

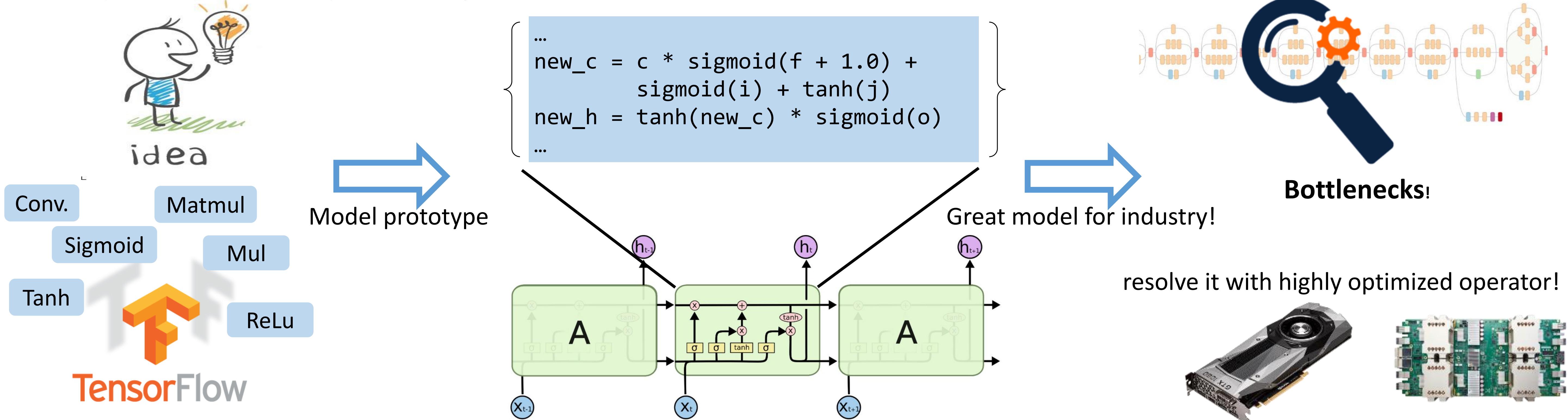
Optimization Mapping for Deep Learning

Wencong Xiao[†], Cheng Chen^{*}, Youshan Miao^{*}, Jilong Xue^{*}, Ming Wu^{*}

[†]Beihang University, ^{*}Microsoft Research

Motivating Scenario

Common process for deep learning



Hardware trend: heterogeneous devices with various accelerated libs

- Custom scenarios: cloud, mobile
- Custom algorithms: CNN, RNN



GPU-based:
nvidia cuDNN

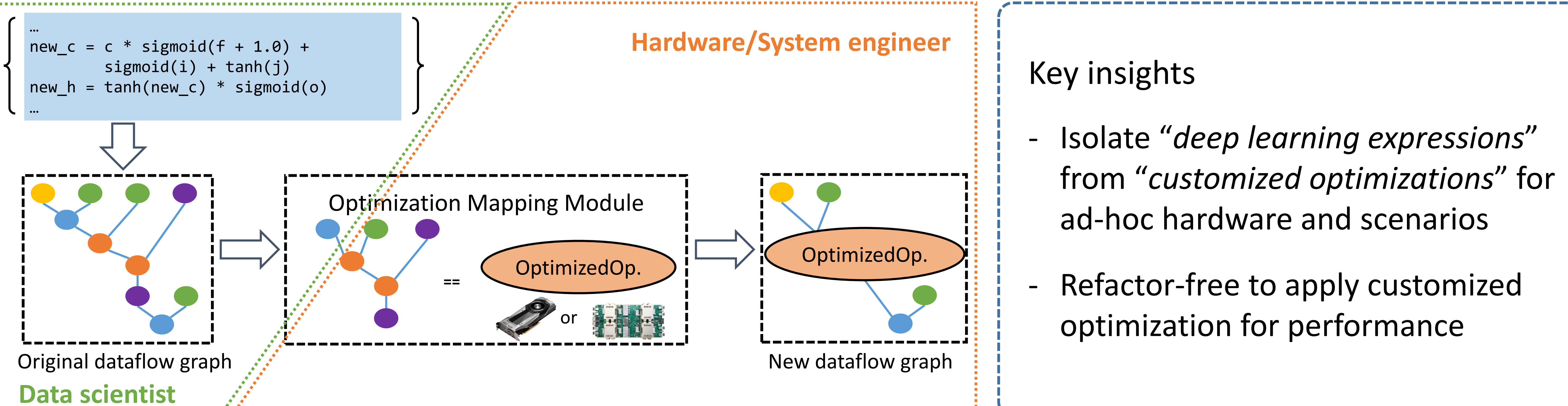


FPGA-based:
Microsoft Brainwave

ASIC-based: Google TPU
Intel Nervana
Cambricon
...

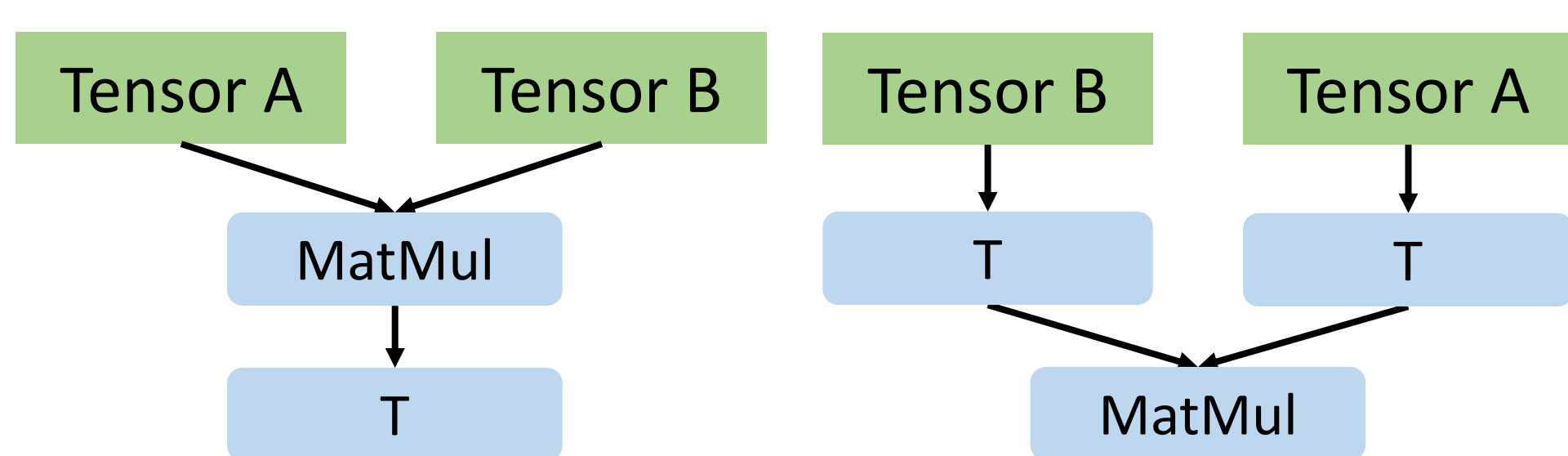


Optimization Mapping

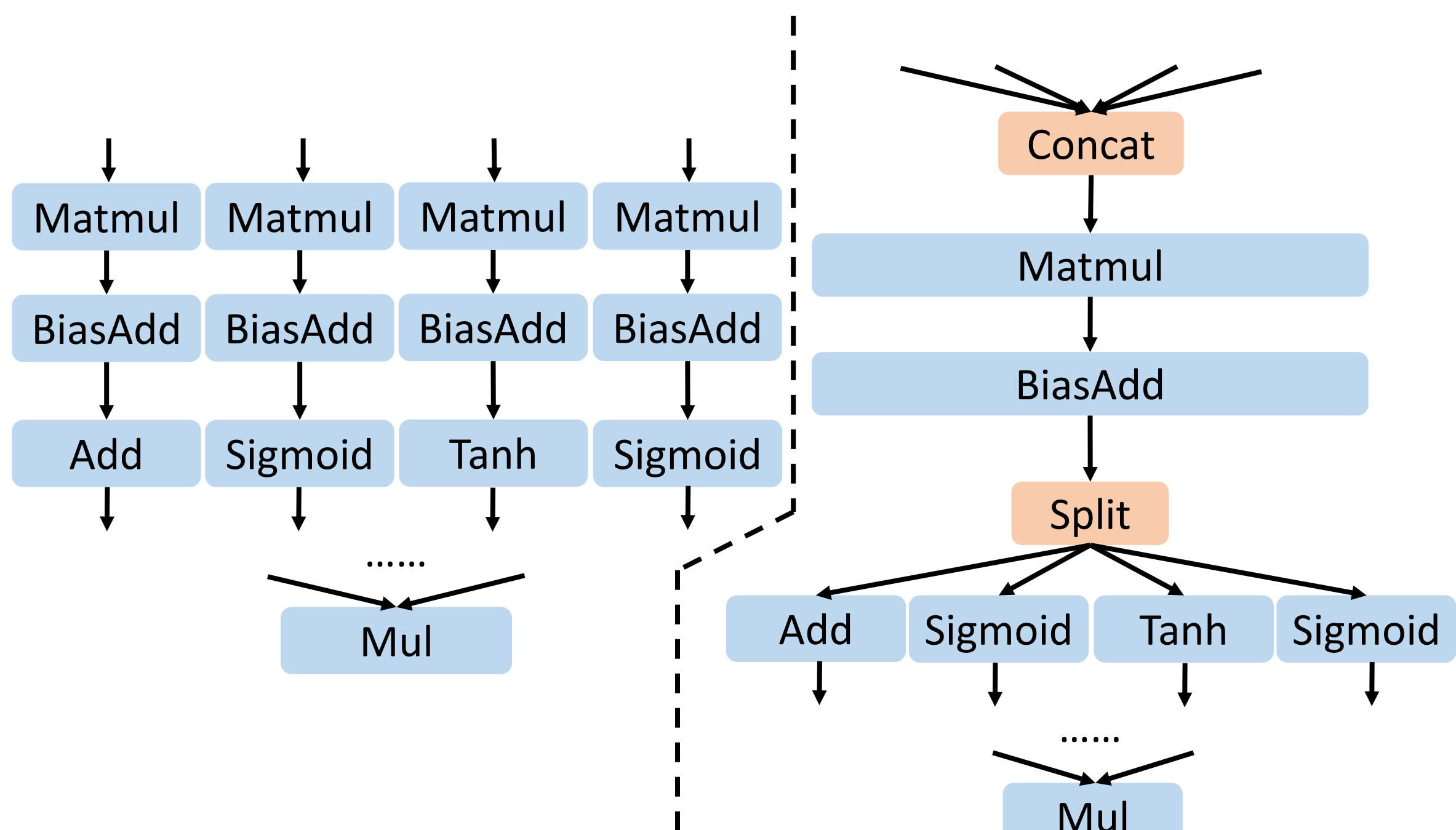


Key challenges, technologies, results

Intermediate representation without ambiguity

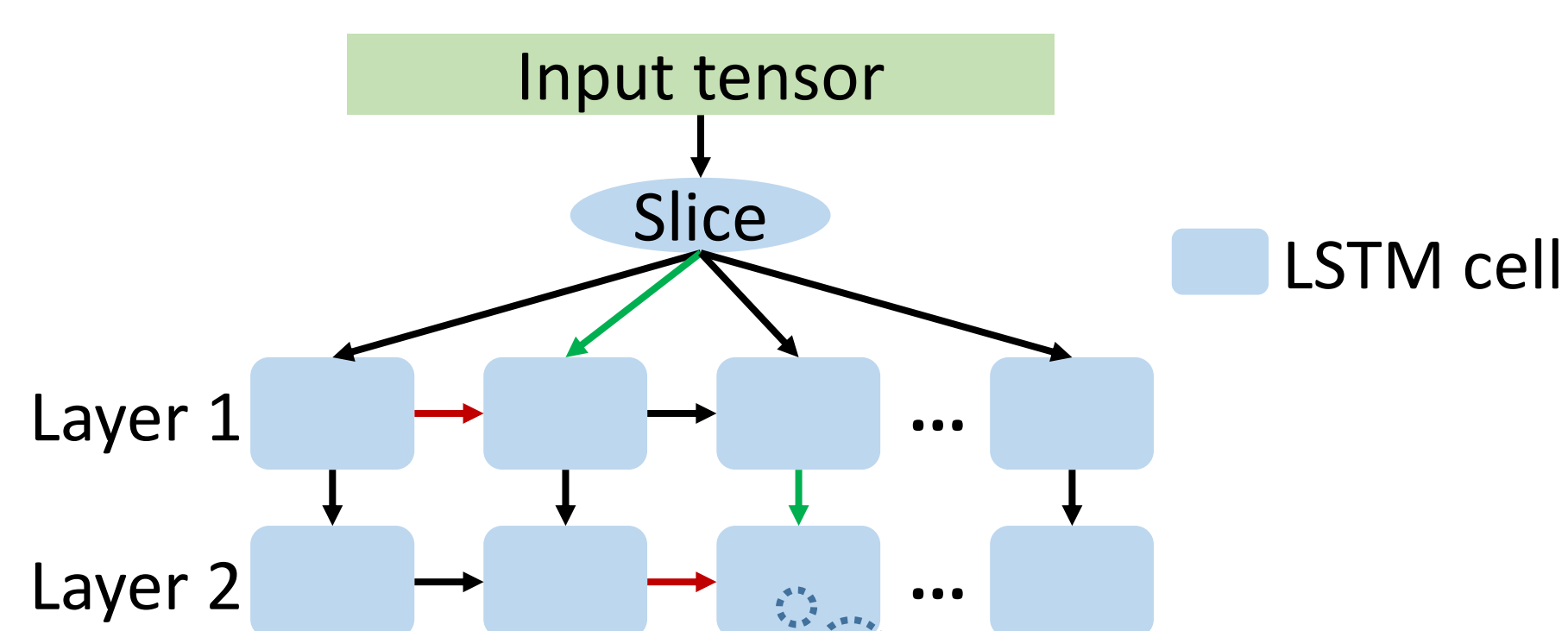


Math expression: $(AB)^T = B^T A^T$



Non-math expression: Native LSTM = Concat-split LSTM

Dynamic subgraph matching



```

cuDNN_LSTM_pattern{
  "cell": { ... },
  "op_type": {
    cuDNN_LSTM,
    #layer,
    #step},
  "cell_input": [
    {external_input | other_cell},
    {other_cell}]
}
    
```

Search space optimization

- Heterogeneous vertices
- Bottom-up search
- Outputs of operator can be used by unlimited operators
- Inputs of operator are limited

Preliminary results

- Implemented as an optimizer in Tensorflow r1.3
- Leverage defined cuDNN LSTM pattern to automatically map cuDNN LSTM operator to native LSTM
- Improve performance by 4.12x with refactor-free

