

# Datacenter Simulation Methodologies

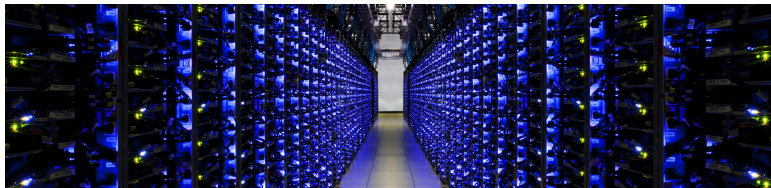
## Introduction

Tamara Silbergleit Lehman, Qiuyun Wang, Seyed Majid Zahedi  
and Benjamin C. Lee



# Future of Computer Architecture

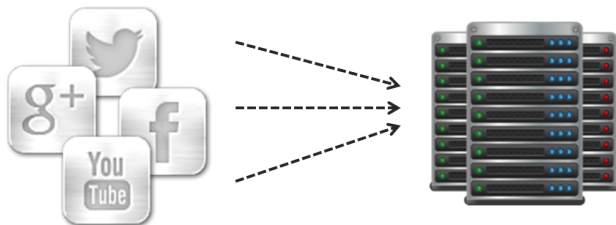
- Methodology supports general-purpose design
- Research bifurcates into mobile, datacenter systems



---

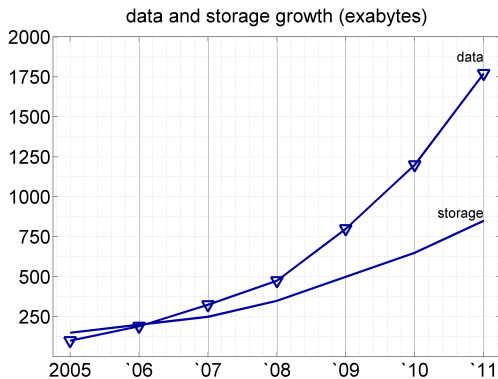
Figures: [Google], [Apple]

# Cloud Computing Applications



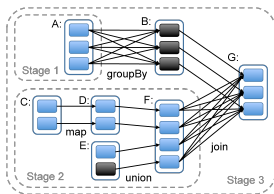
- Data volumes are growing exponentially
- Cloud applications are diversifying rapidly
- Computing capability must grow
- Datacenters dissipate tens of megawatts

# Data Deluge



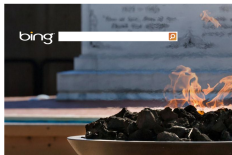
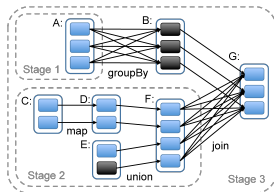
- data  $\uparrow 1.5\times/\text{yr}$
- storage  $\uparrow 1.4\times/\text{yr}$
- datacenter demand  $\uparrow$

Chart: [IDC'08]



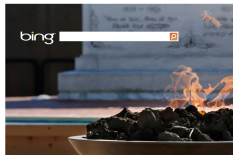
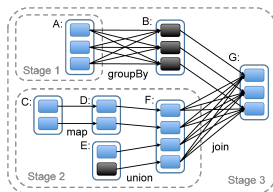
- **Machine Learning**  
MapReduce, Spark,  
GraphX, GraphLab

# Big Data Domains



- **Machine Learning**  
MapReduce, Spark, GraphX, GraphLab
- **Search & Retrieval**  
Bing, Google, Solr

# Big Data Domains



- **Machine Learning**  
MapReduce, Spark, GraphX, GraphLab
- **Search & Retrieval**  
Bing, Google, Solr
- **Cloud Computing**  
EC2, Azure, AppEngine



Sign Up

My Account / Console English

AWS Products & Solutions

AWS Product Information



Developers Support

## Amazon EC2 Details

- EC2 Overview
- EC2 FAQs
- EC2 Pricing
- Amazon EC2 SLA

## Amazon Elastic Compute Cloud (Amazon EC2)

Amazon Elastic Compute Cloud (Amazon EC2) is a web service that provides resizable compute capacity in the cloud. It is designed to make web-scale computing easier for developers.

Amazon EC2's simple web service interface allows you to obtain and

Get Started with  
AWS for Free

Sign Up Now »

AWS Free Trial includes 12M

**Big data demands big computing, yet we face challenges...**

- Architecture Design
- Systems Management
- Research Coordination



## **Pursue performance and energy efficiency**

### **Design processors**

- Microarchitecture (big versus small)
- Heterogeneity (big and small)

### **Design memory systems**

- Technologies (DRAM, PCM, MRAM, etc.)
- Interfaces (DDR<sub>x</sub>, LP-DDR<sub>x</sub>, buffers, etc.)

### **Design communication**

- Inter-processor
- Processor-memory
- Processor-accelerator

## Pursue performance and fairness

### Allocation

- What resources are demanded?
- Analyze utility, preferences

### Scheduling

- When are resources demanded?
- Analyze phase behavior

### Sharing

- Which tasks to co-locate?
- What is each task's share?

## **Anticipate management risk during design**

- Design hardware, manager concurrently

## **Design for manageability**

- Identify hardware to ease allocation
- Organize system to ease scheduling
- Reduce variance in task performance

**By the end of the tutorial, participants will be able to...**

- Deploy a full-system, cycle-accurate simulator
- Simulate processor and memory systems
- Simulate datacenter workload
- Design processor, memory systems for datacenter workloads

## **We describe our practice and experience**

- Describe experience in datacenter research as architects
- Describe strategies for integrating disparate frameworks
- Present a coherent methodology

## **We provide breadth and highlight existing frameworks**

- Draw on related tutorials on simulators, applications
- Specify minimum steps for end-to-end experiments
- Refer to other materials for depth

Time	Topic
09:00 - 09:30	Introduction
<b>09:30 - 10:30</b>	<b>Setting up MARSSx86 and DRAMSim2</b>
10:30 - 11:00	Break
11:00 - 12:00	Spark simulation
12:00 - 13:00	Lunch
13:00 - 13:30	Spark continued
13:30 - 14:30	GraphLab simulation
14:30 - 15:00	Break
15:00 - 16:15	Web search simulation
16:15 - 17:00	Case studies