

High-Performance Deep-Learning Operators on CPU and GPU via Multi-Dimensional Homomorphisms

a.rasch@wwu.de

Ari Rasch, Richard Schulze, Sergei Gorlatch



Generation

Let T and T' be two arbitrary types. A function $h : T[N_1] \dots [N_d] \rightarrow T'$ on d -dimensional arrays is called a *Multi-Dimensional Homomorphism (MDH)* iff there exist *combine operators* $\otimes_1, \dots, \otimes_d : T' \times T' \rightarrow T'$, such that for each $k \in [1, d]$ and arbitrary, concatenated input array $a \mathbin{++}_k b$ in dimension k :

$$h(a \mathbin{++}_k b) = h(a) \otimes_k h(b) \quad \text{[IJPP'18]}$$

MDHs can be uniformly represented via our `md_hom` parallel pattern:

$$\text{md_hom}(f, (\otimes_1, \dots, \otimes_d))(a) = \otimes_1 \dots \otimes_d f(a[i_1] \dots [i_d])$$

$i_1 \in [1, N_1] \quad i_d \in [1, N_d]$

Important Deep-Learning Operators can be expressed as MDHs:

Linear Algebra (BLAS)

GEMM = `md_hom(*, (++, ++, +)) o view(A,B)(i,j,k)(A[i,k],B[k,j])`
 GEMV = `md_hom(*, (++, +)) o view(A,B)(i, k)(A[i,k],B[k])`
 DOT = `md_hom(*, (+)) o view(A,B)(k)(A[k],B[k])`

Convolution

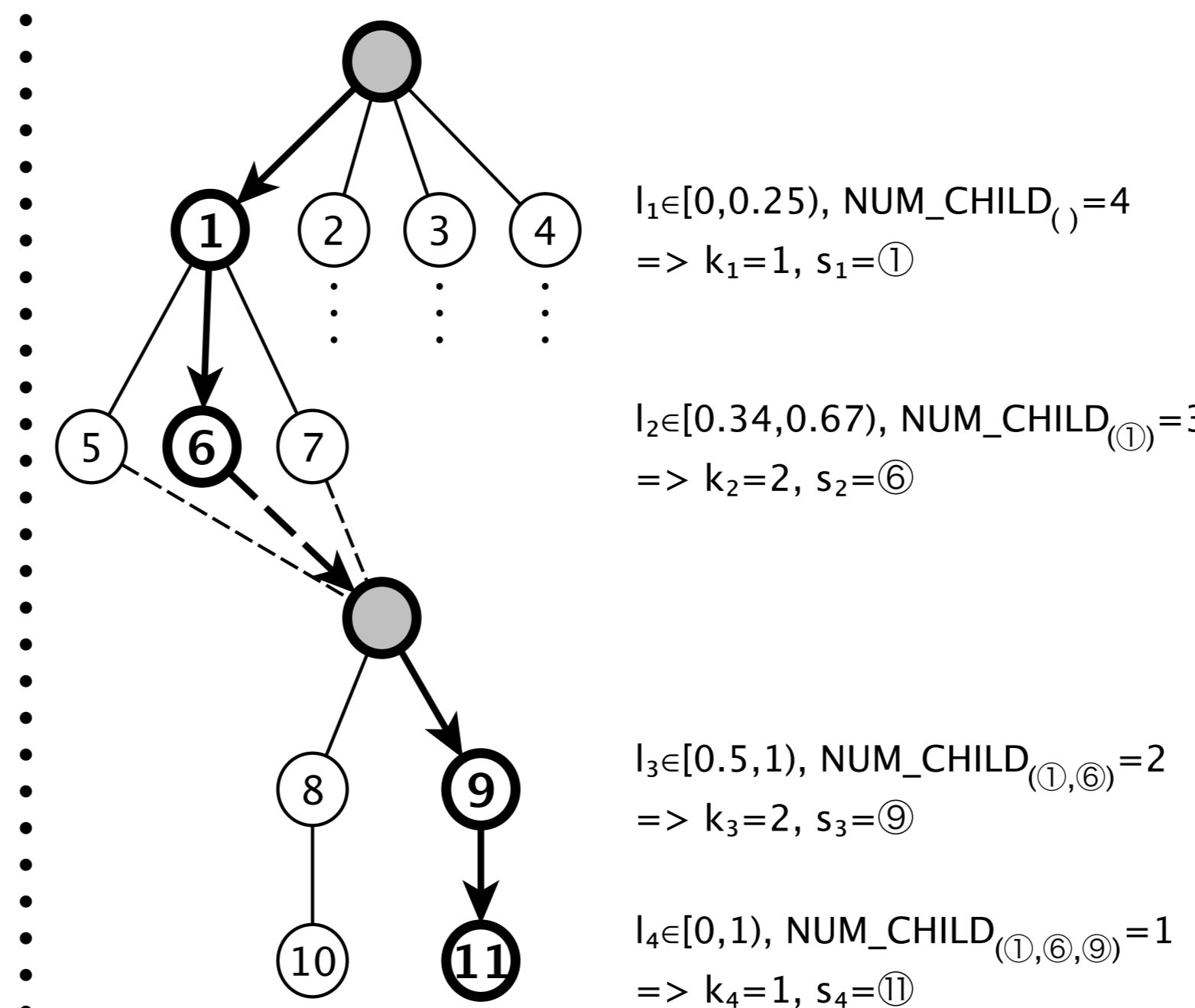
CONV = `md_hom(*, (++, ++, ++, ++, ++, ++, +)) o view(images,filter)(n,k,p,q,c,r,s)(images[n,c,p+r,q+s],filter[k,c,r,s])`

Tensor Contractions

TC = `md_hom(*, (++, ..., ++, ++, +, ..., +)) o view(...)`

Optimization

Our **Auto-Tuning Framework (ATF)** is a **general-purpose approach** that supports auto-tuning of programs with **interdependent tuning parameters**. [CCPE'18]



We provide a novel **chain-of-trees** search space structure for interdependent tuning parameters.

```
#atf::tp name /* name */
range /* range */
constraint /* constraint */
```

We extend the traditional definition of *tuning parameters* by a **parameter constraint**.

ATF efficiently **generates / stores / explores** the search spaces of **interdependent tuning parameters**

2.75x faster than TVM

4x faster than Intel MKL/NVIDIA cuBLAS

Our MDH approach achieves often **better performance** than well-performing **machine- and hand-optimized approaches**.

3.31x faster than NVIDIA cuDNN

2x faster than COGENT & Tensor Comprehensions

Generating **Auto-Tunable OpenCL Code**
[PACT'19]

`md_hom(f, (\otimes_1, \dots, \otimes_k))`

