

Hello!

Tensor  
Vis based machine  
↳ LLVM

# CS 744: TVM

Shivaram Venkataraman

Fall 2020

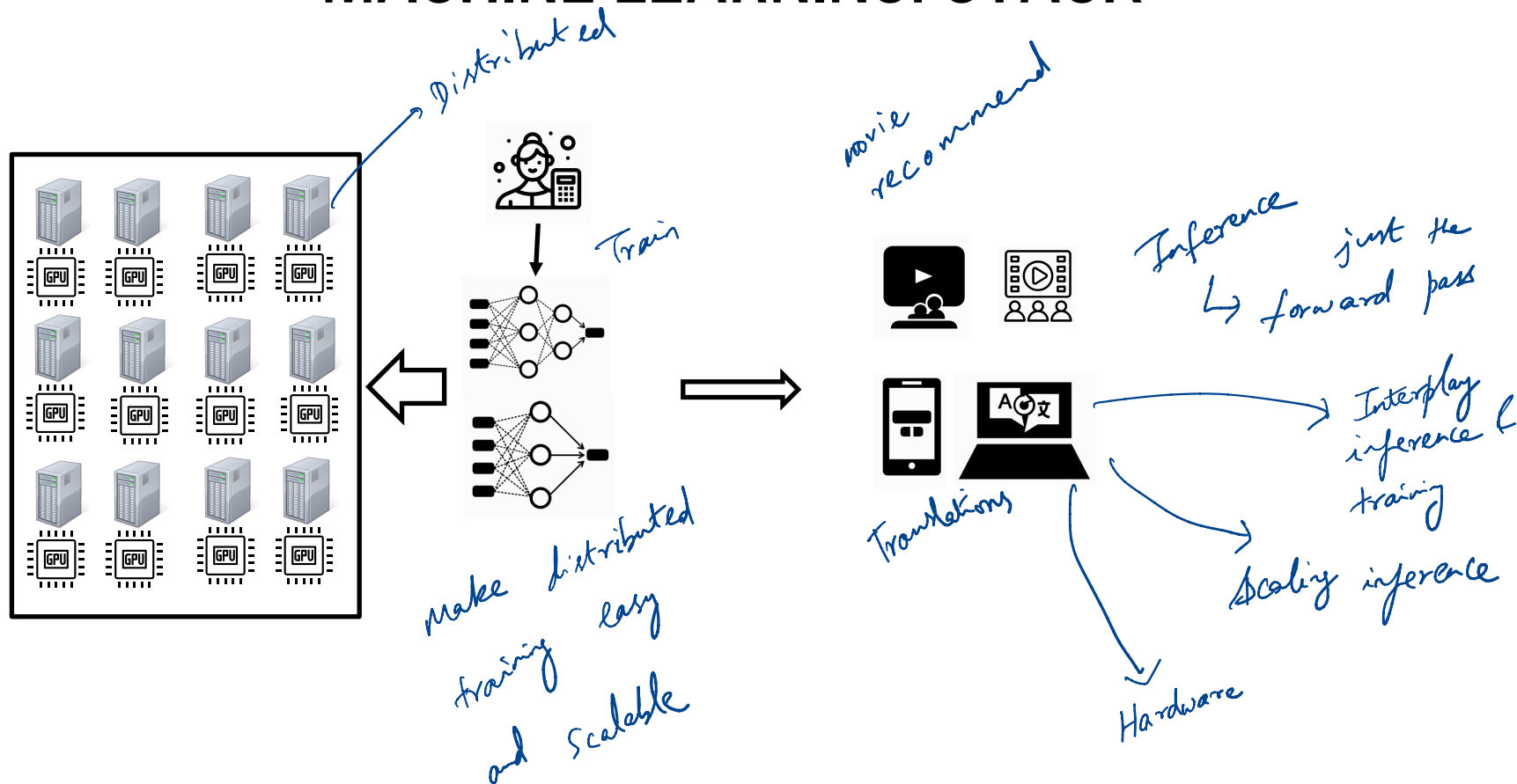
# ADMINISTRIVIA

~~Assignments are done~~

- Course project titles →
- Project proposal aka Introduction (10/16)
  - Introduction
  - Related Work
  - Timeline (with eval plan)
- Midterm: Oct 22

→ 2 page writeup

# MACHINE LEARNING: STACK



# MOTIVATION: PERFORMANCE PORTABILITY

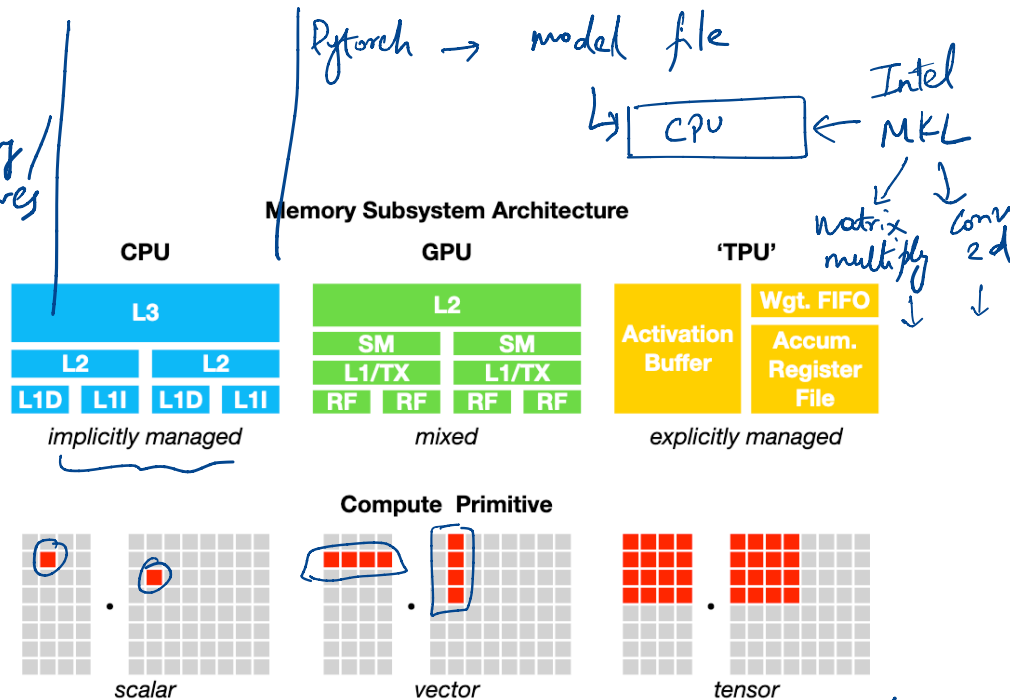
Different hardware

has different mem hierarchy/  
compute primitives

You want high performance  
across hardware backends

Dependence on vendor specific  
libraries

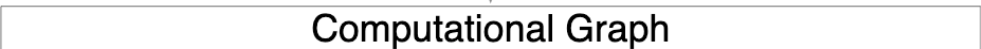
ML models evolve fast  $\Rightarrow$  new operators  
new combination of operators  $\Rightarrow$  not available  
in existing vendor  
libraries



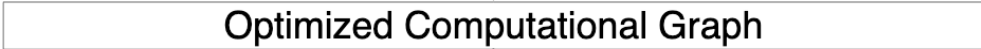
TVM



python code describes ML model

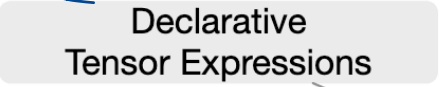


Section 3

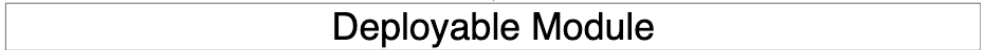
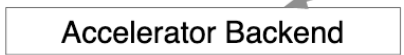
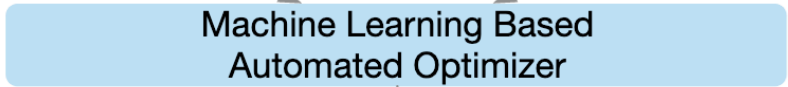


Operator-level Optimization and Code Generation

Section 4



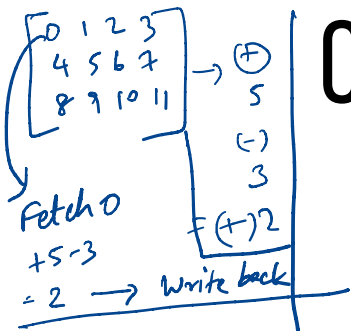
Section 5



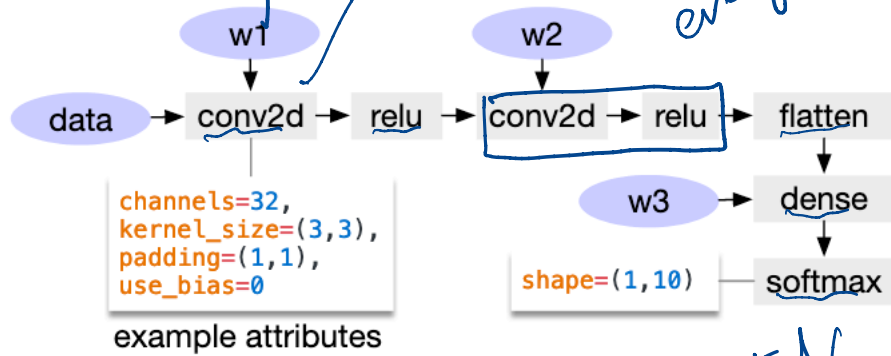
TVM

Binary file that runs on hardware

# OPTIMIZATION COMPUTATION GRAPHS



know the dimensions of the input & output at every vertex



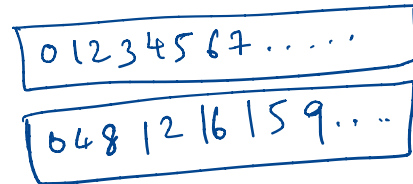
## (1) Operator Fusion

↳ Combine any graph vertexes into one low level op

→ 1-1 operators, "map" → Combine  
→ Sum reduction, scaling after (spark map)

## Data layout

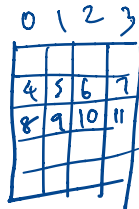
Row Major, Column Major, Blocked



RM  
CM

Invert or change data layout

2-layer net as NN is represented a graph



# TENSOR EXPRESSION LANGUAGE

operators ()

expressed in tensor expression language

```
m, n, h = t.var('m'), t.var('n'), t.var('h')
```

```
A = t.placeholder((m, h), name='A')
```

```
B = t.placeholder((n, h), name='B')
```

```
k = t.reduce_axis((0, h), name='k')
```

```
C = t.compute((m, n), lambda y, x:
```

```
    t.sum(A[k, y] * B[k, x], axis=k))
```

result shape

computing rule

Tensor

math operations

Common Arithmetic, Math operations

Know the shape of the output and the data accessed

Halide  $\rightarrow$  expression

schedule of intrinsics

# CODE GENERATION

OpenMP  $\leftarrow$  gcc

```
for thread_group (by, bx) in cross(64, 64):
  for thread_item (ty, tx) in cross(2, 2):
    local CL[8][8] = 0
    shared AS[2][8], BS[2][8]
    for k in range(1024):
      for i in range(4):
        AS[ty][i*4+tx] = A[k][by*64+ty*8+i*4+tx]
      for each i in 0..4:
        BS[ty][i*4+tx] = B[k][bx*64+ty*8+i*4+tx]
      memory_barrier_among_threads()
```

Co-operative

threads can use AS, BS to do computation

```
def gemm_intrin_lower(inputs, outputs):
  ww_ptr = inputs[0].access_ptr("r")
  xx_ptr = inputs[1].access_ptr("r")
  zz_ptr = outputs[0].access_ptr("w")
  compute = t.hardware_intrin("gemm8x8", ww_ptr, xx_ptr, zz_ptr)
  reset = t.hardware_intrin("fill_zero", zz_ptr)
  update = t.hardware_intrin("fuse_gemm8x8_add", ww_ptr, xx_ptr, zz_ptr)
  return compute, reset, update
```

gemm8x8 = t.decl\_tensor\_intrin(y.op, gemm\_intrin\_lower)

launch 50 threads  
each thread does  
load  $a[i,j]$ ,  $b[i,j]$   
compute  
write it back  
Nested parallelism  $\rightarrow$  Cooperative parallelism

for i in 1:10  
for j in 1:5  
 $a[i,j] = b[i,j] + 2$   
 $\rightarrow$  independent of other loop iterations

load, store, add Tensorization

$\rightarrow$  what is the hardware instruction set

Allows you to register operator intrinsics

Extensible!



# LATENCY HIDING

same as as  
Pytorch etc.

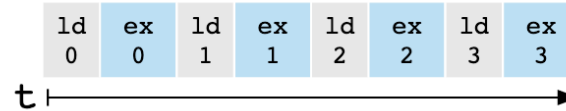
## What is the goal?

→ Overlap computation and communication

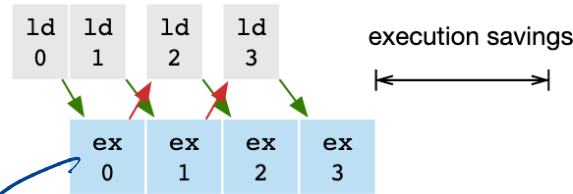
Schedule that utilizes memory bandwidth & compute units

this execution happens after this load

### Monolithic Pipeline



### Decoupled Access-Execute Pipeline



→ read after write (RAW) dependence  
→ write after read (WAR) dependence

### Instruction Stream

```
ld.perform_action(ld0)
ex.perform_action(ex0)
ld.perform_action(ld1)
ex.perform_action(ex1)
...
```

```
ld.perform_action(ld0)
ld.push_dep_to(ex)
ld.perform_action(ld1)
ld.push_dep_to(ex)
ex.pop_dep_from(ld)
ex.perform_action(ex0)
ex.push_dep_to(ld)
ex.pop_dep_from(ld)
ex.perform_action(ex1)
ex.push_dep_to(ld)
ld.pop_dep_from(ex)
ld.perform_action(ld2)
...
```

# AUTOMATING OPTIMIZATION

Goal: Create a specialized operator for input shape and layout

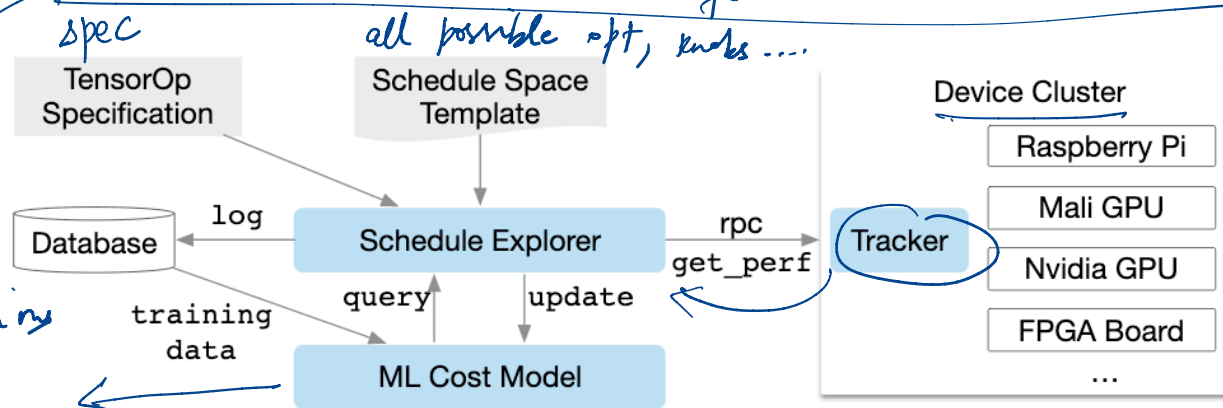
Challenge:

Choose appropriate schedule optimizations

Tiling size, loop unrolling .....

*lots of different choices & also lots of parameters / knobs to choose.*

Automate the optimizer!



*ML on ML?*

*what configurations to try*

# ML-BASED COST MODEL

## Machine Learning Model Design Choices

→ Speed: Faster than time it takes to evaluate a config → *ms to ~ seconds*

Quality: Use a rank objective to predict the relative order of runtime

Gradient tree boosting model ←

*take  
code generate  
features*

memory access count

reuse ratio of each memory buffer at each loop level

one-hot encoding of loop annotations

# ML-BASED COST MODEL

model perf  
when using  
config

$\langle \underline{C_1}, 10ms \rangle$   
 $\langle C_2, 20ms \rangle$

## Iteration

- Run on real hardware
- (1) Select a batch of candidates  $\rightarrow$  each candidate is a configuration, schedule
  - (2) Collect data  $\leftarrow$   $\langle \underline{C_3}, 8ms \rangle$   
 $\langle C_4, \text{fail} \rangle$
  - (3) Use as training data to update the model

## How to select candidates?

$\rightarrow$  Step (1) above

training data  
to GBT

## Parallel Simulated Annealing

Start from a random config

Walk to a nearby config  $\rightarrow$

Successful if cost decreases Else Reject

$C_3 \rightarrow C_3'$  ask model is  $C_3'$  better than  $C_3$   
Yes  $\rightarrow$  go  $\leftarrow$  try  $C_3'$  on cluster  
No  $\rightarrow$  generate another config

# DISTRIBUTED DEVICE POOL

Pool of devices to speed up profiling

RPC interface to run a trial on device

Share device pools for multiple graphs

# SUMMARY

TVM: Compiler for ML inference models

Support high performance for range of models, hardware devices

Key ideas

→ Graph-level optimizations → *operator fusion*

→ Tensor expression language: Code-gen, Latency hiding etc

ML-based Cost Model for automation

# DISCUSSION

<https://forms.gle/WiVgJ3abGXXgfBN99>

Consider that you are building an optimizer for Spark programs instead of ML inference. What would be some configuration knobs that you could similarly tune? What might be different from the TVM optimizer?

Similar logic  $\rightarrow$  latency hiding

overlap comp, communication

In TVM  
 $\hookrightarrow$  dimensions,  
access patterns

operator fusion  $\rightarrow$  map operations

$\hookrightarrow$  lazy!

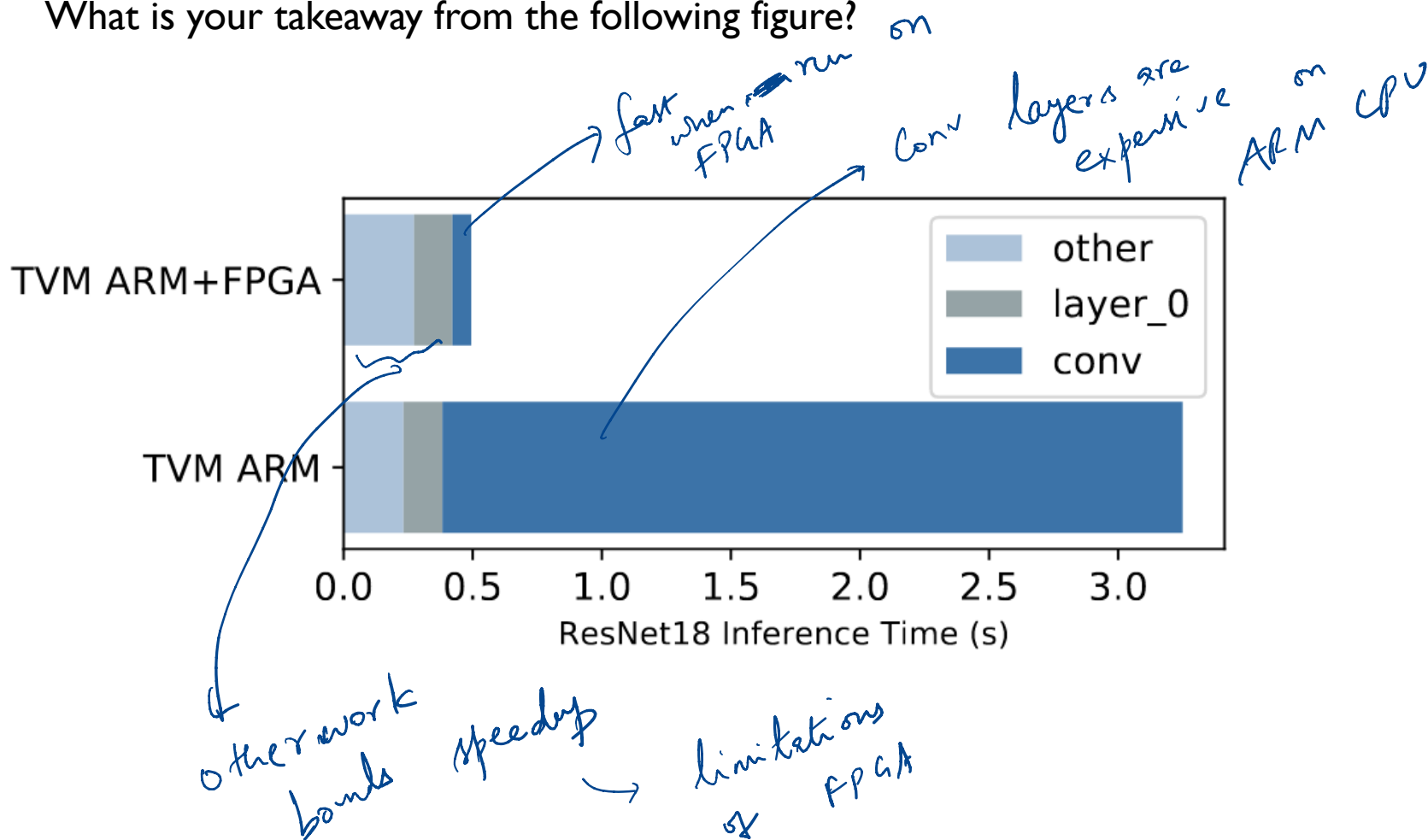
$\hookrightarrow$  operators are user defined  
challenging !?

Partitioning  $\rightarrow$  Can you automate  
 $\hookrightarrow$  number of partitions / co-partitioning  
 $\hookrightarrow$  Performance!

Persistence  $\rightarrow$  manually insert `rdl.cache()`  $\rightarrow$  config space!



What is your takeaway from the following figure?



# NEXT STEPS

Next class: Ray

Course project: Oct 16 (introductions)

Midterm: Oct 22

latency hiding in spark?

$rdd1 \leftarrow \text{map tasks}$  ✓  
 $\leftarrow \text{reduce tasks}$

↓  $rdd2 \leftarrow \text{map tasks} \rightarrow$

---

$rdd1 \text{ map}$   
 $rdd2 \text{ map} \leftrightarrow \text{transfer shuffle files}$   
 $\vdots$   
 $\leftarrow \text{reduce} \rightarrow \leftarrow \text{no more wait}$