



Tencent AI Lab

Efficient Phase-Functioned Real-time Character Control in Mobile Games: A TVM Enabled Approach

**Haidong Lan, Wenxi Zhu, Qian Qiu, Dou Wu, Honglin Zhu,
Jingjing Zhao, Xinghui Fu, Minwen Deng, Jintao Meng**

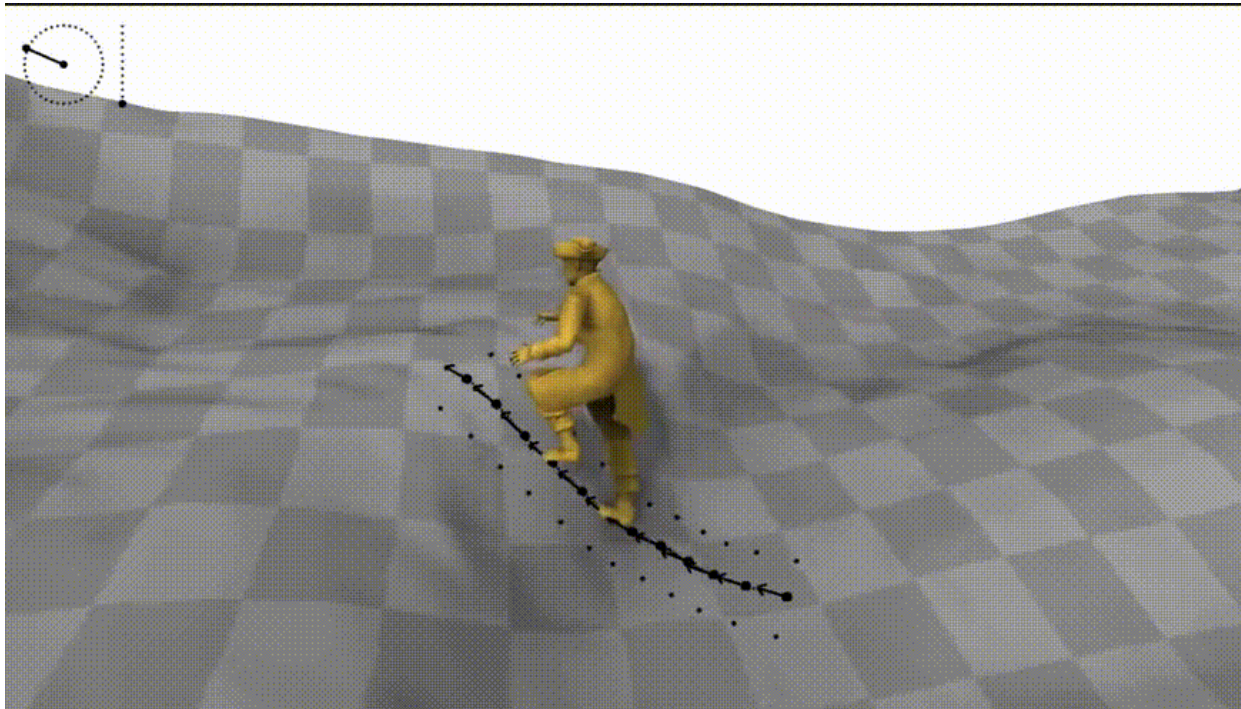
Tencent AI Lab

2022-8-12

Agenda

- **Background**
- **Methods**
- **Performance analysis**
- **Lessons learnt**

Background – Game character control

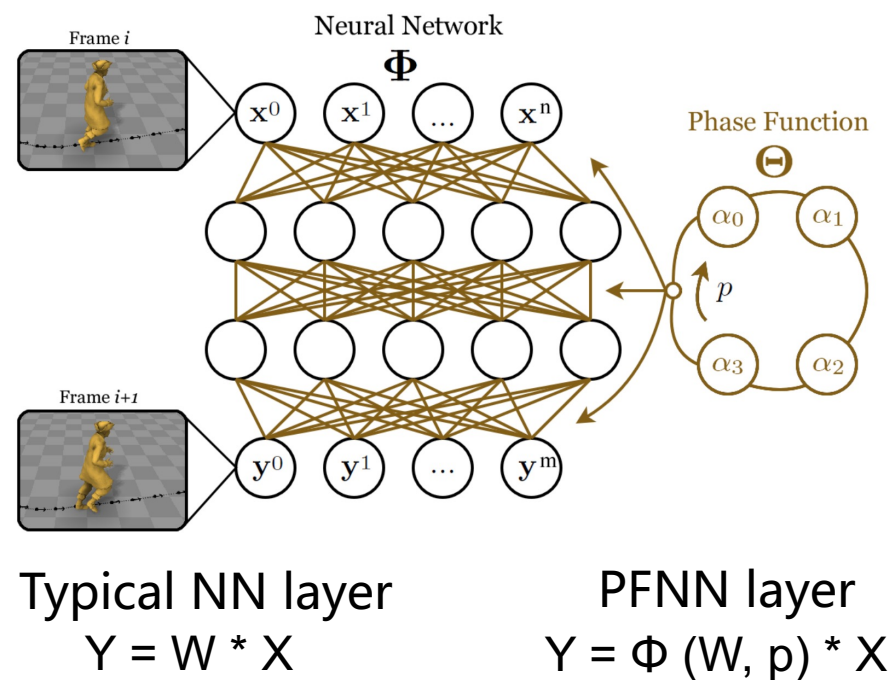


Real-time game character control algorithm

- Predict the anchor point of the game character in the next frame
- Embedded in mobile game system for real-time computation
- Required latency: within 2ms for 8 characters
- Deploy on Android/iOS, verification on macOS/Windows

Video taken from <https://dl.acm.org/doi/10.1145/3072959.3073663>, supplement materials.

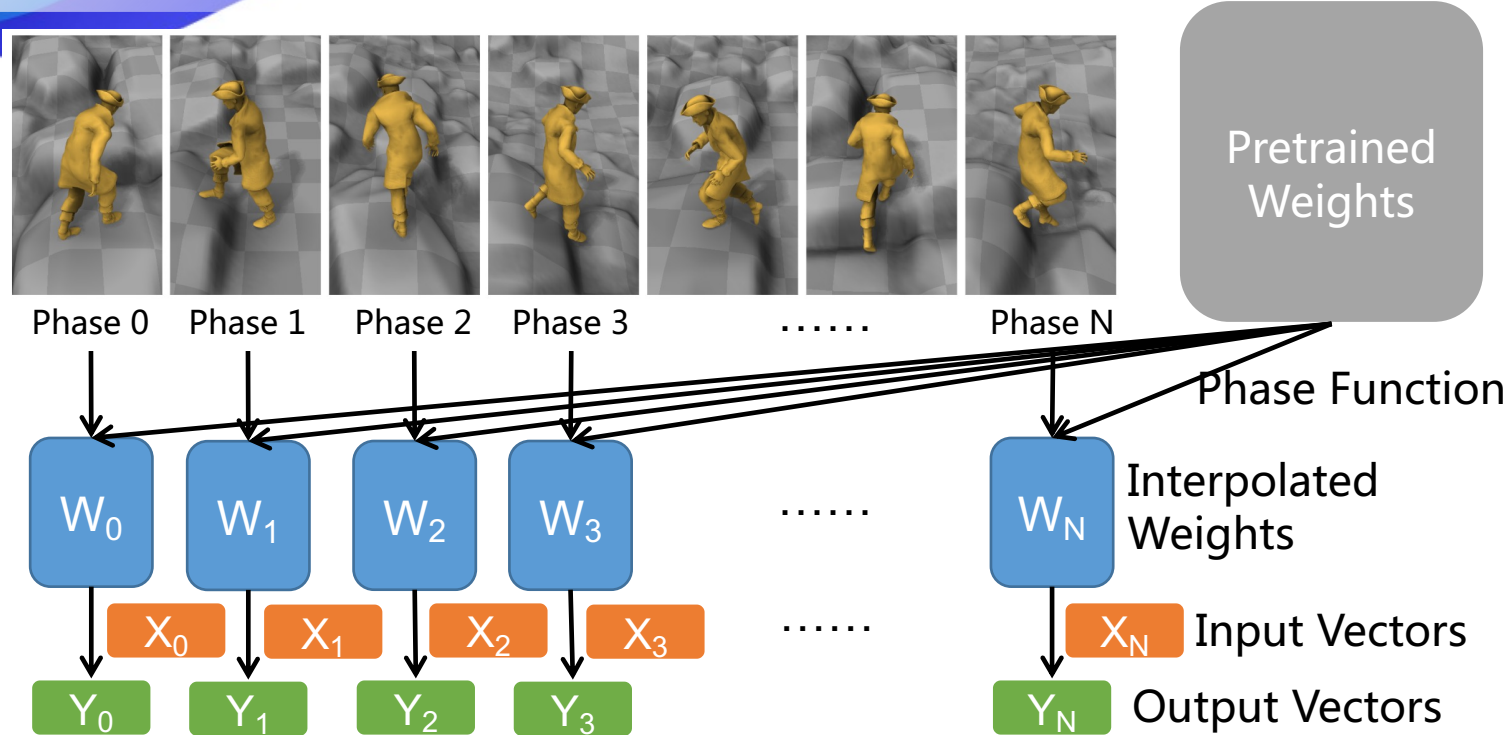
Background – PFNN



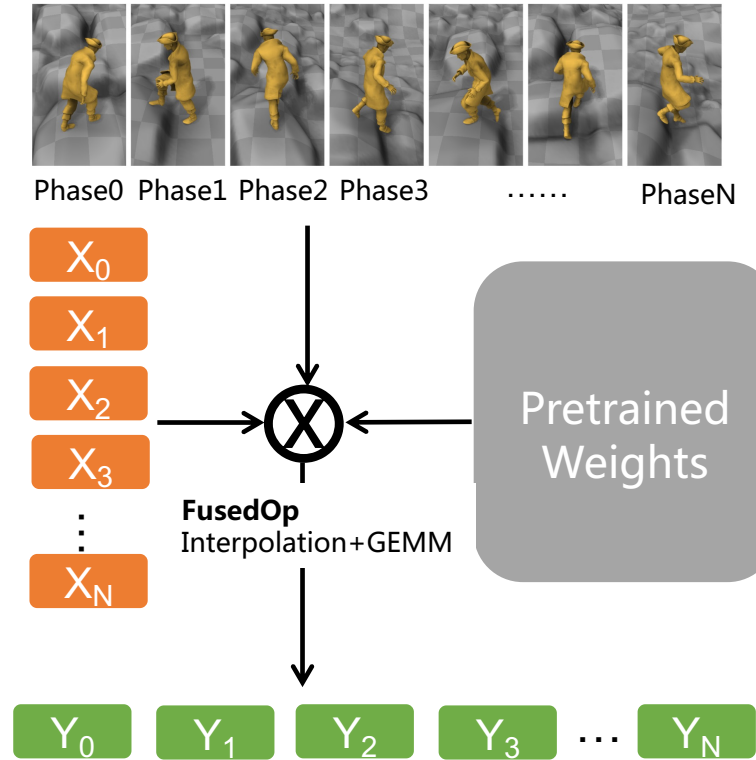
Algorithm

- PFNN (Phased-Functioned Neural Network)
- Phase value corresponds to game character status
- Architecture: fully connected layers where weight matrices are interpolated by a phase function
- The interpolation overhead is noticeable

Background – PFNN in real game scenario



(a) Existing libraries can only compute sequentially for N characters due to distinct phase values



(b) Fused operators can efficiently compute for all characters in a batch

Technical Difficulties

- Different characters cannot share identical phase values and weight matrices.
- Strict requirement on inference latency: within 2ms for 8 characters on mobile.
- Agile development: dev in 2 weeks for desktop and mobile.

Reference: Holden D, Komura T, Saito J. Phase-functioned neural networks for character control[J]. ACM Transactions on Graphics (TOG), 2017, 36(4): 1-13.

Optimization Opportunities

- Interpolation (phase) function is element-wise, fusible with matrix multiplications

Methods – Fused Operator

C: #Characters

Algorithm 1 Non-fused Pseudocode for Equation 5

```

1: for c = 1 to C do
2:   for k = 1 to K do
3:     for n = 1 to N do
4:       S = 0
5:       for i = 1 to 4 do
6:         S = S + Wc,i · αi(k, n)
7:       end for
8:       Ic(k, n) = S
9:     end for
10:  end for
11: end for
12: for c = 1 to C do
13:   Below is equivalent to a GEMV call.
14:   for n = 1 to N do
15:     S = 0
16:     for k = 1 to K do
17:       S = S + Ic(k, n) · Xc(k)
18:     end for
19:     Φ(c, n) = S
20:   end for
21: end for

```

Phase Function

Vector-Matrix Multiplication

Algorithm 2 Fused Pseudocode for Equation 5

```

1: for n = 1 to N do
2:   for c = 1 to C do
3:     S = 0
4:     for k = 1 to K do
5:       T = 0
6:       for i = 1 to 4 do
7:         T = T + Wc,i · αi(k, n)
8:       end for
9:       S = S + T · Xc(k)
10:    end for
11:    Φ(c, n) = S
12:  end for
13: end for

```

PF-GEMM

Observation 1 : The fused operator reduces I_c to local variables.

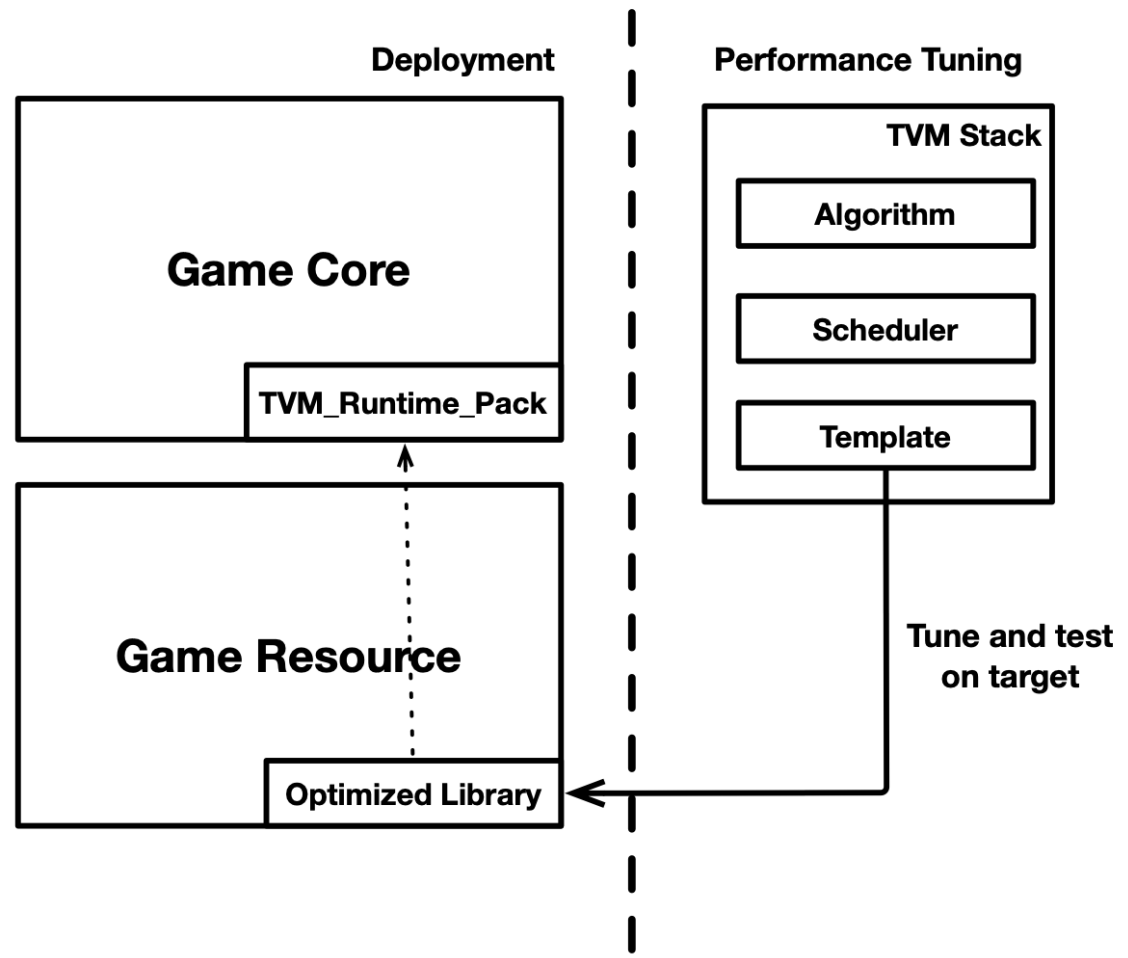
Deduction 1 : No need to separately store interpolated matrices.

Observation 2 : The loop order of fused operator is identical with a matrix multiplication of shape (C, K) x (K, N)

Deduction 2 : The optimization techniques for fused operator can be borrowed from GEMM optimizations.

Methods - Integration

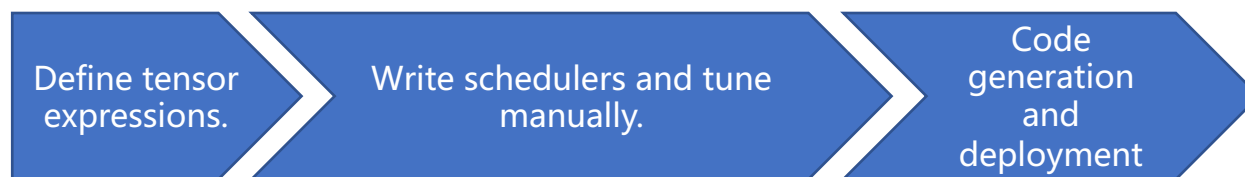
TVM Workflow



Methods – TVM approaches

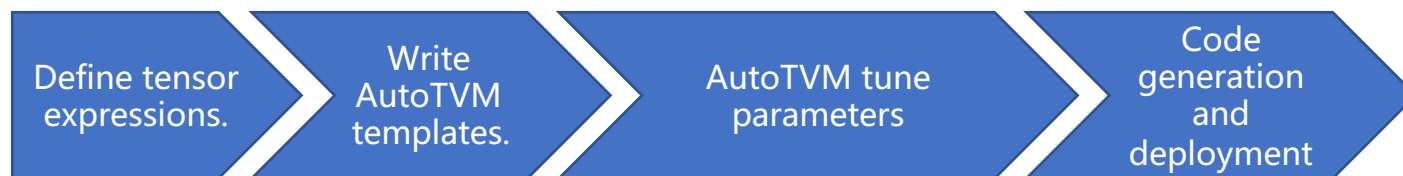
• TVM with manual scheduler

- Manually write algorithms and schedulers.
- TVM acts as cross-platform code generator.
- Requires HPC knowledge but we are cool.



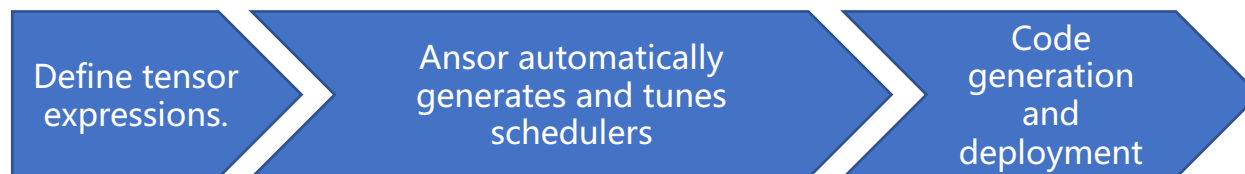
• AutoTVM with external kernels

- Manually written inner kernels (for every hardware) and scheduler templates.
- TVM tunes cache block parameters and generates code for outer loops.
- Best performance with long search time.



• Ansor: end-to-end TVM approach

- We only defined the algorithm.
- Quickly converged to descent results.



Performance analysis

Eigen (baseline BLAS library)

- Cannot fuse operators.
- High latency when multiple game characters are present.

TVM with Hand Schedulers

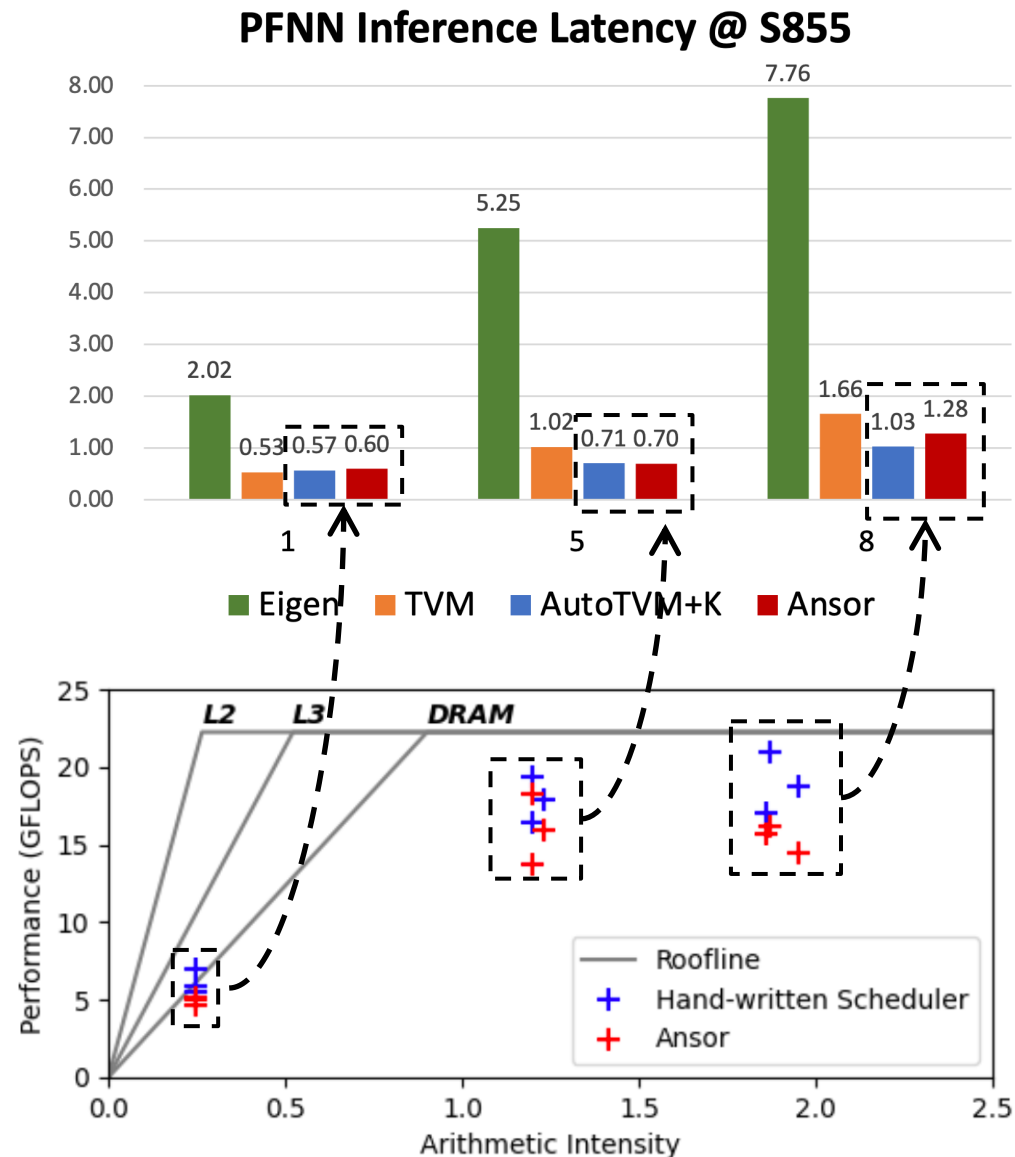
- Our TVM baseline.
- Developed in 2 hours for x86 and Arm processors.
- Parameters are not carefully tuned.
- Met our business goals.

AutoTVM with external inner kernels

- Extremely high performance (30% faster than hand schedulers).
- AutoTVM searches for best cache block size, loop order and register block size.
- We spent three weeks on development and tuning.

Ansor

- Close to AutoTVM approach but converged much faster (~2 hrs vs days).
- Consumes least human efforts.



Performance assessment conducted on a single big core on Qualcomm Snapdragon 855

Lessons learnt – What went right?

Performance

- Operator fusion makes a big difference.
- Great performance even with roughly tuned schedulers.
- AutoTVM can effectively tune hand-written kernels.
- Ansor exhibits amazing performance in an end-to-end workflow.

Usability

- Packed runtime is easy to integrate and easy to customize.
- Cross-platform versions are easy to maintain.

Lessons learnt – What went wrong?

Performance

- LLVM has some corner cases in vectorization (will file an issue for this).
- Better cooperation of hand-written inner kernels and auto scheduling is important to HPC guys.
(TensorIR is magic! We want more documents!)
- AutoTVM cannot converge stably, we made lots of trials to achieve best performance.

Usability

- Static library needs some hack to work. We exported static library in order to accommodate to all IDE versions.
- **Why TVM cannot be installed through pip or anaconda?**
- Shapes are always assumed to be static such that we cannot adjust character number dynamically.



Tencent AI Lab

Thank You – Q&A