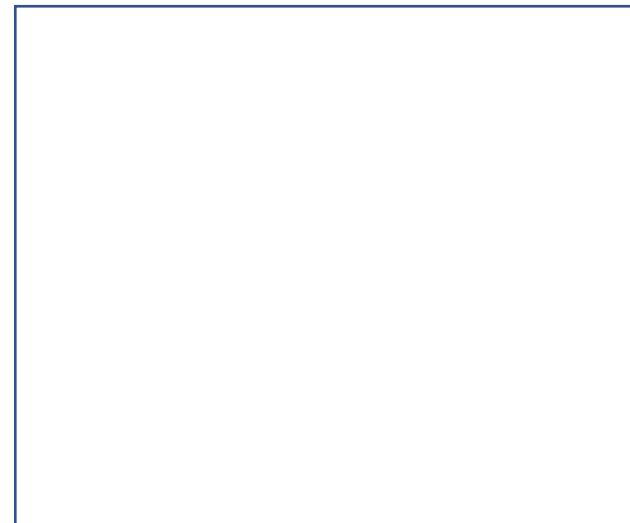
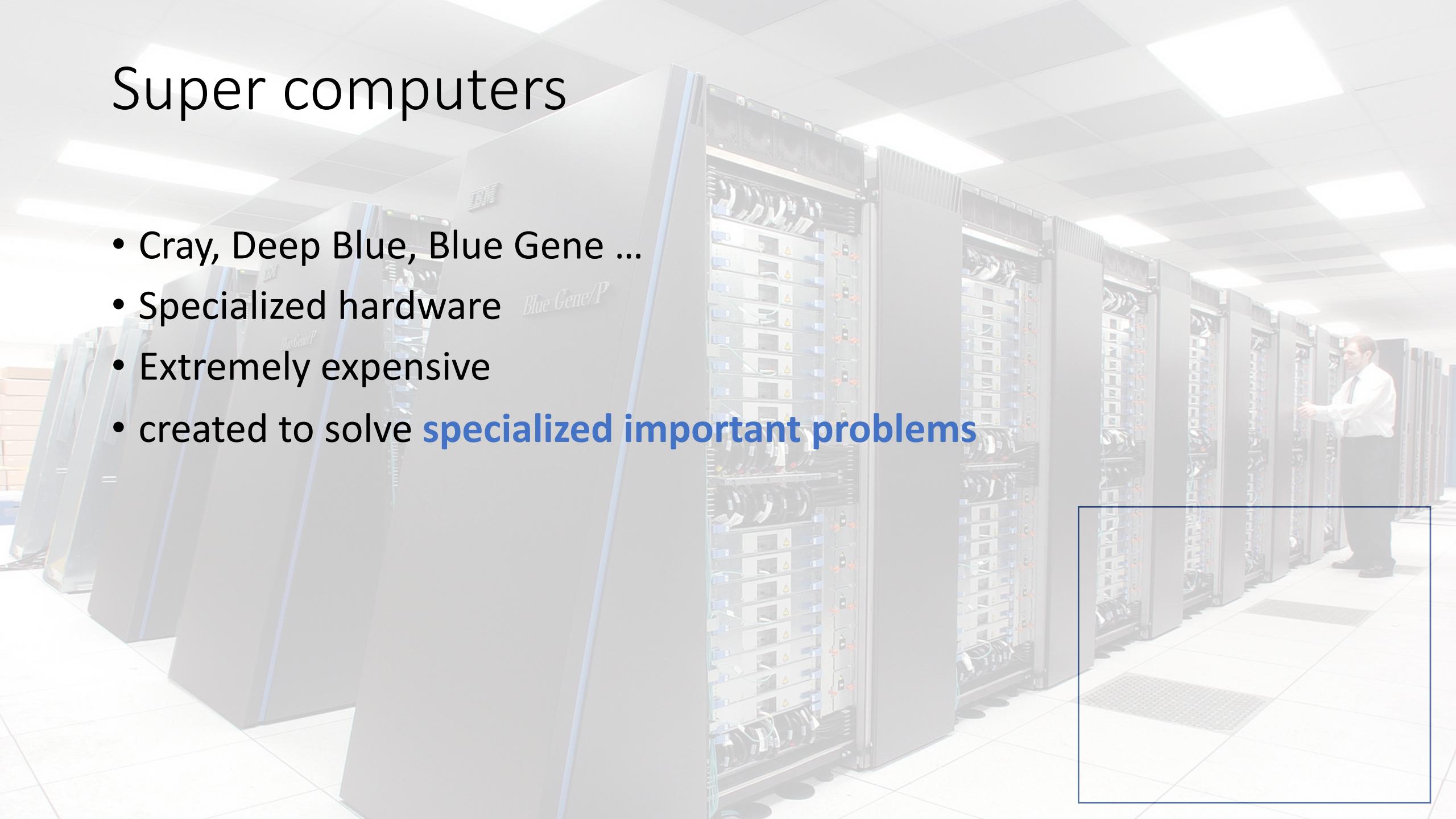


A short history of affordable  
massive computing.



# Super computers

- Cray, Deep Blue, Blue Gene ...
- Specialized hardware
- Extremely expensive
- created to solve **specialized important problems**

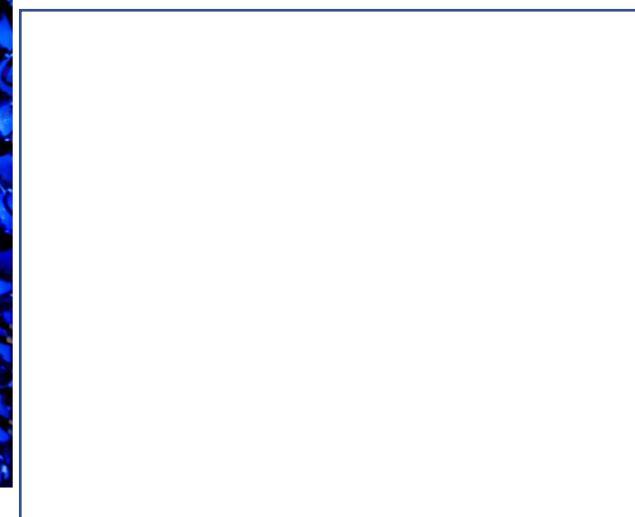
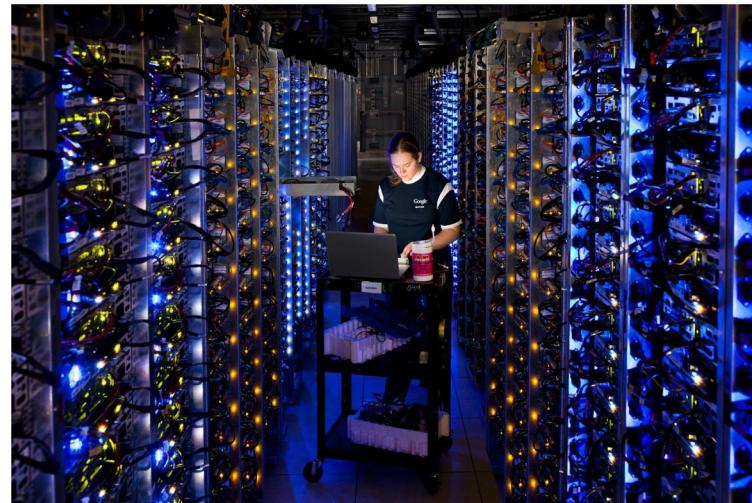


# Data Centers



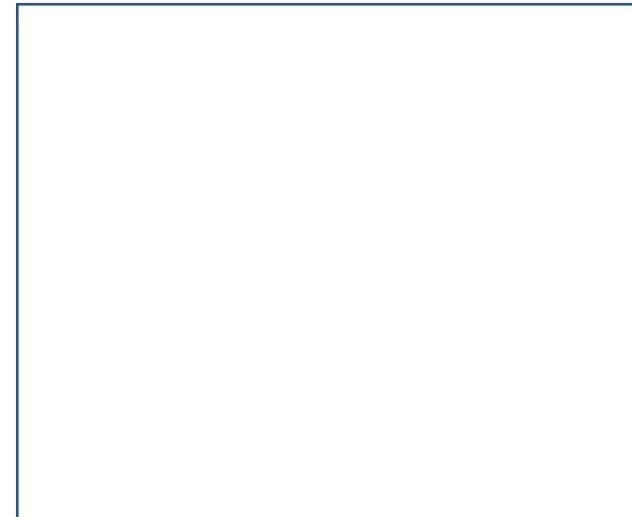
# Data Centers

- The physical aspect of "the cloud"
- Collection of commodity computers
- VAST number of computers (100,000's)
- Created to provide computation for large and small organizations.
- Computation as a commodity.

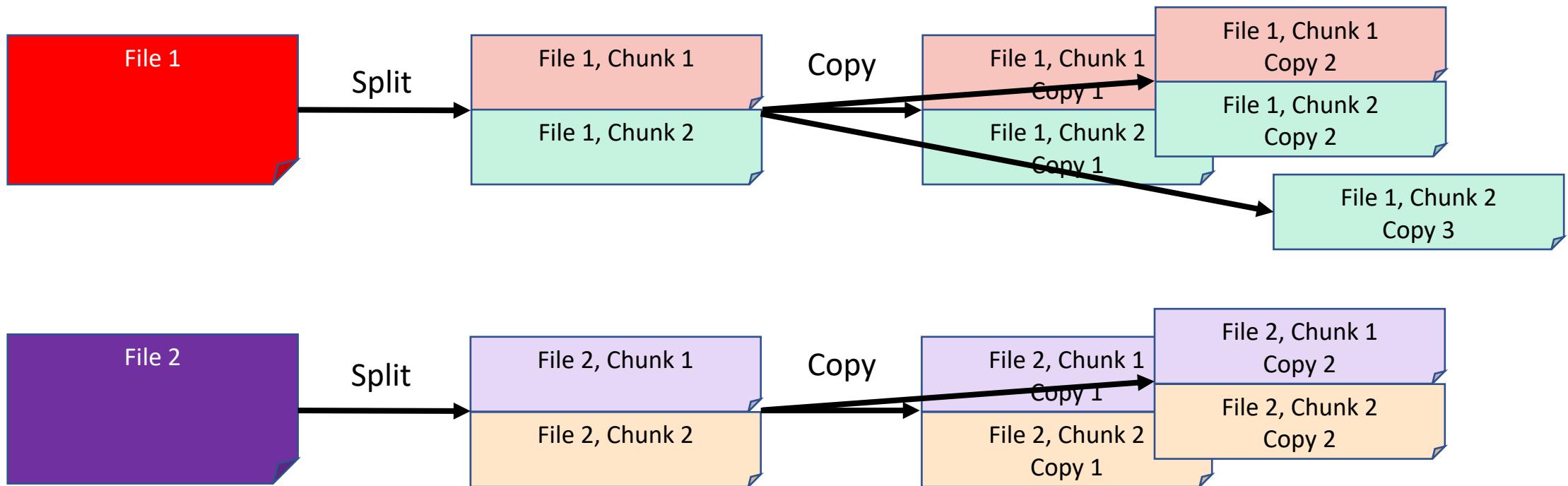


# Making History: Google 2003

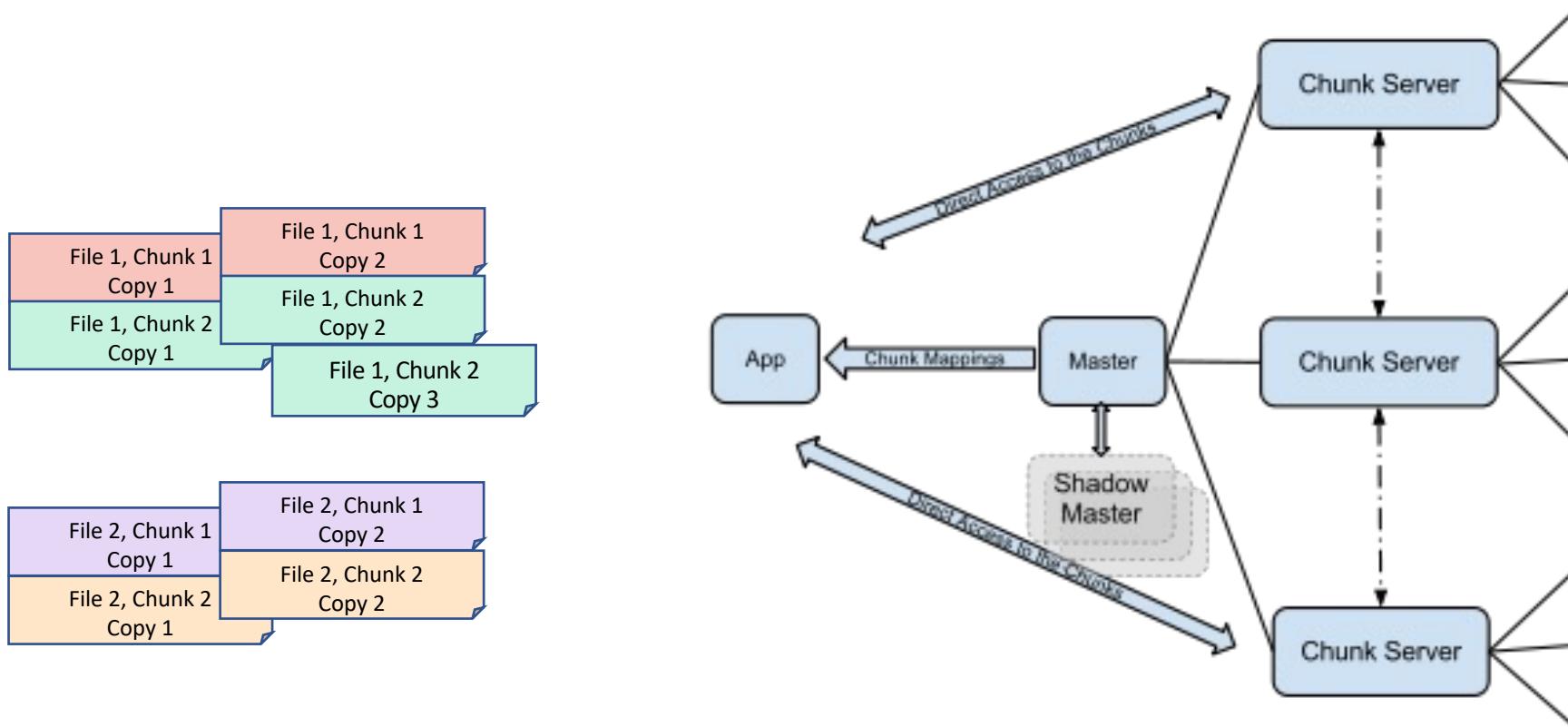
- Larry Page and Sergey Brin develop a method for storing very large files on multiple **commodity** computers.
- Each file is broken into fixed-size **chunks**.
- Each chunk is stored on multiple **chunk servers**.
- The locations of the chunks is managed by the **master**



# HDFS: Chunking files

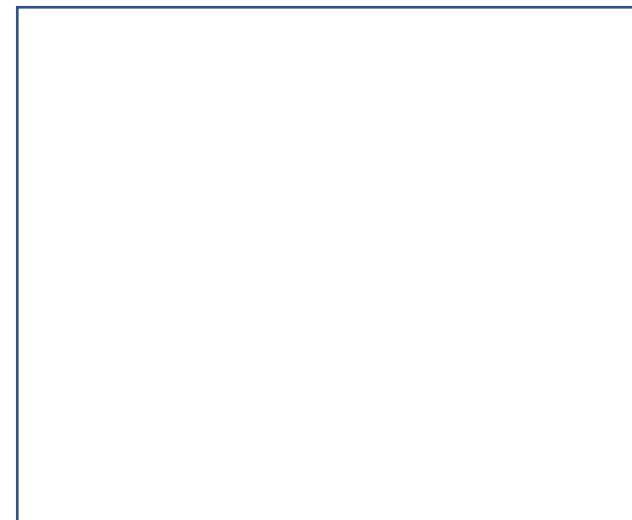


# HDFS: Distributing Chunks

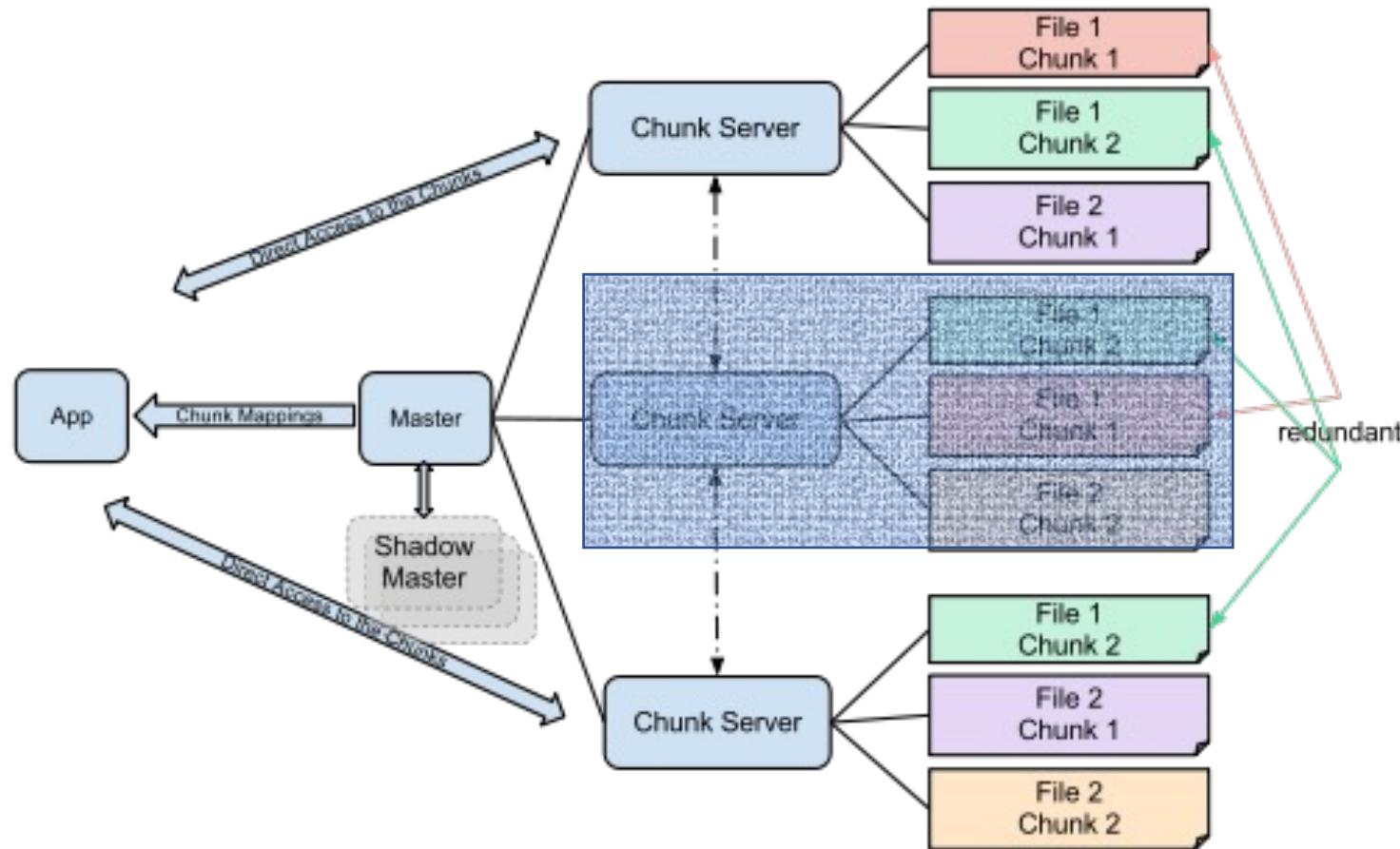


# Properties of GFS/HDFS

- **Commodity Hardware:** Low cost per byte of storage.
- **Locality:** data stored close to CPU.
- **Redundancy:** can recover from server failures.
- **Simple abstraction:** looks to user like standard file system (files, directories, etc.) Chunk mechanism is hidden.

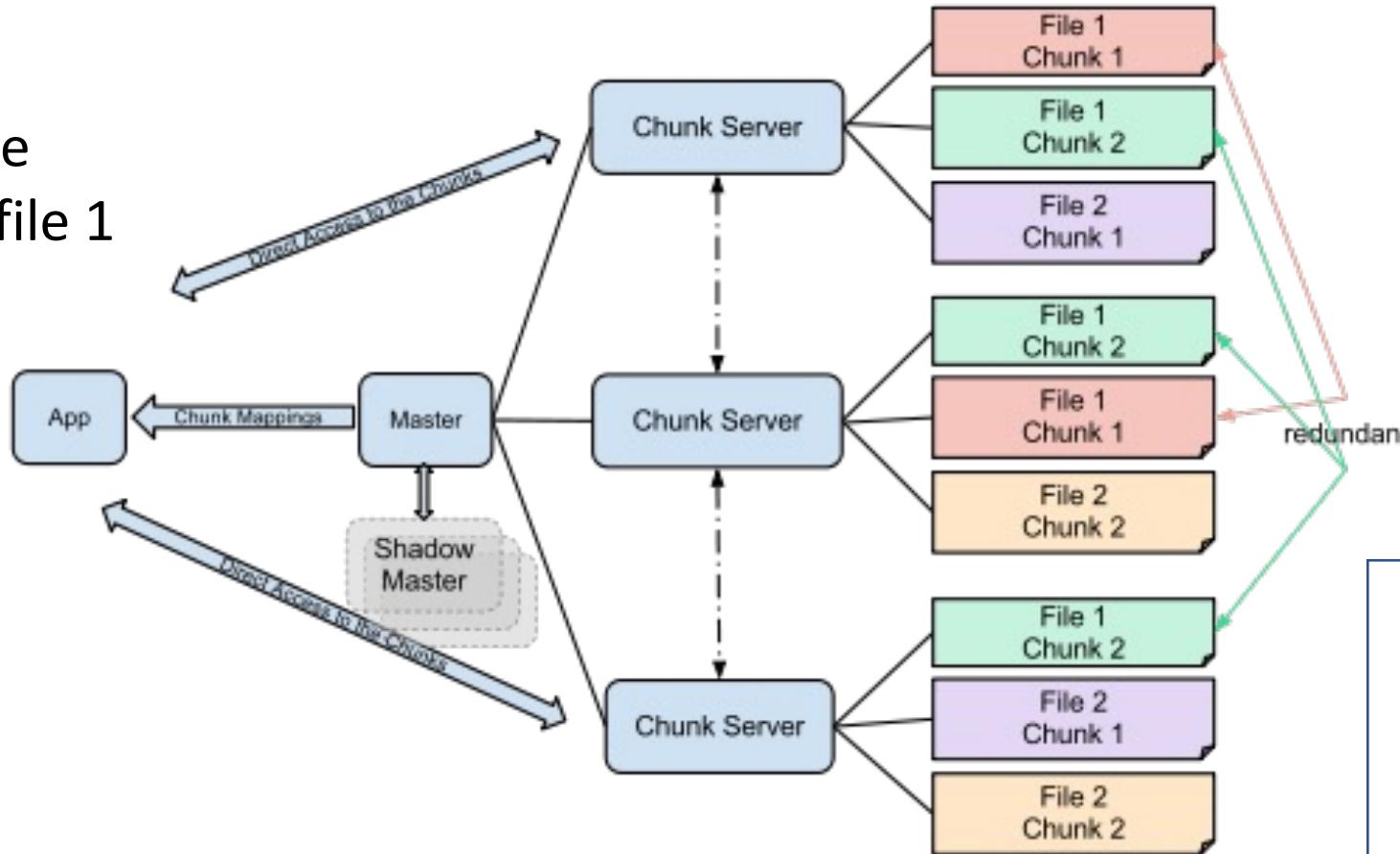


# Redundancy



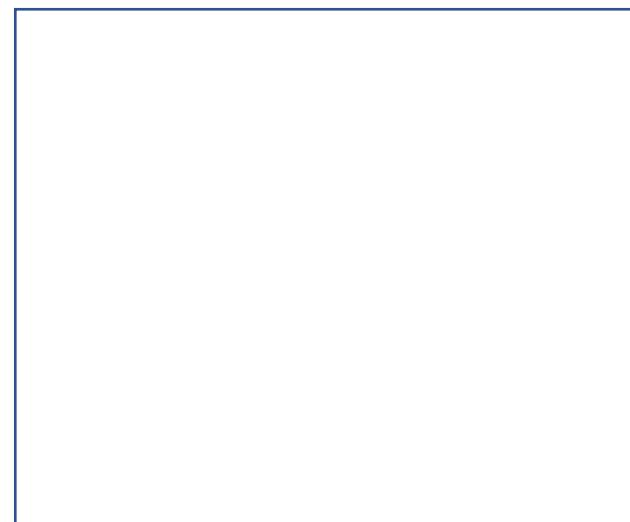
# Locality

Task:  
Sum all of the  
elements in file 1



# Map-Reduce

- HDFS is a **storage abstraction**
- **Map-Reduce** is a **computation abstraction** that works well with HDFS
- Allows programmer to specify parallel computation without knowing how the hardware is organized.
- We will describe Map-Reduce, using Spark, in a later section.



# Spark

- Developed by Matei Zaharia , amplab, 2014
  - Hadoop uses shared **file system** (disk)
  - Spark uses shared **memory** – faster, lower latency.
  - Will be used in this course
- 
- Recall word count by sorting,  
we will redo it using map-reduce!

# The Cloud

- The common name for data centers.
- What is better? Cloud or your local computers?
  - Cloud vs. Local: Rent vs. own: if we want a lot of power for a short time, it is cheaper to rent.
  - Centralized IT: shared staff, shared maintenance, shared upgrade.
  - Storage:
    - Long term – cloud storage (multiple TB) much more expensive than local.
    - Moving TB to/from cloud slow / expensive / physical (snowball)
  - Like a huge supermarket, there are many choices and it is not easy to find the best combination.

# Summary

- Big data analysis is performed on large clusters of commodity computers. – computation as a service.
- HDFS (Hadoop file system): break down files to chunks, make copies, distribute randomly.
- Hadoop Map-Reduce: a computation abstraction that works well with HDFS
- Spark: Sharing **memory** instead of sharing **disk**.