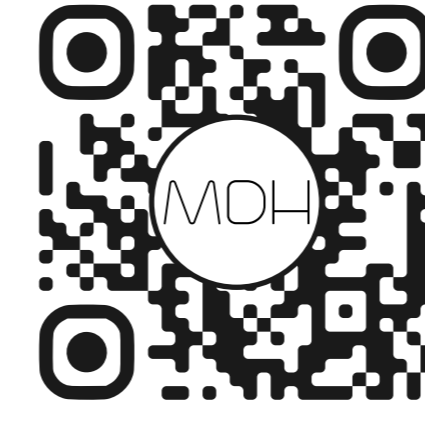




Generation



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 $a \text{++}_k b$ in dimension k :

$$h(a \text{++}_k b) = h(a) \otimes_k h(b)$$

DL computations can be expressed as **MDH functions**, and **GPU/CPU/...** code generated and **optimized** according to **MDH formalism** [1]

[1] Rasch, (De/Re)-Composition of Data-Parallel Computations via Multi-Dimensional Homomorphisms, **TOPLAS'24**

MDHs can be uniformly expressed via our `md_hom` higher-order function:

$$\text{md_hom}(f, (\otimes_1, \dots, \otimes_D))(a) := \otimes_1 \dots \otimes_D f(a[i_1, \dots, i_D])$$

$i_1 \in I_1 \quad i_D \in I_D$

```
CONV<...> = out_view<...>( 0:(n,p,...)->(n,p,q,k) ) o
             md_hom<...>( *, (++, ++, ++, ++, +, +, +) ) o
             inp_view<...>( I:(n,p,...)->(n,p+r,q+s,c) , F:(n,p,...)->(k,r,s,c) )
```

Convolutions

```
MatMul<...> = out_view<...>( C:(i,j,k)->(i,j) ) o
             md_hom<...>( *, (++, ++, +) ) o
             inp_view<...>( A:(i,j,k)->(i,k), B:(i,j,k)->(k,j) )
```

Linear Algebra

```
BiasAdd<NHWC><...> = out_view( 0:(n,h,w,c)->(n,h,w,c) ) o
                    md_hom<...>( +, (++, ++, ++, ++, +) ) o
                    inp_view<...>( I1:(n,h,w,c)->(n,h,w,c), I2:(n,h,w,c)->(c) )
```

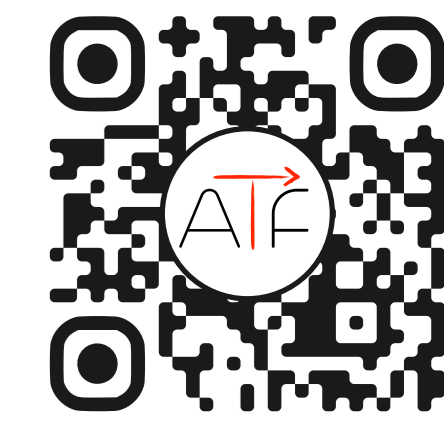
Point-Wise OPs

```
out_view(...) o
md_hom(...) o
inp_view(...)
```

Code Generation
(auto-tunable)

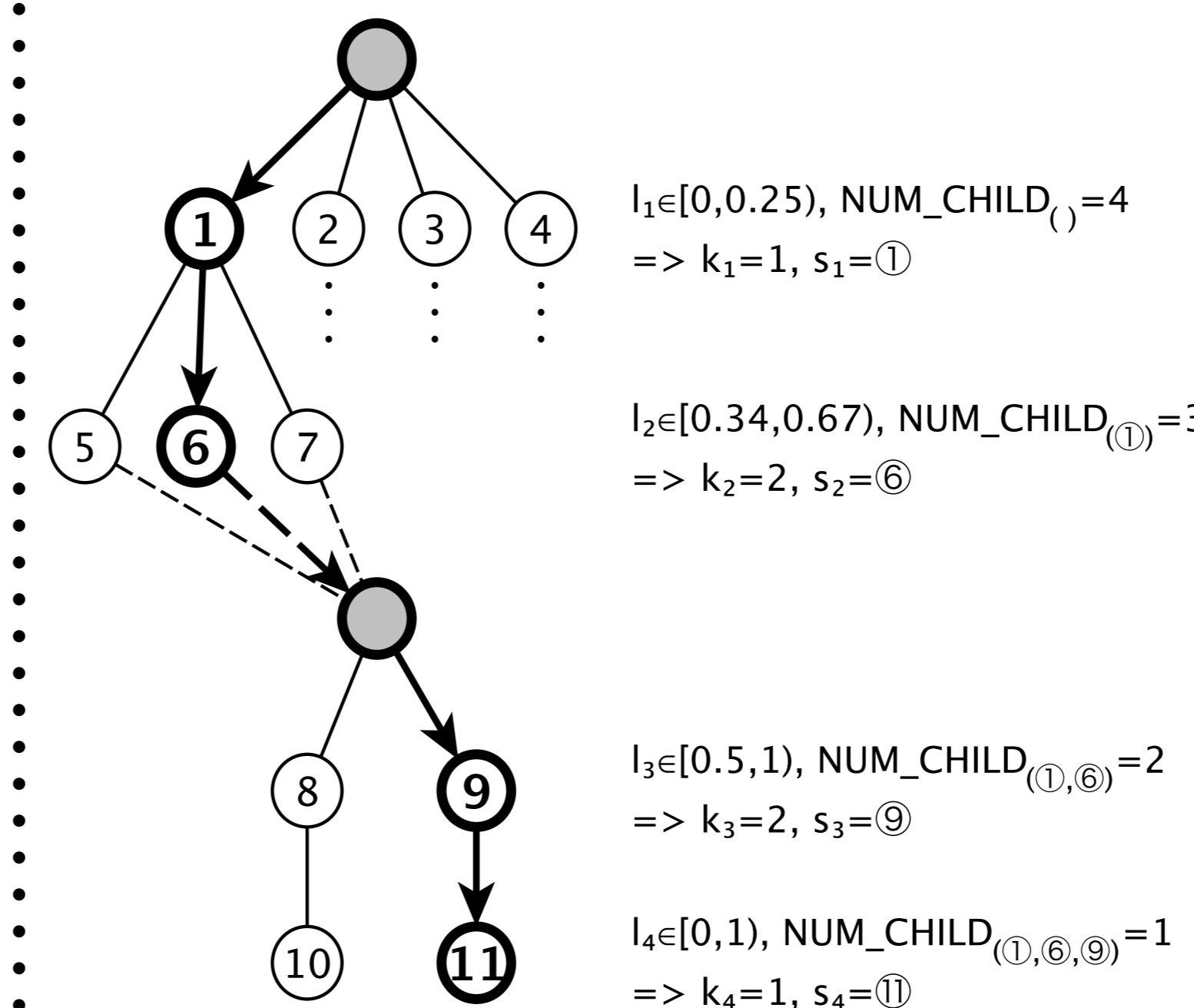


Optimization



Our **Auto-Tuning Framework (ATF)** is a **general-purpose approach** that **automatically optimizes** (*auto-tunes*) programs with **constrained tuning parameters** [2]

[2] Rasch, Schulze, Steuwer, Gorlatch, *Efficient Auto-Tuning of Parallel Programs with Interdependent Tuning Parameters via Auto-Tuning Framework (ATF)*, **TACO'21**



CoT (Chain-of-Trees)

A new search space structure for constrained 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 **constrained tuning parameters**

up to **2.67x** speedups over **NVIDIA cuBLAS**

up to **3.5x** speedups over **NVIDIA cuDNN**

MDH+ATF achieves on **GPUs, CPUs, ...** often higher **Performance & Portability & Productivity** than well-performing **hand- and machine-optimized approaches** [1]

up to **9.01x** faster than **Intel oneDNN**

up to **3.01x** faster than **TVM**

Highlights only!