

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

Richard Schulze, Ari Rasch, Sergei Gorlatch



## Introduction

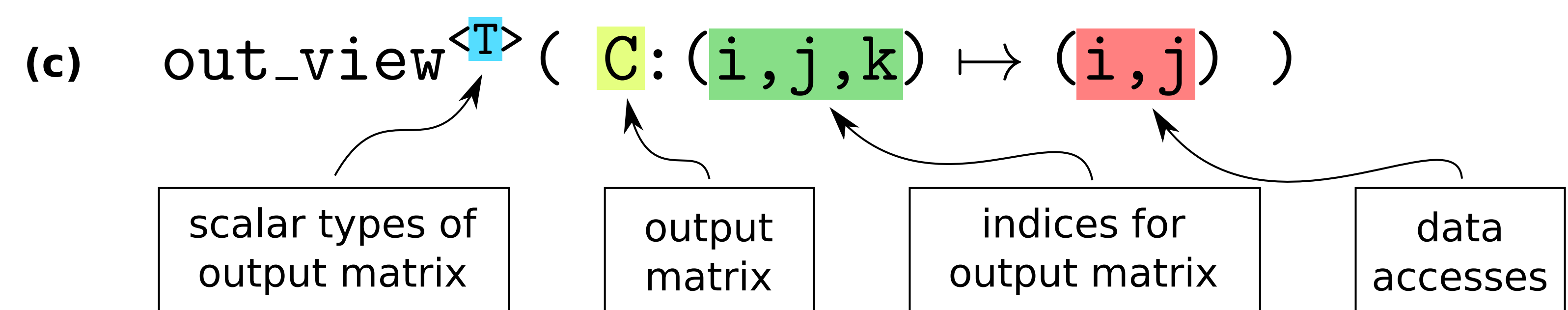
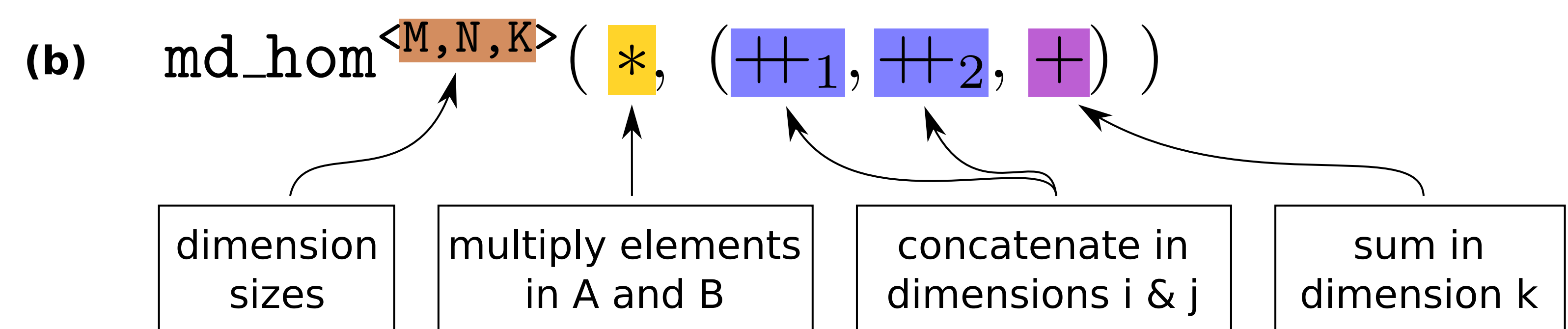
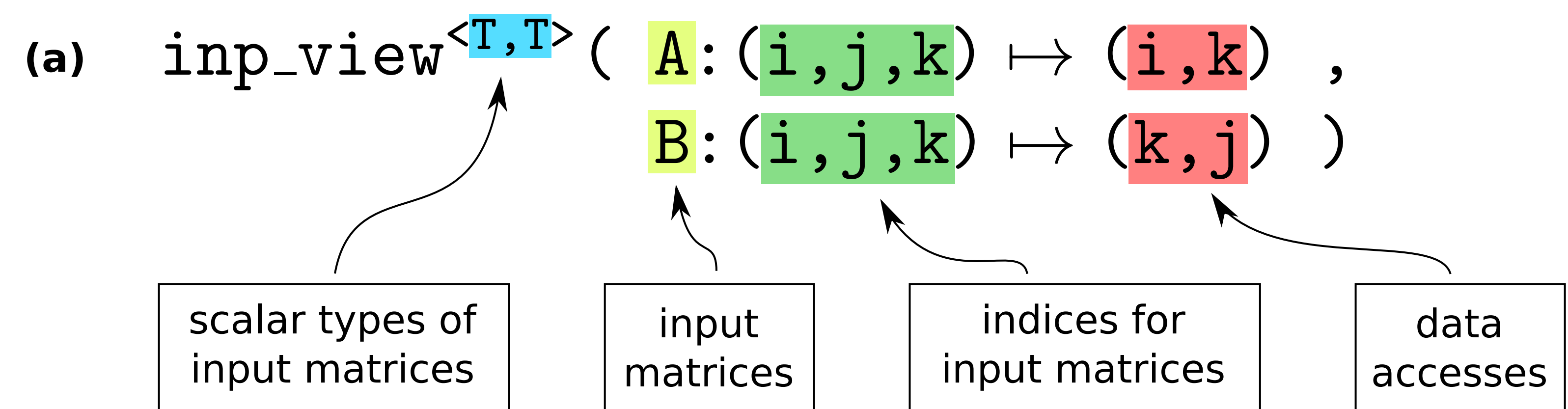
We present our work-in-progress **code generation and optimization approach** for **DL computations**:

- 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

$\text{MatMul}^{\langle T \in \text{Type} \mid M, N, K \in \mathbb{N} \rangle} := \text{out\_view}^{\langle \dots \rangle} (\dots) \circ$  (c)  
 $\text{md\_hom}^{\langle \dots \rangle} (\dots) \circ$  (b)  
 $\text{inp\_view}^{\langle \dots \rangle} (\dots)$  (a)

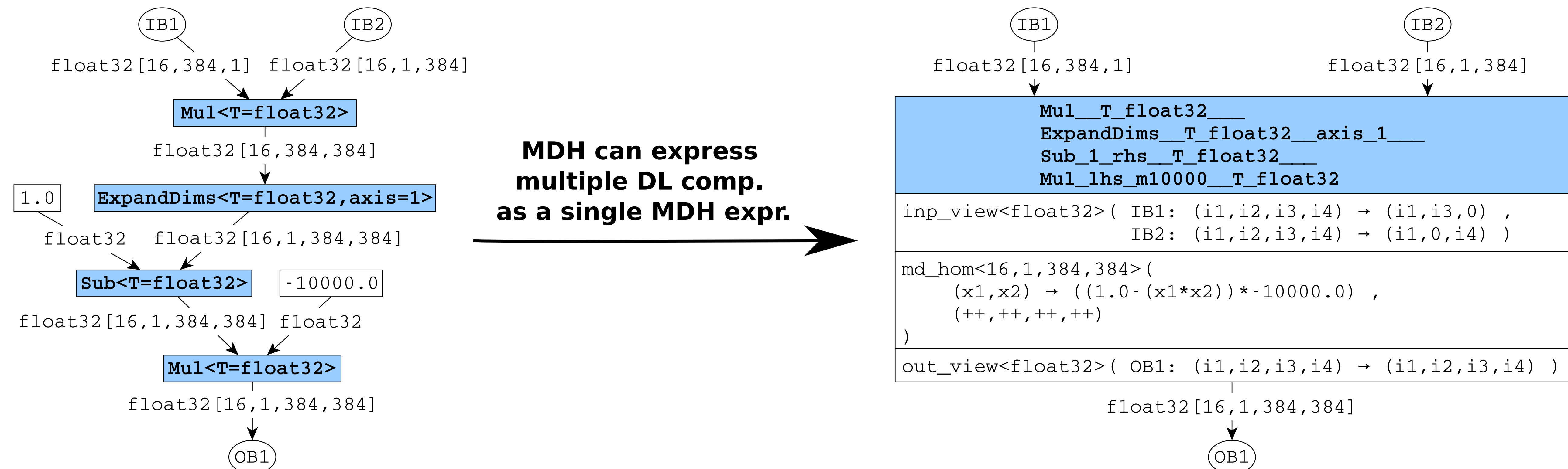


MDH allows us conveniently expressing data-parallel computations and automatically generate CUDA code for them [2, 3].

## DL Computations Expressed in the MDH Formalism

Operator	out_view <sup>&lt;...&gt;</sup>	md_hom <sup>&lt;...&gt;</sup>	inp_view <sup>&lt;...&gt;</sup>
Mul <sup>&lt;...&gt;</sup>	OB1: (i, j) ↦ (i, j)	* , ( ++ <sub>1</sub> , ++ <sub>2</sub> )	IB1: (i, j) ↦ (i, j) , IB2: (i, j) ↦ (i, j)
Sub <sup>&lt;...&gt;</sup>	OB1: (i, j) ↦ (i, j)	- , ( ++ <sub>1</sub> , ++ <sub>2</sub> )	IB1: (i, j) ↦ (i, j) , IB2: (i, j) ↦ (i, j)
ExpandDims <sup>&lt;axis, DEN   ...&gt;</sup>	OB1: (i <sub>1</sub> , ..., i <sub>D</sub> ) ↦ (... , i <sub>axis-1</sub> , 0, i <sub>axis</sub> , ...)	id , ( ++ <sub>1</sub> , ..., ++ <sub>D</sub> )	IB1: (i <sub>1</sub> , ..., i <sub>D</sub> ) ↦ (i <sub>1</sub> , ..., i <sub>D</sub> )
BiasAddGrad <sup>&lt;NHWC   ...&gt;</sup>	OB1: (i, j) ↦ (j)	id , ( +, ++ <sub>2</sub> )	IB1: (i, j) ↦ (i, j)
BatchMatMul <sup>&lt;N, N1   ...&gt;</sup>	OB1: (b1, b2, i, j, k) ↦ (b1, b2, i, j)	* , ( ++ <sub>1</sub> , ..., ++ <sub>4</sub> , + )	IB1: (b1, b2, i, j, k) ↦ (b1, b2, i, k) , IB2: (b1, b2, i, j, k) ↦ (b1, b2, k, j)

Popular DL computations<sup>1</sup> are conveniently expressed in the MDH formalism.



MDH can express multiple DL comp. as a single MDH expr.

<sup>1</sup> 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 competitors 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