





#### Hadoop: A Reality Check







#### **Database Harddrive**



time to find one record  $log_{100}(100,000,000) * 10ms = 40ms$ time to read record 10,000,000 \* 50ms

- $= \log_b N * 10ms$
- = 10ms
  - = **5.8 days**

#### **B-Tree**



#### **Hadoop Harddrive**



throughput time to transfer record 10,000,000 \* 10ms random reads

- = 10MB/s
- = 10ms
- = **1.5 days**
- = (5.8 days)



# Laws of Physics

#### Random (log Scale) Sequential



Adam Jacobs The Pathologies of Big Data

# Hadoop File System



- Files are split into blocks
- Blocks are stored on different servers
- Blocks are replicated

# **Hadoop Computation**



- Map Reduce
- Files are divided into splits
- Each server compute a split
- Data is grouped
- Aggregation from groups

# Hadoop a brief history

#### Nutch: The Free Search Alternative to Google

#### Stefan Krempl 10.06.2004

An open deployable search algorithm is set to allow webmasters to launch their own search engines and transparency into the maturing business

#### Open-Source-Suchmaschine will in die

🕕 uorlesen / MP3-Download

03.09.2003 15:09

Das Projekt Nutch will eine Alternative zu den bekannten I



## **Commits per Month**



Time

### **EMails per Month**





#### Weaknesses





# Strength



#### **Unstructured Data Growth**

**Unstructured: 61.7% growth** 

Structured: 21.8 % growth

**Enterprise data doubles every three years (Forrester)** 

### What is Hadoop really about?





### What is Hadoop really about?



# We're hiring!



# (Lessions Learned)







DB(b-tree/index)

© Datameer, Inc 2010

23

# STORE RAW DATA









#### Iterate fast with known technologies



### Implementation Effort (month)



### Hadoop Cost



#### **Small Data / Latency**



9

#### Pull vs Push

Very slow Local buffer = risk of lost data Monitor many agents Complicated



Simple Pull as often as required Just one system to monitor Easier to secure

### Working with Streams == push





#### **Just Store Structured Data**

#### DBVisualizer



Hive/Pig



Hadoop







![](_page_33_Picture_0.jpeg)

## **Application Design**

#### JSP (View)

Spring (Controller)

Hibernate (Model)

![](_page_34_Picture_4.jpeg)

Database

Data Vis (View)

Flows (Controller)

Cascading Tuples (Model)

![](_page_34_Picture_9.jpeg)

Hadoop

## **Application Design**

![](_page_35_Figure_1.jpeg)

![](_page_35_Picture_2.jpeg)

## **Application Design**

![](_page_36_Figure_1.jpeg)

![](_page_37_Picture_0.jpeg)

#### Pentenho

![](_page_37_Picture_2.jpeg)

Database

![](_page_37_Picture_4.jpeg)

![](_page_37_Picture_5.jpeg)

Hadoop

![](_page_37_Picture_7.jpeg)

![](_page_38_Picture_0.jpeg)

![](_page_38_Picture_1.jpeg)

![](_page_39_Picture_0.jpeg)

www.datameer.com

sg AT datameer.com

@datameer

![](_page_39_Picture_4.jpeg)