

Acquisition & Analysis of Immense Datasets for e-People

Zhiwei Xu

Institute of Computing Technology (ICT)

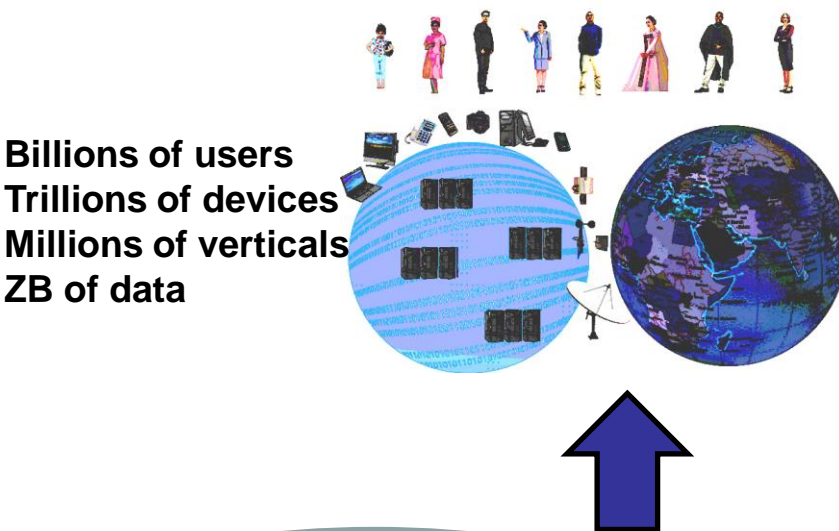
Chinese Academy of Sciences (CAS)

www.ict.ac.cn, zxu@ict.ac.cn

This research is supported in part by the National Basic Research Program of China (Grant 2011CB302502)
and the Strategic Priority Program of Chinese Academy of Sciences (Grant XDA06010400)

A Mega Trend: e-People

- Computing for the Masses
 - IT that directly benefits the masses (billions of individuals), not institutions
 - e-People, not e-Business, e-Science, e-Government
 - CS that utilizes the human-cyber-physical ternary universe
 - e-People is not fully realized if we have to use cyber devices



	2007		2030	
	Per Capita	WW Total	Per Capita	WW Total
Storage	44.7 GB	295 EB	5.23 TB	41.8 ZB
Communication	9.85 GB	65 EB	2.88 TB	23 ZB
GP Computing	1 GIPS	6.39 EIPS	5 TIPS	40 ZIPS
SP Computing	28.6 GIPS	189 EIPS	321 TIPS	2570 ZIPS

2007 data: Hilbert and López, *Science* 2011: 332 (6025), 60-65.

2030 projection: from a conservative estimation by CAS

Institutional Computing
 e-Business
 e-Science
 e-Government

Cyberspace Computing
 IT services
 IT software
 IT hardware

The Chinese Academy of Sciences NICT Project

- New generation ICT
 - 10-year research project (2012-2021)
 - >\$200M for phase one (2012-2016)
 - 19 institutes, over 200 faculty members
 - Aiming at China's needs in 2020-2050
- Human-cyber-physical ternary computing for ZB of data
 - **Functional sensing**
 - **Customizable Internet**
 - **Cloud-sea computing**
 - Billion-thread cloud servers for EB data processing
 - GB-TB terminal devices (human facing)
 - KB-GB sensor nodes (physical world facing)

Potentials for New CS & Gadgets

- Scientific problem
 - Can we realize the EB→ZB transition without increasing energy 1000X?
 - Data assets are different from “money” assets and real estate assets
 - Can I “withdraw” my data from Amazon and deposit them to B&N? (cf. Google’s DLF)
 - Do I own the “smart grid” data from my home?
 - China National Grid is installing 200 million smart meters
- Computer science requirements
 - Personal data asset algebra and normal forms
 - Personal data asset management system
 - XaaS: Everything as a Computer (2020-2030)
 - Cf. Internet of Everything, X as a Service (SaaS, PaaS, IaaS)
 - Home as a Computer, Mobile as a Computer, Building as a Computer
 - TB “smart phones” @2W
 - PB wuTV @20W (home datacenter, physical world facing)
 - Personal Watson (personal intelligent machine) @2000W

A Cloud-Sea Computing Architecture

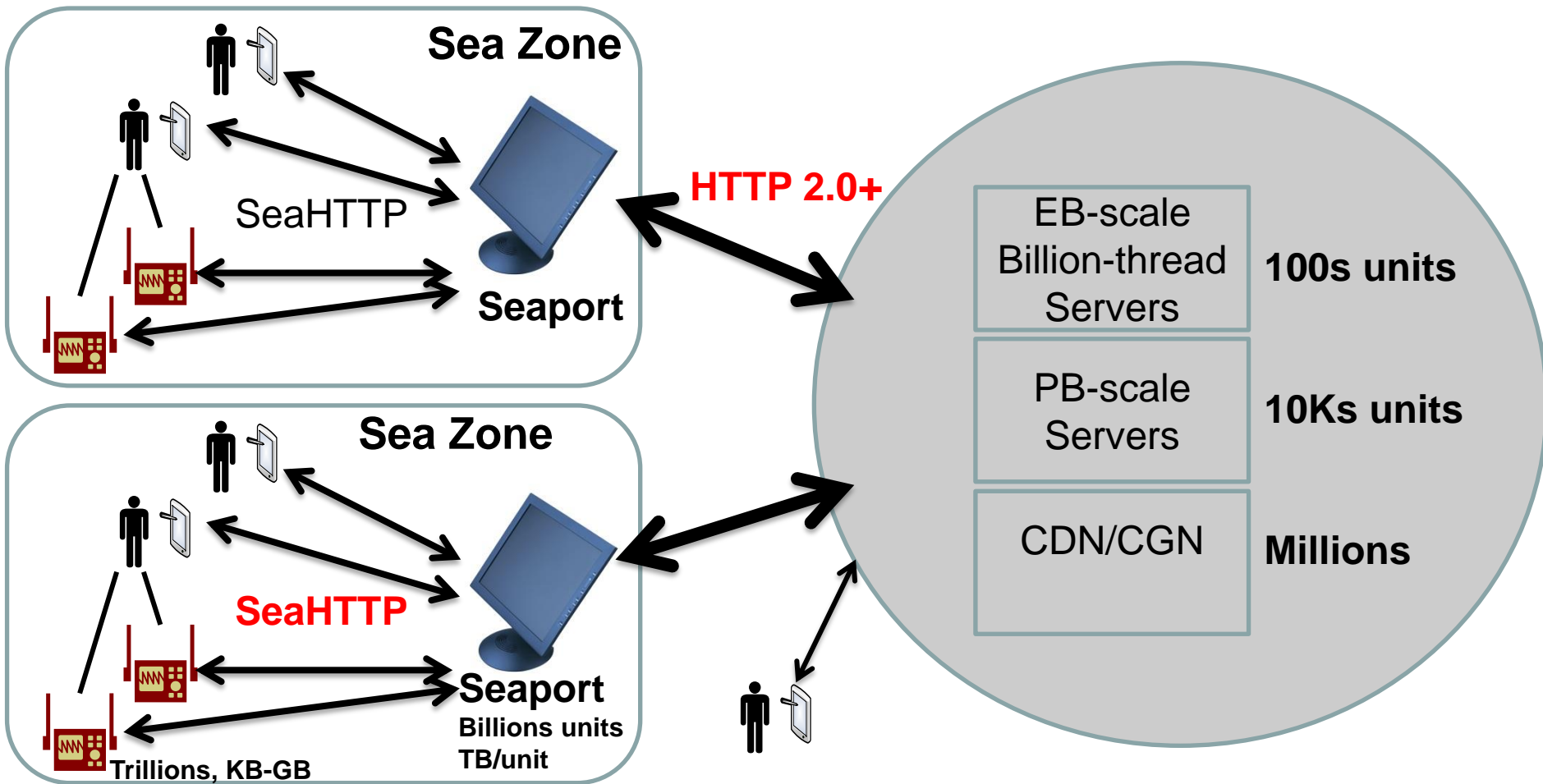
REST 2.0

Sea-side functions

sensing, interaction, local processing

Cloud-side functions

aggregation, request-response, big data



Acquisition and Analysis of Home Appliances Data

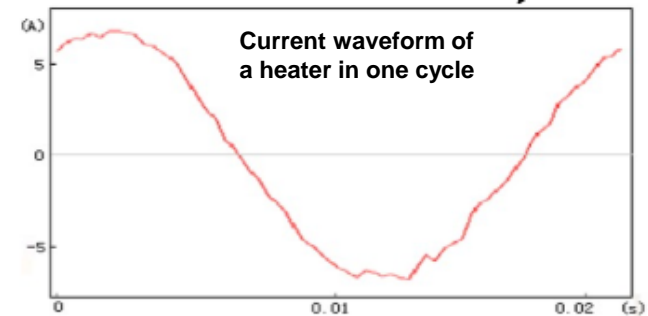
- Application examples (2020-2030)
 - Web search → Grid search
 - “Top 100 green households in Beijing and London”
 - Appliances R&D
 - Utilizing field data for all appliances
- Acquisition challenge
 - Can we timely acquire massive and accurate field data from billions of households, for each and every appliance (lamp, refrigerator, etc.) in every household, with 1-3 sensors per home?
- Analysis challenge
 - Can we timely process and query EB-ZB field data from billions of households, to obtain appliance-specific behavior?

Data Acquisition

- Home is a rich source of personal data
 - Behavior, health, environment, electricity data, etc.

- Electric devices (appliances) data

- ~50 devices per home, 220V@50Hz
- Up to 128th harmonics
 - 256 samples/cycle, 10 bytes/sample
- 6.4 MB/s, or 200TB per year per home
- For China, 200TB x 0.5 billion homes = 100 ZB per year



- Functional sensing of home appliances data

- Function is formalized behavior

- On-off behavior data for each device
- Event behavior data
- Finite behavior data (up to k th harmonics for a given finite k)
- Infinite behavior data

- Data storage needs can be reduced 10,000 times

- 20GB/year per home for aggregated data
- 1TB/year per home for disaggregated data for each device

Data Computing

- Data computing R&D needs workloads

- Close to reality lab data acquisition
- Use Internet services workloads to develop techniques
- Utilize existing ecosystems

- Three examples

- Off-line (back end): RCFile for Hadoop Hive
 - Production use: Facebook, Taobao, Netflix, Twitter, Yahoo, LinkedIn, AOL, Salesforce.com, etc.
- On-line (front end): CCIndex on Hbase
 - Production use in Taobao, Tencent
- High-speed communication: DataMPI

Alexa Top Sites
(2013.06.14)

1. Facebook
2. Google
3. YouTube
4. Yahoo!
- 5. Baidu**
6. Wikipedia
7. Windows Live
8. Twitter
- 9. QQ (Tencent)**
- 10. Taobao**

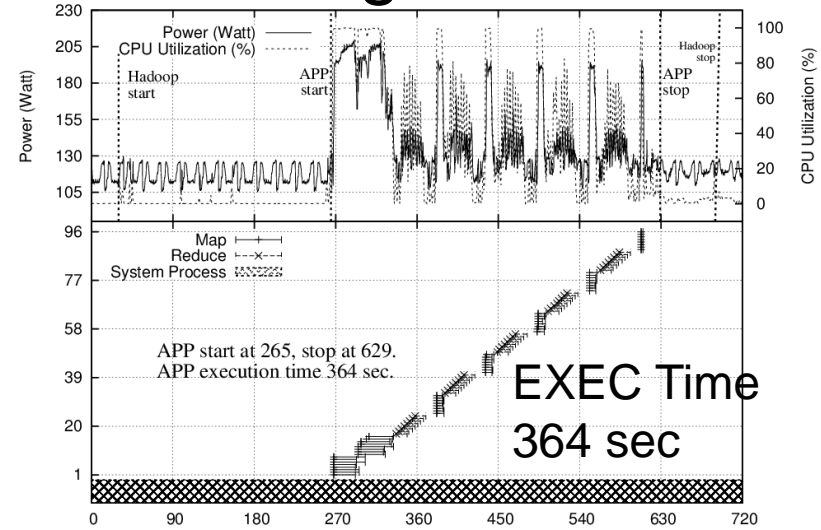
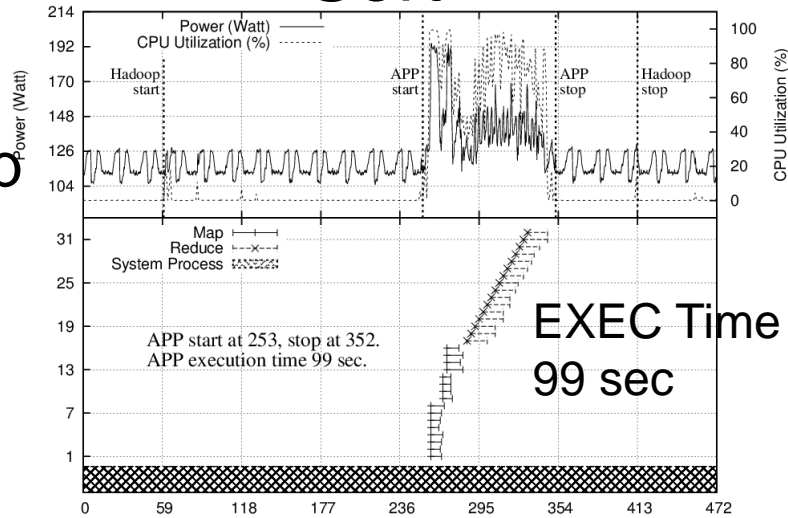
22. eBay

DataMPI open sourced at datampi.org

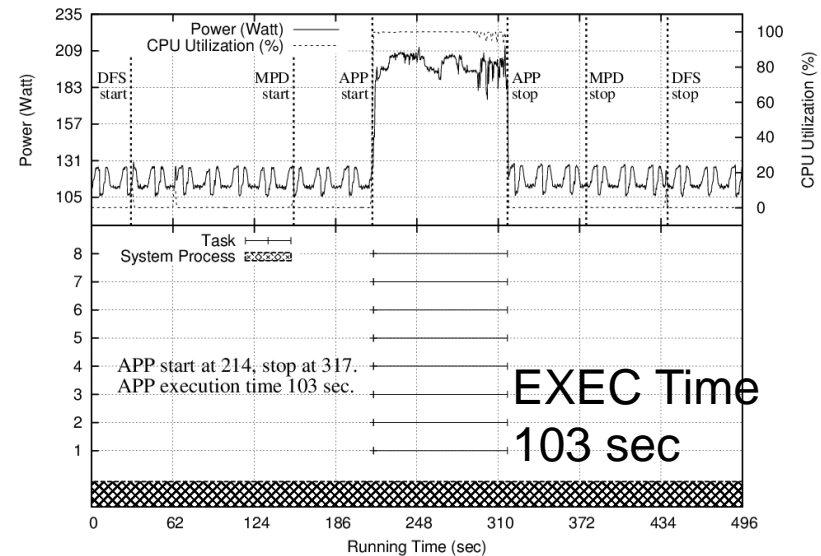
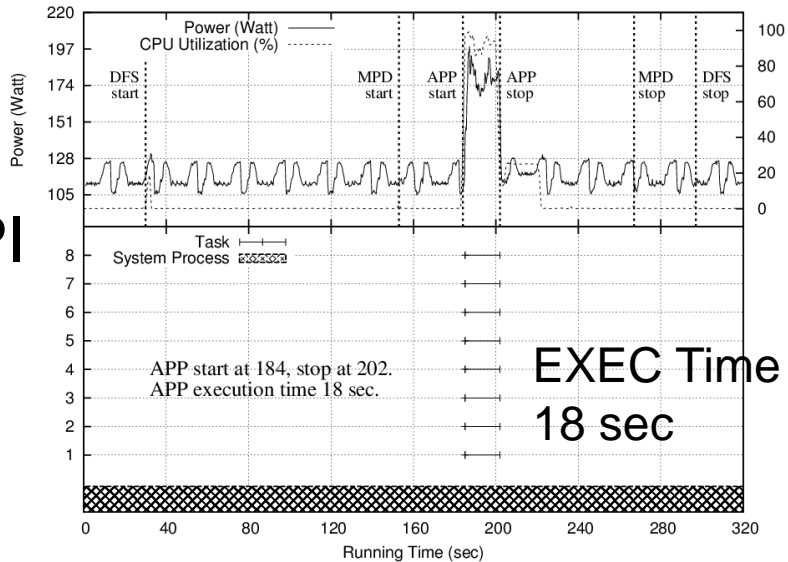
Sort

PageRank

Hadoop



DataMPI



References

- Zhiwei Xu: High-Performance Techniques for Big Data Computing in Internet Services. SC Companion 2012: 1861-1895
- Zhiwei Xu: Measuring Green IT in Society. IEEE Computer 45(5): 83-85 (2012)
- Zhiwei Xu: How much power is needed for a billion-thread high-throughput server? Frontiers of Computer Science 6(4): 339-346 (2012)
- Zhiwei Xu, Guojie Li: Computing for the masses. Commun. ACM 54(10): 129-137 (2011)
- Jingjie Liu, Lei Nie, Zhiwei Xu: The Input-Sensing Problem in Ternary Computing and Its Application in Household Energy-Saving. GreenCom 2011: 131-138
- Yongqiang He, Rubao Lee, Yin Huai, Zheng Shao, Namit Jain, Xiaodong Zhang, Zhiwei Xu: RCFile: A fast and space-efficient data placement structure in MapReduce-based warehouse systems. ICDE 2011: 1199-1208
- Xiaoyi Lu, Bing Wang, Li Zha, Zhiwei Xu: Can MPI Benefit Hadoop and MapReduce Applications? ICPP Workshops 2011: 371-379
- Yongqiang Zou, Jia Liu, Shicai Wang, Li Zha, Zhiwei Xu: CCIndex: A Complementary Clustering Index on Distributed Ordered Tables for Multi-dimensional Range Queries. NPC 2010: 247-261

谢谢!
Thank you!



zxu@ict.ac.cn