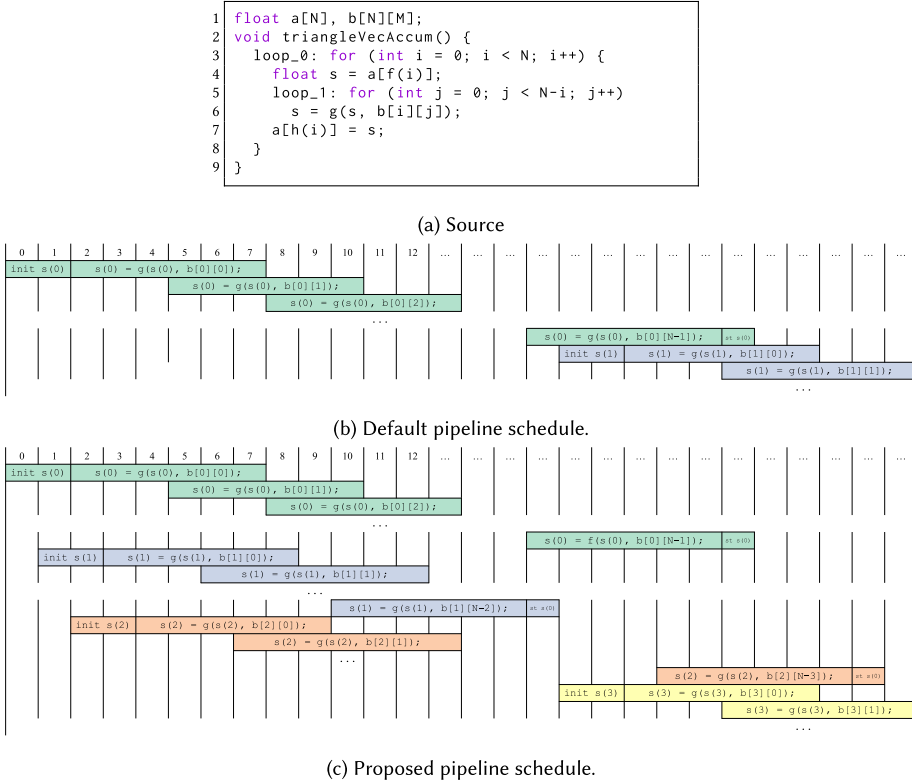


Fig. 9. A motivating example of computing a triangle matrix. Assume there is no inter-iteration dependency in `loop_0`, and each instance of `loop_1` has a minimum `II` of 3. The default pipeline schedule only starts the second iteration of `loop_0` after the last iteration of `loop_1` in the first iteration of `loop_0` starts. Our approach inserts the following iterations of `loop_0` into the empty slots of its first iteration.

## 5.2 Problem Formulation and Dependency Analysis

Here, we first show how to formalise the problem of  $C$ -slow pipelining based on the model in Section 3. We then show how to analyse the correct  $C$  for each loop nest based on the dependency analysis using the generated Boogie program by our tool.

**5.2.1 Problem Formulation.**  $C$ -slow pipelining is amenable for improving the throughput when (1) a hardware design that has an `II` of greater than 1; and (2) it allows out-of-order execution of control flow. To simplify the problem, we restrict the scope of our work to nested loops. A loop in a CFG is defined as a set of consecutive BBs with back edges. Here, we define the following terms for an inner loop in a loop nest:

- $B_L \subseteq B$  denotes the set of all the BBs in a loop,
- $E_L(j, k) \subseteq E_B$  denotes the BB executions in the  $k$ th iteration of the  $j$ th instance of the loop.

Since we only focus on loops, the abstraction has been raised to the loop level for this problem; here, we use  $E_L(j, k)$  to denote a particular loop execution event. This improves the scalability of our analysis. The `II` of each loop may vary between loop iterations as it is dynamically scheduled. For a loop, the maximum `II` of the loop at runtime can be defined as:

$$\forall b, k. b \in B_L \wedge k > 1 \wedge (b, k) \in E_B \wedge (b, k-1) \in E_B \Rightarrow II \geq t(b, k) - t(b, k-1). \quad (23)$$

An  $\Pi$  of greater than one means there are empty pipeline slots in the schedule, which could be attributed to a lack of hardware resources or inter-iteration dependencies. Here, we assume infinite buffering and only analyse the case for the stall caused by inter-iteration dependencies.

Constraint 6 forces each iteration of a loop to start execution sequentially, which can derive the following constraints for the loop:

$$\forall j, k, e, e'. e \in E_L(j, k-1) \wedge e' \in E_L(j, k) \Rightarrow t(e) < t(e'), \quad (24)$$

$$\forall j, k, k', e, e'. e \in E_L(j-1, k) \wedge e' \in E_L(j, k') \Rightarrow t(e) < t(e'). \quad (25)$$

Constraint 24 means that all the BB executions in an iteration of the loop cannot start unless all the BB executions in its last iteration have started. Constraint 25 means that all the BB executions in an instance of the loop cannot start unless all the BB executions in its last instance of the loop have started. The difference of the start times between two iterations when an  $\Pi$  is greater than 1, also known as the empty pipeline slots, cannot be filled with the following iterations because of Constraint 24 and Constraint 25. If it is proven that the  $j$ th instance and the  $(j-1)$ th instance of the loop are independent, Constraint 25 can be relaxed, as the absence of dependency enables early execution of the  $j$ th instance. This allows the following iterations to start early, filling the empty slots as illustrated in Figure 9(c).

The parallelism among iterations depends on the number of consecutive independent instances, also known as the dependency distance of the outer loop. The minimum dependency distance  $Q$  of the outer loop of a loop  $l$  must satisfy the following:

$$d''(j, j') = (\exists k, k', e, e'. e \in E_L(j, k) \wedge e' \in E_L(j', k') \Rightarrow d(e, e')), \quad (26)$$

$$\forall j, q. j > Q \wedge 1 \leq q \leq Q \Rightarrow \neg d''(j - q, j). \quad (27)$$

The constraint for a C-slowed nested loop is that there are always *at most*  $Q$  outer-loop iterations executing concurrently.

Although the outer-loop iterations are parallelised, the execution of inner loops still follows Constraint 24, where the inter-iteration dependencies of the inner loops are respected. Constraint 25 transformed from Constraint 6 is then relaxed to the following combined with Constraint 24.

$$\forall j, k, e, e'. j > Q \wedge e \in E_L(j - Q, k) \wedge e' \in E_L(j, k) \Rightarrow t(e) < t(e') \quad (28)$$

Constraint 7 still holds, as only the independent iterations execute earlier. Constraint 8 and Constraint 9 still hold, as the hardware property remains the same. The starting order of BB executions outside the C-slowed loop remains the same as vanilla Dynamic. The starting order of BB executions inside the C-slowed loop is now in parallel and out-of-order.

The following sections explain how to solve two main problems: (1) How to efficiently determine  $C$  for a correct schedule with better performance and area efficiency? (2) How to transform the hardware to realise C-slow pipelining?

**5.2.2 Exploration for  $C$  Using Dependency Analysis.** The dependency constraint above is equivalent to the minimum dependency distance of the outer loop must be greater than  $C$ , i.e.,  $C \leq Q$ . A loop-carried data dependency always has a dependency distance of 1, therefore, we only need to analyse memory dependencies. Our toolflow automatically generates a Boogie program to describe the memory behaviour of the nested loop and calls the Boogie verifier to prove the absence of memory dependency within a given distance.

For example, Figure 10 illustrates the Boogie program generated for the motivating example in Figure 9. It tries to prove the absence of memory dependency between any two outer-loop iterations with a distance less than  $C$ , which mainly includes two parts. The Boogie procedure in Figure 10(a) arbitrarily picks a memory access from the nested loop during the whole execution and

```

1 procedure pickOneMemoryAccessFromLoop() returns (
2 valid: bool, stmt: int, addr: Index, array: Array,
3 iteration: Index, type: MemoryType) {
4   loop_0: for (int i = 0; i < N; i++) {
5     // s = a[f(i)];
6     if (*) {
7       valid := true; stmt := 0; addr := (f(i));
8       array := a; iteration := (i); type := LOAD;
9       return; }
10    loop_1: for (int j = 0; j < g(i); j++)
11      // s = f(s, b[i][j]);
12      if (*) {
13        valid := true; stmt := 1; addr := (i, j);
14        array := b; iteration := (i); type := LOAD
15        ; return; }
16    // a[h(i)] = s;
17    if (*) {
18      valid := true; stmt := 2; addr := (h(i));
19      array := a; iteration := (i); type := STORE;
20      return; }
21  }
22  valid := false;
23  return;
24 }

```

```

1 // C : Given dependency distance
2 procedure main(C: int) {
3   // assume that all the arrays have arbitrary values
4   havoc a, b;
5
6   // valid: whether the returned memory access is valid
7   // stmt: which statement that executes the memory access
8   // addr: which address the memory access touches
9   // array: which array the memory access touches
10  // iteration: the iteration index of the current outer-loops
11  // type: the type of the memory access, either load or store
12  call valid_0, stmt_0, addr_0, array_0, iteration_0, type_0
13  := pickOneMemoryAccessFromLoop();
14  call valid_1, stmt_1, addr_1, array_1, iteration_1, type_1
15  := pickOneMemoryAccessFromLoop();
16
17  assert !valid_0 || !valid_1 ||
18  array_0 != array_1 ||
19  stmt_0 == stmt_1 ||
20  (type_0 == LOAD && type_1 == LOAD) ||
21  iteration_0 >= iteration_1 ||
22  getDistance(iteration_0, iteration_1) >= C ||
23  addr_0 != addr_1;
24 }

```

(a) Procedure that arbitrarily picks a memory access.

(b) Main procedure that describes absent dependency for a given  $C$ .

Fig. 10. A Boogie program generated for the example in Figure 9. It tries to prove the absence of memory dependency between any two outer-loop iterations with a distance less than  $C$ .

returns its parameters. The returned parameters include the label of the statement being executed, the array and address of the accessed memory, the iteration index of the outer loop, and the type of the memory access. Detailed definitions of these parameters are listed in lines 6–11 in Figure 10(b).

In Figure 10(a), the for loop structures in Boogie are directly generated by the automated tool named EASY [13]. In the loop body, each memory access is replaced with an `if(*)` statement. The `if(*)` statement arbitrarily chooses to return the parameters of the current memory access or continue. The procedure is then able to capture *all* the memory access that may execute during the whole execution.

If all these memory accesses are skipped, then the procedure exits at line 23 with a false `valid` bit, indicating that the returned parameters are invalid. Figure 10(b) describes the main Boogie procedure for dependency analysis. It takes a given  $C$  as an input. In line 4, it assumes that arrays  $a$  and  $b$  hold arbitrary values. This makes the verification results independent from the program inputs. Lines 12 and 13 arbitrarily pick two memory accesses from the nested loop using the procedure in Figure 10(a).

The assertion at line 15 proves Constraint 27 for a given  $C$ . First, the picked two memory accesses from lines 12 and 13 must hold valid parameters (line 15). Second, two accesses touching different arrays cannot have dependency (line 16). Two accesses executed by the same statement are safe (line 17), where the dependency is captured by the hardware logic based on the starting order of BB iteration under Constraint 8. Two loads cannot have dependency (line 18). Two returned memory accesses are arbitrary and have no difference. Here, we assume the memory access with index 0 executes in an earlier outer-loop iteration than the one with index 1 (line 19). In the  $C$ -slow pipelining formulation, only the outer-loop iterations with an iteration distance less than  $C$  can execute concurrently (line 20). Any two memory accesses that exclude the cases at lines 15–20 cannot access the same address. The assertion must hold for *any* two memory accesses for any input values. The Boogie verifier automatically verifies whether the assertion always holds for a given  $C$ . If the assertion always holds, then it is safe to parallelize  $C$  iterations of the outer loop.

Besides the dependency constraints, a resource constraint of  $C$ -slow pipelining is that each path must be able to hold at least  $C$  sets of data. Dynamatic already inserts buffers into the hardware

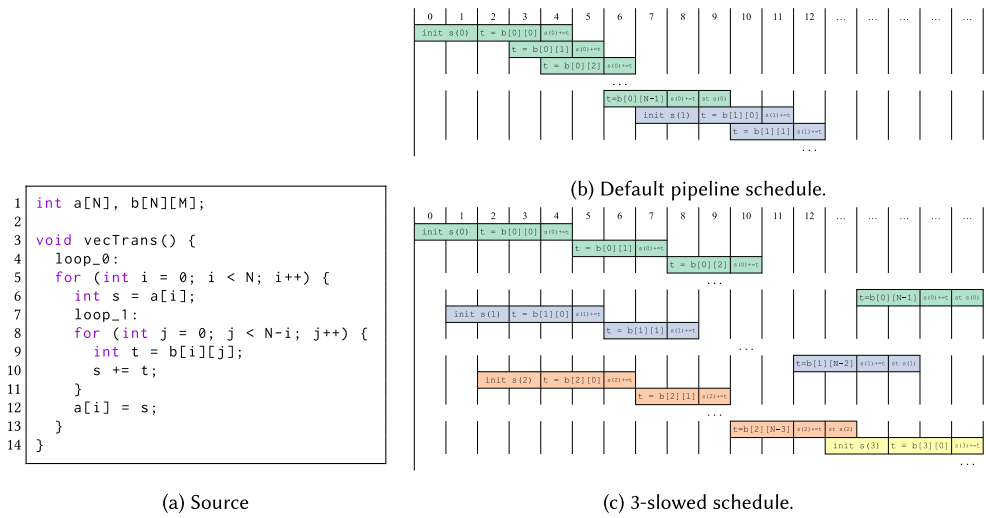


Fig. 11. A large  $C$  may not improve the overall throughput but slow the execution of single instances. A  $C$  of 3 only leads to an additional area for the code example above.

for high throughput. Our hardware transformation pass inserts an additional FIFO with a depth of  $C$  (named  $C$ -slow buffers) in each control path cycle of the inner loop to adapt  $C$ -slow pipelining, which can hold at least  $C$  tokens.

**5.2.3 Exploration for  $C$  Using Throughput Analysis.** The dependency analysis above only defines a set of  $C$ s that are suitable for  $C$ -slow pipelining. We show how to automatically determine an optimised  $C$  among these  $C$ s using throughput analysis. The Boogie program determines an upper bound  $C$  that cannot break any dependency. However, a large  $C$  may not improve the overall throughput but only cause more area overhead.

For example, Figure 11 illustrates an example where  $C$ -slow pipelining does not improve overall throughput. Figure 11(a) shows a function named `vecTrans` that transforms an array named `a`. The function contains a nested loop. In the outer loop, it loads the element in array `a` and accumulates the values in matrix `b` onto the element. In the inner loop, the elements in matrix `b` are accumulated in a triangle form.

The default pipeline schedule of function `vecTrans` is shown in Figure 11(b). In the schedule, both the inner loop and the outer loop are fully pipelined. Although there is a carried dependency in the inner loop on the variable `s`, the integer adder has a latency of one clock cycle, leading to an  $\Pi$  of 1.

This example is not amenable for  $C$ -slow pipelining, however, this cannot be identified by Boogie. In the dependency distance analysis, Boogie found that  $C$  can be any positive integer, as there is no inter-iteration dependency in the outer loop. For example, Figure 11(c) shows a 3-slowed schedule for function `vecTrans`. In the schedule, the start time of iterations in the first inner loop instance is delayed by three cycles, allowing two iterations in the second and third inner loop instances to start early. Such transformation preserves correctness but has no impact on the overall throughput. For this example,  $C$ -slow pipelining only causes more area overhead by adding the additional scheduler for out-of-order execution.

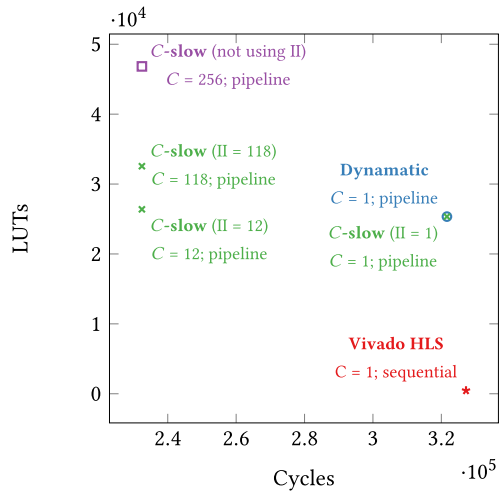
A condition where  $C$ -slow pipelining potentially improves the overall throughput is that the  $\Pi$  of an inner loop is greater than 1. This leaves empty pipeline slots for early execution of the later

```

1 float a[N], b[N][M];
2 void dynamicVecAccum() {
3   loop_0:
4   for (int i = 0; i < N; i++) {
5
6     int s = a[f(i)];
7
8     loop_1:
9     for (int j = 0; j < i; j++) {
10
11       // Dynamic carried dependency on s
12       // causes II = 1 or 118
13       if (b[i][j] != 0)
14         s = g(s, b[i][j]);
15
16     }
17
18     // Dynamic memory dependency distance
19     // causes dynamic II
20     // d >= 256
21     a[h(i)] = s;
22
23   }
24 }

```

(a) Source code.



(b) Area and performance using different approaches.

Fig. 12. An example where both the inner loop and the outer loop have a dynamic carried dependency, leading to dynamic IIs. Choosing  $C$  from average IIs achieves the best performance and has the minimum area among dynamically scheduled hardware. The results are measured from uniformly distributed data.

inner loop instances. When the program is dynamically scheduled, the II of an inner loop may vary at runtime. Here, we use probabilistic analysis to statically infer an optimised  $C$  from the average II.

Figure 12(a) shows a code example where the dependencies in both the inner loop and the outer loop are dynamic. In the outer loop, the function loads and computes an element in array  $a$  at an index of  $f(i)$ . It then updates the element in the same array at an index of  $h(i)$ . In the inner loop, a variable  $s$  updates itself based on the elements in matrix  $b$ , which is used to update array  $a$  in the outer loop.

In the outer loop, the memory dependencies between the loads and stores with array  $a$  are dynamic and depend on the iteration index of the out-loop  $i$ . For simplicity, assume that the minimum dependency distance in the outer loop for given  $f(i)$  and  $g(i)$  is no less than 256. That is,  $C$ -slow pipelining is valid for any  $C$  where  $1 \leq C \leq 256$ .

In the inner loop, the data-dependent condition causes two possible IIs. When the condition at line 11 is false, function  $g$  is skipped, leading to an II of 1. When the condition is true, function  $g$  is executed and the carried dependency on  $s$  causes an II of 118. The II of 118 is contributed by the latency between the input and the output of function  $g$ . The overall throughput then depends on the distribution of elements in matrix  $b$ , which affects the condition.

We then have 256 options to choose  $C$  no greater than 256, limited by the dependency in the outer loop. Here, we discuss three approaches for choosing  $C$  based on the II of the loop. First, choosing a  $C$  based on an optimistic II reduces the area overhead of adding the scheduler for  $C$ -slow pipelining. However, a small  $C$  may limit the parallelism among inner loop instances when there are more empty slots for certain input data. For this example, an extreme case is where the minimum II is 1, which disables  $C$ -slow pipelining ( $C = 1$  is equivalent to a single thread). Second, choosing a  $C$  based on a conservative II enables sufficient pipeline slots for early execution following inner loop instances. However, a large  $C$  may add unnecessary area overhead when there is only a small number or none of the empty pipeline slots, as illustrated in Figure 11(b). Finally, choosing a  $C$  based on an average II may balance the tradeoff between area overhead and performance improvement, achieving a more efficient hardware design. Our toolflow reuses the

results of the throughput analysis during the Dynamic synthesis flow. The buffering process in Dynamic already uses static analysis to estimate the II of each control path.

Figure 12(b) shows the results of the area and performance of the example in Figure 12(a) using different pipeline approaches. First, we evaluate three baselines: vanilla Dynamic [33] for dynamic scheduling, Vivado HLS [51] for static scheduling, naive C-slow pipelining with only dependency constraints [17]. Vanilla Dynamic allows pipelining loop with dynamic dependency, achieving better performance than static scheduling. However, the use of load-store queues that dynamically schedules memory operations causes significant area overhead. However, Vivado HLS that uses static scheduling cannot resolve data-dependent dependencies at compile time and keeps the loop sequential. The resultant hardware design below vanilla Dynamic has poor performance but high area efficiency because of resource sharing at compile time. Finally, the naive C-slow pipelining only analyses the dependency distance in the outer loop for choosing  $C$ . It generates hardware sitting on the top left that has better performance than both, because it allows early execution of later inner loop instances. However, the value of  $C$  is over-approximated to a large value. This results in unnecessary area overhead, since only a small portion of the inserted C-slow buffer in the control path is used.

We also evaluate the three approaches mentioned above for choosing an optimised  $C$ . First, the minimum II indicates that there is no empty slot in the inner loop schedule, preventing the transformation for C-slow pipelining. This leads to the same design as vanilla Dynamic. The case where these two points overlap is only for this particular benchmark, where the minimum II of 1. Otherwise, the result with the minimum II greater than 1 may not overlap with the result by vanilla Dynamic. Second, the maximum II indicates the best case that maximises parallelism for c-slow pipelining. However, this could still cause unnecessary area overhead when the iterations that have the maximum II are rare, such as in the case where only one iteration has the maximum II and the rest have the minimum II. Finally, the average II estimates the overall throughput of the inner loop. Such analysis reduces the C-slow buffer size while preserving sufficient slots for parallelism and maintaining the high performance of the naive C-slow approach. The constraint for an optimised  $C$  is then:

$$C \leq \min(Q, II_{av}), \quad (29)$$

where  $Q$  is the maximum  $C$  that passes the Boogie verification (also known as minimum dependency distance), and  $II_{av}$  is the average II of the inner loop. Detailed algorithm implementation is shown in Algorithm 1.

Here, we take the motivating example as a case study and then discuss the overall results for all the benchmarks. Figure 13 shows the total clock cycles of the hardware for the motivating example with different  $C$ . Only  $C \leq 6$  for this example does not break memory dependency, where  $C_D = 6$ . When  $C$  increases initially, more outer-loop iterations are parallelized, significantly improving the throughput. The average II of the inner loop is 5. When  $C$  is greater than 5, the throughput remains the same, since almost all the empty pipeline slots have been filled. The overhead caused by  $C = 6$  is a larger size of FIFOs used for storing the data.

### 5.3 Hardware Transformation

Once the value of  $C$  for a loop is determined, our toolflow inserts a component named *loop scheduler* between the entry and exit of each loop. Each  $C$  is used as a parameter of the corresponding loop scheduler. The loop scheduler dynamically schedules the control flow and ensures that at most  $C$  iterations can execute concurrently. Any outermost loop or unverified loop has its loop scheduler holding  $C = 1$  and executes control flow sequentially.

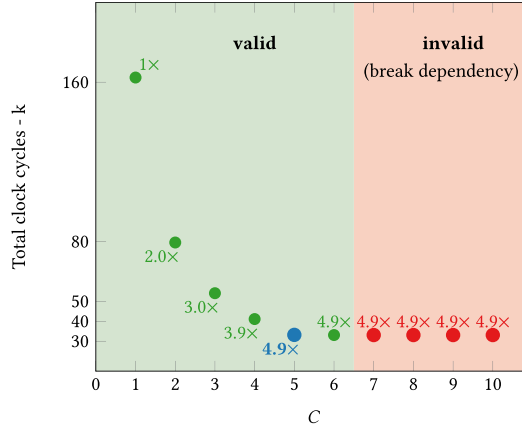


Fig. 13. Speedup by varying  $C$  for the example in Figure 9.  $C > 6$  breaks the memory dependency in the outer loop. Increasing  $C$  initially improves the throughput. However, once all the empty slots are filled, further increasing  $C$  has less effect on the throughput. Our tool automatically determines an optimal  $C = 5$ , shown in blue.  $f(i) = i$ ,  $g(x, y) = x + y$  and  $h(i) = i * i + 7$ .

---

**ALGORITHM 1:** The algorithm for finding  $C$  using average  $II$ .

---

**Require:**  $M$  ▷ the input program module  
 $C_L \leftarrow \{\}$  ▷ a dictionary to return, which indicates an  $C$  for each nested loop  
 $L \leftarrow \text{get\_nested\_loops}(M)$  ▷ a set of nested loops from the input program  
**for** each nested loop  $l$  in  $L$  **do**  
   $L' \leftarrow \text{get\_sub\_loops}(l)$  ▷ a set of sub-loops for loop  $l$   
   $II \leftarrow \text{get\_average\_II}(L')$  ▷ the set of average  $II$ s of the sub-loops in loop  $l$   
   $C \leftarrow \min II$  ▷ the starting  $C$  for searching by picking the minimum value of  $II$ s  
   $B \leftarrow \text{get\_Boogie\_program}(M, l)$  ▷ the Boogie program for checking  $C$  for loop  $l$   
   $D \leftarrow \text{Boogie\_verify}(B, C)$  ▷ the verification state returned by Boogie verifier  
  **while**  $C > 1$  &  $D \neq \text{success}$  **do**  
     $C \leftarrow C - 1$  ▷ a decremented  $C$  for next-round verification  
     $D \leftarrow \text{Boogie\_verify}(B, C)$   
  **end while**  
   $C_L(l) = C$  ▷ the optimised  $C$  for loop  $l$   
**end for**  
**return**  $C_L$

---

Figure 14 shows the proposed loop scheduler integrated into a dynamically scheduled control flow graph. For example, the control flow graph of the code in Figure 9 from vanilla Dynamatic is shown in Figure 14(a). Each dotted block represents a BB. The top BB represents the entry control of the outer loop, which starts the outer iteration and decides whether to execute the inner loop. The middle BB represents the control of the inner loop. The bottom BB represents the exit control of the outer loop, which decides whether to exit the outer loop.

In each BB, a merge is used to accept a control token that triggers the start of the current BB execution. Then, a fork is used to produce other tokens to trigger all the data operations inside this BB, hidden in the ellipsis. The control token flows through the fork to a branch. The branch decides the next BB to trigger based on the BB condition. The control flow in Figure 14(a) follows the following steps:

- (1) A control token enters the top BB to start.
- (2) The token goes through the top BB and enters the middle BB to start the inner loop.



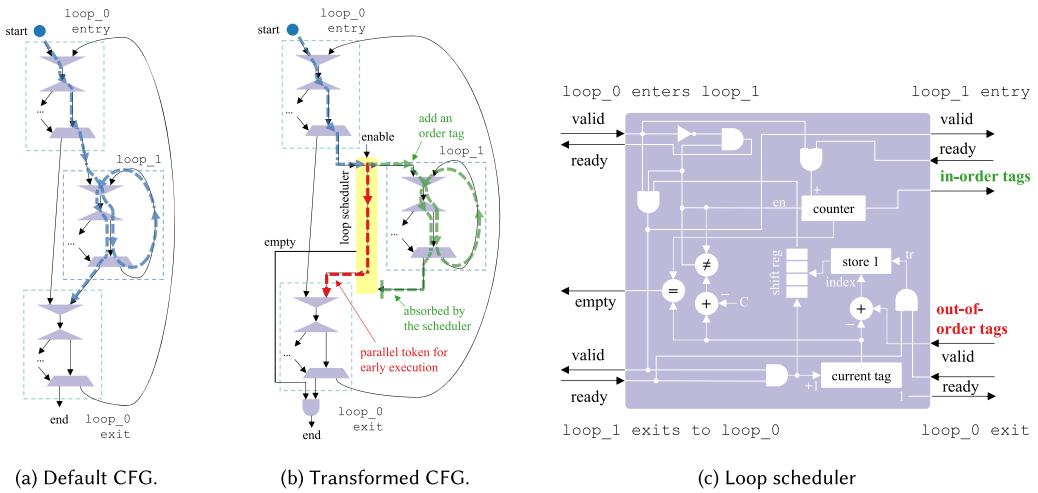


Fig. 14. Our toolflow considers each instance of the innermost loop as a thread and achieves the schedule in Figure 9(c). The dashed arrows represent the token transition in the control flow. The scheduler tags the control tokens in the innermost loop to reorder them at the output after out-of-order execution.

- (3) The token circulates in the middle BB through the back edge until the exit condition is met.
- (4) The token exits the middle BB and enters the bottom BB. It either goes back to the top BB to repeat 2) or exits, depending on the exit condition.

The control flow is sequential, as there is always at most one control token in the control path. Figure 14(b) shows the control flow graph with the proposed loop scheduler integrated into the inner loop. The loop scheduler for the outer loop has  $C$  of 1 and is neglected for simplicity. The control flow is then:

- (1) A control token  $t_1$  enters the top BB to start.
- (2)  $t_1$  goes through the top BB and enters the middle BB with a tag added by the loop scheduler. The loop scheduler checks if there are fewer than  $C(=3$  in this example) tokens in the inner loop. If yes, then it immediately produces another token  $t_2$  and sends it to the bottom BB to execute the control flow early (indicated as the red dashed arrow).
- (3)  $t_1$  circulates in the middle BB.  $t_2$  goes to the top BB and enters the middle BB with another tag. The loop scheduler produces another token  $t_3$  and sends it to the bottom BB.
- (4) The above repeats and  $t_4$  is produced.  $t_1, t_2,$  and  $t_3$  are all circulating in the middle BB.
- (5)  $t_4$  reaches the branch in the top BB but is *blocked* by the loop scheduler until one token exits the middle BB and is consumed by the loop scheduler.
- (6) The above repeats until the last token exits the bottom BB. An AND gate is inserted at the exit of the bottom BB to synchronise the control flow. It requires a token from the exit of the nested loop, and there is no token remaining in the inner loop.

Since the execution order of loop iterations has changed, all the data flows must be strictly scheduled by the control flow. Dynamic uses merges to accept input data at the input of a BB when there is always at most one valid input. There may be multiple valid inputs at the start of BB after the transformation. To preserve correctness, we replace all the merges in the inner loop with muxes, such that the data is always synchronised with the control token and can recover in-order using the tag in the control token. An advantage of this approach is that only the control tokens

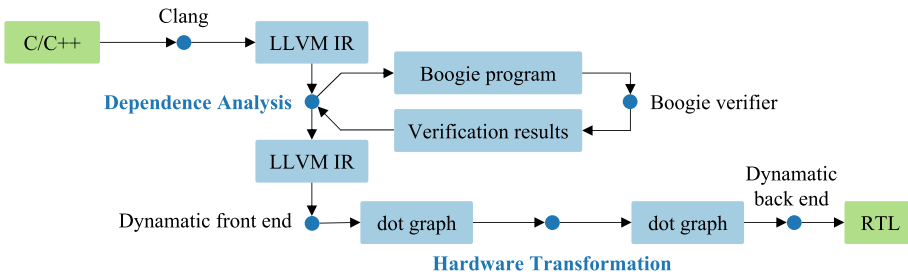


Fig. 15. Our work integrated into Dynamatic. Our contributions are highlighted in bold blue text.

need to be tagged to preserve the original order, where the data flow is always synchronised by the control flow.

The design of the loop scheduler is shown in Figure 14(c). It guards the entry and exit of the inner loop. The ready signal at the bottom right is forced to 1, as the scheduler already controls the input throughput using  $C$ , and there cannot be back pressure at the output. In the scheduler, a counter is used to count the number of executing control tokens in the inner loop. Based on the value of the counter and the specified  $C$ , it decides whether to accept the token from the outer loop and replicates a token to the output to the outer loop for early execution of the next outer-loop iteration. The input token from the exit of the inner loop decrements the counter value by one, allowing the next token from the outer loop to enter the inner loop. An empty bit is used to indicate whether there is no control token in the inner loop.

## 6 TOOLFLOW

Our toolflow is implemented as a set of LLVM passes and integrated into the open-sourced HLS tool Dynamatic for prototyping. As illustrated in Figure 15, the input C program is first lowered into LLVM IR. Then, the dependency in the code is analysed by the generated Boogie program, as explained in Section 4.2.2 and Section 5.2.2. Our Boogie program generator generates Boogie assertions and calls the Boogie verifier to automatically verify the absence of dependency for the extracted instances (subgraphs or outer iterations). The front end of Dynamatic translates the LLVM IR into a dot graph that represents the hardware netlist. Our back-end toolflow inserts additional components for the corresponding transformation, explained in Section 4.3 and Section 5.3, resulting in a new hardware design in the form of a dot graph. Finally, the back end of Dynamatic transforms the new dot graph to RTL code, representing the final hardware design. Our work can also be integrated into other HLS tools, such as CIRCT HLS [21].

## 7 EXPERIMENTS

We compare our work with Xilinx Vivado HLS [51] and the original Dynamatic [33]. To make the comparison as controlled as possible, all the approaches only use scheduling, pipelining, and array partitioning. We use two benchmark sets to evaluate the designs in terms of total circuit area and wall-clock time. Cycle counts were obtained using the Vivado XSIM simulator, and area results were obtained from the post-Place & Synthesis report in Vivado. We used the UltraScale+ family of FPGA devices for experiments, and the version of Xilinx software is 2019.2.

### 7.1 Experiment Setup and Benchmarks

The benchmarks are chosen based on whether our approaches are applicable. Finding suitable benchmarks is a perennial problem for papers that push the limits of HLS, in part because existing

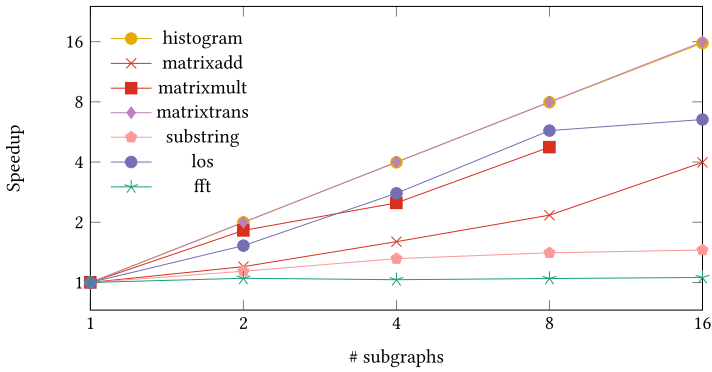


Fig. 16. Speedup, compared to original Dynamic, as more subgraphs in the CFG are parallelised.

benchmark sets such as PolyBench [45] and CHStone [27] tend to be tailored to what HLS tools can already comfortably handle. We use two open-sourced benchmark sets for evaluation. One is the LegUp benchmark set by Chen and Anderson [12] for evaluating multi-threaded HLS. The LegUp benchmark set manually specifies the threads using Pthreads [2]. We inlined all the threads to a sequential program. The other benchmark set is from PolyBench [45] but modified for sparse computation for evaluating dynamic-scheduling HLS. We only include the benchmarks where our approach is applicable. The other benchmarks will have the same results as the original Dynamic. The second benchmark set aims to evaluate dynamic loop pipelining and contains few loop kernels. To create more opportunities for our optimisation to be applied, we unrolled the outermost loops by a factor of 8. This is the largest factor that still led to the designs fitting our target FPGA. We also partitioned the memory in the blocking scheme to increase memory bandwidth. The benchmarks that we used are listed as follows: *histogram* constructs a histogram from an integer array, *matrixadd* sums a float array, *matrixmult* multiplies two float matrices, *matrixtrans* transposes a single matrix, *substring* searches for a pattern in an input string, *los* checks for obstacles on a map, *fft* performs the fast Fourier transformation, *trVecAccum* transforms a triangular matrix, *covariance* computes the covariance matrix, *syr2k* is a symmetric rank-2k matrix update, and *gesummv* is scalar, vector, and matrix multiplication.

## 7.2 Dynamic Inter-block Scheduling

Figure 16 assesses the extent to which more parallelisation of subgraphs leads to more speedups compared to the original Dynamic, using the seven LegUp benchmarks. We see that all the lines except *fft* indicate speedup factors above 1. Placing more subgraphs in parallel leads to more speedup, with *histogram* and *matrixtrans* achieving optimal speedups. In the *fft* benchmark, two reasons for the lack of speedup are: (1) that other parts of the CFG have to be started sequentially and (2) that the memory is naively partitioned in a block scheme, so the memory bandwidth is limited, and there is serious contention between BBs for the LSQs.

Detailed results for both benchmark sets are shown in Table 2. For the LegUp benchmark set, we observe the following:

- (1) Static scheduling (Vivado HLS) is the clear winner in terms of area (see rows “LUTs” and “DSPs”), but in the context of dynamic scheduling, our approach brings only a negligible area overhead, because we only insert small components into the hardware.
- (2) Inter-block scheduling requires substantially fewer cycles than the original Dynamic, thanks to the parallelism it exploits between BBs (see column “Cycles”).

Table 2. Evaluation of Our Work on Two Benchmark Sets

Benchmarks		LegUp benchmarks [12]						C-slow benchmarks [17]					
		histogram	matrixadd	matrixmult	matrixtrans	substring	los	fft	trVecAccum	covariance	syr2k	gesummv	
Code size	loops	9	8	72	8	16	24	24	16	48	24	16	
	bbs	91	17	145	17	54	89	65	49	113	49	33	
	insts	384	112	1,521	81	255	513	457	249	593	385	297	
	graphs	9	8	72	8	8	8	8	8	24	8	8	
	both	156	9.15	101	2.67	14.4	46.5	351	151	56.8	30.6	27.6	
LUTs (1,000s)	vhls	<b>1.67</b>	<b>1.11</b>	<b>6.87</b>	<b>0.103</b>	<b>0.938</b>	<b>2.68</b>	<b>2.65</b>	<b>1.21</b>	<b>3.83</b>	<b>2.04</b>	<b>2.05</b>	
	dhls	156	9.15	79.4	2.67	14.4	46.5	351	149	56.6	30.5	26.9	
	inter-block	156	9.17	79.8	2.69	14.6	46.1	351	150	55.5	30.7	26.9	
	cslow	156	9.15	101	2.67	14.4	46.5	351	151	56.8	30.6	27.6	
	both	156	9.17	100	2.69	14.6	46.1	351	153	65.7	34.4	30.2	
DSPs	vhls	<b>0</b>	<b>2</b>	<b>5</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>16</b>	<b>10</b>	<b>5</b>	<b>5</b>	<b>144</b>	
	dhls	0	30	320	0	0	0	192	40	72	152	144	
	inter-block	0	30	320	0	0	0	192	40	72	152	144	
	cslow	0	30	320	0	0	0	192	40	72	152	144	
	both	0	30	320	0	0	0	192	40	72	152	144	
Cycles (1,000s)	vhls	197	262	4,195	65.6	98.3	48.8	86	1,060	668	647	2,130	
	dhls	317	106	1,090	65.6	217	114	5.39	393	605	602	787	
	inter-block	<b>39.8</b>	<b>48.9</b>	229	<b>8.2</b>	<b>154</b>	<b>19.9</b>	<b>5.15</b>	161	77.1	84.1	327	
	cslow	317	106	<b>164</b>	65.6	217	114	5.39	256	263	255	524	
	both	<b>39.8</b>	<b>48.9</b>	<b>164</b>	<b>8.2</b>	<b>154</b>	<b>19.9</b>	<b>5.15</b>	<b>33.2</b>	<b>33.6</b>	<b>33.9</b>	<b>66.6</b>	
Fmax (MHz)	vhls	<b>464</b>	<b>159</b>	<b>155</b>	<b>562</b>	<b>470</b>	281	<b>155</b>	<b>159</b>	<b>155</b>	<b>155</b>	155	
	dhls	57.7	110	123	227	126	272	81.8	132	132	126	162	
	inter-block	58.3	110	104	210	129	<b>282</b>	103	121	102	98.9	<b>163</b>	
	cslow	57.7	110	83.3	227	126	272	81.8	157	89.7	130	119	
	both	58.3	110	72.3	210	129	<b>282</b>	103	117	102	128	120	
Wall-clock time (ms)	vhls	<b>0.424</b>	1.65	27	0.117	<b>0.209</b>	0.174	0.555	6.65	4.3	4.17	13.7	
	dhls	5.49	0.968	8.86	0.288	1.72	0.419	0.0659	2.97	4.57	4.79	4.87	
	inter-block	0.682	<b>0.446</b>	2.2	<b>0.039</b>	1.2	<b>0.0705</b>	<b>0.0501</b>	1.32	0.759	0.85	2.01	
	cslow	5.49	0.968	<b>1.97</b>	0.288	1.72	0.419	0.0659	1.64	2.93	1.96	4.39	
	both	0.682	<b>0.446</b>	2.27	<b>0.039</b>	1.2	<b>0.0705</b>	<b>0.0501</b>	<b>0.284</b>	<b>0.328</b>	<b>0.264</b>	<b>0.553</b>	
Relative area-delay product	vhls	<b>1</b>	<b>1</b>	1	<b>1</b>	<b>1</b>	<b>1</b>	<b>1</b>	<b>1</b>	<b>1</b>	<b>1</b>	1	
	dhls	1,210	4.83	3.79	64	126	41.8	15.7	55	15.7	94	4.67	
	inter-block	151	2.23	<b>0.946</b>	8.71	88.9	6.96	12	24.7	2.56	16.5	1.93	
	cslow	1,210	4.83	1.07	64	126	41.8	15.7	40.3	9.61	39	5.17	
	both	151	2.23	1.22	8.71	88.9	6.96	12	5.41	1.31	5.12	<b>0.595</b>	

For each benchmark, we highlight the **best** results in each dimension. vhls = Vivado HLS; dhls = original Dynamic; cslow = C-slow pipelining; inter-block = inter-block scheduling; both = inter-block scheduling + C-slow pipelining. The code size lists the number of loops, BBs, instructions, and extracted subgraphs.

- (3) Inter-block scheduling achieves up to 8.05× speedup (see row “Wall clock time” for histogram). We further observe that the area-delay products we obtain are down to 0.125× of the original Dynamic.
- (4) Although Vivado HLS has low performance in cycles, its high clock frequency makes it win for histogram and substring. Also, it uses if-conversion to simplify BBs (unlike our work), which results in fewer BBs. The BBs in the innermost subgraphs still start sequentially, leading to large cycle counts.
- (5) The performance of fft does not change, because the performance bottleneck is the memory bandwidth, where the parallelised loops always access the same memory block concurrently. This could be further improved by optimising the array partitioning scheme, but it is orthogonal to this work.

### 7.3 C-slow Pipelining

For the C-slow pipelining benchmark set, we make the following observations for C-slow pipelining:

- (1) The observations on the area are similar to the results for inter-block scheduling, as the C-slow scheduler is small.

Table 3. Verification Time in Seconds for Dependence Check in Boogie

Benchmarks	histogram	matrixadd	matrixmult	matrixtrans	substring	los	fft	triangleVecAccum	covariance	syr2k	gesummv
Inter-Block	1.128	0.002	1.080	1.102	1.122	1.197	1.182	1.080	1.098	1.210	1.103
C-Slow	-	-	-	-	-	-	-	1.215	1.098	1.088	1.050

- (2) C-slow pipelining requires substantially fewer cycles than the original Dynamatic, thanks to the parallelism it exploits between outer-loop iterations (see rows “Cycles”).
- (3) C-slow pipelining enables a speedup up to 4.5× with only 8% area overhead (see rows “Wall clock time” for *matrixmult*). We further observe that the area-delay products we obtain are down to 0.28× of the original Dynamatic.

#### 7.4 Combining Inter-block Scheduling and C-slow Pipelining

For the C-slow pipelining benchmark set, we also make the following observations for combining both approaches:

- (1) The area-delay product of Dynamatic is significantly worse than Vivado HLS, because the version of Dynamatic we used does not support resource sharing, leading to significant area overhead (although it is now supported [34]).
- (2) Unrolling alone is not enough to obtain substantial speedups, because the BBs still have to start sequentially (see column “Cycles → unroll”).
- (3) By applying both techniques simultaneously on the unrolled programs, we achieve a 14.3× average speedup with a 10% area overhead. That significant speedup can be attributed in part to the reordering of BBs.

#### 7.5 Verification Times

Table 3 lists the verification time taken by the Boogie verifier to check the absence of dependency. It shows that for all the benchmarks, the verification time takes no more than two seconds in a run. The verification time scales exponentially with the number of memory statements in the input programs and the complexity of the memory access pattern. However, optimisations such as profiling and affine analysis could simplify the formal verification problem in Boogie for better scalability.

## 8 CONCLUSIONS

Existing dynamic-scheduling HLS tools require all BBs to start in strict program order to respect any dependencies between BBs, regardless of whether dependencies are actually present. This leads to missed opportunities for performance improvements by having BBs start simultaneously. In this article, we propose a general dependency framework that analyses the inter-BB dependencies and identifies the existing restrictions of the dynamic scheduling approach for HLS. Our model helps guide the researchers to push the limit of the state-of-the-art scheduling technique by lifting these restrictions and exploring more optimisations. This could further improve the performance of the generated HLS hardware by identifying the absence of certain dependencies in the constraints in static analysis and exploiting hardware parallelism.

We also illustrate two examples of using the proposed model for optimising dynamic scheduling. First, we show how to parallelise two consecutive loops if they are proven independent. Our tool takes an arbitrary program and automatically generates a Boogie program to verify the absence of dependency between these two loops. Second, we show how to enable C-slow pipelining for a nested loop if both the II of the inner loop and the minimum dependency distance of the outer loop are greater than one. Our tool takes a nested loop and automatically generates a Boogie

program to verify the dependency constraints and uses static throughput analysis for determining an optimised  $C$  for high-performance and area-efficient C-slow pipelining.

These two proposed optimisation techniques optimise the control flow of the generated hardware from different sides and have shown that they can be composited. One aims to optimise the schedules between different loop statements, and the other aims to optimise the schedules between different iterations of the same loop statement. However, both techniques achieve hardware parallelism by statically proving the absence of dependency between two runtime events. Particularly, both use Microsoft Boogie verifier as the formal verification tool in the back end for verifying the dependency constraints statically. The performance gain is significant, while the area overhead is negligible because of more dependency information obtained using the static analysis. Such static analysis is still conservative compared to the theoretical dynamic scheduling model, because it over-approximates the dependency check for two runtime events to the dependency check for two statements. The proposed approaches close the performance gap between the dynamic scheduling implementation and the theoretical dynamic scheduling. Our future work will explore the fundamental limits of the proposed model, both theoretically and practically.

## ACKNOWLEDGMENT

For the purpose of open access, the authors have applied a Creative Commons Attribution (CC BY) license to any Accepted Manuscript version arising.

## REFERENCES

- [1] Amazon. 2022. Amazon EC2 F1 Instances. Retrieved from <https://aws.amazon.com/ec2/instance-types/f1/>.
- [2] B. Barney. 2021. POSIX Threads Programming. Retrieved from <https://computing.llnl.gov/tutorials/pthreads>.
- [3] Uday Bondhugula, Albert Hartono, J. Ramanujam, and P. Sadayappan. 2008. A practical automatic polyhedral parallelizer and locality optimizer. In *29th ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI'08)*. Association for Computing Machinery, New York, NY, 101–113. DOI: <https://doi.org/10.1145/1375581.1375595>
- [4] David R. Butenhof. 1997. *Programming with POSIX Threads*. Addison-Wesley Professional.
- [5] Daniel Cabrera, Xavier Martorell, Georgi Gaydadjiev, Eduard Ayguade, and Daniel Jiménez-González. 2009. OpenMP extensions for FPGA accelerators. In *International Symposium on Systems, Architectures, Modeling, and Simulation*. IEEE, 17–24.
- [6] A. Canis, S. D. Brown, and J. H. Anderson. 2014. Modulo SDC scheduling with recurrence minimization in high-level synthesis. In *24th International Conference on Field Programmable Logic and Applications (FPL'14)*. 1–8. DOI: <https://doi.org/10.1109/FPL.2014.6927490>
- [7] Andrew Canis, Jongsok Choi, Mark Aldham, Victor Zhang, Ahmed Kammoona, Jason H. Anderson, Stephen Brown, and Tomasz Czajkowski. 2011. LegUp: High-level synthesis for FPGA-based processor/accelerator systems. In *19th ACM/SIGDA International Symposium on Field Programmable Gate Arrays (FPGA'11)*. ACM, 33–36.
- [8] L. P. Carloni, K. L. McMillan, and A. L. Sangiovanni-Vincentelli. 2001. Theory of latency-insensitive design. *IEEE Trans. Comput.-aid. Des. Integ. Circ. Syst.* 20, 9 (2001), 1059–1076. DOI: <https://doi.org/10.1109/43.945302>
- [9] V. G. Castellana, A. Tumeo, and F. Ferrandi. 2014. High-level synthesis of memory bound and irregular parallel applications with Bambu. In *IEEE Hot Chips 26 Symposium (HCS'14)*. IEEE.
- [10] Catapult High-Level Synthesis. 2021. Retrieved from <https://www.mentor.com/hls-lp/catapult-high-level-synthesis>.
- [11] Celoxica. 2005. Handel-C. Retrieved from <http://www.celoxica.com>.
- [12] Yu Ting Chen and Jason H. Anderson. 2017. Automated generation of banked memory architectures in the high-level synthesis of multi-threaded software. In *27th International Conference on Field Programmable Logic and Applications (FPL'17)*. 1–8. DOI: <https://doi.org/10.23919/FPL.2017.8056841>
- [13] J. Cheng, S. T. Fleming, Y. T. Chen, J. Anderson, J. Wickerson, and G. A. Constantinides. 2022. Efficient memory arbitration in high-level synthesis from multi-threaded code. *IEEE Trans. Comput.* 71, 4 (2022), 933–946. DOI: <https://doi.org/10.1109/TC.2021.3066466>
- [14] Jianyi Cheng, Lana Josipović, George A. Constantinides, Paolo Jenne, and John Wickerson. 2020. Combining dynamic & static scheduling in high-level synthesis. In *ACM/SIGDA International Symposium on Field-Programmable Gate Arrays (FPGA'20)*. Association for Computing Machinery, New York, NY, 288–298. DOI: <https://doi.org/10.1145/3373087.3375297>

- [15] Jianyi Cheng, Lana Josipović, George A. Constantinides, and John Wickerson. 2022. Dynamic inter-block scheduling for HLS. In *32nd International Conference on Field-Programmable Logic and Applications (FPL'22)*.
- [16] Jianyi Cheng, John Wickerson, and George A. Constantinides. 2021. Exploiting the correlation between dependence distance and latency in loop pipelining for HLS. In *31st International Conference on Field-Programmable Logic and Applications (FPL'21)*. 341–346. DOI: <https://doi.org/10.1109/FPL53798.2021.00066>
- [17] Jianyi Cheng, John Wickerson, and George A. Constantinides. 2022. Dynamic C-slow pipelining for HLS. In *IEEE 30th Annual International Symposium on Field-Programmable Custom Computing Machines (FCCM'22)*. 1–10. DOI: <https://doi.org/10.1109/FCCM53951.2022.9786096>
- [18] Jongsok Choi, Stephen Brown, and Jason Anderson. 2013. From software threads to parallel hardware in high-level synthesis for FPGAs. In *International Conference on Field-Programmable Technology (FPT'13)*. 270–277. DOI: <https://doi.org/10.1109/FPT.2013.6718365>
- [19] N. Y. S. Chong. 2014. *Scalable Verification Techniques for Data-parallel Programs*. Doctoral Thesis. Imperial College London, London, UK.
- [20] Jason Cong and Zhiru Zhang. 2006. An efficient and versatile scheduling algorithm based on SDC formulation. In *43rd Annual Design Automation Conference (DAC'06)*. Association for Computing Machinery, New York, NY, 433–438. DOI: <https://doi.org/10.1145/1146909.1147025>
- [21] CIRCT contributors. 2023. CIRCT-based HLS Compilation Flows, Debugging, and Cosimulation Tools. Retrieved from <https://github.com/circt-hls/circt-hls>.
- [22] Philippe Coussey, Daniel D. Gajski, Michael Meredith, and Andres Takach. 2009. An introduction to high-level synthesis. *IEEE De. Test Comput.* 26, 4 (2009), 8–17. DOI: <https://doi.org/10.1109/MDT.2009.69>
- [23] L. Dagum and R. Menon. 1998. OpenMP: An industry standard API for shared-memory programming. *IEEE Computat. Sci. Eng.* 5, 1 (1998), 46–55. DOI: <https://doi.org/10.1109/99.660313>
- [24] Leonardo de Moura and Nikolaj Bjørner. 2008. Z3: An efficient SMT solver. In *Tools and Algorithms for the Construction and Analysis of Systems*, C. R. Ramakrishnan and Jakob Rehof (Eds.). Springer Berlin, 337–340.
- [25] Tobias Grosser, Hongbin Zheng, Raghesh Aloor, Andreas Simbürger, Armin Größlinger, and Louis-Noël Pouchet. 2011. Polly—Polyhedral optimization in LLVM. In *1st International Workshop on Polyhedral Compilation Techniques (IMPACT'11)*.
- [26] S. Gupta, N. Dutt, R. Gupta, and A. Nicolau. 2003. SPARK: A high-level synthesis framework for applying parallelizing compiler transformations. In *16th International Conference on VLSI Design*. 461–466. DOI: <https://doi.org/10.1109/ICVD.2003.1183177>
- [27] Yuko Hara, Hiroyuki Tomiyama, Shinya Honda, Hiroaki Takada, and Katsuya Ishii. 2008. CHStone: A benchmark program suite for practical C-based high-level synthesis. In *IEEE International Symposium on Circuits and Systems (ISCAS'08)*. 1192–1195. DOI: <https://doi.org/10.1109/ISCAS.2008.4541637>
- [28] Ian Page and Wayne Luk. 1991. Compiling occam into field-programmable gate arrays. In *FPGAs, Oxford Workshop on Field Programmable Logic and Applications*, Vol. 15. Abingdon EE&CS Books, Abingdon, 271–283.
- [29] Intel Compiler. 2022. Retrieved from <https://www.intel.com/content/www/us/en/developer/tools/oneapi/dpc-compiler.html#gs.sa60u7>.
- [30] Intel FPGA SDK for OpenCL Software Technology. 2021. Retrieved from <https://www.intel.co.uk/content/www/uk/en/software/programmable/sdk-for-opencl/overview.html>.
- [31] Intel HLS Compiler. 2022. Retrieved from <https://www.intel.co.uk/content/www/uk/en/software/programmable/quartus-prime/hls-compiler.html>.
- [32] Lana Josipović, Philip Brisk, and Paolo Ienne. 2017. An out-of-order load-store queue for spatial computing. *ACM Trans. Embed. Comput. Syst.* 16, 5s (Sept. 2017).
- [33] Lana Josipović, Radhika Ghosal, and Paolo Ienne. 2018. Dynamically scheduled high-level synthesis. In *ACM/SIGDA International Symposium on Field-Programmable Gate Arrays (FPGA'18)*. ACM, 127–136.
- [34] Lana Josipović, Axel Marmet, Andrea Guerrieri, and Paolo Ienne. 2022. Resource sharing in dataflow circuits. In *IEEE 30th Annual International Symposium on Field-Programmable Custom Computing Machines (FCCM'22)*. 1–9. <https://doi.org/10.1109/FCCM53951.2022.9786084>
- [35] Lana Josipović, Andrea Guerrieri, and Paolo Ienne. 2022. From C/C++ code to high-performance dataflow circuits. *IEEE Trans. Comput.-aid Des. Integ. Circ. Syst.* 41, 7 (2022), 2142–2155. DOI: <https://doi.org/10.1109/TCAD.2021.3105574>
- [36] K. Rustan M. Leino. 2008. This is Boogie 2. Retrieved from <https://www.microsoft.com/en-us/research/publication/this-is-boogie-2-2/>.
- [37] Charles E. Leiserson, Flavio M. Rose, and James B. Saxe. 1983. Optimizing synchronous circuitry by retiming (preliminary version). In *Third Caltech Conference on Very Large Scale Integration*. Springer, 87–116.
- [38] Y. Y. Leow, C. Y. Ng, and W. F. Wong. 2006. Generating hardware from OpenMP programs. In *IEEE International Conference on Field Programmable Technology*. 73–80. DOI: <https://doi.org/10.1109/FPT.2006.270297>

- [39] Rui Li, Lincoln Berkley, Yihang Yang, and Rajit Manohar. 2021. Fluid: An asynchronous high-level synthesis tool for complex program structures. In *27th IEEE International Symposium on Asynchronous Circuits and Systems (ASYNC'21)*. 1–8. DOI : <https://doi.org/10.1109/ASYNC48570.2021.00009>
- [40] J. Liu, J. Wickerson, and G. A. Constantinides. 2016. Loop splitting for efficient pipelining in high-level synthesis. In *IEEE 24th Annual International Symposium on Field-Programmable Custom Computing Machines (FCCM'16)*. 72–79. DOI : [urlhttps://doi.org/10.1109/FCCM.2016.27](https://doi.org/10.1109/FCCM.2016.27)
- [41] Q. Liu, G. A. Constantinides, K. Masselos, and P. Y. K. Cheung. 2007. Automatic on-chip memory minimization for data reuse. In *15th Annual IEEE Symposium on Field-Programmable Custom Computing Machines (FCCM'07)*. IEEE, 251–260.
- [42] Yury Markovskiy and Yatish Patel. 2002. Simple symmetric multithreading in xilinx FPGAs. (2002).
- [43] Microsoft. 2022. Project Catapult. Retrieved from <https://www.microsoft.com/en-us/research/project/project-catapult/>.
- [44] Sayuri Ota and Nagisa Ishiura. 2019. Synthesis of distributed control circuits for dynamic scheduling across multiple dataflow graphs. In *34th International Technical Conference on Circuits/Systems, Computers and Communications (ITC-CSCC'19)*. 1–4. DOI : <https://doi.org/10.1109/ITC-CSCC.2019.8793453>
- [45] Louis-Noël Pouchet et al. 2012. PolyBench: The polyhedral benchmark suite. Retrieved from <http://www.cs.ucla.edu/pouchet/software/polybench>.
- [46] Stratus High-Level Synthesis. 2021. Retrieved from [https://www.cadence.com/en\\_US/home/tools/digital-design-and-signoff/synthesis/stratus-high-level-synthesis.html](https://www.cadence.com/en_US/home/tools/digital-design-and-signoff/synthesis/stratus-high-level-synthesis.html).
- [47] Qiuyue Sun, Amir Taherin, Yawo Siatitse, and Yuhao Zhu. 2020. Energy-efficient 360-Degree video rendering on FPGA via algorithm-architecture co-design. In *ACM/SIGDA International Symposium on Field-Programmable Gate Arrays (FPGA'20)*. Association for Computing Machinery, New York, NY, 97–103. DOI : <https://doi.org/10.1145/3373087.3375317>
- [48] Girish Venkataramani, Mihai Budiu, Tiberiu Chelcea, and Seth Copen Goldstein. 2004. C to asynchronous dataflow circuits: An end-to-end toolflow. In *IEEE 13th International Workshop on Logic Synthesis (IWLS'04)*.
- [49] Jie Wang, Licheng Guo, and Jason Cong. 2021. AutoSA: A polyhedral compiler for high-performance systolic arrays on FPGA. In *ACM/SIGDA International Symposium on Field-Programmable Gate Arrays (FPGA'21)*. Association for Computing Machinery, New York, NY, 93–104. DOI : <https://doi.org/10.1145/3431920.3439292>
- [50] Nicholas Weaver, Yury Markovskiy, Yatish Patel, and John Wawrzynek. 2003. Post-placement C-slow retiming for the Xilinx Virtex FPGA. In *ACM/SIGDA 11th International Symposium on Field Programmable Gate Arrays (FPGA'03)*. Association for Computing Machinery, New York, NY, 185–194. DOI : <https://doi.org/10.1145/611817.611845>
- [51] Xilinx Vivado HLS. 2022. Retrieved from <https://www.xilinx.com/support/documentation/navigation/design-hubs/dh0012-vivado-high-level-synthesis-hub.html>.
- [52] Tanner Young-Schultz, Lothar Lilge, Stephen Brown, and Vaughn Betz. 2020. Using OpenCL to enable software-like development of an FPGA-accelerated biophotonic cancer treatment simulator. In *ACM/SIGDA International Symposium on Field-Programmable Gate Arrays (FPGA'20)*. Association for Computing Machinery, New York, NY, 86–96. DOI : <https://doi.org/10.1145/3373087.3375300>
- [53] Z. Zhang and B. Liu. 2013. SDC-based modulo scheduling for pipeline synthesis. In *IEEE/ACM International Conference on Computer-Aided Design (ICCAD'13)*. 211–218. DOI : <https://doi.org/10.1109/ICCAD.2013.6691121>
- [54] Yuan Zhou, Khalid Musa Al-Hawaj, and Zhiru Zhang. 2017. A new approach to automatic memory banking using trace-based address mining. In *ACM/SIGDA International Symposium on Field-Programmable Gate Arrays (FPGA'17)*. Association for Computing Machinery, New York, NY, 179–188. DOI : <https://doi.org/10.1145/3020078.3021734>
- [55] Wei Zuo, Yun Liang, Peng Li, Kyle Rupnow, Deming Chen, and Jason Cong. 2013. Improving high level synthesis optimization opportunity through polyhedral transformations. In *ACM/SIGDA International Symposium on Field Programmable Gate Arrays (FPGA'13)*. Association for Computing Machinery, New York, NY, 9–18. DOI : <https://doi.org/10.1145/2435264.2435271>

Received 3 November 2022; revised 1 April 2023; accepted 10 May 2023