



Optimizing DNN Algorithms on Heterogenous Computing

Supervisors: Professor Jia, Jiaya & Professor Yu, Bei

Mentor: Mr. Qi, Sun

LJ2001

Department of Computer Science and Engineering

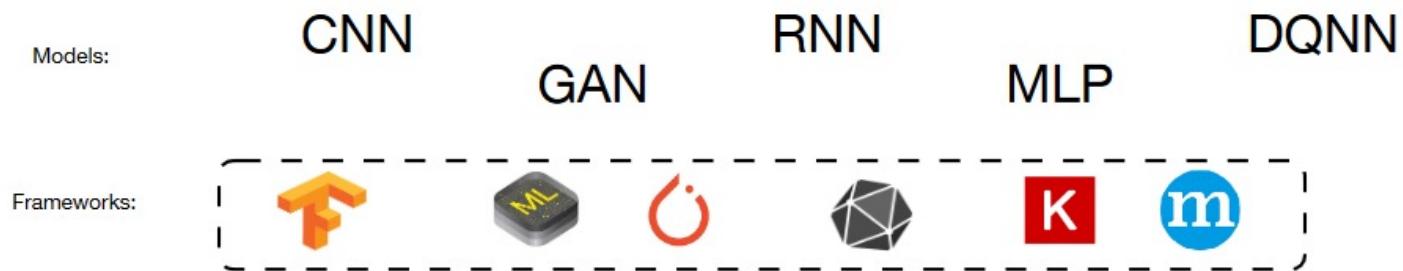
April 26, 2021

Presented by Choi, Jang
Hyeon

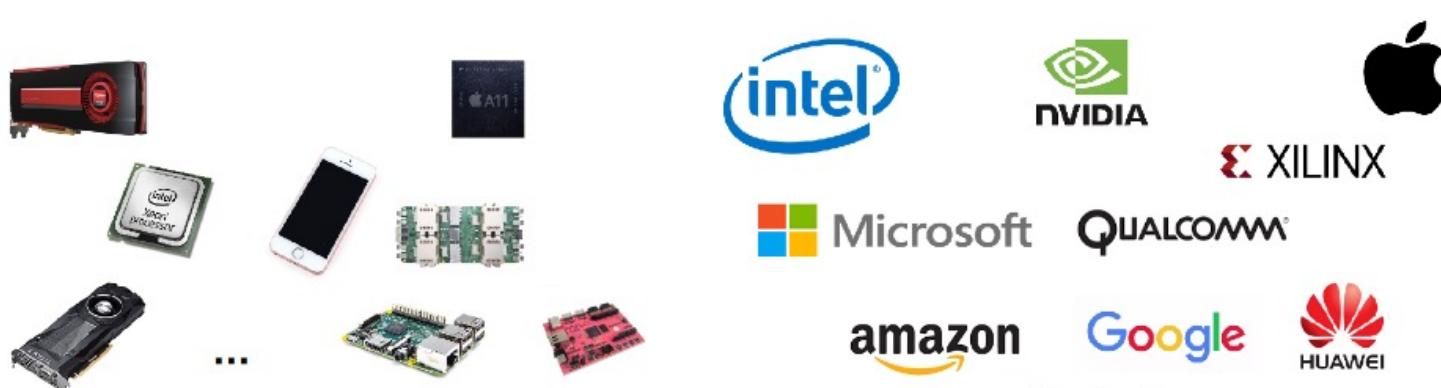
Outline

- ▶ Enlarging the configuration space
 - ▶ Adding parameter candidates.
 - ▶ Introducing new knobs into a design space.
- ▶ Attempts on merging the configuration space

Brief Recap

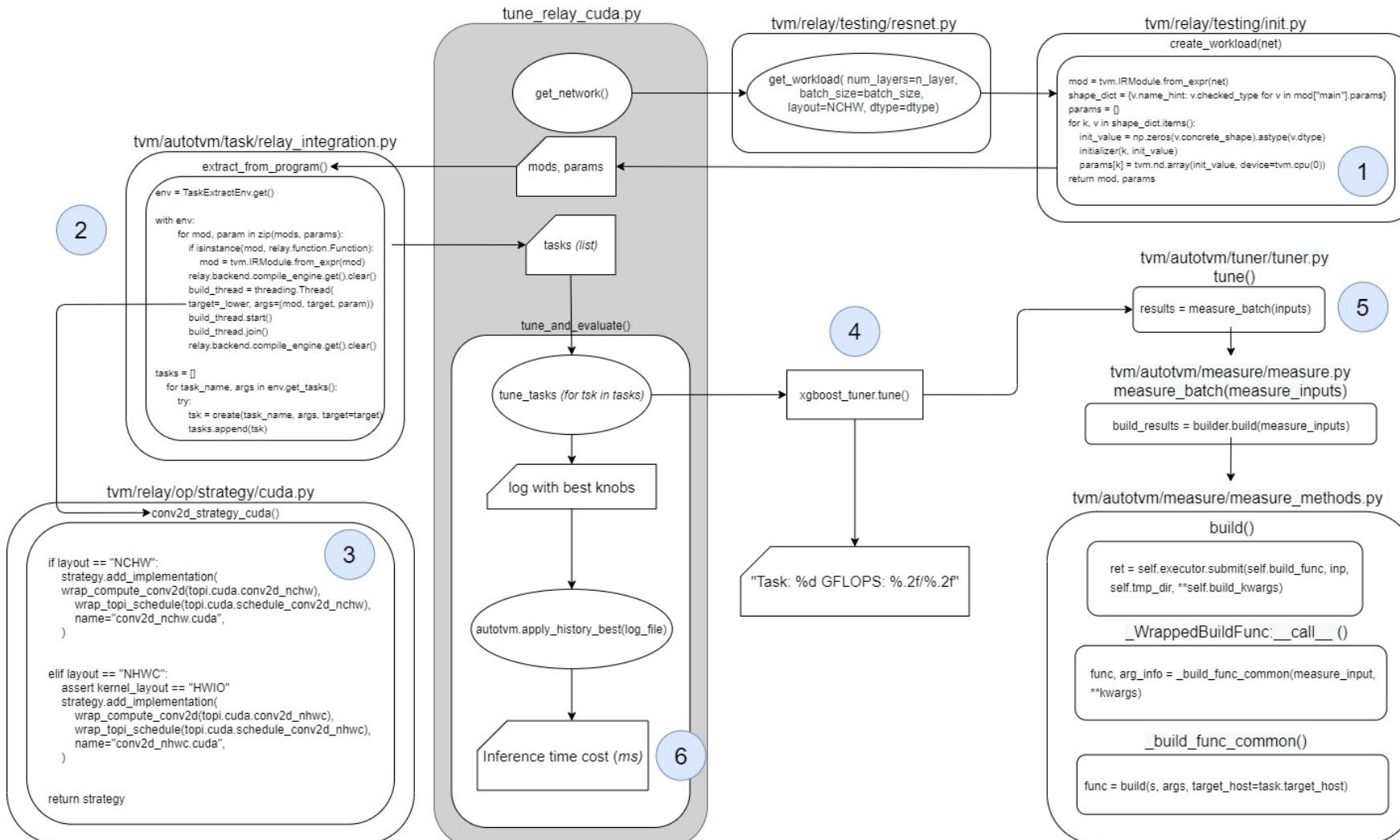


Challenge: Efficiently deploying deep learning everywhere



<https://sampl.cs.washington.edu/tvmconf/#about-tvmconf>

General Workflow of AutoTVM



<https://github.com/apache/tvm/tree/main/python/tvm>

https://tvm.apache.org/docs/tutorials/autotvm/tune_relay_cuda.html

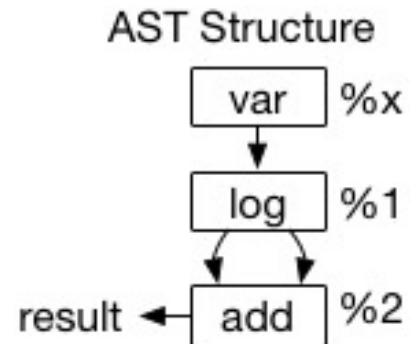
IR Module

Python Code

```
x = relay.var("x")
v1 = relay.log(x)
v2 = relay.add(v1, v1)
f = relay.Function([x], v2)
```

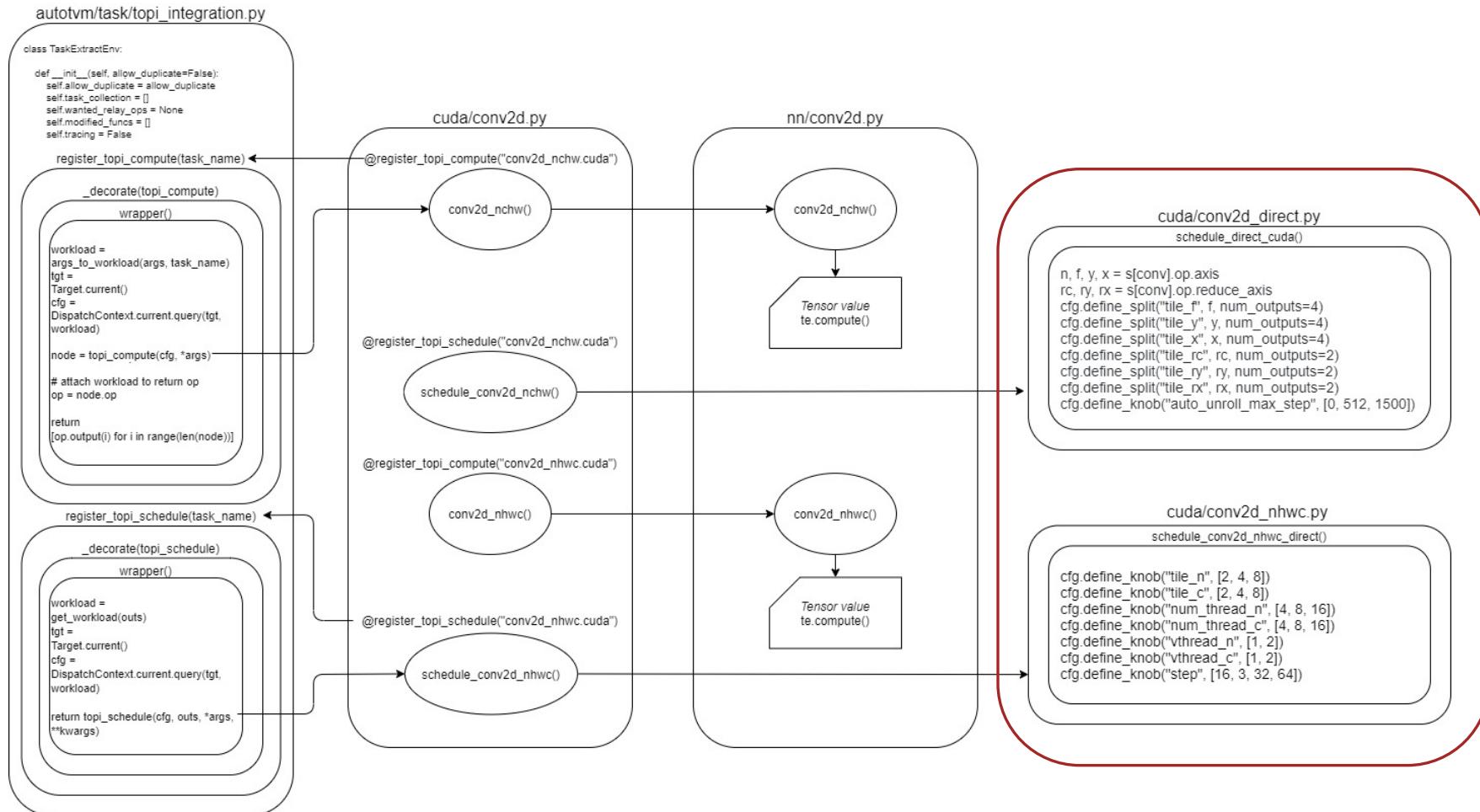
Text Form

```
fn (%x) {
    %1 = log(%x)
    %2 = add(%1, %1)
    %2
}
```



https://tvm.apache.org/docs/dev/relay_intro.html

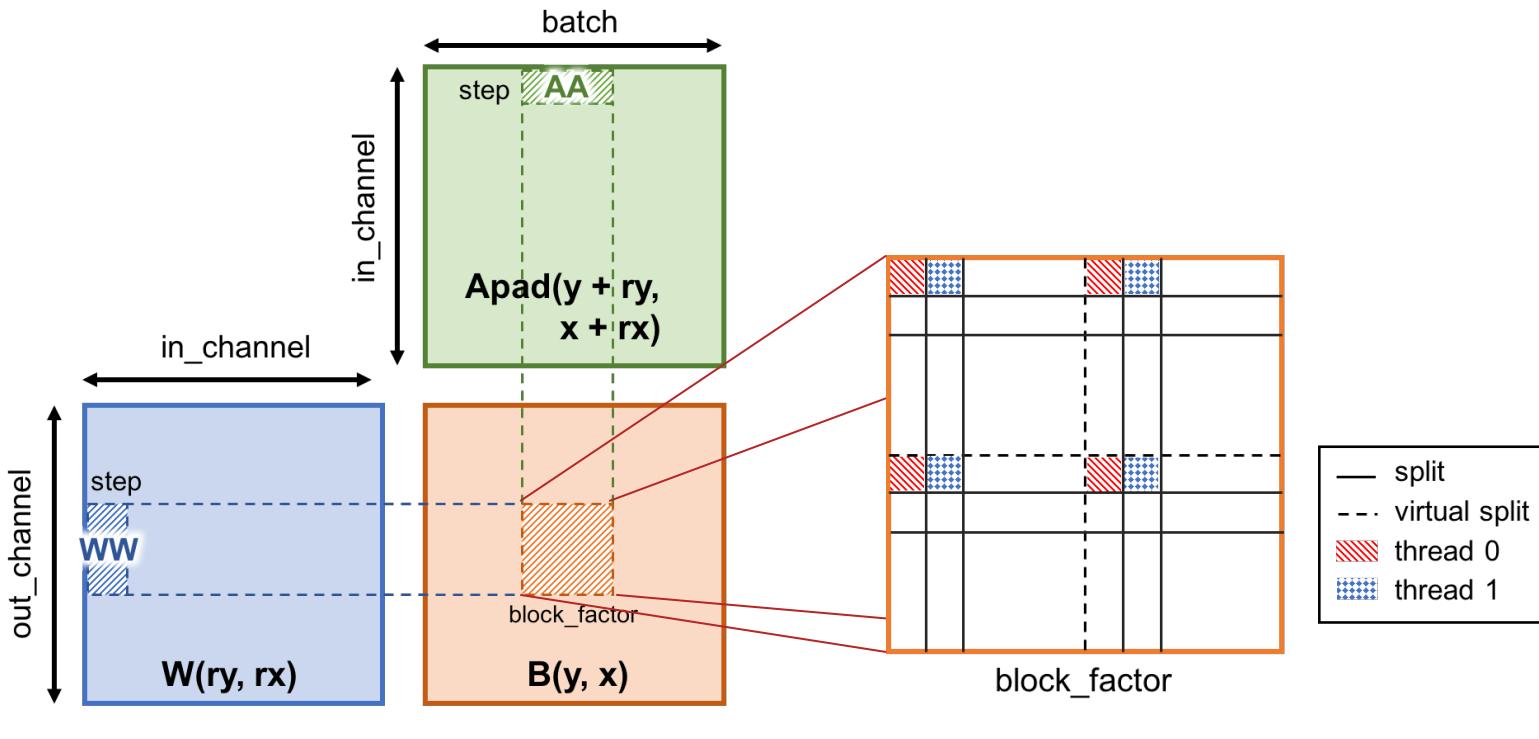
Design Space Templates



https://github.com/apache/tvm/blob/main/python/tvm/topi/cuda/conv2d_direct.py

https://github.com/apache/tvm/blob/main/python/tvm/topi/cuda/conv2d_nhwc.py

Defining Splits vs. Knobs



```

A = t.placeholder((1024, 1024))
B = t.placeholder((1024, 1024))
k = t.reduce_axis((0, 1024))
C = t.compute((1024, 1024), lambda y, x:
    t.sum(A[k, y] * B[k, x], axis=k))
s = t.create_schedule(C.op)

for y in range(1024):
    for x in range(1024):
        C[y][x] = 0
        for k in range(1024):
            C[y][x] += A[k][y] * B[k][x]

+ Loop Tiling
yo, xo, ko, yi, xi, ki = s[C].tile(y, x, k, 8, 8, 8)

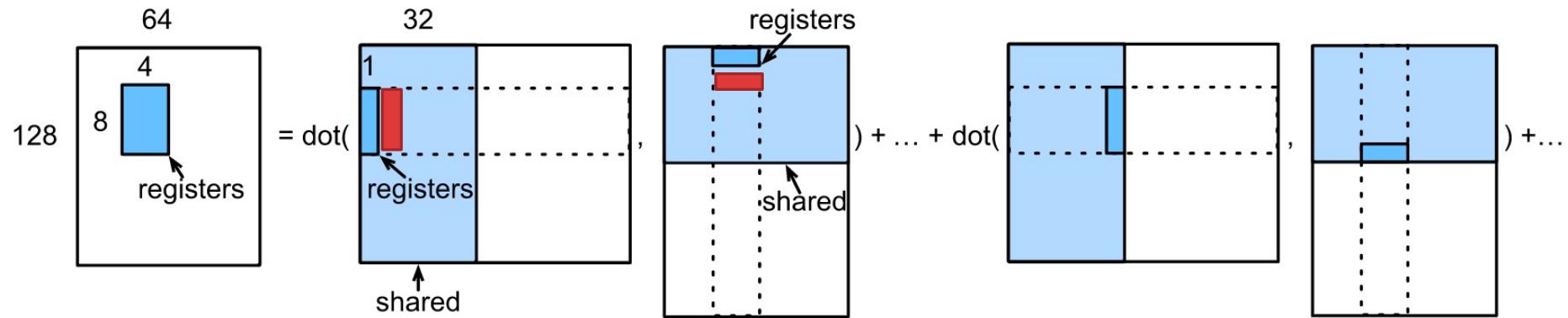
for yo in range(128):
    for xo in range(128):
        C[yo*8:yo*8+8][xo*8:xo*8+8] = 0
        for ko in range(128):
            for yi in range(8):
                for xi in range(8):
                    for ki in range(8):
                        C[yo*8+yi][xo*8+xi] +=
                            A[ko*8+ki][yo*8+yi] * B[ko*8+ki][xo*8+xi]

```

https://tvm.apache.org/docs/tutorials/optimize/opt_conv_cuda.html

<https://www.programmersought.com/article/38296967626/>

Matrix Computation



https://tvm.d2l.ai/chapter_gpu_schedules/matmul.html

Enlarging the Configuration Space

cuda/conv2d_direct.py
schedule_direct_cuda()

```
n, f, y, x = s[conv].op.axis
rc, ry, rx = s[conv].op.reduce_axis
cfg.define_split("tile_f", f, num_outputs=4)
cfg.define_split("tile_y", y, num_outputs=4)
cfg.define_split("tile_x", x, num_outputs=4)
cfg.define_split("tile_rc", rc, num_outputs=2)
cfg.define_split("tile_ry", ry, num_outputs=2)
cfg.define_split("tile_rx", rx, num_outputs=2)
cfg.define_knob("auto_unroll_max_step", [0, 512, 1500])
```

cuda/conv2d_direct.py
schedule_direct_cuda()

```
n, f, y, x = s[conv].op.axis
rc, ry, rx = s[conv].op.reduce_axis
cfg.define_split("tile_f", f, num_outputs=4)
cfg.define_split("tile_y", y, num_outputs=4)
cfg.define_split("tile_x", x, num_outputs=4)
cfg.define_split("tile_rc", rc, num_outputs=2)
cfg.define_split("tile_ry", ry, num_outputs=2)
cfg.define_split("tile_rx", rx, num_outputs=2)
cfg.define_knob("auto_unroll_max_step", [0, 512, 1024, 1500])
```

cuda/conv2d_nhwc.py
schedule_conv2d_nhwc_direct()

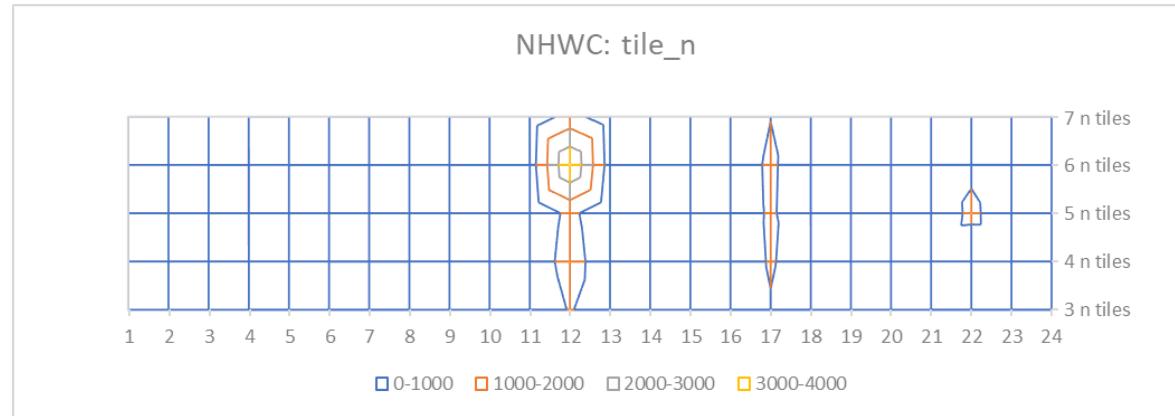
```
cfg.define_knob("tile_n", [2, 4, 8])
cfg.define_knob("tile_c", [2, 4, 8])
cfg.define_knob("num_thread_n", [4, 8, 16])
cfg.define_knob("num_thread_c", [4, 8, 16])
cfg.define_knob("vthread_n", [1, 2])
cfg.define_knob("vthread_c", [1, 2])
cfg.define_knob("step", [16, 3, 32, 64])
```

cuda/conv2d_nhwc.py
schedule_conv2d_nhwc_direct()

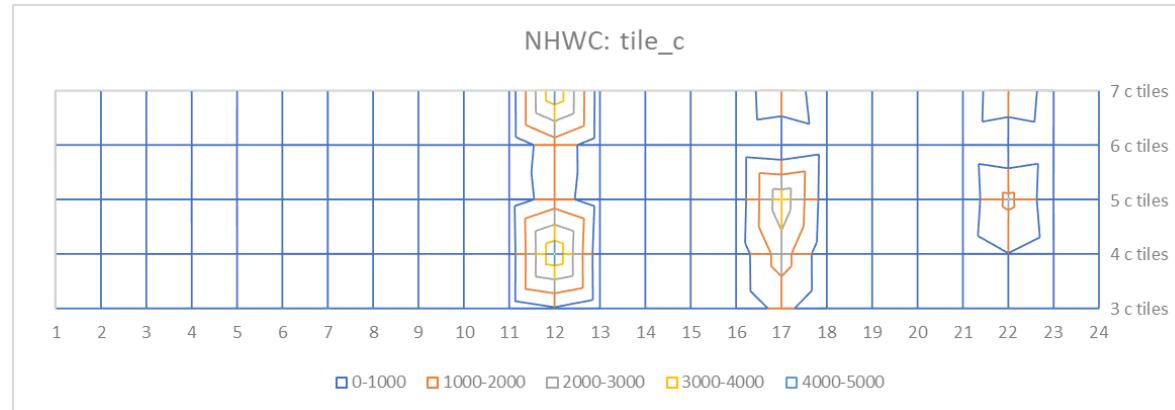
```
cfg.define_knob("tile_n", [2, 4, 8, 16, 32, 64, 128])
cfg.define_knob("tile_c", [2, 4, 8, 16, 32, 64, 128])
cfg.define_knob("num_thread_n", [4, 8, 16])
cfg.define_knob("num_thread_c", [4, 8, 16])
cfg.define_knob("vthread_n", [1, 2])
cfg.define_knob("vthread_c", [1, 2])
cfg.define_knob("step", [16, 3, 32, 64])
```

Enlarging the Configuration Space: Results

NHWC: tile_n



NHWC: tile_c

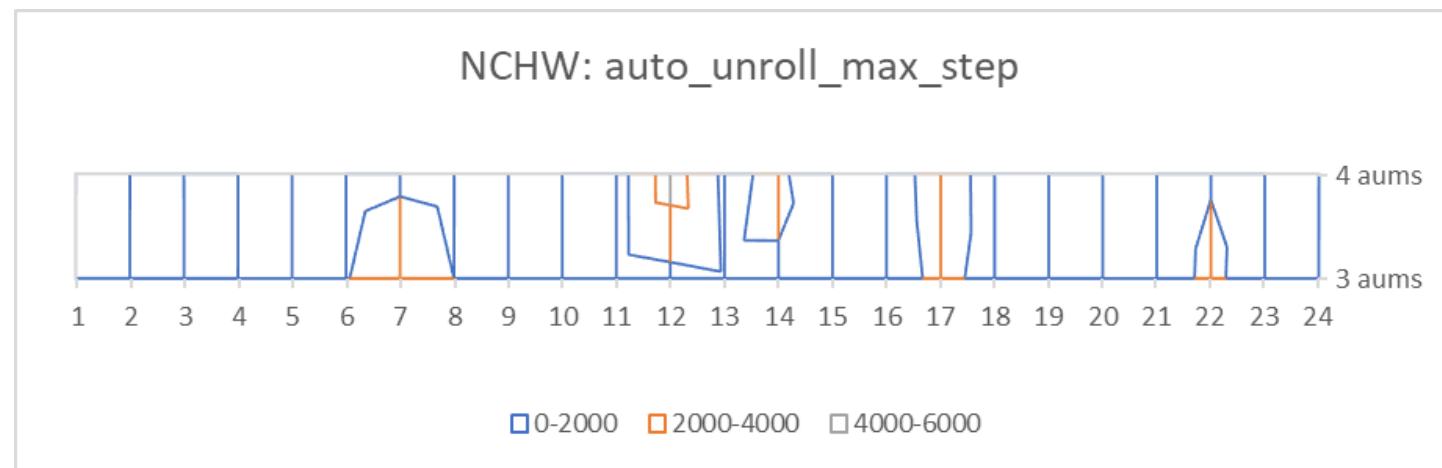




Enlarging the Configuration Space: Results

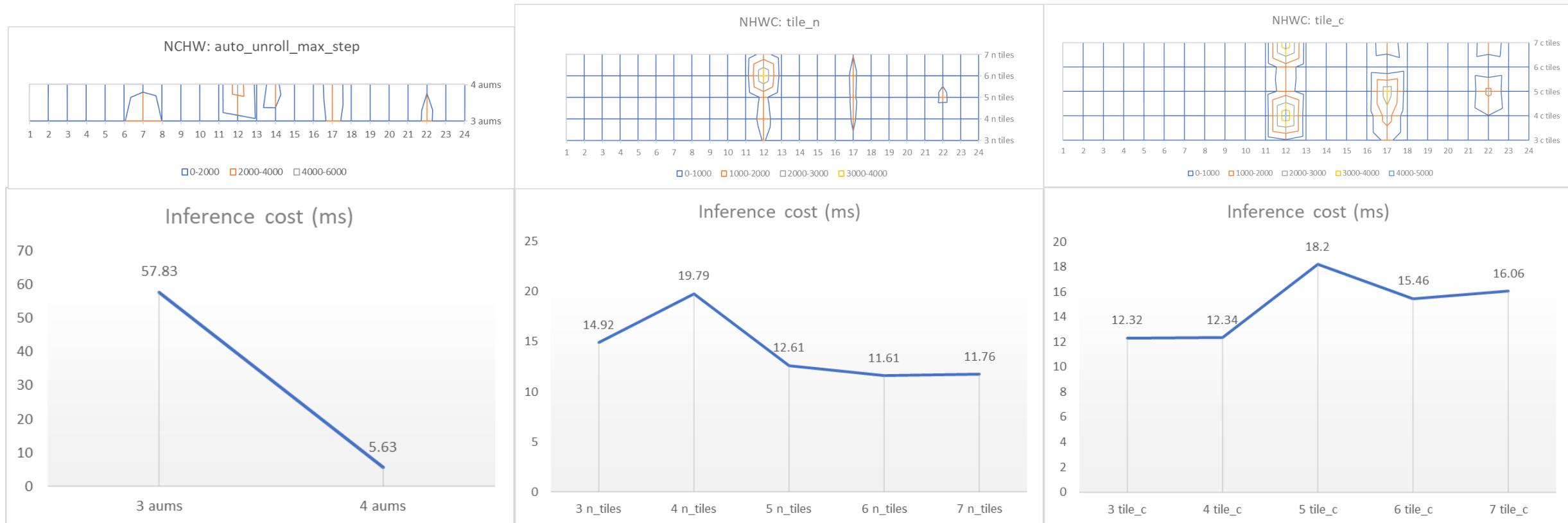


NCHW: auto_unroll_max_step





Enlarging the Configuration Space: Results



Enlarging the Configuration Space II

```
cuda/conv2d_direct.py —————> cuda/conv2d_direct.py

schedule_direct_cuda()
n, f, y, x = s[conv].op.axis
rc, ry, rx = s[conv].op.reduce_axis
cfg.define_split("tile_f", f, num_outputs=4)
cfg.define_split("tile_y", y, num_outputs=4)
cfg.define_split("tile_x", x, num_outputs=4)
cfg.define_split("tile_rc", rc, num_outputs=2)
cfg.define_split("tile_ry", ry, num_outputs=2)
cfg.define_split("tile_rx", rx, num_outputs=2)
cfg.define_knob("auto_unroll_max_step", [0, 512, 1500])

cuda/conv2d_nhwc.py
schedule_conv2d_nhwc_direct()
cfg.define_knob("tile_n", [2, 4, 8])
cfg.define_knob("tile_c", [2, 4, 8])
cfg.define_knob("num_thread_n", [4, 8, 16])
cfg.define_knob("num_thread_c", [4, 8, 16])
cfg.define_knob("vthread_n", [1, 2])
cfg.define_knob("vthread_c", [1, 2])
cfg.define_knob("step", [16, 3, 32, 64])

schedule_direct_cuda()
n, f, y, x = s[conv].op.axis
rc, ry, rx = s[conv].op.reduce_axis
cfg.define_split("tile_f", f, num_outputs=4)
cfg.define_split("tile_y", y, num_outputs=4)
cfg.define_split("tile_x", x, num_outputs=4)
cfg.define_split("tile_rc", rc, num_outputs=2)
cfg.define_split("tile_ry", ry, num_outputs=2)
cfg.define_split("tile_rx", rx, num_outputs=2)
cfg.define_knob("num_thread_n", [4, 8, 16])
cfg.define_knob("num_thread_c", [4, 8, 16])
cfg.define_knob("vthread_n", [1, 2])
cfg.define_knob("vthread_c", [1, 2])
cfg.define_knob("auto_unroll_max_step", [0, 512, 1500])
target = tvm.target.Target.current()
if target.kind.name == "nvptx", "rocm":
    cfg.define_knob("unroll_explicit", [1])
else:
    cfg.define_knob("unroll_explicit", [0, 1])
num_thread_n = cfg["num_thread_n"].val
num_thread_c = cfg["num_thread_c"].val
vthread_n = cfg["vthread_n"].val
vthread_c = cfg["vthread_c"].val
thread_x = te.thread_axis(0, num_thread_c, "threadIdx.x")
thread_y = te.thread_axis(0, num_thread_n, "threadIdx.y")
thread_z = te.thread_axis((0, vthread_c), "vthread", name="vx")
thread_yz = te.thread_axis((0, vthread_n), "vthread", name="vy")

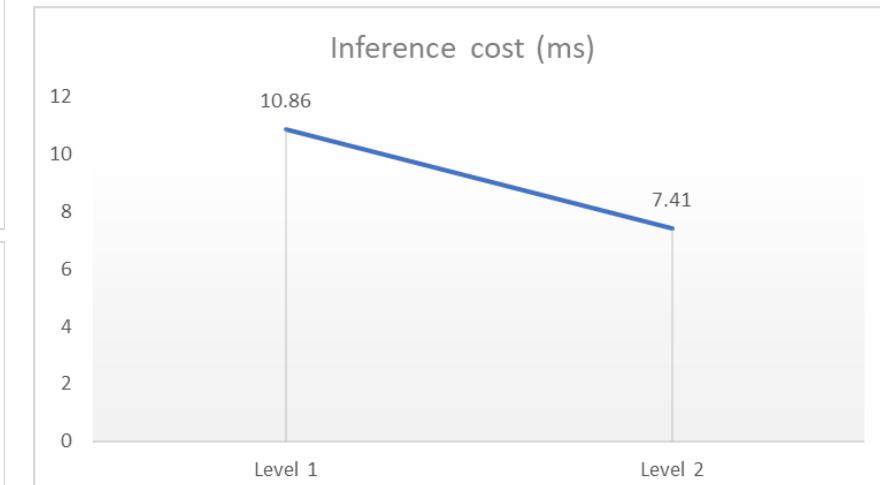
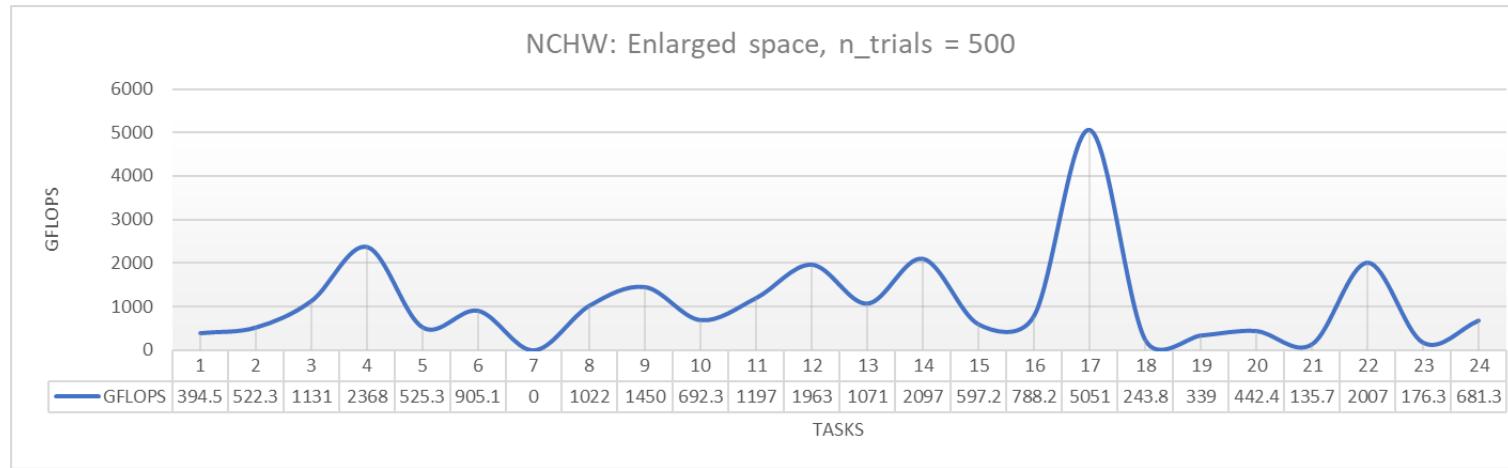
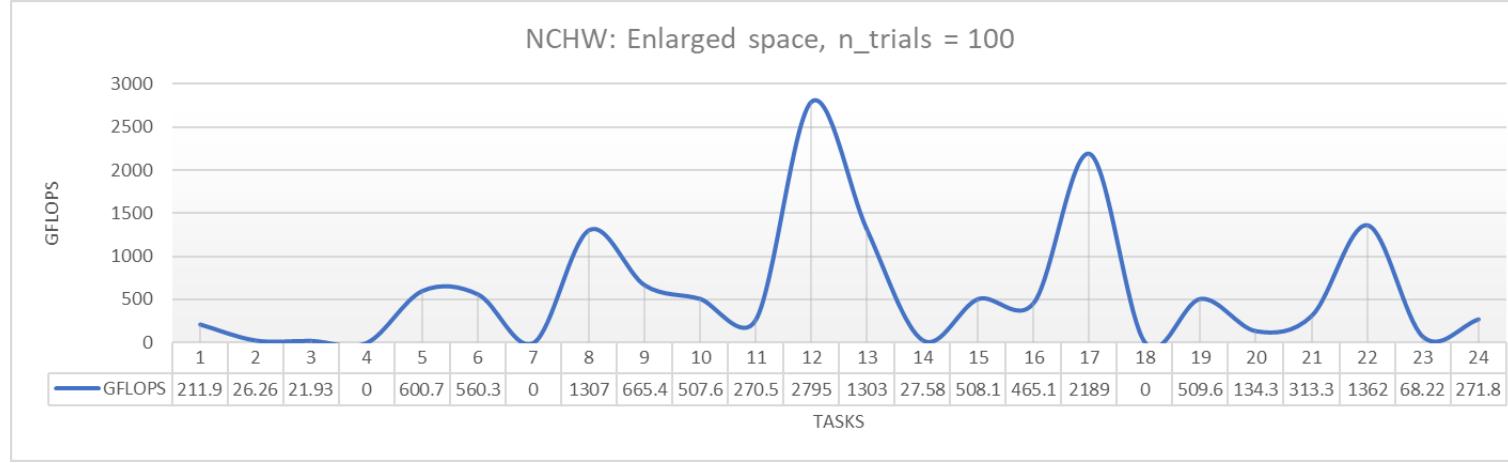
# tile and bind spatial axes
n, f, y, x = s[output].op.axis
kernel_scope, n = s[output].split(n, nparts=1)
bf, vf, tf, fi = cfg["tile_f"].apply(s, output, f)
by, vy, ty, yi = cfg["tile_y"].apply(s, output, y)
bx, vx, tx, xi = cfg["tile_x"].apply(s, output, x)

bf = s[output].fuse(n, bf)
s[output].bind(bf, te.thread_axis("blockIdx.z"))
s[output].bind(by, te.thread_axis("blockIdx.y"))
s[output].bind(bx, te.thread_axis("blockIdx.x"))
s[output].bind(vf, te.thread_axis("vthread"))
s[output].bind(vy, thread_yz)
s[output].bind(vx, thread_xz)
s[output].bind(tx, thread_y)
s[output].bind(xi, thread_z)
s[output].reorder(bf, by, vf, vy, vx, tf, ty, bx, fi, yi, xi)
s[OL].compute_at(s[output], bx)

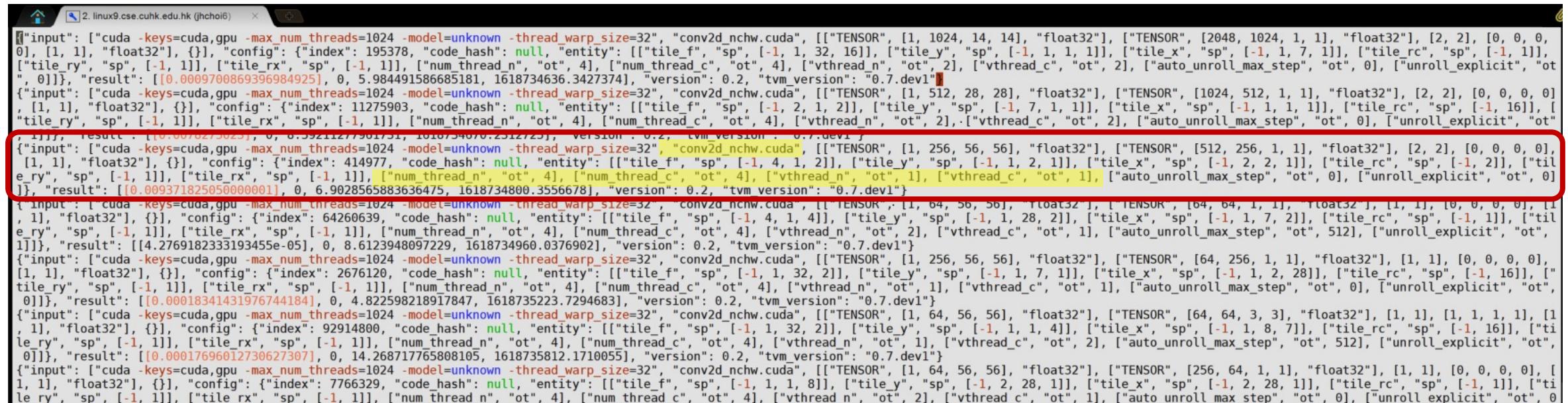
s[AA].compute_at(s[OL], rxo)
s[WW].compute_at(s[OL], rxo)
```



Enlarging the Configuration Space II: Results



Enlarging the Configuration Space II: Results



```
[{"input": ["cuda -keys=cuda,gpu -max_num_threads=1024 -model=unknown -thread_warp_size=32", "conv2d_nchw.cuda", [{"TENSOR": [1, 1024, 14, 14], "float32"}, {"TENSOR": [2048, 1024, 1, 1], "float32"}, [2, 2], [0, 0, 0, 0], [1, 1], "float32"], {}], "config": {"index": 195378, "code_hash": null, "entity": [{"tile_f", "sp", [-1, 1, 32, 16]}, {"tile_y", "sp", [-1, 1, 1, 1]}, {"tile_x", "sp", [-1, 1, 7, 1]}, {"tile_rc", "sp", [-1, 1]}, {"tile_ry", "sp", [-1, 1]}, {"tile_rx", "sp", [-1, 1]}, {"num_thread_n", "ot", 4}, {"num_thread_c", "ot", 4}, {"vthread_n", "ot", 2}, {"vthread_c", "ot", 2}, {"auto_unroll_max_step", "ot", 0}, {"unroll_explicit", "ot", 0}], "result": [[0.0009700869396984925], 0, 5.984491586685181, 1618734636.3427374], "version": 0.2, "tvm_version": "0.7.dev1"}], {"input": ["cuda -keys=cuda,gpu -max_num_threads=1024 -model=unknown -thread_warp_size=32", "conv2d_nchw.cuda", [{"TENSOR": [1, 512, 28, 28], "float32"}, {"TENSOR": [1024, 512, 1, 1], "float32"}, [2, 2], [0, 0, 0, 0], [1, 1], "float32"], {}], "config": {"index": 11275903, "code_hash": null, "entity": [{"tile_f", "sp", [-1, 2, 1, 2]}, {"tile_y", "sp", [-1, 7, 1, 1]}, {"tile_x", "sp", [-1, 1, 1, 1]}, {"tile_rc", "sp", [-1, 16]}, {"tile_ry", "sp", [-1, 1]}, {"tile_rx", "sp", [-1, 1]}, {"num_thread_n", "ot", 4}, {"num_thread_c", "ot", 4}, {"vthread_n", "ot", 2}, {"vthread_c", "ot", 2}, {"auto_unroll_max_step", "ot", 0}, {"unroll_explicit", "ot", 0}], "result": [[0.0009700869396984925], 0, 5.984491586685181, 1618734636.3427374], "version": 0.2, "tvm_version": "0.7.dev1"]], {"input": ["cuda -keys=cuda,gpu -max_num_threads=1024 -model=unknown -thread_warp_size=32", "conv2d_nchw.cuda", [{"TENSOR": [1, 256, 56, 56], "float32"}, {"TENSOR": [512, 256, 1, 1], "float32"}, [2, 2], [0, 0, 0, 0], [1, 1], "float32"], {}], "config": {"index": 414977, "code_hash": null, "entity": [{"tile_f", "sp", [-1, 4, 1, 2]}, {"tile_y", "sp", [-1, 1, 2, 1]}, {"tile_x", "sp", [-1, 2, 2, 1]}, {"tile_rc", "sp", [-1, 2]}, {"tile_ry", "sp", [-1, 1]}, {"tile_rx", "sp", [-1, 1]}, {"num_thread_n", "ot", 4}, {"num_thread_c", "ot", 4}, {"vthread_n", "ot", 1}, {"vthread_c", "ot", 1}, {"auto_unroll_max_step", "ot", 0}, {"unroll_explicit", "ot", 0}], "result": [[0.009371825050000001], 0, 6.9028565883636475, 1618734800.3556678], "version": 0.2, "tvm_version": "0.7.dev1"]], {"input": ["cuda -keys=cuda,gpu -max_num_threads=1024 -model=unknown -thread_warp_size=32", "conv2d_nchw.cuda", [{"TENSOR": [1, 64, 56, 56], "float32"}, {"TENSOR": [64, 64, 1, 1], "float32"}, [1, 1], [0, 0, 0, 0], [1, 1], "float32"], {}], "config": {"index": 64260639, "code_hash": null, "entity": [{"tile_f", "sp", [-1, 4, 1, 4]}, {"tile_y", "sp", [-1, 1, 28, 2]}, {"tile_x", "sp", [-1, 1, 7, 2]}, {"tile_rc", "sp", [-1, 1]}, {"tile_ry", "sp", [-1, 1]}, {"tile_rx", "sp", [-1, 1]}, {"num_thread_n", "ot", 4}, {"num_thread_c", "ot", 4}, {"vthread_n", "ot", 2}, {"vthread_c", "ot", 1}, {"auto_unroll_max_step", "ot", 512}, {"unroll_explicit", "ot", 1}], "result": [[4.2769182333193455e-05], 0, 8.6123948097229, 1618734960.0376902], "version": 0.2, "tvm_version": "0.7.dev1"]], {"input": ["cuda -keys=cuda,gpu -max_num_threads=1024 -model=unknown -thread_warp_size=32", "conv2d_nchw.cuda", [{"TENSOR": [1, 256, 56, 56], "float32"}, {"TENSOR": [64, 256, 1, 1], "float32"}, [1, 1], [0, 0, 0, 0], [1, 1], "float32"], {}], "config": {"index": 2676120, "code_hash": null, "entity": [{"tile_f", "sp", [-1, 1, 32, 2]}, {"tile_y", "sp", [-1, 1, 7, 1]}, {"tile_x", "sp", [-1, 1, 2, 28]}, {"tile_rc", "sp", [-1, 16]}, {"tile_ry", "sp", [-1, 1]}, {"tile_rx", "sp", [-1, 1]}, {"num_thread_n", "ot", 4}, {"num_thread_c", "ot", 4}, {"vthread_n", "ot", 1}, {"vthread_c", "ot", 1}, {"auto_unroll_max_step", "ot", 0}, {"unroll_explicit", "ot", 0}], "result": [[0.00018341431976744184], 0, 4.822598218917847, 1618735223.7294683], "version": 0.2, "tvm_version": "0.7.dev1"]], {"input": ["cuda -keys=cuda,gpu -max_num_threads=1024 -model=unknown -thread_warp_size=32", "conv2d_nchw.cuda", [{"TENSOR": [1, 64, 56, 56], "float32"}, {"TENSOR": [64, 64, 3, 3], "float32"}, [1, 1], [1, 1, 1, 1], [1, 1], "float32"], {}], "config": {"index": 92914800, "code_hash": null, "entity": [{"tile_f", "sp", [-1, 1, 32, 2]}, {"tile_y", "sp", [-1, 1, 1, 4]}, {"tile_x", "sp", [-1, 1, 8, 7]}, {"tile_rc", "sp", [-1, 16]}, {"tile_ry", "sp", [-1, 1]}, {"tile_rx", "sp", [-1, 1]}, {"num_thread_n", "ot", 4}, {"num_thread_c", "ot", 4}, {"vthread_n", "ot", 1}, {"vthread_c", "ot", 2}, {"auto_unroll_max_step", "ot", 512}, {"unroll_explicit", "ot", 0}], "result": [[0.00017696012730627307], 0, 14.268717765808105, 1618735812.1710055], "version": 0.2, "tvm_version": "0.7.dev1"]], {"input": ["cuda -keys=cuda,gpu -max_num_threads=1024 -model=unknown -thread_warp_size=32", "conv2d_nchw.cuda", [{"TENSOR": [1, 64, 56, 56], "float32"}, {"TENSOR": [256, 64, 1, 1], "float32"}, [1, 1], [0, 0, 0, 0], [1, 1], "float32"], {}], "config": {"index": 7766329, "code_hash": null, "entity": [{"tile_f", "sp", [-1, 1, 1, 8]}, {"tile_y", "sp", [-1, 2, 28, 1]}, {"tile_x", "sp", [-1, 2, 28, 1]}, {"tile_rc", "sp", [-1, 1]}, {"tile_ry", "sp", [-1, 1]}, {"tile_rx", "sp", [-1, 1]}, {"num_thread_n", "ot", 4}, {"num_thread_c", "ot", 4}, {"vthread_n", "ot", 2}, {"vthread_c", "ot", 1}, {"auto_unroll_max_step", "ot", 0}, {"unroll_explicit", "ot", 0}], "result": [[0.00017696012730627307], 0, 14.268717765808105, 1618735812.1710055], "version": 0.2, "tvm_version": "0.7.dev1"]]
```



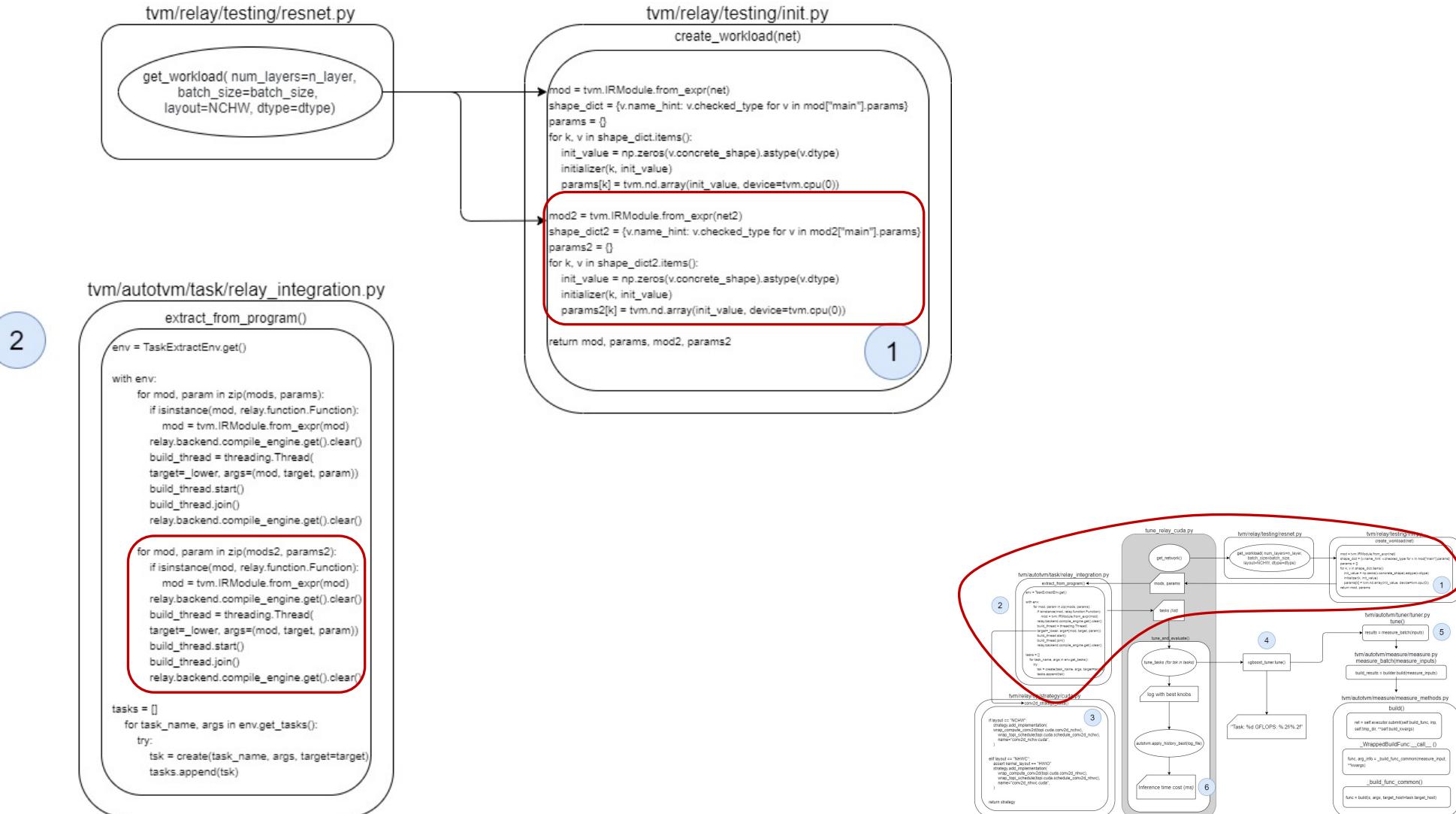
Attempts to Merge the Configuration Space



- ▶ Creating two IR modules from ResNet models of NCHW and NHWC layouts.
 - ▶ Formatting the log file.
- ▶ Including NHWC tasks and schedule when tuning a NCHW module.



Attempt 1: Creating two IR modules



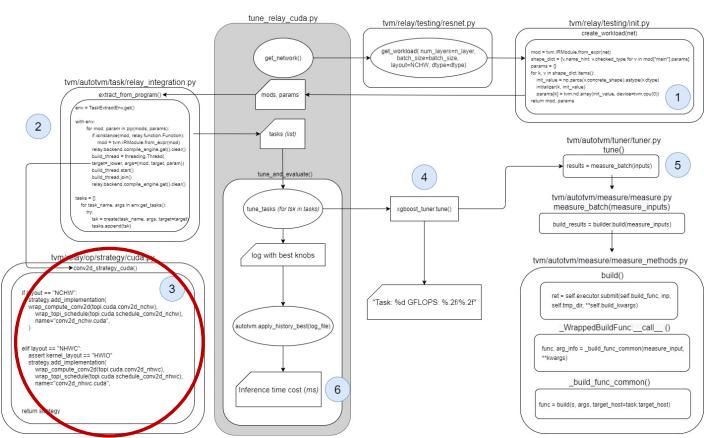
Attempt 1: Formatting the log file

```
1 import json
2
3 res = []
4 with open('resnet-50_NCHW.log', 'r') as f:
5     lines = f.readlines()
6     for line in lines:
7         dic = json.loads(line)
8         keylist = dic.get("input")
9         if keylist[1] in ["conv2d_nhwc.cuda", "conv2d_nhwc_winograd_direct.cuda"]:
10             if keylist[1] == "conv2d_nhwc.cuda":
11                 keylist[1] = "conv2d_nchw.cuda"
12             elif keylist[1] == "conv2d_nhwc_winograd_direct.cuda":
13                 keylist[1] = "conv2d_nchw_winograd.cuda"
14
15             a_list = keylist[2][0][1] #target tensor to rearrange
16             order = [0, 3, 1, 2] #nhwc to nchw order
17             a_list = [a_list[i] for i in order]
18             keylist[2][0][1] = a_list
19
20             keylist[2][1][1] = keylist[2][1][1][::-1] #reversed TENSOR
21
22             dic["input"] = keylist
23             res.append(dic)
24 f.close()
25
26 with open("resnet-50_NCHWnew.log", "a") as w:
27     for item in res:
28         w.write(json.dumps(item))
29         w.write('\n')
30
31 w.close()
```

Attempt 2: Binding two types of tasks

tvm/relay/op/strategy/cuda.py
conv2d_strategy_cuda()

```
if layout == "NCHW":  
    strategy.add_implementation(  
        wrap_compute_conv2d(topi.cuda.conv2d_nchw),  
        wrap_topi_schedule(topi.cuda.schedule_conv2d_nchw),  
        name="conv2d_nchw.cuda",  
    )  
  
    strategy.add_implementation(  
        wrap_compute_conv2d(topi.cuda.conv2d_nchw),  
        wrap_topi_schedule(topi.cuda.schedule_conv2d_nhwc),  
        name="conv2d_nhwc.cuda",  
    )  
  
elif layout == "NHWC":  
    assert kernel_layout == "HWIO"  
    strategy.add_implementation(  
        wrap_compute_conv2d(topi.cuda.conv2d_nhwc),  
        wrap_topi_schedule(topi.cuda.schedule_conv2d_nhwc),  
        name="conv2d_nhwc.cuda",  
    )  
  
return strategy
```





If there was a next attempt...



- ▶ Differentiate between “conv2d_nchw.cuda” and “conv2d_nhwc.cuda” tasks.
 - ▶ Devise a method to format tensors and other parameters according to the tasks and schedule of the other memory layout.

```
task name conv2d_nhwc.cuda
task name conv2d_nhwc.cuda
task name conv2d_nhwc_winograd_direct.cuda
task name conv2d_nhwc.cuda
task name conv2d_nhwc.cuda
task name conv2d_nhwc.cuda
task name conv2d_nhwc.cuda
task name conv2d_nhwc_winograd_direct.cuda
task name conv2d_nhwc.cuda
task name conv2d_nhwc_winograd_direct.cuda
task name conv2d_nhwc.cuda
task name conv2d_nhwc_winograd_direct.cuda
task name conv2d_nhwc.cuda
```

```
kernel shape: 2048 512 1 1
kernel shape: 512 512 3 3
kernel shape: 512 2048 1 1
kernel shape: 2048 512 1 1
kernel shape: 512 1024 1 1
kernel shape: 1024 256 1 1
kernel shape: 256 256 3 3
kernel shape: 256 1024 1 1
kernel shape: 1024 256 1 1
kernel shape: 256 512 1 1
kernel shape: 512 128 1 1
kernel shape: 128 128 3 3
kernel shape: 128 512 1 1
kernel shape: 512 128 1 1
kernel shape: 128 256 1 1
kernel shape: 256 64 1 1
kernel shape: 64 64 3 3
kernel shape: 64 256 1 1
kernel shape: 256 64 1 1
kernel shape: 64 64 1 1
kernel shape: 64 3 7 7
kernel shape: 256 64 1 1
kernel shape: 512 256 1 1
kernel shape: 1024 512 1 1
kernel shape: 2048 1024 1 1
```

```
kernel shape: 1 1 512 2048
kernel shape: 3 3 512 512
kernel shape: 1 1 2048 512
kernel shape: 1 1 512 2048
kernel shape: 1 1 1024 512
kernel shape: 1 1 256 1024
kernel shape: 3 3 256 256
kernel shape: 1 1 1024 256
kernel shape: 1 1 256 1024
kernel shape: 1 1 512 256
kernel shape: 1 1 128 512
kernel shape: 3 3 128 128
kernel shape: 1 1 512 128
kernel shape: 1 1 128 512
kernel shape: 1 1 256 128
kernel shape: 1 1 64 256
kernel shape: 3 3 64 64
kernel shape: 1 1 256 64
kernel shape: 1 1 64 256
kernel shape: 1 1 64 64
kernel shape: 7 7 3 64
kernel shape: 1 1 64 256
kernel shape: 1 1 256 512
kernel shape: 1 1 512 1024
kernel shape: 1 1 1024 2048
```

Conclusion

- ▶ Enlarged the configuration space
 - ▶ Expands the exploration area of the search algorithm
 - ▶ A more comprehensive log file can be formed
 - ▶ Low inference costs and a higher peak GFLOPS could be achieved.
- ▶ Merging the design spaces requires more work
 - ▶ Cross-compatibility of the tensor data

Auto-scheduler: Anso

- ▶ AutoTVM = template-based search (manual)
- ▶ Autoscheduler = automatic space construction
 - ▶ Less search time + higher performances.

	AutoTVM Workflow	Auto-scheduler Workflow
Step 1: Write a compute definition (relatively easy part)	# Matrix multiply C = te.compute((M, N), lambda x, y: te.sum(A[x, k] * B[k, y], axis=k))	# The same
Step 2: Write a schedule template (difficult part)	# 20-100 lines of tricky DSL code # Define search space cfg.define_split("tile_x", batch, num_outputs=4) cfg.define_split("tile_y", out_dim, num_outputs=4) ... # Apply config into the template bx, txz, tx, xi = cfg["tile_x"].apply(s, C, C.op.axis[0]) by, tyz, ty, yi = cfg["tile_y"].apply(s, C, C.op.axis[1]) s[C].reorder(by, bx, tyz, txz, ty, tx, yi, xi) s[CC].compute_at(s[C], tx) ...	# Not required
Step 3: Run auto-tuning (automatic search)	tuner.tune(...)	task.tune(...)

<https://tvm.apache.org/2021/03/03/intro-auto-scheduler>



Thank you!



References

1. Chen, T., Moreau, T., Jiang, Z., Zheng, L., & Yan, E. (2018). TVM: An Automated End-to-End Optimizing Compiler for Deep Learning. In *13th USENIX Symposium on Operating Systems Design and Implementation* (pp. 579-594). Carlsbad: USENIX Association. Retrieved from <https://www.usenix.org/system/files/osdi18-chen.pdf>
2. Common architectures in convolutional neural networks. (2020). Retrieved 26 November 2020, from <https://www.jeremyjordan.me/convnet-architectures/#resnet>
3. He, K., Zhang, X., Ren, S., & Sun, J. (2015). Deep Residual Learning for Image Recognition. *Microsoft*, 1-12. Retrieved from <https://arxiv.org/pdf/1512.03385.pdf>
4. NVIDIA TITAN Xp Graphics Card with Pascal Architecture. (2020). Retrieved 25 November 2020, from <https://www.nvidia.com/en-us/titan/titan-xp/>
5. ResNet-50 convolutional neural network - MATLAB resnet50. (2020). Retrieved 24 November 2020, from <https://www.mathworks.com/help/deeplearning/ref/resnet50.html>
6. Zheng, L., & Yan, E. (n.d.). Auto-tuning a convolutional network for NVIDIA GPU. Retrieved September 18, 2020, from https://tvm.apache.org/docs/tutorials/autotvm/tune_relay_cuda.html
7. <https://github.com/apache/tvm/tree/main/python/tvm>
8. <https://tvm.apache.org/2021/03/03/intro-auto-scheduler>