

Auto-scheduler for TVM

Lianmin Zheng
Shanghai Jiao Tong University

Towards “full automation”

Towards “full automation”

Tensor Expression

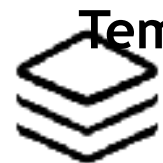
$$C_{ij} = \sum_k A_{ik} B_{kj}$$

Towards “full automation”

Tensor Expression

$$C_{ij} = \sum_k A_{ik} B_{kj}$$

AutoTVM Schedule



Template

Knob 1 ...

Knob 2 ...

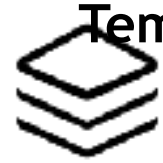
Knob 3 ...

Towards “full automation”

Tensor Expression

$$C_{ij} = \sum_k A_{ik} B_{kj}$$

AutoTVM Schedule



Template

Knob 1 ...
Knob 2 ...
Knob 3 ...



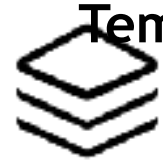
Optimizer

Towards “full automation”

Tensor Expression

$$C_{ij} = \sum_k A_{ik} B_{kj}$$

AutoTVM Schedule

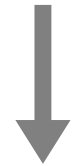


Template

Knob 1 ...
Knob 2 ...
Knob 3 ...



Optimizer



Final Optimized
Code

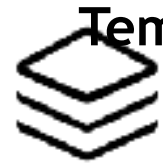


Towards “full automation”

Tensor Expression

$$C_{ij} = \sum_k A_{ik} B_{kj}$$

AutoTVM Schedule



Template

Knob 1 ...
Knob 2 ...
Knob 3 ...



Some
human
effort



Optimizer



Final Optimized
Code

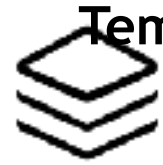


Towards “full automation”

Tensor Expression

$$C_{ij} = \sum_k A_{ik} B_{kj}$$

AutoTVM Schedule



Template

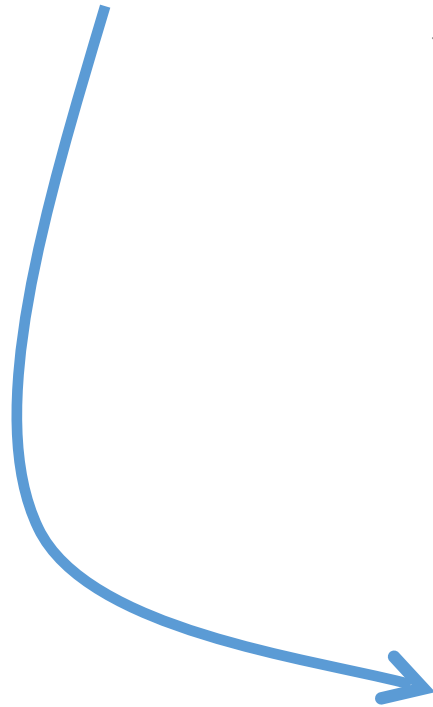
Knob 1 ...
Knob 2 ...
Knob 3 ...



Some
human
effort

Optimizer

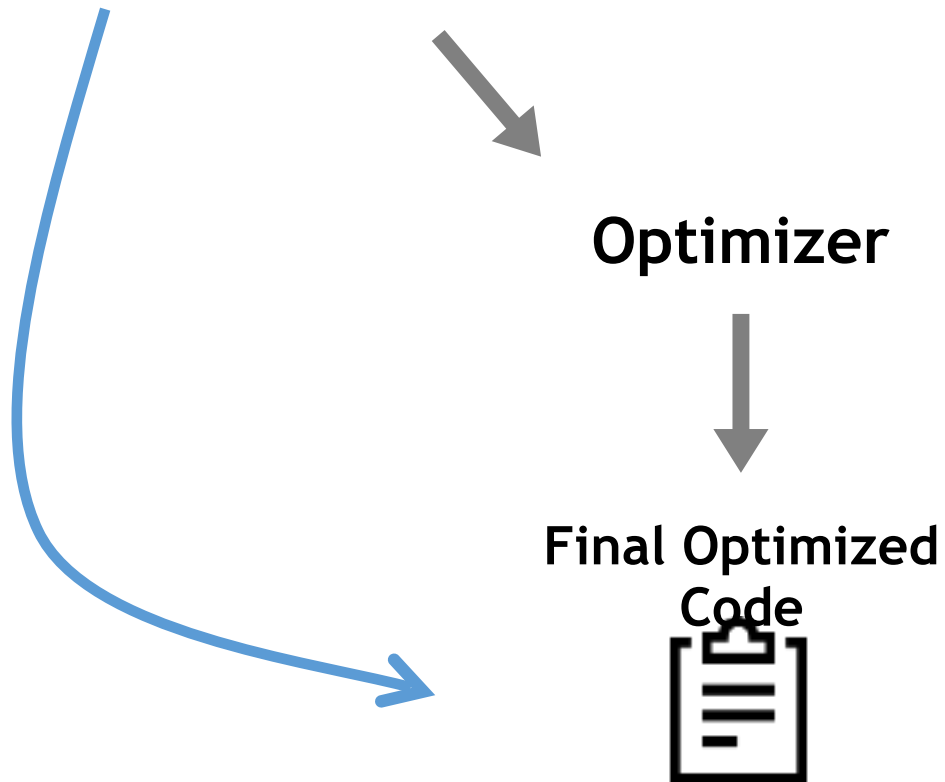
Final Optimized
Code



Towards “full automation”

Tensor Expression

$$C_{ij} = \sum_k A_{ik} B_{kj}$$



Auto-Scheduler

Tensor Expression

$$C_{ij} = \sum_k A_{ik} B_{kj}$$



Deduce
Tuning
Space

**Auto-
Scheduler**



Generate
Schedule
Template



Optimizer



**Final Optimized
Code**



Auto-scheduler

Tensor Expression

$$C_{N,CO,H,W} = \sum_{CI, KH, KW} A_{N,CO,H+KH, W+KW} B_{CO, CI, KH, KW}$$

Auto-scheduler

Tensor Expression

$$C_{N,CO,H,W} = \sum_{CI, KH, KW} A_{N,CO,H+KH, W+KW} B_{CO, CI, KH, KW}$$

Auto-packing

- for vectorization
- for tensorization
- other packing strategies

Auto-scheduler

Tensor Expression

$$C_{N,CO,H,W} = \sum_{CI, KH, KW} A_{N,CO,H+KH, W+KW} B_{CO, CI, KH, KW}$$

Auto-packing

- for vectorization
- for tensorization
- other packing strategies

Tile-knob on axes

Tune multi-level tiling on N, H, W, CO, CI

Auto-scheduler

Tensor Expression

$$C_{N,CO,H,W} = \sum_{CI, KH, KW} A_{N,CO,H+KH, W+KW} B_{CO, CI, KH, KW}$$

Auto-packing

- for vectorization
- for tensorization
- other packing strategies

Tile-knob on axes

Tune multi-level tiling on N, H, W, CO, CI

Location-knob on buffers

Tune *compute_at* location for A, B, C and packing buffers

Auto-scheduler

Tensor Expression

$$C_{N,CO,H,W} = \sum_{CI, KH, KW} A_{N,CO,H+KH, W+KW} B_{CO, CI, KH, KW}$$

Auto-packing

- for vectorization
- for tensorization
- other packing strategies

Tile-knob on axes

Tune multi-level tiling on N, H, W, CO, CI

Location-knob on buffers

Tune *compute_at* location for A, B, C and packing buffers

Annotation-knob on axes

Tune *parallel, unroll*

Auto-scheduler

Tensor Expression

$$C_{N,CO,H,W} = \sum_{CI, KH, KW} A_{N,CO,H+KH, W+KW} B_{CO, CI, KH, KW}$$

Auto-packing

- for vectorization
- for tensorization
- other packing strategies

Tile-knob on axes

Tune multi-level tiling on N, H, W, CO, CI

Location-knob on buffers

Tune *compute_at* location for A, B, C and packing buffers

Annotation-knob on axes

Tune *parallel, unroll*

Map to different backends

Map threading, memory scoping, tensorization

Goal and progress

Replace manual schedules in
current TOPI without
regression

Provide schedules for
new operators:
training ops, conv3d,
...



One month