

Affordable Supercomputing for Data Mining Applications

Computation intensive

TASKS

ARCHITECTURES

Data centers of low-end servers

Many-core

Moore's Law

Web Data:
The ECML/PKDD Discovery Challenge 2010 on Web Quality

Web 3.0

Web Data Federation

Data and Computation intensive tasks → ARCHITECTURE?

Network Algorithms: Bulk Synchronous Parallel

HAMA: incubatory project

Google Pregel clone

Plethora of Solutions?

Conclusions

LAWA: Longitudinal Analytics of Web Archive Data

SCIIMS: Strategic crime and Immigration management System

Image Classification on GPGPU

Feature generation: Histogram of Oriented Gradients

Training images → **HOG descriptors**

Feature generation:

- Histogram of Oriented Gradients (grayscale)
- Results in 200-400 dimensions
- Lowered with PCA

Gaussian Mixture Modeling

GMM soft assignment → **Probability based kernel**

Linear logistic regression

Research data continuous growth in volume

Source	Number of Images
CLEF 2006: IAPR TC-12	20,000
CLEF 2008, PASCAL VOC 2007: MIR Flickr	25,000
CLEF 2011, PASCAL VOC 2010: MIR Flickr	1,000,000

Blockers in Scalability

Number of transactions/second vs **Number of CPUs**

GPU onboard memory

- Global 4-8 GB
- Block shared 10+ KB

„Numbers Everyone Should Know” Jeff Dean, Google

Intra-process communication

- Mutex lock/unlock 100 ns
- Read 1 MB sequentially from network 10,000,000 ns

Disk

- Disk seek 10,000,000 ns
- Read 1 MB sequentially from disk 30,000,000 ns

RAM

- L1 cache reference 0.5 ns
- L2 cache reference 7 ns
- Main memory reference 100 ns
- Read 1 MB sequentially from memory 250,000 ns

Longitudinal Analytics of Web Archive Data in FIRE LAWA

Euromed JUST

SCIIMS

Conclusion

- New hw/sw architectures for scaling
- Interdisciplinary area
- Enables novel apps (e.g. Web 3.0)
- Computation & Data Intensive problems need breakthrough!
- No mature open architecture yet for distributed network analysis (BSP?)