Optimization of neural computations using functional data-parallel language

The problem



Modern neural networks are large and require a lot of computational power, which sequential systems are unable to provide.



Parallel systems that are used for neural training are hard to optimize due to an infinitely large optimizational space.



Hardware configurations vary and evolve rapidly, which harms code portability: manual approach is expensive, and machine learning libraries are not performance portable.

The approach



Encode nets using functional data-parallel language Lift www.lift-lang.org



Rewrite rules encode fine-grained structural optimizations...



Introduce neural networkspecific rewrite rules...



Generate hundreds of code variants



THE UNIVERSITY of EDINBURGH Informatics



Choose the target hardware platform



...and parameters for autotuning



...that approximate calculations and optimize network configuration



Choose the best code variant for a target platform using machine learning

Naums Mogers naums.mogers@ed.ac.uk www.naumsmogers.me Kenneth Heafield

Generic Lift stack

Functional language

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

Generic rewrite rules

- Preserve semantics
- ▲ Improve performance
- Vectorization
- Memory coalescing, tiling
- Blocking
- Expression simplification
- Mapping to ND threads
- Split kernels
- Share 32-bit registers
- Use proprietary subroutines



OpenCl

- Low-level hardware management
- Parametrized yields to autotuning
- Cross-platform



EPSRC Centre for Doctoral Training in **Pervasive Parallelism**

- Learning rate autotuning
- Training batch size autotuning
- Layer size autotuning
- Layer number autotuning
- Gradient quantization
- Vary precision among layers
- Float ops approximation

- ▲ Alter accuracy
- Alter semantics

Neural rewrite rules

- Encode information about neural network
- Specify required accuracy

Additional constructs

- **Domain-specific Lift stack**

Supervisors: Christophe Dubach Mike O'Boyle