

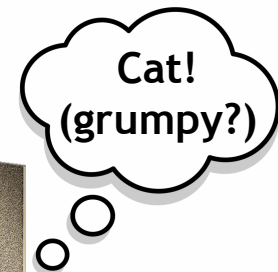
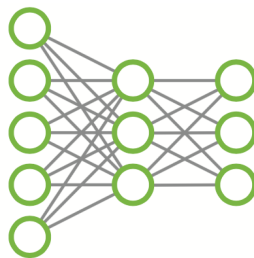


# Vulkan ML

Piers Daniell  
July 2019

# Why machine learning in Vulkan?

- Research showcases potential use of machine learning in interactive and high frame rate applications
  - Character animation (phase function neural network, etc.)
  - Full screen image processing (antialiasing, upscaling, inpainting, DLSS, etc.)
  - Non-Player Character bots (AlphaStar, OpenAI Five, etc.)
  - Image generations (GAN, fire & smoke & clouds, etc.)
- Current machine learning solutions have relatively high interop overhead
  - Interop with third party framework (Python TensorFlow, PyTorch, OpenCL, etc.) introduces bubbles where the CPU/GPU are not doing useful work
  - Sharing data with external APIs can be challenging due to difference in memory models and may require additional copies



# How to do machine learning in Vulkan?

- This is possible already today
  - Just use compute shaders to implement the various algebra operations.
  - Examples: Tencent/ncnn, Alibaba/MNN, Unity ML, etc.
- Or use compilers which will generate SPIR-V code for you
  - Examples: TVM.AI
- But
  - Writing high efficiency layered matrix multiplications, with various activation functions requires some advanced GPU programming skills, with different solutions for different hardware



# The Vulkan ML TSG (Technical Subgroup)

- A new technical subgroup at Khronos has been formed to improve the solution space for machine learning in Vulkan
- Includes representatives from many companies
- Goals
  - Investigate proprietary extensions for inclusion into core Vulkan (VK\_NV\_cooperative\_matrix, etc.)
  - Improvements to compute shaders specific to ML needs
  - New cross vendor extensions (meta-commands, etc.)
- If you are interested, please reach out to us: [pboudier@nvidia.com](mailto:pboudier@nvidia.com)



# VK\_NV\_cooperative\_matrix

- Exposes NVIDIA's Turing Tensor Cores to Vulkan/SPIR-V
- Accelerates large, low-precision matrix multiplies
- Core compute function for deep learning
- FP16 supported today, UINT8/SINT8 coming soon
- Sample code: [https://github.com/jeffbolznm/vk\\_cooperative\\_matrix\\_perf](https://github.com/jeffbolznm/vk_cooperative_matrix_perf)
- Performance on NVIDIA TITAN RTX
  - fp16 matrix math with fp16 accumulation: 100 TFLOPS

```
C:\WINDOWS\system32\cmd.exe
P:\git\github.com\jeffbolznm\vk_cooperative_matrix_perf>build\Release\vk_cooperative_matrix_perf.exe
usage: vk_cooperative_matrix_perf.exe [--correctness]

shader: shaders/shmemfp16.spv

cooperativeMatrixProps = 16x16x16  A = fp16 B = fp16 C = fp16 D = fp16 scope = subgroup
TILE_M=128 TILE_N=128, TILE_K=32 BColMajor=0  56.299752 TFlops
TILE_M=128 TILE_N=128, TILE_K=32 BColMajor=1  50.764185 TFlops
TILE_M=128 TILE_N=256, TILE_K=32 BColMajor=0  76.122378 TFlops
TILE_M=128 TILE_N=256, TILE_K=32 BColMajor=1  74.003313 TFlops
TILE_M=256 TILE_N=128, TILE_K=32 BColMajor=0  79.110662 TFlops
TILE_M=256 TILE_N=128, TILE_K=32 BColMajor=1  73.402560 TFlops
TILE_M=256 TILE_N=256, TILE_K=32 BColMajor=0  100.239919 TFlops ←
TILE_M=256 TILE_N=256, TILE_K=32 BColMajor=1  60.687488 TFlops
```