



WISSENSCHAFTSRAT

# Dataflow and higher level abstractions for parallel programming

Jeronimo Castrillon

Chair for Compiler Construction (CCC)

TU Dresden, Germany

Cyber-Physical Systems Summer School 2019

Alghero, Sardinia, Italy. September 25, 2019

cfaed.tu-dresden.de





**DFG** 

WR

## **Evolution of Hardware**



CONSTRUCTION



2

#### **Evolution of Hardware: Great complexity**



Heterogeneous many-cores, scalable platforms, complex memory hierarchies, n C domain-specific accelerators, ... 2019





Panda, P. R., Dutt, N. D., & Nicolau, A. Memory issues in embedded systems-on-chip: optimizations and exploration. Springer Science & Business Media. 1999

Single thread Smart Smart health transport MEM Communication A15 A15 subsystem support HW queues VLIW DSP A15 A15 L1 L1 Infotainmen DMAs, sema-phores PMU Network L1.L2 system Processor NoC Packet DMA Peripherals Programmer Virtual 3D i/p Control

systems  $\sqrt{n}$ 

=

NLP

AI, ML

HCI

omputina, 2011

CHAIRFOR

COMPILER

CONSTRUCTION

Robotics

© Prof. J. Castrillon. CPS Summer School. Alghero, Italy, 2019

#### **Models: Von Neumann**



#### Program: sequence of instructions in memory



### **Models: Von Neumann and Modern HW**



- Different sources of parallelism
  - SIMD, threads, GPUs, ...
- Fundamentally different resources
  - Dataflow accelerators, Tensor processing units, ...

#### → Lots of effort in auto-parallelization and vectorization



**Theorem** (Allen/Kennedy): Any reordering transformation that preserves every dependence in a program preserves the meaning of that program

© Prof. J. Castrillon. CPS Summer School. Alghero, Italy, 2019



#### Static/dynamic control/data analysis: Sample results



Step	Speedup	No. of PEs	Parallel Efficiency
1	3.61x	16	22.58%
2	5.48x	17	32.3%
3	5.48x	16	34.3%
manual	9.43x	19	49.6%

Table 1: Summary of JPEG Encoder Parallelization by MAPS

#### 2008

Experimental multi-core from Tokyo Institute of Technology (Prof. Isshiki)

[DAC'08]



#### **Problems for auto-parallelizing compilers**





#### Problems for auto-parallelizing compilers (2)



1) Find all dependencies?

2) Coding style and the illusion of infinite shared memory

3) Dependencies cansometimes be violated!(they are artifacts of style)

```
while(!queue.empty())
19
   {
20
      // Dequeue a vertex from queue
21
      s = queue.front();
22
      queue.pop_front();
23
24
      // Apply function f to s. accumulate values
25
26
      result += f(s):
27
      // Get all adjacent vertices of s.
28
     // If an adjacent node hasn't been visited.
29
      // then mark it as visited and enqueue it
30
      for(i=adj[s].begin(); i!=adj[s].end(); ++i)
31
32
      {
        if(!visited[*i])
33
34
        {
          visited[*i] = true;
35
          queue.push_back(*i);
36
37
        3
38
      }
39
    }
40
    return result;
41
42 }
```

[Edler'18]



### **Models: Von Neumann and Modern HW**



- Different sources of parallelism
  - SIMD, threads, GPUs, ...
- Fundamentally different resources
  - Dataflow accelerators, Tensor processing unit, ...



#### Von Neumann and sequential programming

- Difficult to automatically extract parallelism
- Difficult to recognize high-level operations
- Difficult to ensure correctness (functional, timing, ...)



#### **Models and Programming**



Smart health

Aerospace

CI

CHAIRFOR COMPILER

CONSTRUCTION

AI, ML

- □ Specially important for CPS!
  - Domain specific

- Interconnected systems
- Interaction with physical world (Reactive)
- Real-time (WCET, predictable HW)
- Dynamic changing conditions

#### In this lecture

- □ Classic dataflow modeling
- Current work on extensions: dynamic, adaptable... (CPS)

=

Smart

transport

Infotainment

system

NLP

Raise-level of abstraction: Domain-specific languages





## **Dataflow programming**

© Prof. J. Castrillon. CPS Summer School. Alghero, Italy, 2019

#### **Dataflow programming**

- Graph representation of applications: Long history
  - Implicit repetitive execution of tasks
  - Good model for streaming applications
  - Good match for signal processing & multi-media

#### □ The why

- Explicit parallelism
- Often: Determinism
- Better analyzability (scheduling, mapping, optimization)





\_\_\_\_ CHAIRFOR

COMPILER

CONSTRUCTION



#### Example: Synchronous dataflow graph (SDF)



Def.: An **SDF** is an annotated multi-graph G=(V,E,W). V is the set of actors and E  $\subseteq VxV$  the set of channels.  $W = \{w_1,...,w_{|E|}\} \subset N^3$  is a set of annotations for every channel  $e = (a_1,a_2), w_e = (p_e,c_e,d_e)$ :  $p_e$  is the number of tokens produced by a firing of  $a_1$ ,  $c_e$  the tokens consumed by  $a_2$  and  $d_e$  the amount of delay tokens

Originally in [Lee'87]



#### Example: Synchronous dataflow graph (SDF)



Def.: An **SDF** is an annotated multi-graph G=(V,E,W). V is the set of actors and E  $\subseteq VxV$  the set of channels.  $W = \{w_1,...,w_{|E|}\} \subset N^3$  is a set of annotations for every channel  $e = (a_1,a_2), w_e = (p_e,c_e,d_e)$ :  $p_e$  is the number of tokens produced by a firing of  $a_1$ ,  $c_e$  the tokens consumed by  $a_2$  and  $d_e$  the amount of delay tokens



#### **Topology matrix**



Def.: Given an SDF G=(V,E,W), its **topology matrix**  $\Gamma$  has a row for every channel and a column for every actor.  $[\Gamma_{ik}] = p_{ik} - c_{ik}$  (i.e. the amount of tokens produced minus those consumed by actor k on/from channel i)

**Example** 





### **Topology matrix (2)**



CONSTRUCTION

- **State** of an SDF: tokens in the queues
- □ The topology matrix can be used to update the state after a firing



### **Topology matrix (3)**



CONSTRUCTION

- **State** of an SDF: tokens in the queues
- □ The topology matrix can be used to update the state after a firing



#### **Topology matrix and complete cycles**



Complete cycle: sequence of firing that brings the SDF to the original state

$$s_n = s_0 + \Gamma \cdot (v_1 + v_2 + \dots + v_n) = s_0$$
$$\Gamma \cdot (v_1 + v_2 + \dots + v_n) = \Gamma \cdot \vec{r} = 0$$



#### Maximum cycle mean (MCM)

Cfaed Center for Advancing Designed

- Unroll SDF is an HSDF
- □ The MCM relates to the maximum throughput of an HSDF
- Given an HSDF G=(V,E)
  - □ Let O(G) be the set of all cycles

□ A cycle C=((
$$v_i, v_{i+1}$$
), ( $v_{i+2}, v_{i+3}$ ), ..., ( $v_{i+n-1}, v_i$ ))

$$MCM(G) = \max_{C \in O(G)} \left( \frac{\sum_{v:(u,v) \lor (v,u) \in C} \rho(v)}{\sum_{e \in C} w(e)} \right)$$

x delay tokens

WCET



#### **Other models**





#### **Dynamic models and language support**



- Kahn Process Networks (KPNs)
  - Simple syntax for programming
  - Less model information: more analysis effort
- Via tracing compilers can

21

Identify general communication patterns



#### **Programming flow: Overview**



#### **Programming flows**

- □ Many: CAL, DOL, CPN, Daedalus, Ptolemy, CAPH, ...
- Information
  - Dataflow model: Rates, states, actions, traces, ...
  - □ Architecture model: Resources, interconnect, costs, ...
- Optimization: Mapping application to hardware
  - Multi-objective optimization
  - High-level simulation / cost models



CHAIRFOR

COMPILER

CONSTRUCTION



#### Automatic mapping to heterogeneous many-cores



CONSTRUCTION

- Lots of research: Heuristics and meta-heuristics
- Sample results: Effort-quality trade-off



© Prof. J. Castrillon. CPS Summer School. Alghero, Italy, 2019





## **Dataflow and CPSs**

Adaptivity, predictable, multi-app, reactive, ...

© Prof. J. Castrillon. CPS Summer School. Alghero, Italy, 2019

## System dynamics: Hybrid mappings



COMPILER

CONSTRUCTION



#### **Data-level parallelism: Scalable and adaptive**



- Change parallelism from the application specification
- □ Static code analysis to identify possible transformations (or via annotations)
- Implementation in FIFO library (semantics preserving)
- Similar to AdaPNet or parameterized SDFs

27







### Sample symmetries for MJPEG



- Multimedia benchmarks on Adapteva (16 Epiphany cores)
- Not always very intuitive



### Flexible mappings: Run-time analysis



- Linux kernel: symmetry-aware
- Target: Odroid XU4 (big.LITTLE)
- Multi-application scenarios: audio filter (AF) and MIMO
  - 1x AF
  - $4 \times AF$
  - $2 \times AF + 2 \times MIMO$
- 3 mappings to two processors
  - T1: Best CPU time
  - T2: Best wall-clock time
  - T3: GBM heuristic [Castrill12]

Single AF

[SCOPES'17b]





30

#### Flexible mappings: Multi-application results (1)



#### Robustness



- Static mappings, transformed or not, provide good predictability
- However: Many things out of control
  - Application data, unexpected interrupts, unexpected OS decisions



Can we reason about robustness of mapping to external factors?



#### **Design centering**

- Design centering: Find a mapping that can better tolerate variations while staying feasible
- Studied field, in e.g., biology, circuit design or manufacturing systems.
- Currently
  - Using a bio-inspired algorithm
  - Robust against OS changes to the mapping





33



#### **Design centering: Algorithmic**



Intuition: Find the center and the form of a region, in which parameters deliver a correct solution



#### **Evaluation**



- □ Analyze how robust the center really is
  - Perturbate mappings and check how often the constraints are missed
  - Signal processing applications on clustered ARM manycore and NoC manycore (16)



#### **Ongoing work: Improve representations**



- $\Box$  Work on embeddings: Architectures  $\rightarrow$  Real numbers
- Novel mapping representations: faster heuristics & more efficient heuristics?
- Example: T-SNE Visualization for mappings space (8 tasks on Odroid XU4)







## **Domain-specific abstractions**

© Prof. J. Castrillon. CPS Summer School. Alghero, Italy, 2019

#### **Higher-level Abstractions**

38



Dataflow: Nice abstraction across domains (maybe too overloaded)

Need to match domain-specific accelerators with domain-specific abstractions



#### **Example: CFDlang**

39



Origin: Spectral element methods in Computational Fluid Dynamics

$$\mathbf{v}_e = (\mathbf{A} \otimes \mathbf{A} \otimes \mathbf{A}) \mathbf{u}_e$$

Simple expression language, formal semantics, domain expert optimizations



## **Tensor languages (hype): Machine learning**



- Many existing frameworks, e.g., TVM, Tensor Comprehensions, TensorFlow, ...
- Recent work
  - Cross-domain tensor transformations [GPCE'18]
  - Formal semantics for safety checks (particularly important for CPS!)
  - Optimization for emerging memories [LCTES'19]



CHAIRFOR

COMPILER

CONSTRUCTION

[Array'19]

#### **Correct by construction**

- Similar to formal MoCs
  - Correct by design
  - No abstraction leaks



A = placeholder((m,h), name='A') B = placeholder(h,h), name='A')  $k = reducC_{ij} = \sum_{k=1}^{n} A_{ki}B_{kj}$  C = compute((m,k=1, lambda i, j; sum(A[k, i] \* B[k, j], axis=k))

#### Current efforts in formal semantics for safe code generation

$$\begin{bmatrix} \cdot \end{bmatrix} : Context \rightarrow Memory \rightarrow (list of Nat) \rightarrow \mathbb{D}$$

$$\begin{bmatrix} x \end{bmatrix} \Gamma \mu \overline{i} = \mu x \overline{i}$$

$$\begin{bmatrix} x \end{bmatrix} \Gamma \mu \overline{i} = \mu x \overline{i}$$

$$\begin{bmatrix} (e) \end{bmatrix} \Gamma \mu \overline{i} = \begin{bmatrix} e \end{bmatrix} \Gamma \mu \overline{i}$$

$$\begin{bmatrix} (e) \end{bmatrix} \Gamma \mu \overline{i} = \begin{bmatrix} e \end{bmatrix} \Gamma \mu \overline{i}$$

$$\begin{bmatrix} (e) \end{bmatrix} \Gamma \mu \overline{i} = \begin{bmatrix} e \\ e \\ 0 \end{bmatrix} \Gamma \mu \overline{i} = \begin{bmatrix} e \\ 0 \end{bmatrix} \Gamma \mu \overline{i} + \begin{bmatrix} e \\ 1 \end{bmatrix} \Gamma \mu \overline{i}$$

$$\begin{bmatrix} mul \ e_0 \ e_1 \end{bmatrix} \Gamma \mu \overline{i} = \begin{bmatrix} e_0 \end{bmatrix} \Gamma \mu [1 \cdot \begin{bmatrix} e_1 \end{bmatrix} \Gamma \mu \overline{i},$$

$$\begin{bmatrix} mul \ e_0 \ e_1 \end{bmatrix} \Gamma \mu \overline{i} = \begin{bmatrix} e_0 \end{bmatrix} \Gamma \mu [1 \cdot \begin{bmatrix} e_1 \end{bmatrix} \Gamma \mu \overline{i},$$

$$\begin{bmatrix} mul \ e_0 \ e_1 \end{bmatrix} \Gamma \mu \overline{i} = \begin{bmatrix} e_0 \end{bmatrix} \Gamma \mu [1 \cdot \begin{bmatrix} e_1 \end{bmatrix} \Gamma \mu \overline{i},$$

$$\begin{bmatrix} mul \ e_0 \ e_1 \end{bmatrix} \Gamma \mu \overline{i} = \begin{bmatrix} e_0 \end{bmatrix} \Gamma \mu \overline{i} \cdot \begin{bmatrix} e_1 \end{bmatrix} \Gamma \mu \overline{i},$$

$$\begin{bmatrix} mul \ e_0 \ e_1 \end{bmatrix} \Gamma \mu \overline{i} + \begin{bmatrix} e_0 \end{bmatrix} \Gamma \mu \overline{i} \cdot \begin{bmatrix} e_1 \end{bmatrix} \Gamma \mu \overline{i},$$

$$\begin{bmatrix} mul \ e_0 \ e_1 \end{bmatrix} \Gamma \mu \overline{i} + \begin{bmatrix} e_0 \end{bmatrix} \Gamma \mu \overline{i} \cdot \begin{bmatrix} e_1 \end{bmatrix} \Gamma \mu \overline{i},$$

$$\begin{bmatrix} e_1 \ \Gamma \mu \begin{bmatrix} j_1, \dots, j_{i-1}, j_i, \dots, j_k \end{bmatrix} = \\ \begin{bmatrix} e_1 \ \Gamma \mu \begin{bmatrix} j_1, \dots, j_{i-1}, j_i, \dots, j_k \end{bmatrix} = \\ \begin{bmatrix} e_1 \ \Gamma \mu \begin{bmatrix} j_1, \dots, j_{i-1}, j_i, \dots, j_k \end{bmatrix} = \\ \begin{bmatrix} e_1 \ \Gamma \mu \begin{bmatrix} j_1, \dots, j_{i-1}, j_i, \dots, j_k \end{bmatrix} = \\ \begin{bmatrix} e_1 \ \Gamma \mu \begin{bmatrix} j_1, \dots, j_{i-1}, m, j_i, \dots, j_k \end{bmatrix} = \\ \begin{bmatrix} e_1 \ \Gamma \mu \begin{bmatrix} j_1, \dots, j_{i-1}, m, j_i, \dots, j_k \end{bmatrix} = \\ \begin{bmatrix} e_1 \ \Gamma \mu \begin{bmatrix} j_1, \dots, j_{i-1}, m, j_i, \dots, j_k \end{bmatrix} = \\ \begin{bmatrix} e_1 \ \Gamma \mu \begin{bmatrix} j_1, \dots, j_{i-1}, m, j_i, \dots, j_k \end{bmatrix} = \\ \begin{bmatrix} e_1 \ \Gamma \mu \begin{bmatrix} j_1, \dots, j_{i-1}, m, j_i, \dots, j_k \end{bmatrix} = \\ \begin{bmatrix} e_1 \ \Gamma \mu \begin{bmatrix} j_1, \dots, j_{i-1}, m, j_i, \dots, j_k \end{bmatrix} = \\ \begin{bmatrix} e_1 \ \Gamma \mu \begin{bmatrix} j_1, \dots, j_{i-1}, m, j_i, \dots, j_k \end{bmatrix} = \\ \begin{bmatrix} e_1 \ \Gamma \mu \begin{bmatrix} j_1, \dots, j_{i-1}, m, j_i, \dots, j_k \end{bmatrix} = \\ \begin{bmatrix} e_1 \ \Gamma \mu \begin{bmatrix} j_1, \dots, j_{i-1}, m, j_i, \dots, j_k \end{bmatrix} = \\ \begin{bmatrix} e_1 \ \Gamma \mu \begin{bmatrix} j_1, \dots, j_{i-1}, m, j_i, \dots, j_k \end{bmatrix} = \\ \begin{bmatrix} e_1 \ \Gamma \mu \begin{bmatrix} j_1, \dots, j_{i-1}, m, j_i, \dots, j_k \end{bmatrix} = \\ \begin{bmatrix} e_1 \ \Gamma \mu \begin{bmatrix} j_1, \dots, j_{i-1}, m, j_i, \dots, j_k \end{bmatrix} = \\ \begin{bmatrix} e_1 \ \Gamma \mu \begin{bmatrix} j_1, \dots, j_{i-1}, m, j_i, \dots, j_k \end{bmatrix} = \\ \begin{bmatrix} e_1 \ \Gamma \mu \begin{bmatrix} j_1, \dots, j_{i-1}, m, j_i, \dots, j_k \end{bmatrix} = \\ \begin{bmatrix} e_1 \ \Gamma \mu \begin{bmatrix} j_1, \dots, j_{i-1}, m, j_i, \dots, j_k \end{bmatrix} = \\ \begin{bmatrix} e_1 \ \Gamma \mu \begin{bmatrix} j_1, \dots, j_{i-1}, m, j_i, \dots, j_k \end{bmatrix} = \\ \begin{bmatrix} e_1 \ \Gamma \mu \begin{bmatrix} m \begin{bmatrix} m \mu \begin{bmatrix} m$$

41

#### **Emerging architectures**

- Lots of research on tensor computations on tensor accelerators!
- Recently looking into emerging memory architectures
  - Hybrid architectures STT/Re-RAM [TVLSI'18, TVLSI'19]
  - Racetrack memories

42

- Racetrack memories: Predicted extreme density at low latency
  - Multiple bits stored in nano-wires with magnetic domains



![](_page_41_Picture_8.jpeg)

e.g., [Kim'19]

## Architecture and data layout optimization

![](_page_42_Picture_1.jpeg)

□ Architecture – software co-optimization

43

Embedded system for inference: RTM as scratchpad with pre-shifting and other optimizations
RT0

![](_page_42_Figure_4.jpeg)

#### Architecture and data layout optimization

![](_page_43_Picture_1.jpeg)

Data-layout: Reduce the number of shifts

![](_page_43_Figure_3.jpeg)

#### Latency comparison vs SRAM

![](_page_44_Picture_1.jpeg)

- Un-optimized and naïve mapping: Even worse latency than SRAM
- □ 24% average improvement (even with very conservative circuit simulation)

![](_page_44_Figure_4.jpeg)

#### **Energy comparison vs SRAM**

![](_page_45_Picture_1.jpeg)

- Higher savings due to less leakage power
- □ 74% average improvement

![](_page_45_Figure_4.jpeg)

![](_page_46_Picture_0.jpeg)

![](_page_46_Picture_1.jpeg)

![](_page_46_Picture_2.jpeg)

© Prof. J. Castrillon. CPS Summer School. Alghero, Italy, 2019

#### **Summary**

![](_page_47_Picture_1.jpeg)

- Principled methodologies for programming are a must! preaching to the choir
  - Advances in auto-parallelization (polyhedral, dynamic analysis)
  - Explicit parallel MoC-based programming models
  - Quest for more adaptivity but retaining time predictability
  - Higher-level abstractions: Example for tensor-based computations for accelerators
- Lots of challenges remain (thankfully)
  - Cost models and characterization of trade-offs (vs. blind searches)
  - Understand impact of emerging technologies
  - Syntax and semantics (for correctness): Lots of open questions
  - Time-semantics in programming and execution environments

![](_page_47_Picture_13.jpeg)

#### References

![](_page_48_Picture_1.jpeg)

[DAC'08] Ceng, J., et al. "MAPS: an integrated framework for MPSoC application parallelization.", Proceedings of the 45th annual Design Automation Conference, DAC'08. [Aguilar'18] Aguilar, M. A., Software Parallelization and Distribution for Heterogeneous Multi-Core Embedded Systems, RWTH Aachen University, RWTH Aachen University, 2018, 181pp [Edler'18] T. J.K. Edler von Koch. S. Manilov, C. Vasiladiotis, M. Cole, and B.Franke, "Towards a Compiler Analysis for Parallel Algorithmic Skeletons", CC'18. [Lee'87] E.A. Lee and D. G. Messerschmitt. "Synchronous data flow." Proceedings of the IEEE 1987 [Springer'14] J. Castrillon, R. Leupers, "Programming Heterogeneous MPSoCs: Tool Flows to Close the Software Productivity Gap", Springer, pp. 258, 2014. [MCSoC'16] A. Goens, R. Khasanov, J. Castrillon, S. Polstra, A. Pimentel, "Why Comparing System-level MPSoC Mapping Approaches is Difficult: a Case Study", MCSoC-16. [Schor'14] Schor, L.; Bacivarov, I.; Yang, H. & Thiele, L. "AdaPNet: Adapting Process Networks in Response to Resource Variations", ESWEEK'14, 2014 [Stuijk'11] Stuijk, S.; Geilen, M.; Theelen, B. & Basten, T. "Scenario-aware dataflow: Modeling, analysis and implementation of dynamic applications", SAMOS'11 404-411 [Desnos'13] Desnos, K., et al. "Pimm: Parameterized and interfaced dataflow meta-model for mpsocs runtime reconfiguration." SAMOS), 2013. [PARMA-DITAM'18] R. Khasanov, et al, "Implicit Data-Parallelism in Kahn Process Networks: Bridging the MacQueen Gap", PARMA-DITAM 18, ACM, pp. 20-25, Jan 2018. [ACM TACO'17] Goens, A. et al. "Symmetry in Software Synthesis". In: ACM Transactions on Architecture and Code Optimization (TACO) (2017).

[SCOPES'17a] G. Hempel, et al, "Robust Mapping of Process Networks to Many-Core Systems Using Bio-Inspired Design Centering" SCOPES'17.

[SCOPES'17b] Goens, A. et al. "TETRIS: a Multi-Application Run-Time System for Predictable Execution of Static Mappings", SCOPES'17.

[DAC'12] J. Castrillon, A. Tretter, R. Leupers, G. Ascheid. "Communication-aware Mapping of KPN Applications Onto Heterogeneous MPSoCs" in DAC 2012, 1266-1271

[Schwarzer'17] T. Schwarzer, et al. "Symmetry-eliminating design space exploration for hybrid application mapping on many-core architectures." IEEE TCAD 37.2 (2017): 297-310.

[MCSoC'18] A. Goens, C. Menard, J. Castrillon, "On the Representation of Mappings to Multicores", MCSoC-18, pp. 184–191, Hanoi, Vietnam, Sep 2018.

[CyPhy'19] M. Lohstroh, et al. "Reactors: A Deterministic Model for Composable Reactive Systems" Workshop on Design, Modeling and Evaluation of Cyber Physical Systems (CyPhy 2019) Oct 2019.

[**RWDSL'18**] N. A. Rink, et al. "CFDlang: High-level code generation for high-order methods in fluid dynamics". RWDSL'18.

[GPCE'17] A. Susungi, N. A. Rink, J. Castrillon, et al. "Towards Compositional and Generative Tensor Optimizations". GPCE'17, Oct. 2017, pp. 169–175.

[Chen'18] Chen, Tianqi, et al. "TVM: An Automated End-to-End Optimizing Compiler for Deep Learning." 13th OSDI, 2018.

[GPCE'18] A. Susungi, N. A. Rink, A. Cohen, J. Castrillon, and C. Tadonki. "Meta-programming for cross-domain tensor optimizations". GPCE'18, Nov. 2018, pp. 79–92.

[Array'19] N.A. Rink, N. A. and J. Castrillon. "TelL: a type-safe imperative Tensor Intermediate Language", ARRAY, ACM, 2019, pp. 57-68

[Kim'19] J. Kim, et al. "A code generator for high-performance tensor contractions on GPUs." Proceedings of the Symposium on Code Generation and Optimization. IEEE Press, 2019.

[IEEE CAL'19] Asif Ali Khan, Fazal Hameed, Robin Bläsing, Stuart Parkin, Jeronimo Castrillon, "RTSim: A Cycle-accurate Simulator for Racetrack Memories", In IEEE Computer Architecture Letters, Feb 2019.

[LCTES'19] Khan, A. A., Rink, N. A., Hameed, F., Castrillon, J. "Optimizing Tensor Contractions for Embedded Devices with Racetrack Memory Scratch-Pads", Proceedings of LCTES, Jun 2019, pp. 5-18

![](_page_48_Picture_17.jpeg)

© Prof. J. Castrillon. CPS Summer School. Alghero, Italy, 2019