Acceleration of Cloud Data Centers: Challenges and future directions

Christoforos Kachris, NTUA, InAccel

Dimitrios Soudris, NTUA



Making Software Projects Easier to Understand

- Challenges
 - Heterogeneous computing
 - The need of accelerators
 - Domain specific architectures



Data Center traffic





Data Center applications



Accelerators can increase performance at lower TCO for targeted workloads



Financial market

•In the US alone, NYSE processes more than **1.4 billion messages** every day....

• ...That's 3.5 times more than the number of internet searches Google handles daily.



Number of Messages/Transactions (in Millions)

Christoforos Kachris, Microlab@NTUA



Data Science: need for high computing power





How Big are Data Centers

Data Center Site	Sq ft			
Facebook (Santa Clara)	86,000			
Google (South Carolina)	200,000			
HP (Atlanta)	200,000			
IBM (Colorado)	300,000			
Microsoft (Chicago)	700,000			



[Source: "How Clean is Your Cloud?", Greenpeace 2011]



Wembley Stadium: **172,000** square ft





Google data center





Apple's Data Center





Data Centers Power Consumption

 Data centers consumed 330 Billion KWh in 2007 and is expected to reach 1012 Billion KWh in 2020

	2007 (Billion KWh)	2020 (Billion KWh)	
Data Centers	330	1012	
Telecoms	293	951	
Total Cloud	623	1963	

2007 electricity consumption. Billion kwH





Soon we are going to need a power plant next to the Data Centers



Carbon Emissions of ICT Sector



0003 Answer a typical The CO2 emissions of 10,000 Google Google search KILOWATT-Google, Inc.] HOURS OF ELECTRICITY REQUIRED

searches is equal to a five mile trip in the average U.S. automobile [Source:



End of Growth of single program speed

40 years of Processor Performance





Hardware acceleration

- Hardware acceleration is the use of specialized hardware components to perform some functions faster (10x-100x) than is possible in software running on a more general-purpose CPU.
- Hardware acceleration can be performed either by specialized chips (AICS) or
- By programmable specialized chips (FPGAs) that can be configured for specific applications





Hardware accelerators





Accelerators

A GPU is effective at processing the <u>same set of operations</u> in parallel – single instruction, multiple data (SIMD). A GPU has a well-defined instruction-set, and fixed word sizes – for example single, double, or half-precision integer and floating point values.



An FPGA is effective at processing the <u>same or different operations</u> in parallel – multiple instructions, multiple data (MIMD). An FPGA does not have a predefined instruction-set, or a fixed data width.



Processing Platforms

 HW acceleration can be used to reduce significantly the execution time and the energy consumption of several applications (10x-100x)

The Dilemma: Flexibility vs. Efficiency



Programmable Processing

Source: "High-performance Energy-Efficient Reconfigurable Accelerator Circuits for the Sub-45nm Era" July 2011 by Ram K. Krishnamurthy, Circuits Research Labs, Intel Corp.

Hardware accelerators





CPU: High flexibility, low throughput





Chip Accelerators: High throughput, high NRE Initial one-off costs

Programmable accelerators (FPGAs) can provide massive parallelism and higher efficiency than CPUs for certain categories of applications

Best of 2 worlds

InAccel offers novel accelerators as IP cores to boost the performance of their applications (similar to specialized containers)



FPGAs in the cloud

• Altera's acquisition by Intel



• Microsoft's catapult for Bing search



• IBM open power – CAPI interface with FPGAs





Intel Xeon + FPGAs

Software Development for Accelerating Workloads using Xeon and coherently attached FPGA in-socket





Programming Interface





Catapult FPGA acceleration Card





VINEYARD Framework

- Accelerators stored in an AppStore
- Cloud users request accelerators based on applications requirements
- Decouple Hardware

 Software
 designers



Heterogeneous Data Center



FPGA as a Service

• Amazon EC F1's Xilinx FPGA







Commercialization





Heterogeneous made easy





лоторесиир



4x Lower Cost







Acceleration of Spark

- > Spark is the most widely used framework for Data Analytics
- Develop hardware components as IP cores for widely used applications
 - >> Spark
 - Logistic regression
 - Recommendation
 - K-means
 - Linear regression
 - PageRank
 - Graph computing



Seamless integration with Zero-code changes:

Without inaccel > spark-submit With inaccel > spark-submit --inaccel



Acceleration made easy

Now we're ready to build a pipeline and fit it. This puts the data through all of the feature processing, model tuning & training we described in a single call.

```
In [*]: from pyspark.ml import Pipeline
```

pipeline = Pipeline(stages=[labelIndexer, featuresScaler, cv, indexToLabel])

```
%time model = pipeline.fit(train)
```

▼ A	pache Spark:	1 EXECUTORS	8 CORES	Jobs:	1 RUNNING	6 COMPLETED	
	27 March	16:45					
	38	39		10	41		42
Jobs:					3	4	5

The next cell converts the test set from LibSVM to **Parquet** memory format, in order to serve as our streaming source.























- Future Data Center will have to sustain huge amount of network traffic
- However the power consumption will have to remain almost the same
- FPGA acceleration as a promising solution in providing
 - high throughput,
 - low latency and
 - energy efficient processing





Thank you!

Christoforos Kachris chris@inaccel.com