



Graphics and ML SIG Meeting

Jan 6, 2022

10:05am PDT

<https://github.com/riscv-admin/graphics>



@risc_v

Only RISC-V Members May Attend

- Non-members are asked to please leave except for Joint Working Groups (JWG).
- Members share IP protection by virtue of their common membership agreement. Non-members being present jeopardizes that protection. [Joint working groups](#) (JWG) agree that any IP discussed or worked on is fully open source and unencumbered as per the policy.
- It is easy to become a member. Check out riscv.org/membership
- If you need work done between non-members or or other orgs and RISC-V, please use a joint working group (JWG).
 - used to allow non-members in SIGs but the SIGs purpose has changed.
- Please put your name and company (in parens after your name) as your zoom name. If you are an individual member just use the word “individual” instead of company name.
- Non-member guests may present to the group but should only stay for the presentation. Guests should leave for any follow on discussions.

Antitrust Policy Notice

RISC-V International meetings involve participation by industry competitors, and it is the intention of RISC-V International to conduct all its activities in accordance with applicable antitrust and competition laws. It is therefore extremely important that attendees adhere to meeting agendas, and be aware of, and not participate in, any activities that are prohibited under applicable US state, federal or foreign antitrust and competition laws.

Examples of types of actions that are prohibited at RISC-V International meetings and in connection with RISC-V International activities are described in the RISC-V International Regulations Article 7 available here: <https://riscv.org/regulations/>

If you have questions about these matters, please contact your company counsel.

Collaborative & Welcoming Community

RISC-V is a free and open ISA enabling a new era of processor innovation through open standard collaboration. Born in academia and research, RISC-V ISA delivers a new level of free, extensible software and hardware freedom on architecture, paving the way for the next 50 years of computing design and innovation.

We are a transparent, collaborative community where all are welcomed, and all members are encouraged to participate. We are a continuous improvement organization. If you see something that can be improved, please tell us. help@riscv.org

We as members, contributors, and leaders pledge to make participation in our community a harassment-free experience for everyone.

<https://riscv.org/community/community-code-of-conduct/>

Conventions



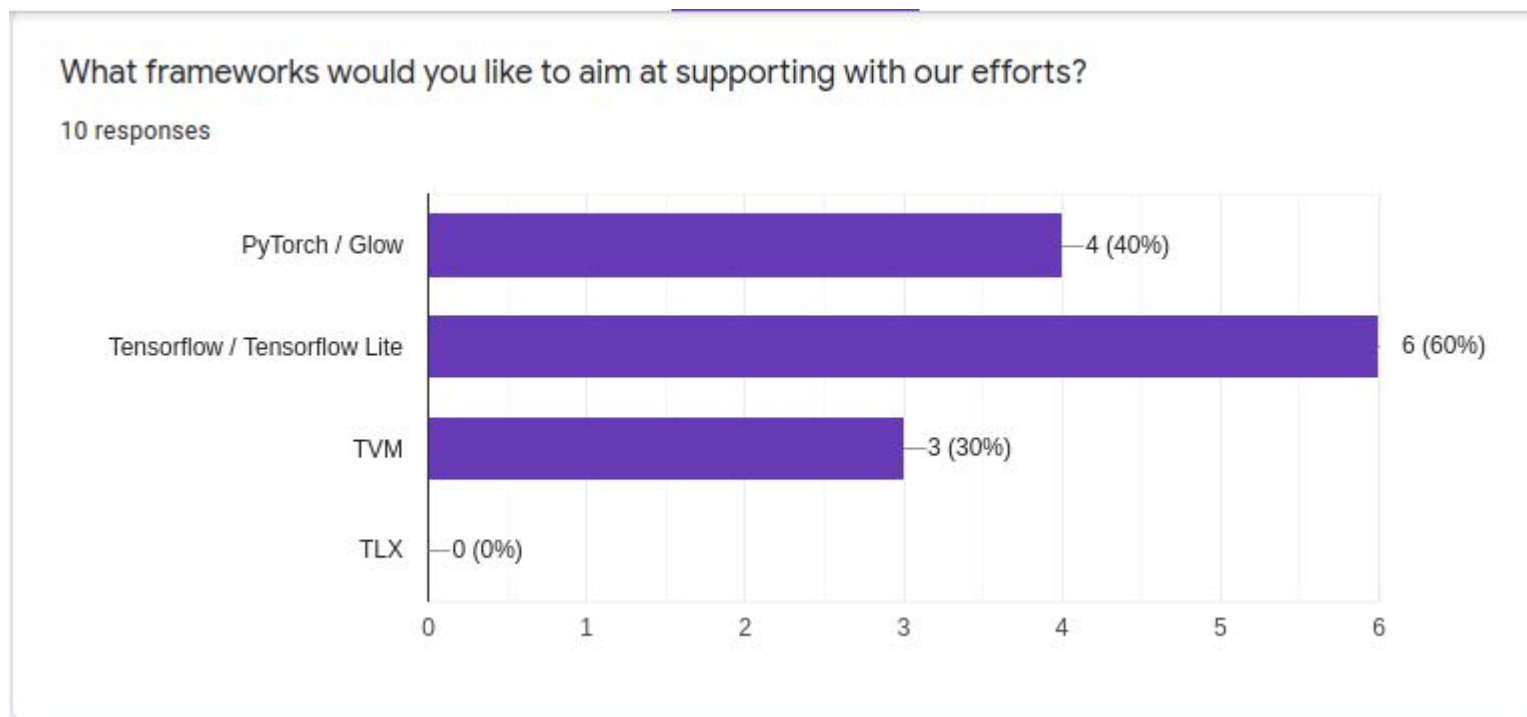
- **For one hour meetings, please start at 5 after the start time** in order to allow people going to other meetings have time for a short break between meetings. 30 minute meetings start on time.
- Unless it is a scheduled agenda topic, we don't solve problems or detailed topics in most meetings unless specified in the agenda because we don't often have enough time to do so and it is more efficient to do so offline and/or in email. We identify items and send folks off to do the work and come back with solutions or proposals.
- If some policy, org, extension, etc. can be doing things in a better way, help us make it better. Do not change or not abide by the item unilaterally. Instead let's work together to make it better.
- Please conduct meetings that accommodates the virtual and broad geographical nature of our teams. This includes meeting times, repeating questions before you answer, at appropriate times polling attendees, guide people to interact in a way that has attendees taking turns speaking, ...
- Where appropriate and possible, meeting minutes will be added as speaker notes within the slides for the Agenda

Agenda

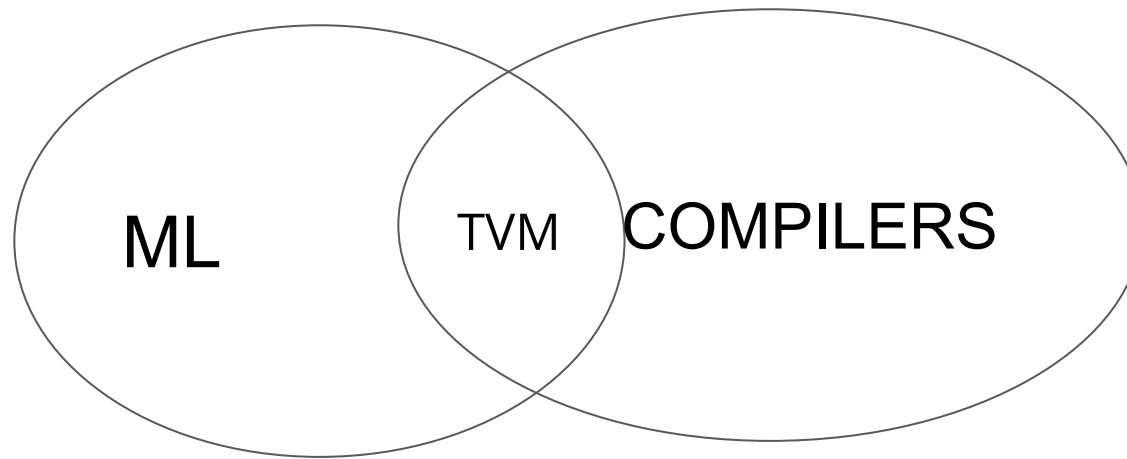


- ML frameworks: the case for TVM support (20 min)
- Status of the RISC-V backend on TVM (5 min)
- Gap analysis for ML: a practical approach (10 min)

ML Frameworks: the case for TVM support (1/6)



ML Frameworks: the case for TVM support (2/6)



ML Frameworks: the case for TVM support (3/6)

Example: kernel that computes element-wise $A*B+C$ for two 2x2 matrices

In C one could write that kernel in at least two ways:

```
for (int i = 0; i < 2; ++i)
  for (int j = 0; j < 2; ++j)
    MUL[i][j] = A[i][j] * B[i][j];
```

```
for (int i = 0; i < 2; ++i)
  for (int j = 0; j < 2; ++j)
    ADD[i][j] = MUL[i][j] + C[i][j];
```

Or the same thing fused:

```
for (int i = 0; i < 2; ++i)
  for (int j = 0; j < 2; ++j)
    MAC[i][j] = A[i][j] * B[i][j] + C[i][j];
```

ML Frameworks: the case for TVM support (3/6)

Example: kernel that computes element-wise $A*B+C$ for two 2x2 matrices

In C one could write that kernel in at least two ways:

```
for (int i = 0; i < 2; ++i)
  for (int j = 0; j < 2; ++j)
    MUL[i][j] = A[i][j] * B[i][j];
```

```
for (int i = 0; i < 2; ++i)
  for (int j = 0; j < 2; ++j)
    ADD[i][j] = MUL[i][j] + C[i][j];
```

Or the same thing fused:

```
for (int i = 0; i < 2; ++i)
  for (int j = 0; j < 2; ++j)
    MAC[i][j] = A[i][j] * B[i][j] + C[i][j];
```

Your “stash of assumptions”:

- One CPU
- Row major strided order traversal is better
- Fused is feasible and better

ML Frameworks: the case for TVM support (4/6)

```
A = te.placeholder((m,n), name="A")
B = te.placeholder((m,n), name="B")
C = te.placeholder((m,n), name="C")
```

```
MUL = te.compute(
    shape=[m,n],
    fcompute=lambda i,j: A[i,j]*B[i,j],
    name="MUL")
```

```
|
ADD = te.compute(
    shape=[m,n],
    fcompute=lambda i,j: MUL[i,j]+C[i,j],
    name="ADD")
```

Computation

```
schedule = te.create_schedule(ADD.op)
```

Schedule

```
mac = tvn.build(schedule, [A, B, C, ADD], target, name="mac")
```

Code generation

ML Frameworks: the case for TVM support (5/6)

- Default scheduler: **NOT FUSED**

```
schedule = te.create_schedule(ADD.op)
```

- Custom scheduler: **FUSED**

```
schedule = te.create_schedule(ADD.op)  
schedule[MUL].compute_at(schedule[ADD], ADD.op.axis[0])
```

ML Frameworks: the case for TVM support (6/6)

Advantages:

- Kernel computation made orthogonal, not only from backend, also from scheduling (like in Halide!)
- Optimal scheduling:
 - Extremely dependent on underlying hardware. To be provided by an expert user
 - Tools for auto-tuning available.
 - Intermediate representation suitable for polyhedral compilers that will help (*)
- More backends available than for any other available framework

(*) In the future, subject to the availability of a cost function

Status of the RISC-V backend on TVM

Volunteers needed for evaluating the tool!

Heads up: no Spike support, but some RISC-V chips from SiFive chips are supported.

Gap analysis for ML: a practical approach

Cherry-picking models vs. running industry standard benchmarks

My view:

The industry standard benchmark is MLPerf (except for RecSys), but supporting every single model in MLPerf is difficult.

Cherry-picking something like ResNet50 that is on MLPerf is a good start.

Backup Slides