

CS133

Parallel & Distributed Computing

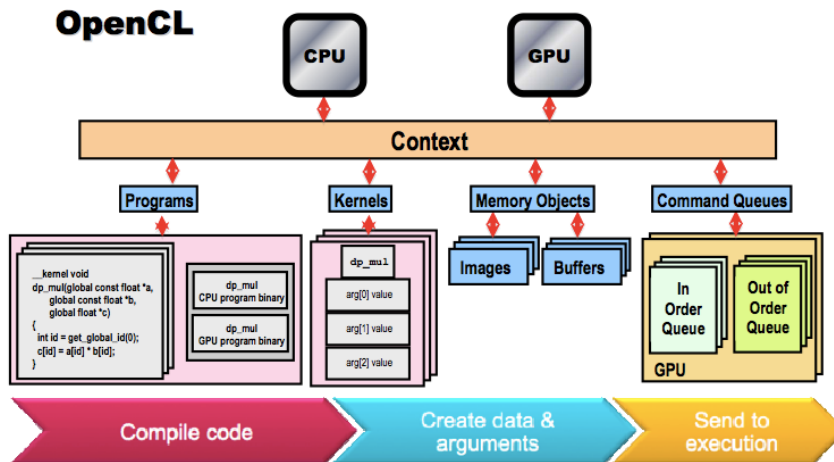
Introduction to Convolutional Neural Networks

Instructor: Jason Cong

cong@cs.ucla.edu

1

Review of OpenCL Basics



© Copyright Khronos Group, 2009 - Page 15

P. Mistry & D. Schaa (Northeastern Univ.), with B. Gaster (AMD), AMD © 2011

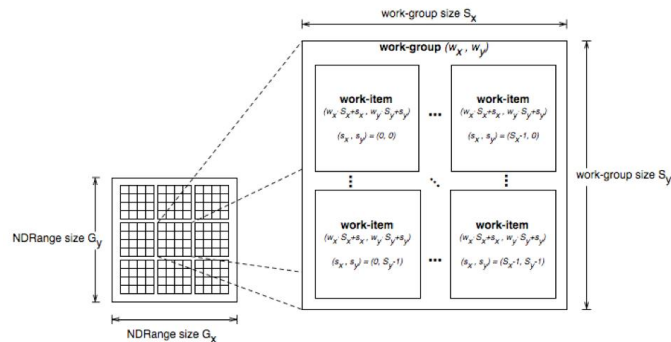
2

2

Correction from Last Lecture -- Thread Structure

□ Work-items can uniquely identify themselves based on:

- A global id (unique within the index space)
- A work-group ID and a local ID within the work-group



P. Mistry & D. Schaa (Northeastern Univ.), with B. Gaster (AMD), AMD © 2011

3

3

Impacts of Deep Learning for Many Applications

Unmanned Vehicle



Speech & Audio



Text & Language



Genomics



Image & Video



Multi-Media



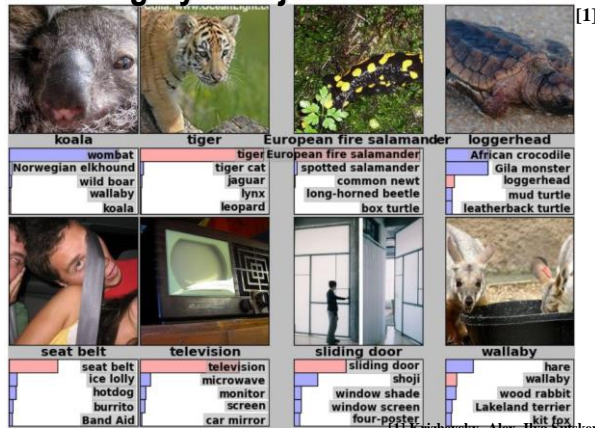
4

All images are from internet search

4

ImageNet Competition

- u 1,200,000 Training Images
 - With 50,000 Validation & 100,000 Test Images
- u 1000 Category of objects



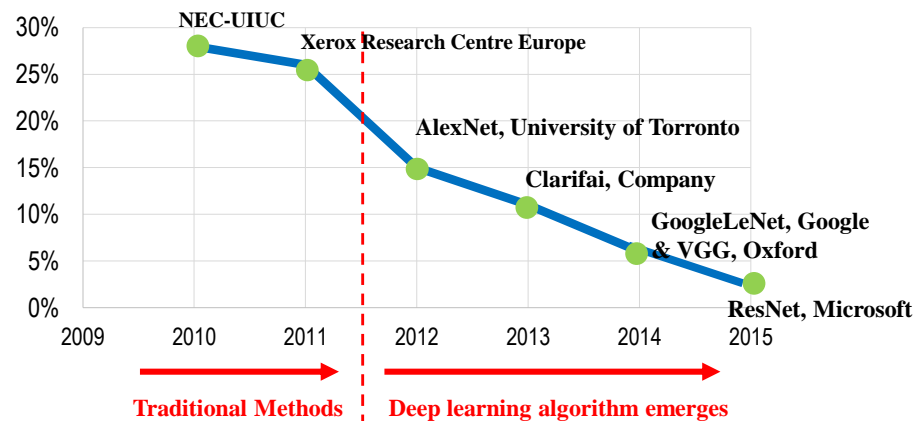
5

H. Krizhevsky, Alex. Ilya. Sutskever, and Geoffrey E. Hinton, NIPS2012

5

ImageNet Competition Results

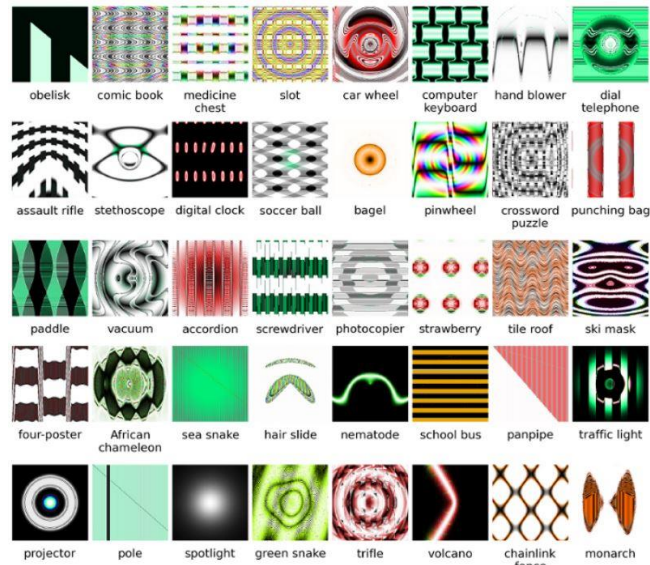
Winning % Error



6

6

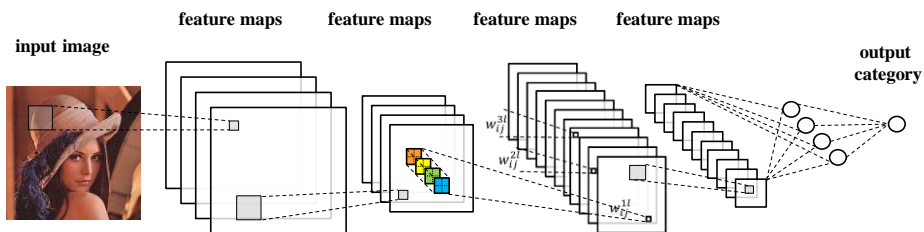
Still not Perfect: Images that Fools CNN/DNN



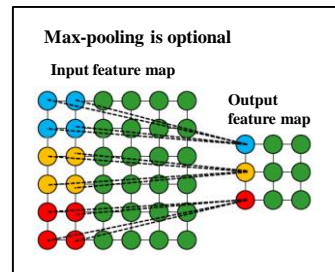
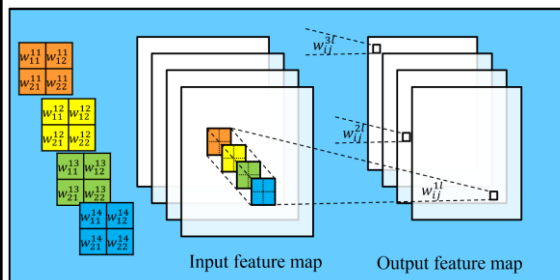
7

7

Inference: Convolutional Neural Network

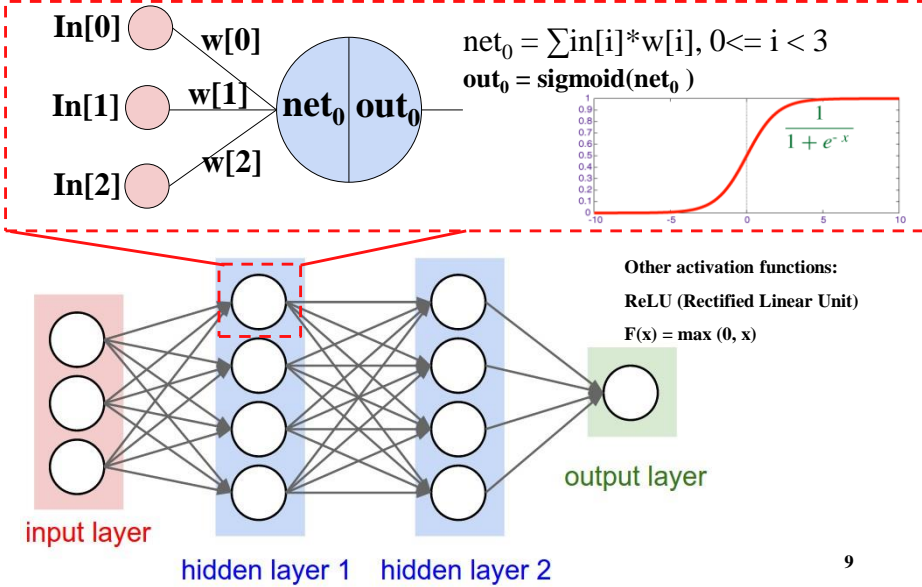


Inference: A feedforward computation



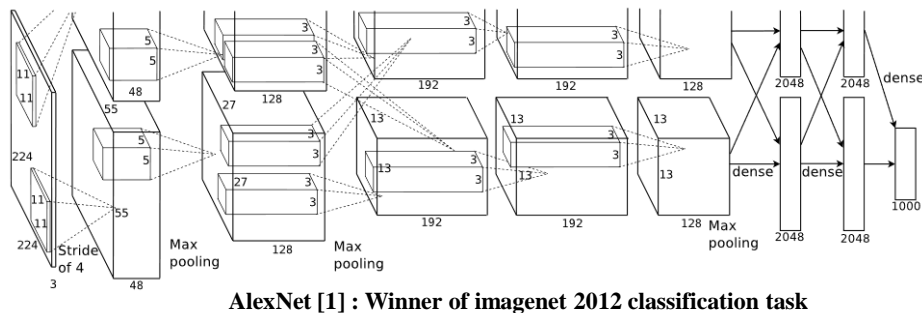
8

Inference: Deep Neural Network (Fully Connected Layers)



9

Real-life CNNs: AlexNet, 8 layers



[1] Krizhevsky, Alex, Ilya Sutskever, and Geoffrey E. Hinton. "Imagenet classification with deep convolutional neural networks." *Advances in neural information processing systems*. 2012.

10

Computation& Storage Analysis

	CONV	POOL	ReLU	Fully connect
Computation Complexity(10^6)	30600	6.12	13.5	122.7
Percentage	99.5%	0.0%	0.1%	0.4%
Storage Complexity (MB)	113	0	0	471.6
Percentage	19.3%	0.0%	0.0%	80.6%
Time% in pure software	96.3%	0.0%	0.0%	3.7%

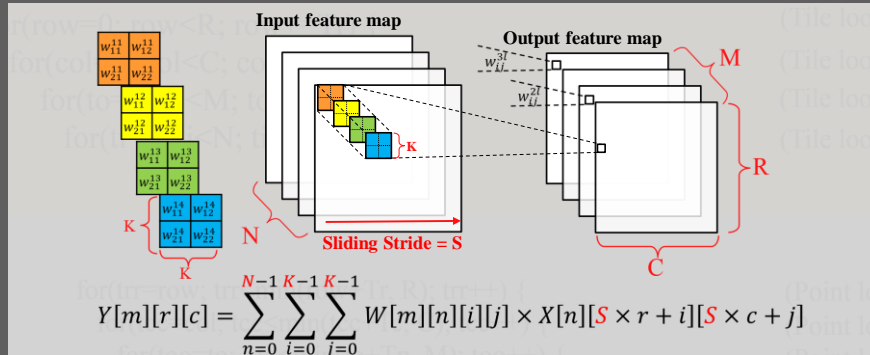
An example of VGG-16 network

A Lot of Interest in Hardware Acceleration !

11

11

Convolution Kernel for Inferencing



```

1 for(row=0; row<R; row++) {
2   for(col=0; col<C; col++) {
3     for(to=0; to<M; to++) {
4       for(ti=0; ti<N; ti++) {
5         for(i=0; i<K; i++) {
6           for(j=0; j<K; j++) {
              output_fm[to][row][col] +=
              weights[to][ti][i][j]*input_fm[ti][S*row+i][S*col+j];
            }
          }
        }
      }
    }
  }
}

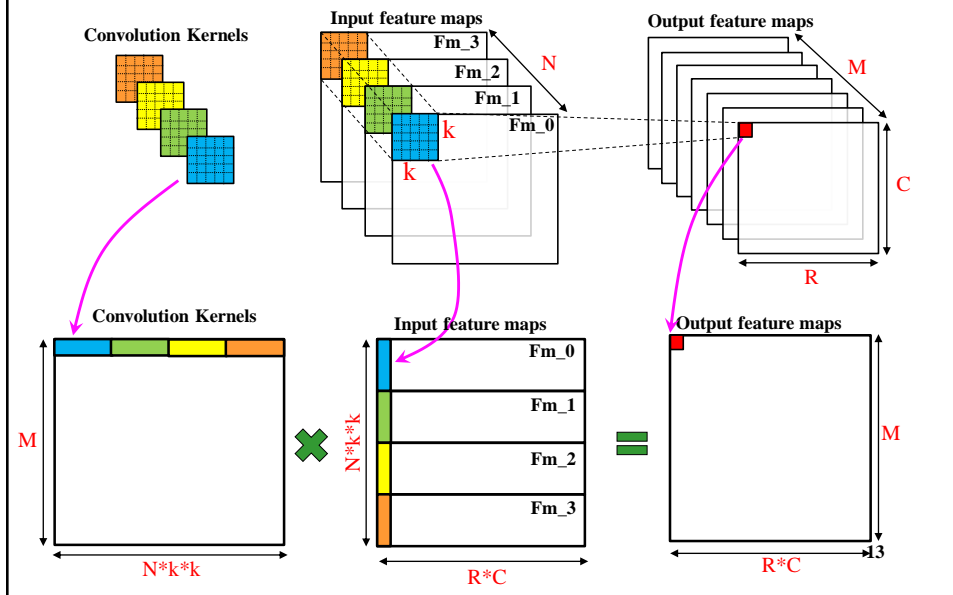
```

R, C, M, N, K, S are all configuration parameters of the convolutional layer

12

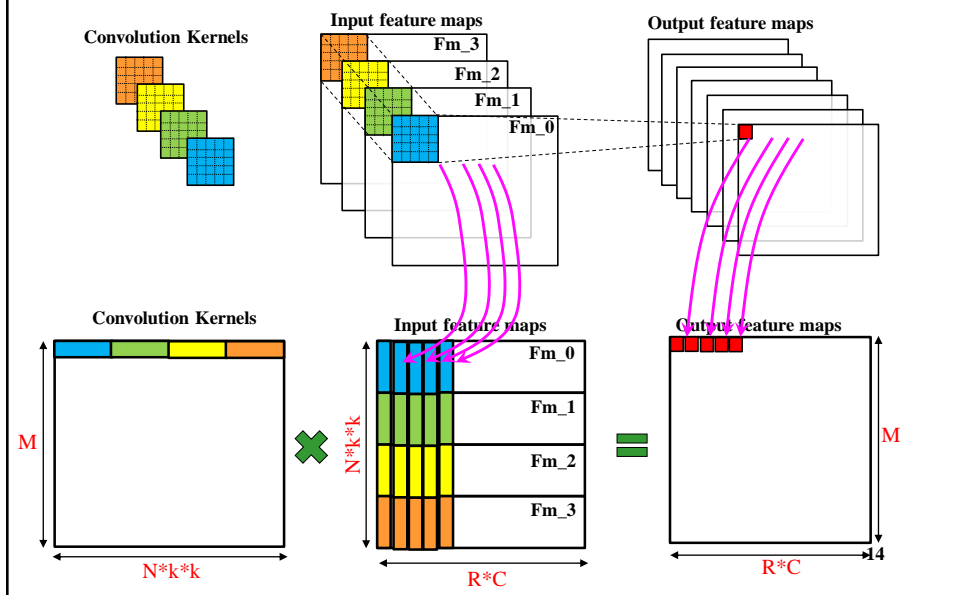
12

Transforming Convolution to MM



13

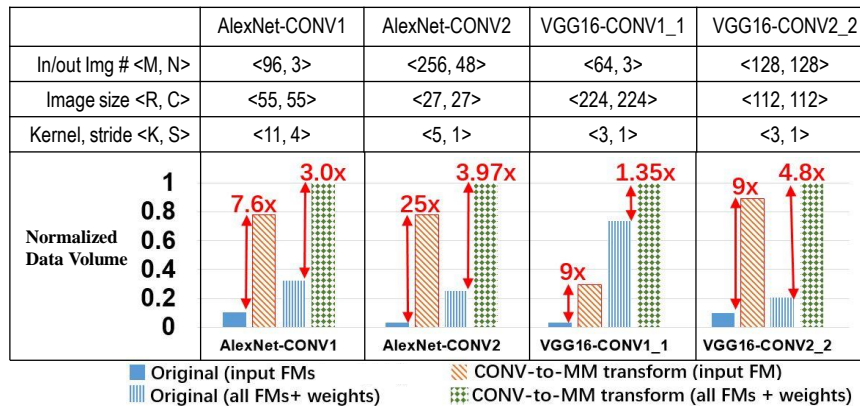
Transforming Convolution to MM



14

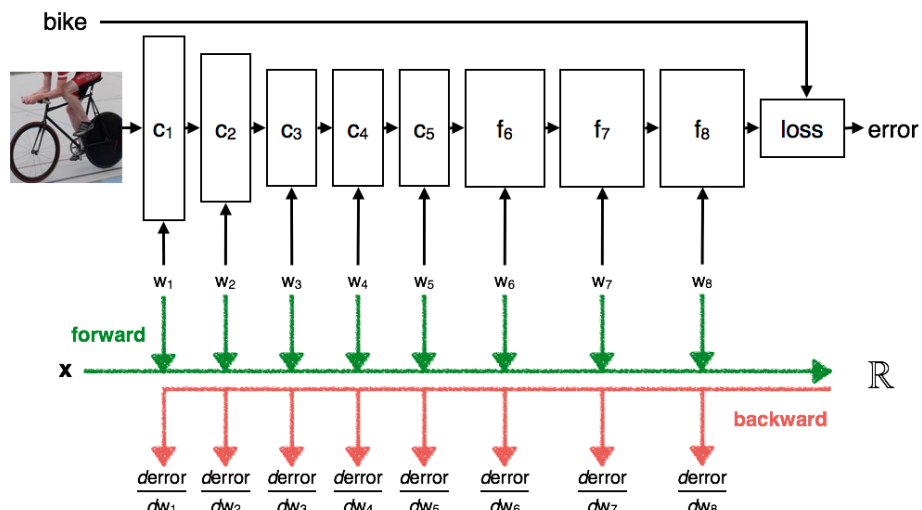
Data Duplications

Data Size	Weight	Input	Output
Convolution	$M \times N \times K \times K$	$\sim N \times (R/S) \times (C/S)$	$M \times R \times C$
Convolution MM	$M \times N \times K \times K$	$N \times K \times K \times R \times C$	$M \times R \times C$
MM/CONV Ratio	1	$\sim (K/S) \times (K/S)$	1



15

CNN Training



16

16

Feedforward Computation on FPGA

1	for(row=0; row<R; row+=Tr) {	(Tile loop)
2	for(col=0; col<C; col+=Tc) {	(Tile loop)
3	for(to=0; to<M; to+=Tm) {	(Tile loop)
4	for(ti=0; ti<N; ti+=Tn) {	(Tile loop)
Off-chip Data Transfer: Memory Access Optimization		
On-chip Data: Computation Optimization		
5	for(trr=row; trr<min(row+Tr, R); trr++) {	(Point loop)
6	for(tcc=col; tcc<min(tcc+Tc, C); tcc++) {	(Point loop)
7	for(too=to; too<min(to+Tm, M); too++) {	(Point loop)
8	for(tii=ti; tii<(ti+Tn, N); tii++) {	(Point loop)
9	for(i=0; i<K; i++) {	(Point loop)
10	for(j=0; j<K; j++) {	(Point loop)
	output_fm[to][row][col] +=	
	weights[to][ti][i][j]*input_fm[ti][S*row+i][S*col+j];	
	}}}}}	
	}}}}}	17
	}}}}	

17

CS 133 Worksheet

- Given that L1 cache is 32KB and L2 cache is 1MB on your machine, please analyze if each of the following set of matrices in your Lab 3 can fit into L1 or L2 cache.
- All input feature maps
 - All output feature maps
 - All weight matrices

18

18

Acknowledgements

- **Some slides on CNN are compiled with help of Chen Zhang, a visiting PhD student from Peking Univ. , now at Microsoft Research Asia (MSRA)**