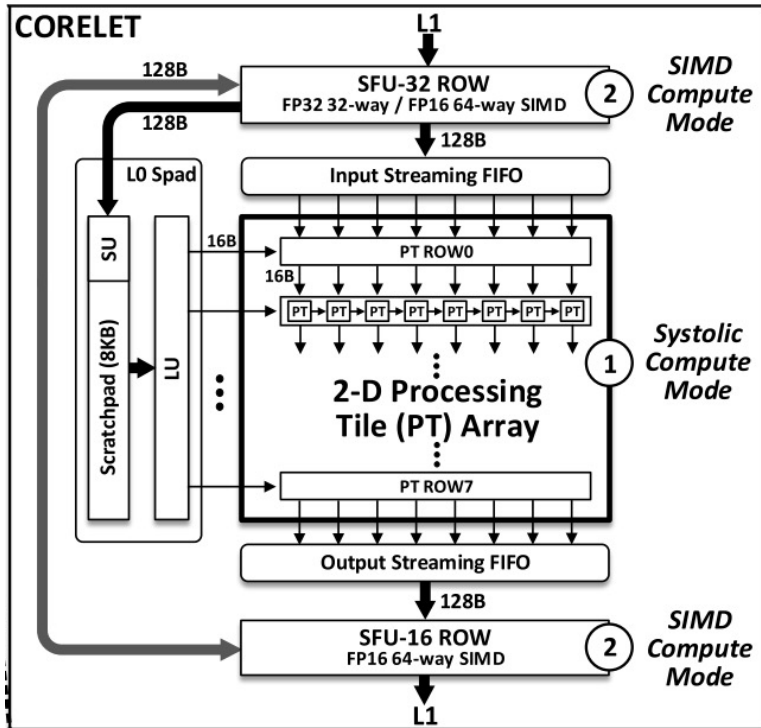# ECE260B Winter 22

## Project Description

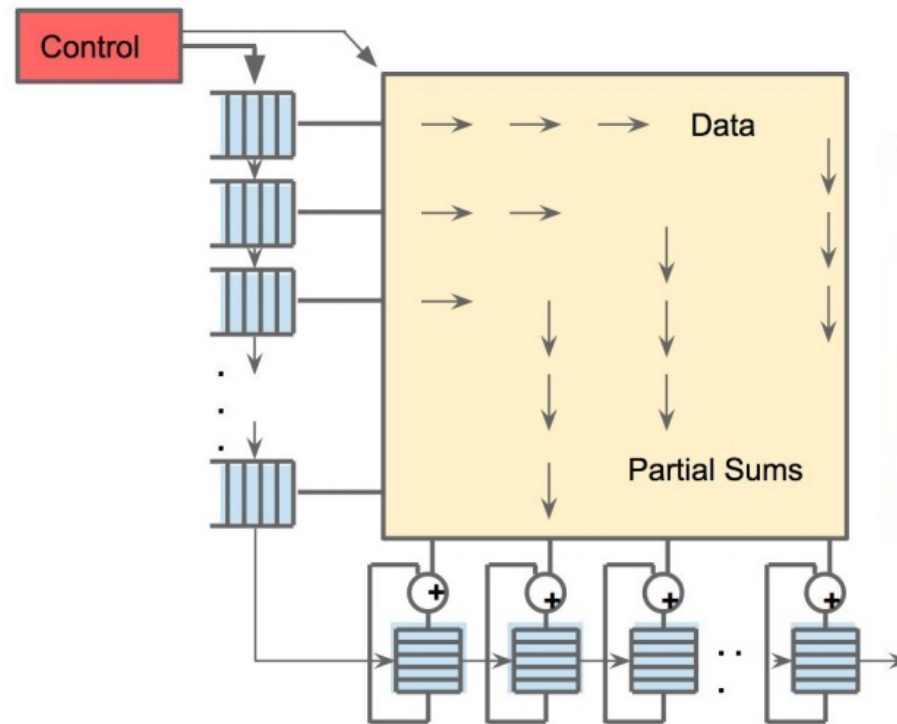**Prof. Mingu Kang**

**UCSD Computer Engineering**

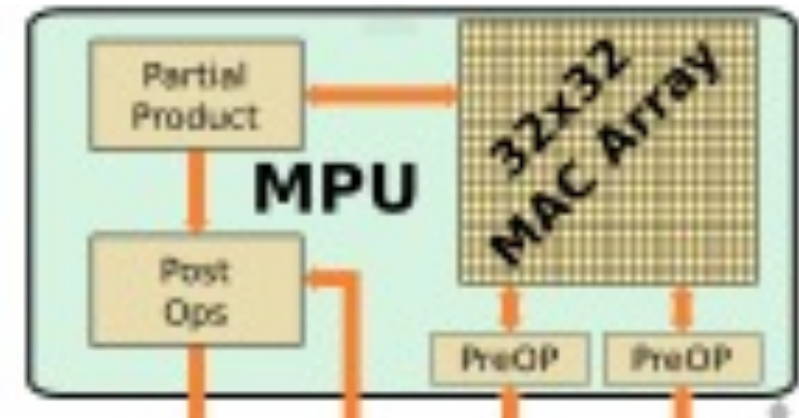# 2D Systolic-based Architecture



**IBM Rapid**

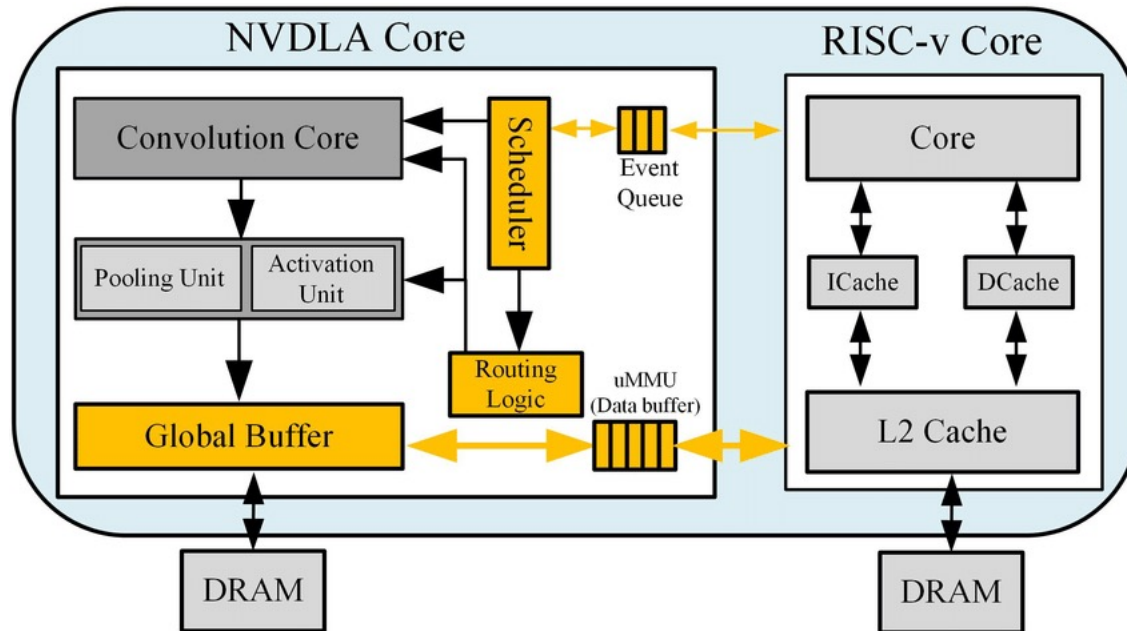**Google TPU**

**Intel Nervana**

J.Oh, "A 3.0 TFLOPS 0.62V Scalable Processor Core for High Compute Utilization AI Training and Inference"
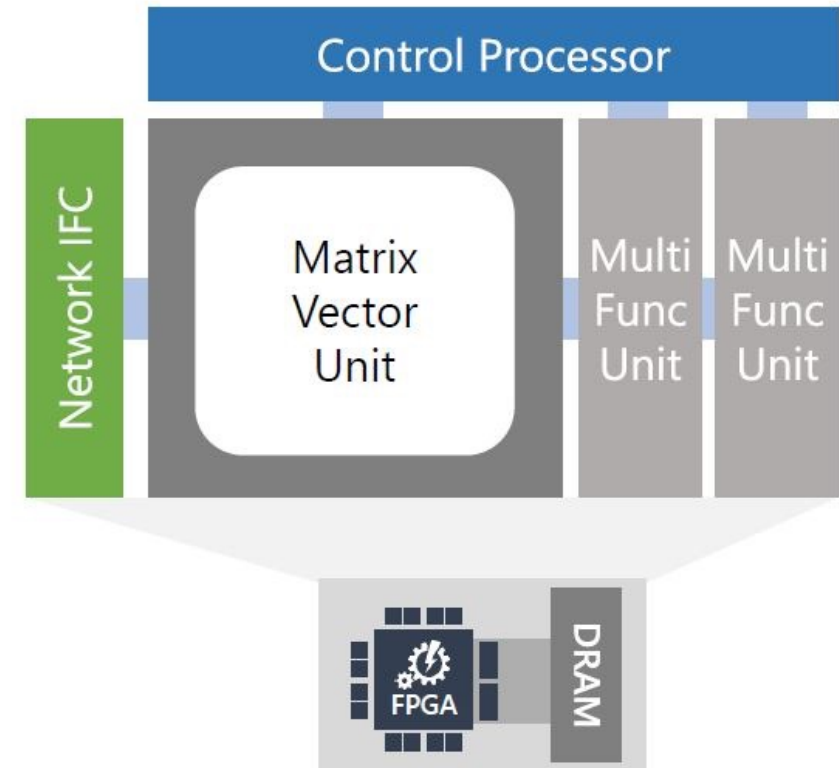https://arxiv.org/pdf/1704.04760.pdf
https://fuse.wikichip.org/news/3270/intel-axes-nervana-just-two-months-after-launch/

# 1-D Vector Processor based Architecture

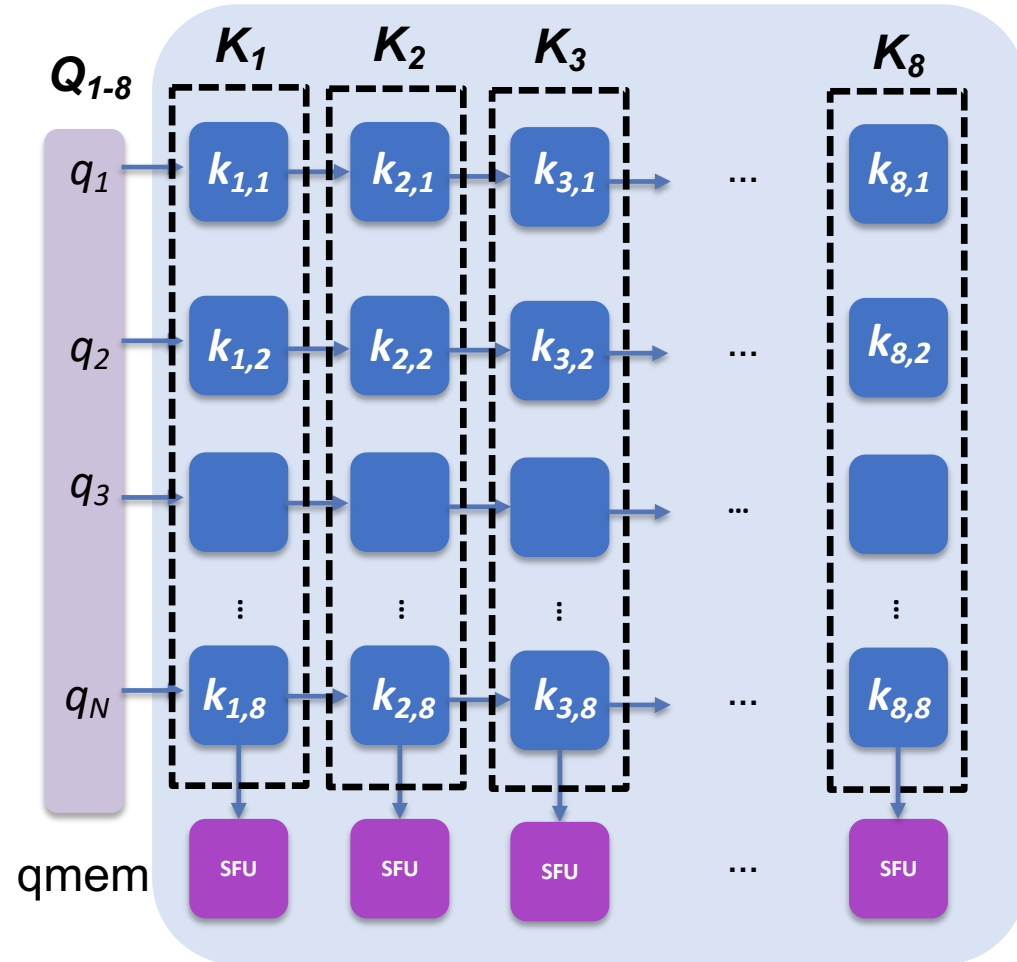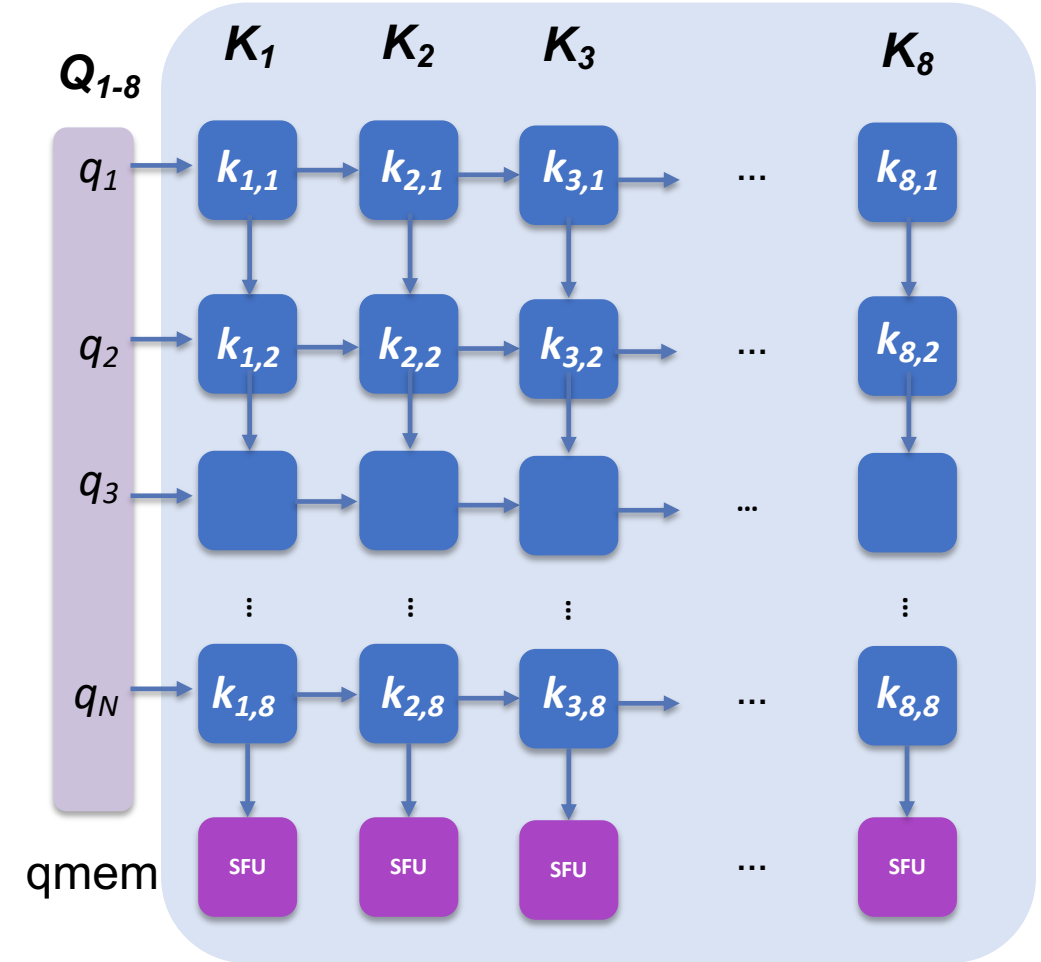

**NVIDA NVDLA**
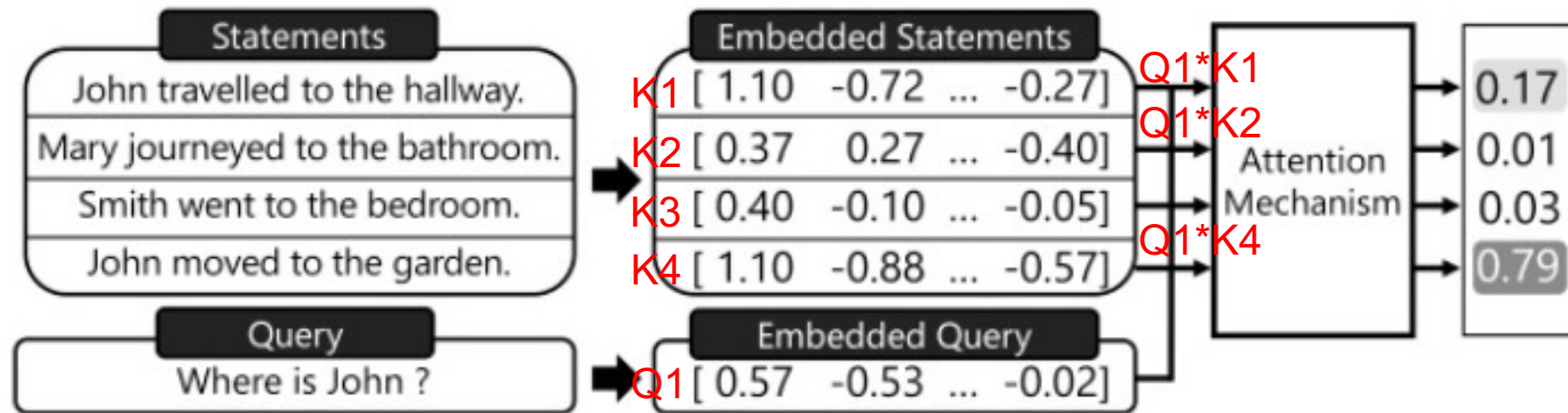
**Microsoft Brainwave**

# Architecture



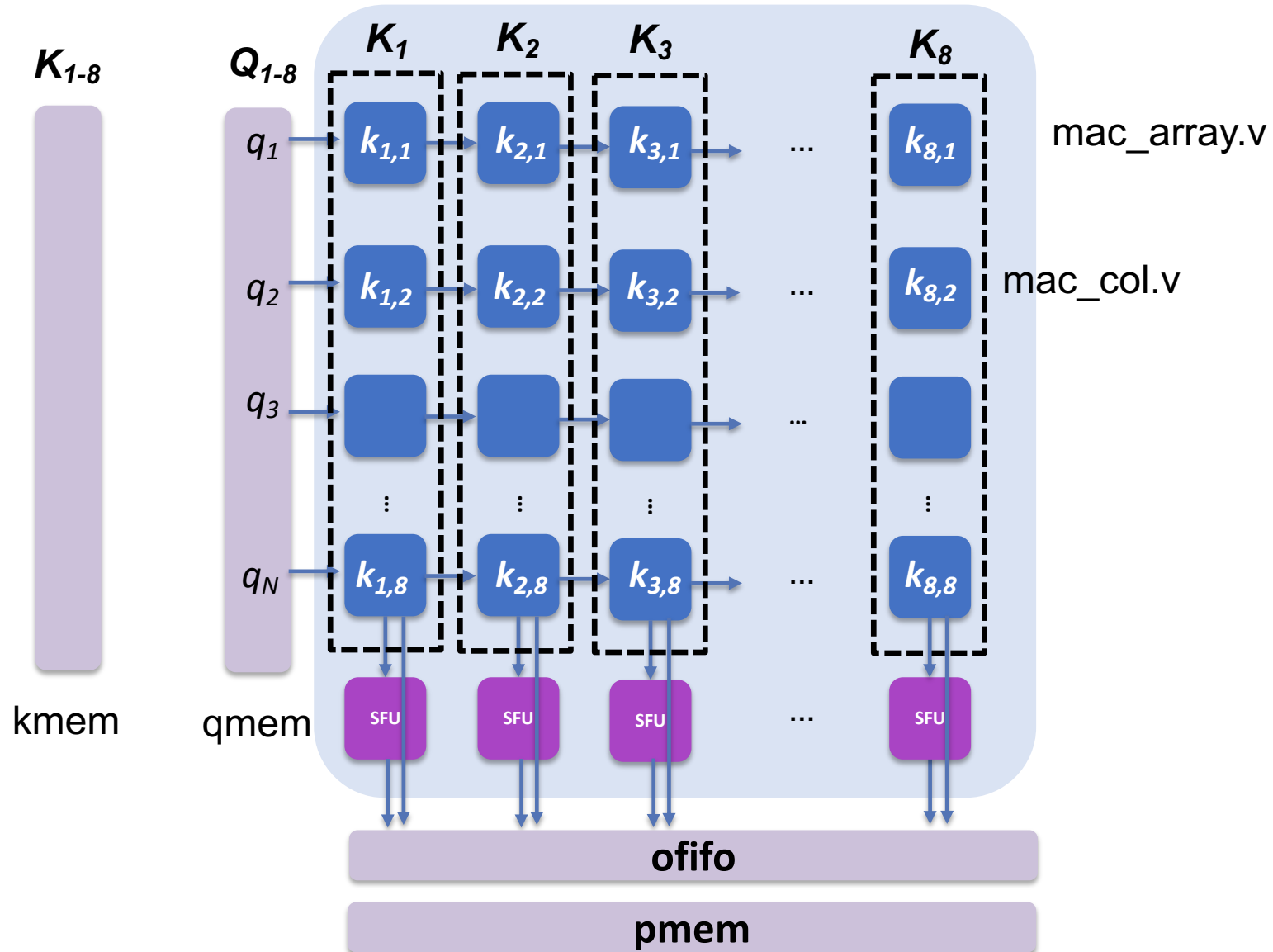1-D vector processor architecture

2D-systolic array architecture

# Front-end Computation in Transformer Model for Natural Language Processing



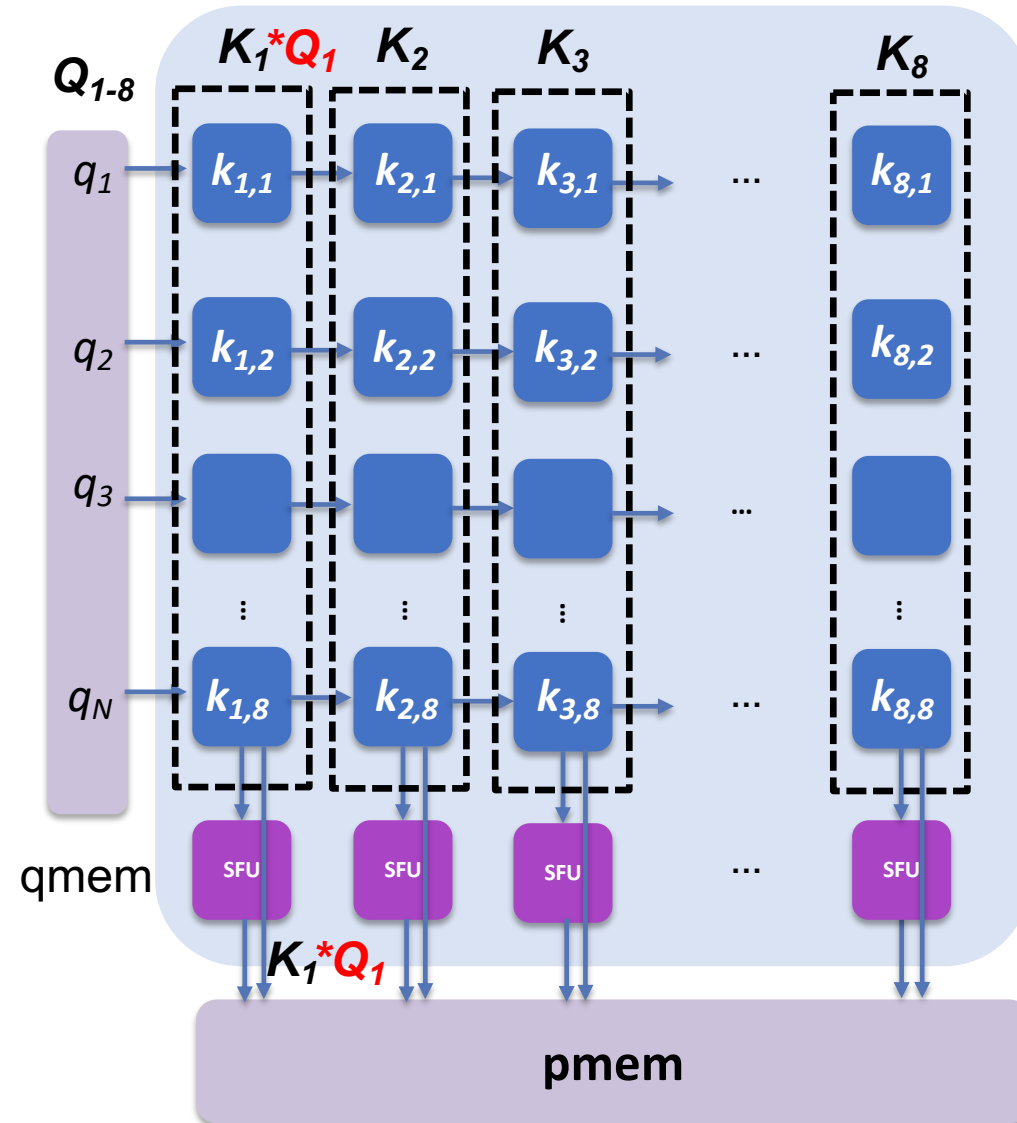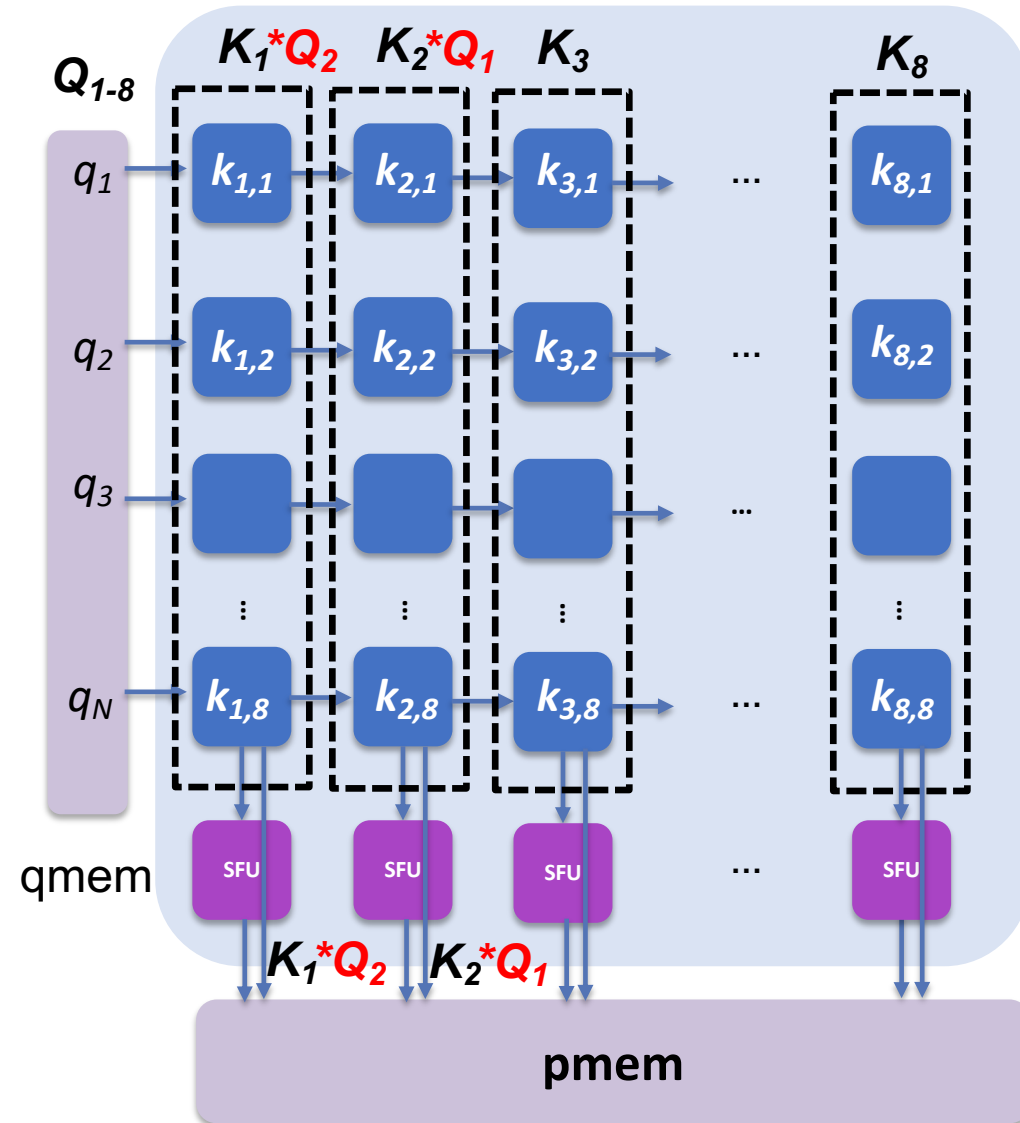Facebook bAbi QA task processing

# Architecture

# Computation

# Computation

# Computation

# Computation

# Processing Stage

- Stage1: memory write stage

  - qmem & kmem: fill both memory with "qdata.txt" and "kdata.txt"


- Stage2: k-data loading stage

  - move the data from kmem to process's local register


- Stage3: execute

  - move the data from qmem to input of the processor

# More Detail

# OFIFO Operation



WR = 1

cyc = 1

WR = 1   WR = 1   WR = 1

cyc = 3

WR = 1   WR = 1

cyc = 2

WR = 1   WR = 1   WR = 1   WR = 1   WR = 1

ready to
be read

...

pmem

cyc = 8

# Normalization

1st time
[Q0K0, Q0K1, ..., Q0K7] / sum[ abs(Q0K0), abs(Q0K1), ..., abs(Q0K7) ]

2nd time
[Q1K0, Q1K1, ..., Q1K7] / sum[ abs(Q1K0), abs(Q1K1), ..., abs(Q1K7) ]
...

sfp

Write back

fetch

[Q0K0, Q0K1, ..., Q0K7] @ add = 0
[Q1K0, Q1K1, ..., Q1K7] @ add = 1
[Q2K0, Q2K1, ..., Q2K7] @ add = 2
....

pmem

# Extension for Dual Core



Normalize by "sum( abs(Q0K0), abs(Q0K1), ..., abs(Q0K7 ) ) from core0
+ sum( abs(Q0K8), abs(Q0K9), ..., abs(Q0K15) ) from core1"

# Optimization

- Optimize the implementation to maximize the frequency and minimize power

    - However, the project is not a spec competition

- Apply techniques you have learned from this course

- Potential optimization (does not mean you need to do everything below)

    - Pipeline, multi-cycle path, memory double-buffering, clock gating, loop unrolling, more # of

      cores, any other creative idea

    - Async interface handling (synchronizer, hand shaking, ….)

    - Maximize parallelism between processing stages

    - Design your own controller

    - better verification (with random input)

    - Minimize WNS and DRC errors

# How to prove ?

- Prove the efficacy of your idea (why ? how ? how much ?)

  - measure power and check functionality with VCD

    (SDF is not required, but it will be plus alpha)

  - clear comparison before vs. after

  - clear description in the report regarding your design choice

    e.g., why you chose the decision, why you cut the pipeline boundary, ….

  - report is VERY important

# Project Policy

- TA & Instructor **will not**

  - solve your project instead of you

  - debug your code (e.g., I got this error with screenshot, can you check my code ?)

  - analyze your issue (e.g., my performance is not so good, do you know why ?)

  - comment on your result (e.g., Can you check my result is correct ?, Can I get A grade with this optimization ?)

  - "fighting against errors" is the biggest part of this project


- However, TA & Instructor will

  - clarify the project problem

  - help any environment issue (e.g., it seems there is no license file, tool does not load, …)

  - provide general guideline for a certain issue


- You can still discuss between students through discussion tab