

Optimization of neural computations in a functional data-parallel language

Naums Mogers

GPU code optimization: **portability** versus **performance**

- Manual optimization → **good performance**
 - **expensive** to do
 - **not portable**
 - lack **usability**
 - does not support **new devices**
- Autotuners (PetaBricks, CLTune) → Functionally portable
 - **not performance-portable**
 - **no** structural optimizations

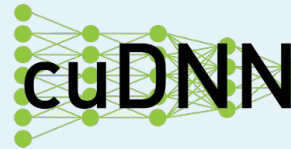
High-level Neural Network syntax

Caffe TensorFlow theano torch

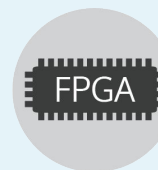
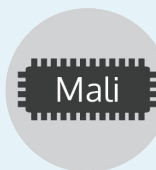
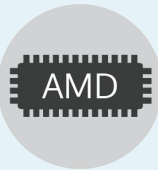
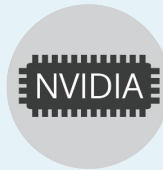
Optimization & compilation methods



Manual
implementation



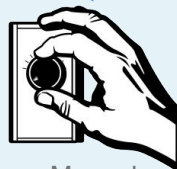
Hardware



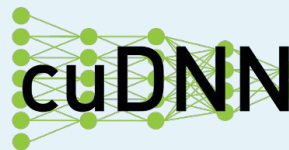
High-level Neural Network syntax

Caffe TensorFlow theano torch

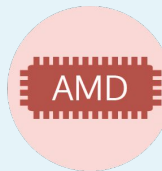
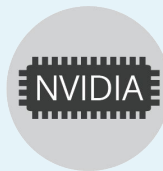
Optimization & compilation methods



Manual
implementation



Hardware



The Lift language

- Functional
 - Abstracted from hardware
 - Algorithm-centred
 - Pure and safe
 - High-level, easy to use
- Data-parallel

The Lift language

- Functional
 - Abstracted from hardware
 - Algorithm-centred
 - Pure and safe
 - High-level, easy to use
- Data-parallel
- Chooses the best OpenCL derivation for the target hardware
 - Both functionally and performance portable
 - Doesn't require hardware knowledge

Lift's rewrite rules

- Semantics-preserving transformations encoding fine-grained optimizations

Lift's rewrite rules

- Semantics-preserving transformations encoding fine-grained optimizations

Vectorization

Lift's rewrite rules

- Semantics-preserving transformations encoding fine-grained optimizations

Vectorization

Memory tiling

Lift's rewrite rules

- Semantics-preserving transformations encoding fine-grained optimizations

Vectorization

Memory coalescing

Memory tiling

Lift's rewrite rules

- Semantics-preserving transformations encoding fine-grained optimizations

Vectorization

Blocking

Memory coalescing

Memory tiling

Lift's rewrite rules

- Semantics-preserving transformations encoding fine-grained optimizations

Vectorization

Blocking

ND mapping

Memory coalescing

Memory tiling

Lift's rewrite rules

- Semantics-preserving transformations encoding fine-grained optimizations

Vectorization

Blocking

ND mapping

Memory coalescing

Memory tiling

Simplification

The method:

Neural Network-specific extension
of the Lift language

The extension:

- Neural Network (NN)-specific optimizations

The user

- Encodes the NN in Lift
- Specifies the minimum required accuracy



The user

- Encodes the NN in Lift
- Specifies the minimum required accuracy



Lift

- Applies generic optimizations
- Optimizes the NN code without preserving semantics
- Abides by the required accuracy

Proposed optimizations:

- Approximations
 - Floating operations
 - Different layer precisions
 - Gradient quantization

Proposed optimizations:

- Approximations
 - Floating operations
 - Different layer precisions
 - Gradient quantization
- NN configuration autotuning
 - Layer number
 - Layer size
 - Training batch size
 - Learning rate

The extension will be evaluated by

- Implementing Convolutional Neural Network (CNN) forward-propagation and training in Lift

The extension will be evaluated by

- Implementing Convolutional Neural Network (CNN) forward-propagation and training in Lift
- Comparing CNN performance in domain-specific Lift vs generic Lift

The extension will be evaluated by

- Implementing Convolutional Neural Network (CNN) forward-propagation and training in Lift
- Comparing CNN performance in domain-specific Lift vs generic Lift
- Comparing CNN performance in domain-specific Lift vs Caffe

Evaluation metrics

- Forward-propagation runtime
- Training runtime
- The range of platforms supported

- Current GPU optimization methods are **not performance portable**
- Lift approach is **performance portable**
- Extend Lift to Neural Network-specific optimizations