# **GDC** 2019

🚓 unity



# Math for Game Developers Inside Neural Networks

Michael Buttner, Unity Labs

Principal Research Engineer michaelbu@unity3d.com

### Image descriptions



Man in black shirt is playing guitar.

Construction worker in orange safety vest is working on road.

Two young girls are playing with lego toy.

Boy is doing backflip on wakeboard.

[Karpathy and Fei-Fei 2015]



### **Style transfer**







[Gatys, Leon A and Ecker, Alexander S and Bethge, Matthias 2015]



### Image enhancement





https://github.com/alexjc/neural-enhance





### Beyond images and classification

Vast amount of neural network applications are "offline" Majority of examples focus on image-related problems NN inference is usually implemented inefficiently

- Neural Networks in performance critical code
- Inference and Training without frameworks
- Data-oriented design



### **Data oriented NN Inference**



Multilayer Perceptron (MLP) – DenseLayer - FullyConnectedLayer

Optimize data layout for memory access patterns

Continuous memory access to avoid cache misses

Minimize memory access

No dynamic memory allocation

128-bit SSE Load & Store instructions

$$a_l = \sigma(W_l a_{l-1} + b_l)$$

$$\hat{y} = \sigma \left( W_2 \left( \sigma \left( W_1 + \left( \sigma (W_0 x + b_0) \right) b_1 \right) \right) + b_2 \right)$$



### **Instruction Timings**

Instruction	Latency
SQRTSS/PS	9-10
VSQRTPS	9-10
SQRTSD/PD	14-15
VSQRTPD	14-15
RSQRTSS/PS	5
VRSQRTPS	5

Instruction	Latency
FSIN	50-170
FCOS	50-115
FSINCOS	60-120
FPTAN	~90
FPATAN	50-160

https://www.agner.org/optimize/instruction\_tables.pdf

By Agner Fog. Technical University of Denmark. Copyright © 1996 – 2018. Last updated 2018-09-15.

AMD Ryzen

### **Memory Caching**



CppCon 2014: Mike Acton "Data-Oriented Design and C++"







CppCon 2014: Mike Acton "Data-Oriented Design and C++"

### **DirectX Math**



"An all inline SIMD C++ linear algebra library for use in games and graphics apps"

https://github.com/Microsoft/DirectXMath

### **Alignment & Padding**

Row major layout & 128-bit padding for vectors and matrices





### MatrixMN & VectorN

### struct MatrixMN

};

float\* ptr; uint16\_t numRows; uint16\_t numColumns;

INLINE MatrixMN(const uint32\_t rows, const uint32\_t columns);

INLINE MemoryRequirements GetMemoryRequirements() const; INLINE MemoryRequirements operator()() const; INLINE MemoryBlock& operator()(MemoryBlock& memoryBlock);

INLINE uint32\_t GetNumRows() const; INLINE uint32\_t GetAlignedNumRows() const; INLINE uint32\_t GetNumColumns() const; INLINE uint32\_t GetAlignedNumColumns() const;

```
INLINE float* GetPtr();
INLINE const float* GetPtr() const;
INLINE float* GetRowPtr(const uint32_t row);
INLINE const float* GetRowPtr(const uint32_t row) const;
INLINE XMFLOAT4A* GetPtrSSE();
INLINE const XMFLOAT4A* GetPtrSSE() const;
```

### struct VectorN

```
float* ptr;
uint16_t numValues;
```

INLINE VectorN(const uint32\_t values);

INLINE MemoryRequirements GetMemoryRequirements() const; INLINE MemoryRequirements operator()() const; INLINE MemoryBlock& operator()(MemoryBlock& memoryBlock);

```
INLINE uint32_t GetNumValues() const;
INLINE uint32_t GetAlignedNumValues() const;
```

```
INLINE float* GetPtr();
INLINE const float* GetPtr() const;
INLINE XMFLOAT4A* GetPtrSSE();
INLINE const XMFLOAT4A* GetPtrSSE() const;
```

Matrix and Vector reference memory Decoupled memory management

### **MemoryBlocks & MemoryRequirements**



### **Stack allocation**

Requires tight control of scopes

Reasonably sized memory blocks can be allocated on the stack

sub rsp, <size>

#define ALIGNED\_ALLOCA(size, alignment) align\_ptr( \
 reinterpret\_cast<size\_t>(alloca((size) + alignment)), alignment)
#define STACK\_MEMORY(requirements) \
 ((void\*)ALIGNED\_ALLOCA((requirements.Size()), \
 (requirements.Alignment())))

void\* Allocate(size\_t size) { return alloca(size); }

### SSE MatMul



const XMVECTOR v(XMLoadFL

const XMVECTOR const XMVECTOR const XMVECTOR

XMStoreFloat4A

void Mul(const MatrixMN& lhs, const VectorN& rhs, VectorN& result)

ASSERT(lhs.GetNumRows() == result.GetNumValues()); ASSERT(lhs.GetNumColumns() == rhs.GetNumValues());

const uint32 t advance(lhs.GetAlignedNumColumns() >> 2); const uint32 t blockAdvance(advance \* 3);

XMFLOAT4A\* writePtr(result.GetPtrSSE());

const XMFLOAT4A\* r0(lhs.GetPtrSSE()); const XMFLOAT4A\* r1(r0 + advance); const XMFLOAT4A\* r2(r1 + advance); const XMFLOAT4A\* r3(r2 + advance);

const uint32\_t numRows(lhs.GetAlignedNumRows() >> 2); for (uint32 t row = 0; row < numRows; ++row)</pre>

const XMFLOAT4A\* readPtr(rhs.GetPtrSSE());

XMVECTOR d0(XMVectorZero()); XMVECTOR d1(XMVectorZero()); XMVECTOR d2(XMVectorZero()); XMVECTOR d3(XMVectorZero());

for (uint32 t i = 0; i < advance; ++i)</pre>

const XMVECTOR v(XMLoadFLoat4A(readPtr++));

d0 = XMVectorAdd(d0, XMVector4Dot(XMLoadFLoat4A(r0++), v)); d1 = XMVectorAdd(d1, XMVector4Dot(XMLoadFLoat4A(r1++), v)); d2 = XMVectorAdd(d2, XMVector4Dot(XMLoadFLoat4A(r2++), v)); d3 = XMVectorAdd(d3, XMVector4Dot(XMLoadFLoat4A(r3++), v));

const XMVECTOR d4(XMVectorPermute<0, 1, 4, 5>(d0, d1)); const XMVECTOR d5(XMVectorPermute<0, 1, 4, 5>(d2, d3)); const XMVECTOR d(XMVectorPermute<0, 2, 4, 6>(d4, d5));

r0 += blockAdvance: r1 += blockAdvance: r2 += blockAdvance; r3 += blockAdvance;

ector4 ector4 ector4 ector4 х i>(d2, →(d4, d

oat4A(r0++), v)) oat4A(r1++), v)) oat4A(r2++), v)) oat4A(r3++), v))



### **Composite operations**

All this pain just for a simple matrix multiplication?

"Layer Fusion" | Composite operations

 $\hat{y} = \sigma(Wx + b)$  can be combined into a single operation

Reduce memory access & hide memory latency

$$\sigma(Wx+b) \longrightarrow \begin{array}{c} u = Wx \\ v = u + b \\ s = \sigma(v) \end{array}$$

```
void Sigmoid(MatrixMN& lhs)
```

```
const XMVECTOR one(XMVectorReplicate(1.0f));
```

```
XMFLOAT4A* ptr(lhs.GetPtrSSE());
```

```
const uint32_t numValues(lhs.GetAlignedNumColumns() >> 2);
for (uint32_t i = 0; i < numValues; ++i)</pre>
```

```
const XMVECTOR x(XMLoadFLoat4A(ptr));
```

```
// f(x) = 1 / (1 + e^-x)
const XMVECTOR result(
    XMVectorMax(x, XMVectorZero()) +
        XMVectorExpE(XMVectorMin(x, XMVectorZero())) - one);
```

// f(x) = x / (1 + abs(x)) | "Fast sigmoid"
//const XMVECTOR result(XMVectorDivide(x,
// XMVectorAdd(one, XMVectorAbs(x))));

```
XMStoreFLoat4A(ptr++, result);
```

### SSE FullyConnected Layer

const XMVECTOR d4(XMVectorPermute<0, 1, 4, 5>(d0, d1)); const XMVECTOR d5(XMVectorPermute<0, 1, 4, 5>(d2, d3)); const XMVECTOR d(XMVectorPermute<0, 2, 4, 6>(d4, d5));

XMStoreFloat4A(writePtr++, d);

const XMVECTOR d4(XMVectorPermute<0, 1, 4, 5>(d0, d1)); const XMVECTOR d5(XMVectorPermute<0, 1, 4, 5>(d2, d3)); const XMVECTOR d(XMVectorPermute<0, 2, 4, 6>(d4, d5));

const XMVECTOR x(XMVectorAdd(d, XMLoadFloat4A(biasPtr++)));

```
XMStoreFloat4A(writePtr++, y);
```

ASSERT(lhs.GetNumRows() == result.GetNumValues()): ASSERT(lhs.GetNumColumns() == rhs.GetNumValues()); const uint32 t advance(lhs.GetAlignedNumColumns() >> 2); const uint32 t blockAdvance(advance \* 3); XMFLOAT4A\* writePtr(result.GetPtrSSE()); const XMFLOAT4A\* r0(lhs.GetPtrSSE()); const XMFLOAT4A\* r1(r0 + advance); const XMFLOAT4A\* r2(r1 + advance); const XMFLOAT4A\* r3(r2 + advance); const uint32\_t numRows(lhs.GetAlignedNumRows() >> 2); for (uint32 t row = 0; row < numRows; ++row)</pre> const XMFLOAT4A\* readPtr(rhs.GetPtrSSE()); XMVECTOR d0(XMVectorZero()); XMVECTOR d1(XMVectorZero()): XMVECTOR d2(XMVectorZero()); XMVECTOR d3(XMVectorZero()); for (uint32 t i = 0; i < advance; ++i)</pre> const XMVECTOR v(XMLoadFLoat4A(readPtr++)); d0 = XMVectorAdd(d0, XMVector4Dot(XMLoadFLoat4A(r0++), v)); d1 = XMVectorAdd(d1, XMVector4Dot(XMLoadFLoat4A(r1++), v)); d2 = XMVectorAdd(d2, XMVector4Dot(XMLoadFLoat4A(r2++), v)); d3 = XMVectorAdd(d3, XMVector4Dot(XMLoadFLoat4A(r3++), v)); const XMVECTOR d4(XMVectorPermute<0, 1, 4, 5>(d0, d1)); const XMVECTOR d5(XMVectorPermute<0, 1, 4, 5>(d2, d3)); const XMVECTOR d(XMVectorPermute<0, 2, 4, 6>(d4, d5)); r0 += blockAdvance: r1 += blockAdvance; r2 += blockAdvance:

void Mul(const MatrixMN& lhs, const VectorN& rhs, VectorN& result)

### **Performance gain**

Performance gain depends on network architecture and feature vector size

Complex architectures gain the most

- Experiment with naïve implementation
- Lock down architecture
- Side-by-side implementation
- Profile and optimize

### Complex manifold example

- Feature vectors size 400+ values
- From 3.0 ms to 0.7 ms (76% gain)



float\* writePtr(result.GetPtr());

const XMFLOAT4A\* w0(m\_weights[index][index0].GetPtrSSE()); const XMFLOAT4A\* w1(m weights[index][index1].GetPtrSSE()); const XMFLOAT4A\* w2(m\_weights[index][index2].GetPtrSSE()); const XMFLOAT4A\* w3(m\_weights[index][index3].GetPtrSSE());

XMVector4Dot(Blend(a, b, c, d, t),

XMLoadFloat4A(readPtr++));

< 🕄 unitv

 $\chi$  –

## **Neural Network Training**

The goal is to adjust all  $W_l$ ,  $b_l$  such that for any training sample (x, y) the loss *L* becomes 0

 $\frac{\partial L}{\partial W_l}$ ,  $\frac{\partial L}{\partial b_l}$  tell us how to adjust  $W_l$ ,  $b_l$  for any given (x, y)

Loss function can include additional constraints For example maintain unit quaternion property  $N(x) = \sigma(Wx + b)$ 

$$\hat{y} = N_2 \left( N_1 \big( N_0(x) \big) \right)$$

$$L = \frac{1}{2}(\hat{y} - y))^2$$



Calculate gradients

### **Neural Network Training**

```
void Train(const Tensor& input, const Tensor& output,
           const uint32 t numEpochs, const Scalar eta)
    for (uint32 t i = 0; i < numEpochs; ++i)</pre>
        Epoch(input, output, eta);
void Epoch(const Tensor& input, const Tensor& output, const Scalar eta)
    ASSERT(input.GetNumRows() == output.GetNumRows());
    const uint32_t numLayers(uint32_t(weights.Length()));
    const uint32 t numTrainingSamples(uint32 t(input.GetNumRows()));
    for (uint32 t i = 0; i < numTrainingSamples; ++i)</pre>
        CalculateGradients(input.GetRow(i), output.GetRow(i));
```

```
void UpdateWeightsAndBiases(const Scalar factor)
```

```
ASSERT(weights.Length() == gradient_weights.Length());
const uint32_t numWeights(uint32_t(weights.Length()));
for (uint32_t i = 0; i < numWeights; ++i)</pre>
```

ElementwiseSub(weights[i], gradient\_weights[i], factor);

```
ASSERT(biases.Length() == gradient_biases.Length());
const uint32_t numBiases(uint32_t(biases.Length()));
for (uint32_t i = 0; i < numBiases; ++i)</pre>
```

ElementwiseSub(biases[i], gradient\_biases[i], factor);

### **Calculate Gradients**



void CalculateGradients(const Tensor& x, const Tensor& v)

### **Beyond Backpropagation**

What if we are not using a simple MLP network?

What if we want to experiment with different architectures?

 $y(\theta) = \sigma(u(\theta)x + v(\theta))$  $u(\theta) = \dots$  $v(\theta) = \dots$ 

Scalar x1(2); Scalar x2(3); Scalar loss = cos(x1) + x1 \* exp(x2); const float dydx1(x1.derivative()); const float dydx2(x2.derivative()); Backpropagation is just a special case of reverse-mode automatic differentiation, applied to a neural network

"AD allows exact and efficient calculation of derivatives, by systematically invoking the chain rule of calculus at the elementary operator level during program execution"

Neural networks are more than just a series of "layers"

### **Algorithmic Differentiation**

 $f(x_1, x_2) = \cos(x_1) + x_1 \exp(x_2)$ 

Expressions can be divided into elementary operations

Computational graphs are a subset of abstract syntax trees

$$\begin{array}{ll} f(x_1, x_2) = w_1 & \dot{f}(x_1, x_2) = \dot{w}_1 \\ w_1 = w_2 + w_3 & \dot{w}_1 = \dot{w}_2 + \dot{w}_3 \\ w_2 = \cos(w_5) & \dot{w}_2 = -\sin(w_5) \dot{w}_5 \\ w_3 = w_4 w_5 & \dot{w}_3 = \dot{w}_4 w_5 + \dot{w}_5 w_4 \\ w_4 = e^{w_6} & \dot{w}_4 = e^{w_6} \dot{w}_6 \\ w_5 = x_1 & \dot{w}_5 = \operatorname{seed} \in \{0, 1\} \\ w_6 = x_2 & \dot{w}_6 = \operatorname{seed} \in \{0, 1\} \end{array}$$



Forward primal trace calculates rate of change of *f* with respect to *one* of its inputs

### **Adjoint Mathematics**

We call  $f_i$  the result of the operation on  $n_i$ , we denote  $f = f_N$  the final result on the op node  $n_N$  and  $C_i$  the set of arguments to node  $n_i$ .

The adjoint  $a_i$  of the result  $f_i$  on node  $n_i$  is the partial derivative of the final result  $f = f_N$  to  $f_i$ 

$$a_i = \frac{df}{df_i}$$

Since obviously  $a_N = \frac{\partial a_N}{\partial a_N} = 1$  and directly from the derivative chain rule, the fundamental adjoint equation is:

$$a_j = \sum_{i/n_j \in C_i} \frac{\partial f_i}{\partial f_j} a_i$$



 $\overline{w}_1 = \frac{\partial f}{\partial f} = 1$ +  $w_1$  $f(x_1, x_2) = \cos(x_1) + x_1 \exp(x_2)$  $f(x_1, x_2) = w_1$  $\overline{w}_2 = 0$  $\overline{w}_3 = 0$ COS  $w_1 = w_2 + w_3$  $w_2$  $W_3$ \*  $w_2 = \cos(w_5)$  $w_3 = w_4 w_5$  $w_4 = e^{w_6}$  $\overline{w}_5 = 0$  $\overline{w}_4 = 0$  $w_4 exp$  $W_5$  $\chi_1$  $w_5 = x_1$  $w_6 = x_2$  $\overline{w}_6 = 0$  $W_6$  $\chi_2$  $\frac{\partial f}{\partial x_1} = -\sin(x_1) + e^{x_2} \qquad \text{Let } x_1 = 2 \text{ and } x_2 = 3,$  $\frac{\partial f}{\partial x_2} = e^{x_2} x_1 \qquad \text{then } \frac{\partial f}{\partial x_1} = 19.19 \text{ and } \frac{\partial f}{\partial x_2} = 40.17$ 

+  $\overline{w}_1 = \frac{\partial f}{\partial f} = 1$  $w_1$  $f(x_1, x_2) = \cos(x_1) + x_1 \exp(x_2)$  $f(x_1, x_2) = w_1$  $\overline{w}_2 = 1$  $\overline{w}_3 = 1$ COS  $w_1 = w_2 + w_3$  $w_2$  $W_3$ \*  $w_2 = \cos(w_5)$  $w_3 = w_4 w_5$  $w_4 = e^{w_6}$  $\overline{w}_5 = 0$  $\overline{w}_4 = 0$  $w_4 exp$  $W_5$  $\chi_1$  $w_5 = x_1$  $w_6 = x_2$  $\overline{w}_6 = 0$  $W_6$  $\chi_2$  $\frac{\partial f}{\partial x_1} = -\sin(x_1) + e^{x_2} \qquad \text{Let } x_1 = 2 \text{ and } x_2 = 3,$  $\frac{\partial f}{\partial x_2} = e^{x_2} x_1 \qquad \text{then } \frac{\partial f}{\partial x_1} = 19.19 \text{ and } \frac{\partial f}{\partial x_2} = 40.17$ 



 $\overline{w}_1 = \frac{\partial f}{\partial f} = 1$  $W_1$ +  $f(x_1, x_2) = \cos(x_1) + x_1 \exp(x_2)$  $f(x_1, x_2) = w_1$  $\overline{w}_2 = 1$  $\overline{w}_3 = 1$ COS  $w_1 = w_2 + w_3$  $w_2$  $W_3$ \*  $w_2 = \cos(w_5)$  $w_3 = w_4 w_5$  $w_4 = e^{w_6}$  $\overline{w}_5 = -0.91 \quad w_4 \quad exp$  $\overline{w}_4 = 0$  $W_5$  $\chi_1$  $w_5 = x_1$  $w_6 = x_2$  $\overline{w}_6 = 0$  $W_6$  $\chi_2$  $\frac{\partial f}{\partial x_1} = -\sin(x_1) + e^{x_2} \qquad \text{Let } x_1 = 2 \text{ and } x_2 = 3,$  $\frac{\partial f}{\partial x_2} = e^{x_2} x_1 \qquad \text{then } \frac{\partial f}{\partial x_1} = 19.19 \text{ and } \frac{\partial f}{\partial x_2} = 40.17$ 









## **Algorithmic Adjoint Differentiation**

Reverse adjoint trace calculates rate of change of *f* with respect to *all* of its inputs

```
Evaluating \nabla f(x) is as fast evaluating f(x)
```

Forward pass evaluates the expression

Reverse pass evaluates  $\nabla f(x)$ 

AAD allows us to use arbitrary control structures like loops or branches, intermediate variables and functions

Data-oriented design & SSE

```
void Train(const Tensor& input, const Tensor& output,
           const uint32 t numEpochs, const Scalar eta)
    for (uint32 t i = 0; i < numEpochs; ++i)</pre>
        Epoch(input, output, eta);
void Epoch(const Tensor& input, const Tensor& output, const Scalar eta)
    ASSERT(input.GetNumRows() == output.GetNumRows());
    const uint32 t numLayers(uint32 t(weights.Length()));
    const uint32 t numTrainingSamples(uint32 t(input.GetNumRows()));
    for (uint32 t i = 0; i < numTrainingSamples; ++i)</pre>
        CalculateGradients(input.GetRow(i), output.GetRow(i));
        UpdateWeightsAndBiases(eta);
```

Feature vectors size 360 values Multiple hidden layers & 1d Convolution From 0.5 ms to 0.1 ms (~80% gain)

> ~70.000 training samples Stochastic gradient descent Single training update < 0.3 ms Single epoch ~20 seconds No framework dependencies

Michael Buttner, Unity Labs Principal Research Engineer michaelbu@unity3d.com

### **Algorithmic Adjoint Differentiation**

