

https://mdh-lang.org

Code Generation & Optimization for Deep-Learning Computations on GPUs via Multi-Dimensional Homomorphisms

Richard Schulze, Ari Rasch, Sergei Gorlatch



Introduction

We present our work-in-progress **code generation and optimization approach** for **DL computations**:

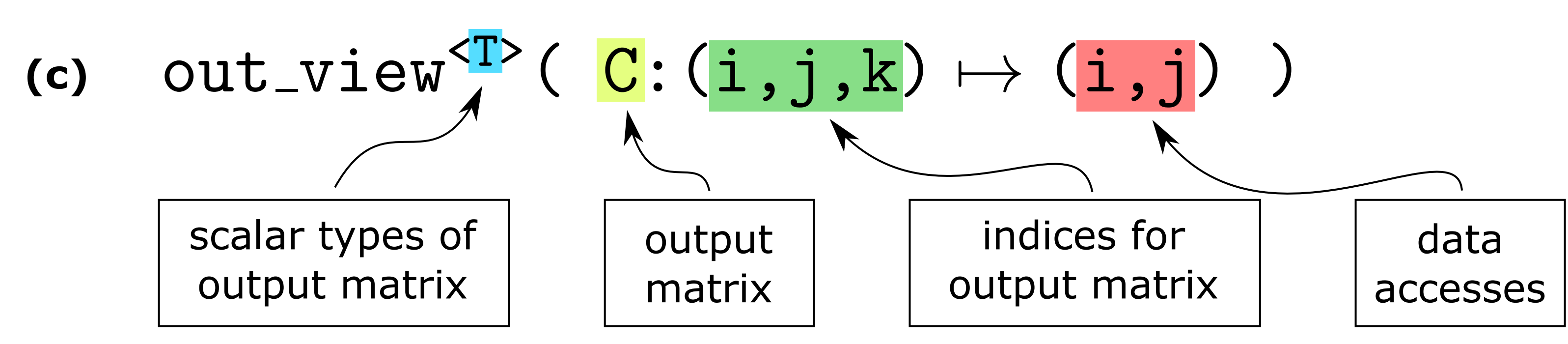
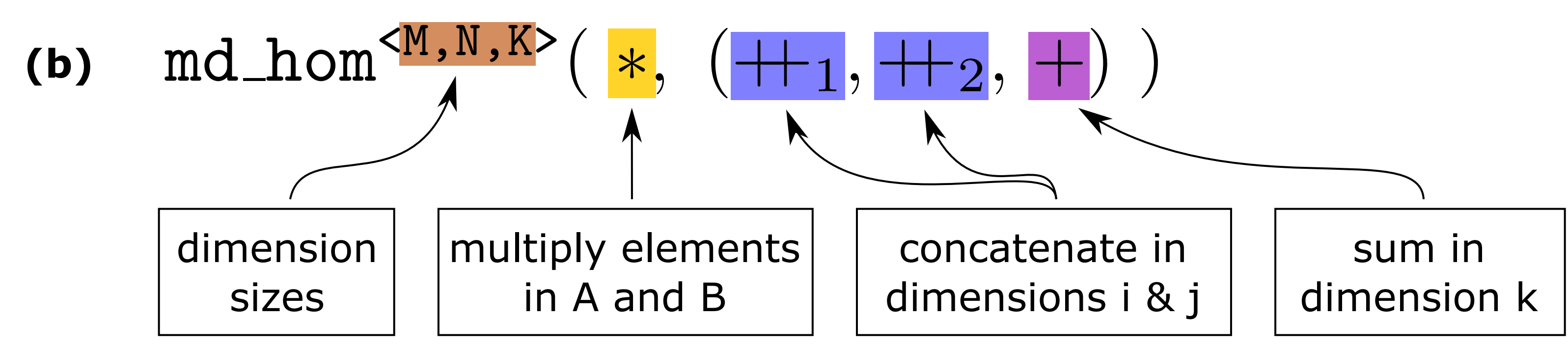
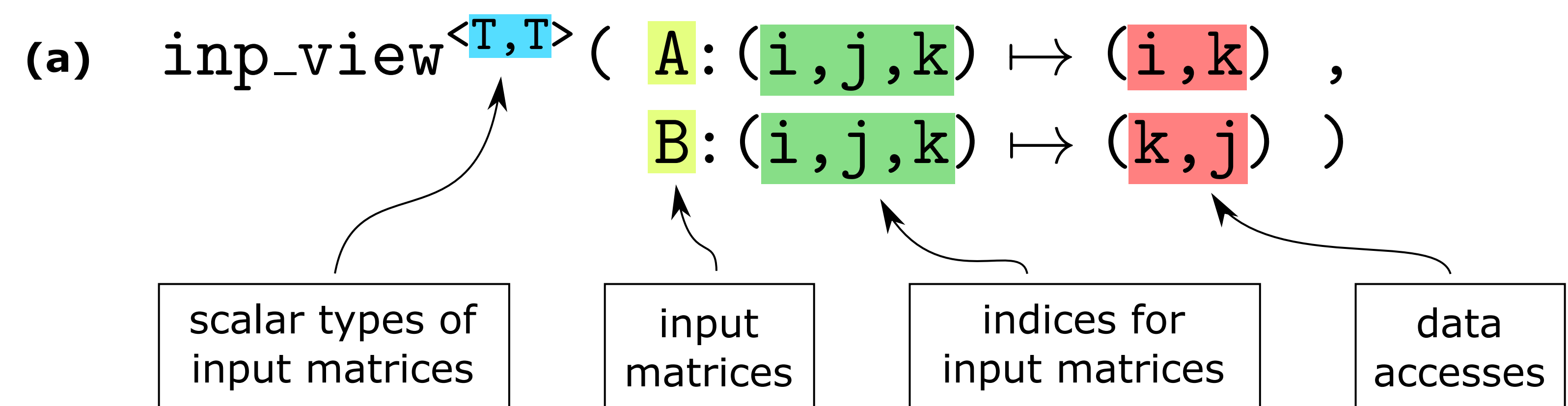
- based on the formalism of **Multi-Dimensional Homomorphisms (MDH)** [1]
- achieves **high-performance** for popular DL computations by exploiting the already existing MDH GPU code generation and optimization approach
- **more expressive** than the state-of-the-art DL abstractions (e.g., as provided by TensorFlow): we are capable of expressing multiple DL computations as a single MDH expression

[1] Rasch, Gorlatch, "Multi-Dimensional Homomorphisms and Their Implementation in OpenCL", IJPP'18

The MDH Formalism

$$\text{MatMul}^{\langle T \in \text{Type} \mid M, N, K \in \mathbb{N} \rangle} := \text{out_view}^{\langle \dots \rangle} (\dots) \circ \text{md_hom}^{\langle \dots \rangle} (\dots) \circ \text{inp_view}^{\langle \dots \rangle} (\dots)$$

(c) (b) (a)

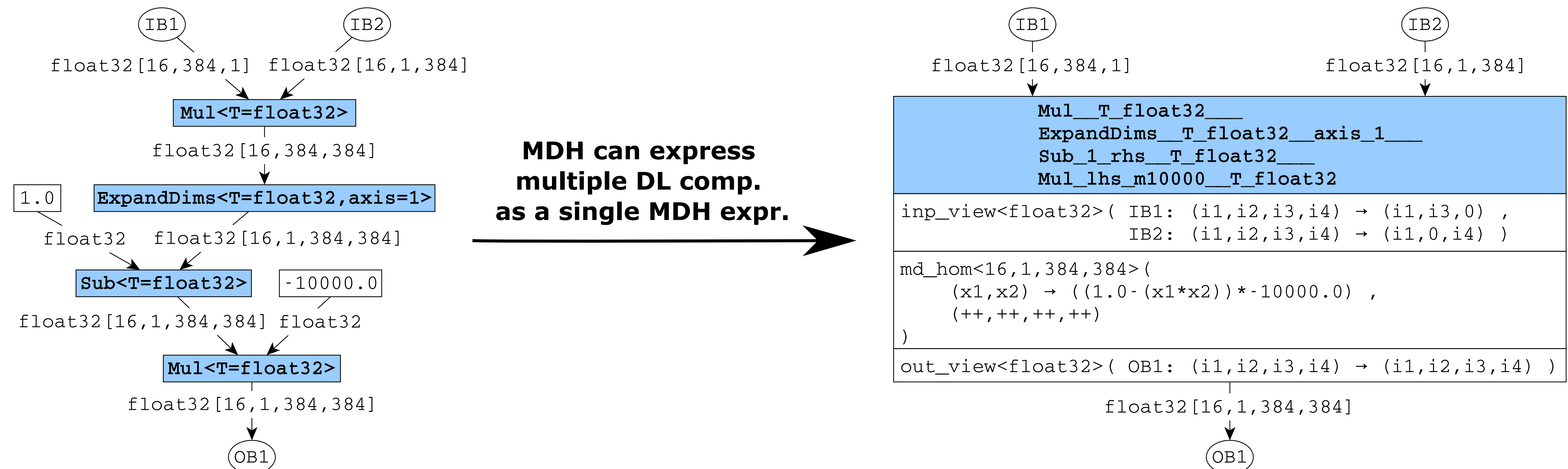


MDH allows us conveniently expressing data-parallel computations and automatically generate CUDA code for them [2, 3].

DL Computations Expressed in the MDH Formalism

Operator	out_view ^{<...>}	md_hom ^{<...>}	inp_view ^{<...>}
Mul ^{<...>}	OB1: (i, j) ↦ (i, j)	* , (++ ₁ , ++ ₂)	IB1: (i, j) ↦ (i, j) , IB2: (i, j) ↦ (i, j)
Sub ^{<...>}	OB1: (i, j) ↦ (i, j)	- , (++ ₁ , ++ ₂)	IB1: (i, j) ↦ (i, j) , IB2: (i, j) ↦ (i, j)
ExpandDims ^{<axis, D ∈ ℕ ...>}	OB1: (i ₁ , ..., i _D) ↦ (... , i _{axis-1} , 0, i _{axis} , ...)	id , (++ ₁ , ..., ++ _D)	IB1: (i ₁ , ..., i _D) ↦ (i ₁ , ..., i _D)
BiasAddGrad ^{<NHWC ...>}	OB1: (i, j) ↦ (j)	id , (+, ++ ₂)	IB1: (i, j) ↦ (i, j)
BatchMatMul ^{<N, N ...>}	OB1: (b1, b2, i, j, k) ↦ (b1, b2, i, j)	* , (++ ₁ , ..., ++ ₄ , +)	IB1: (b1, b2, i, j, k) ↦ (b1, b2, i, k) , IB2: (b1, b2, i, j, k) ↦ (b1, b2, k, j)

Popular DL computations¹ are conveniently expressed in the MDH formalism.



¹ Taken from the TensorFlow implementation of the real-world BERT neural network.

Experimental Results

2.9x faster than TVM for BiasAddGrad

1.5x faster than TensorFlow for BiasAddGrad

1.1x faster than TVM for BatchMatMul

3.8x faster than TVM for a subgraph of BERT

Our preliminary experimental results on NVIDIA V100 GPU show that we can achieve **better performance** than well-performing machine- and hand-optimized approaches on real-world data sizes.

4.9x faster than TensorFlow for a subgraph of BERT

1.9x faster than TC for BatchMatMul

1.7x faster than TC for a subgraph of BERT

1.7x faster than TC for BiasAddGrad

[2] Rasch, Schulze, Gorlatch, "Generating Portable High-Performance Code via Multi-Dimensional Homomorphisms", PACT'19

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