

# GFS (SOSP '03) <sup>MR OSDI '04</sup>

## Goal

Shared : M/R HTML

Performance .

Fault tolerance .

1 failure / year \* machine

300TB: 1000 machines  $\Rightarrow$  3 failure / day

## GFS Approach .

• FS-like API . (POSIX) mount .

"lib"  $\left\{ \begin{array}{l} \text{write / append} \\ \text{read} \end{array} \right.$

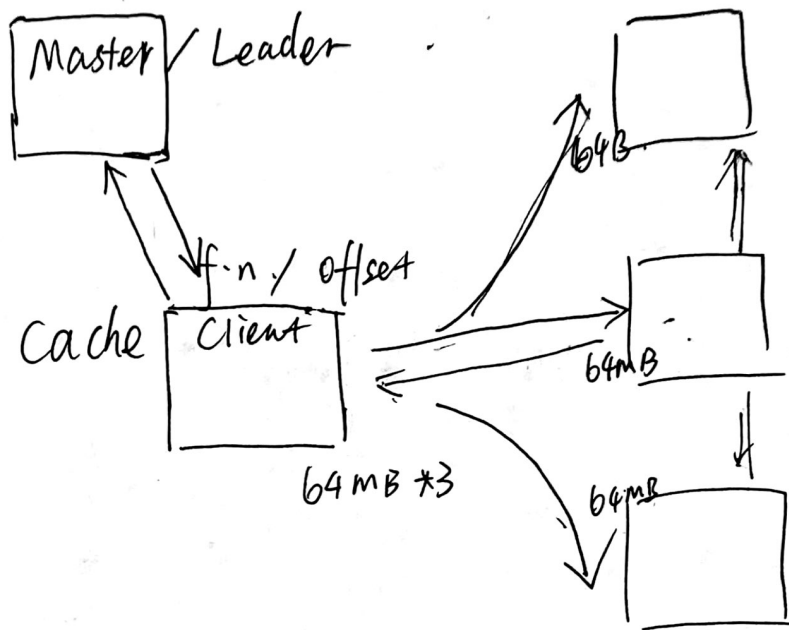
• Single Master . / Metadata servers .

filename  $\rightarrow$  chunks  
offset

•  $2^2$  64 MB  $\frac{116000}{\text{US}}$   $2^2$  4 KB  $\left. \begin{array}{l} \text{page size} \\ \text{block size disk} \end{array} \right\}$

• 3-way Replication .

# How GFS works



## Performance

Why single master + 64 MB sufficient?  
large files. sequential read/write

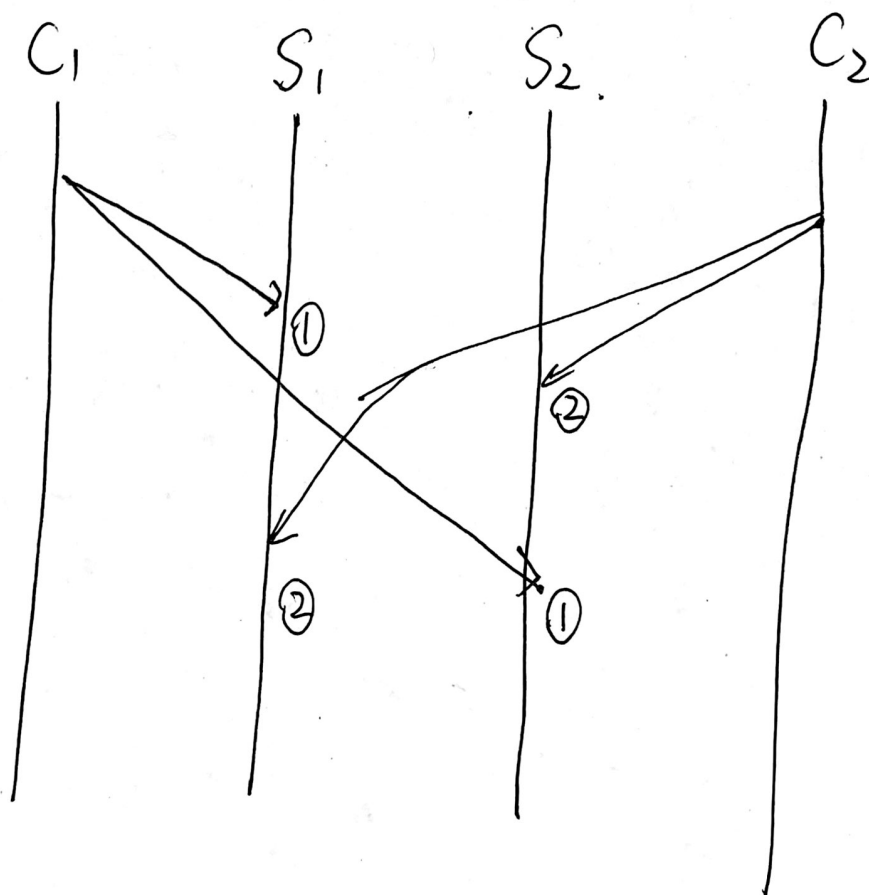
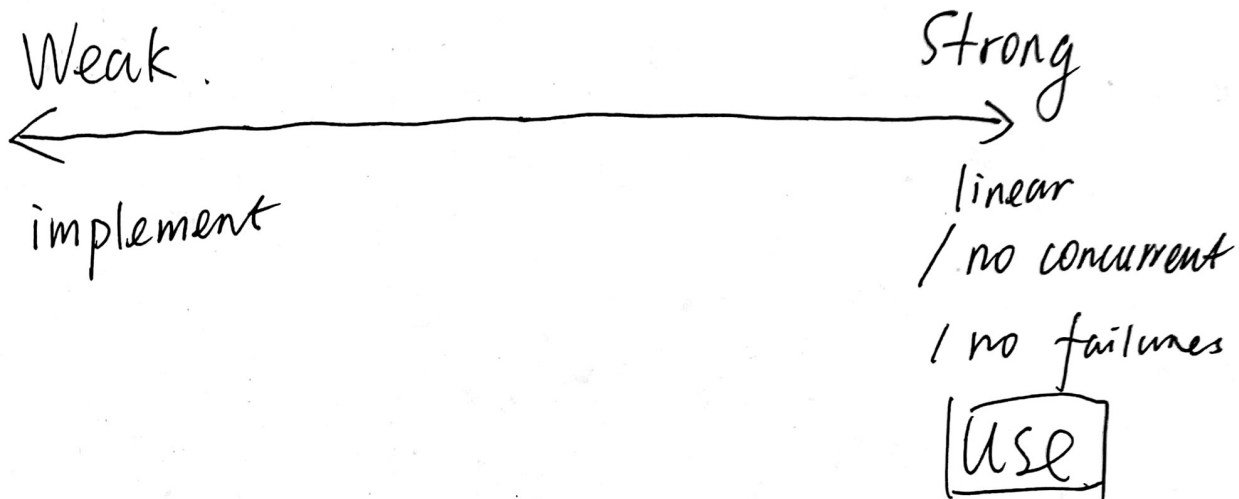
Not good for  
small files. → aggregate

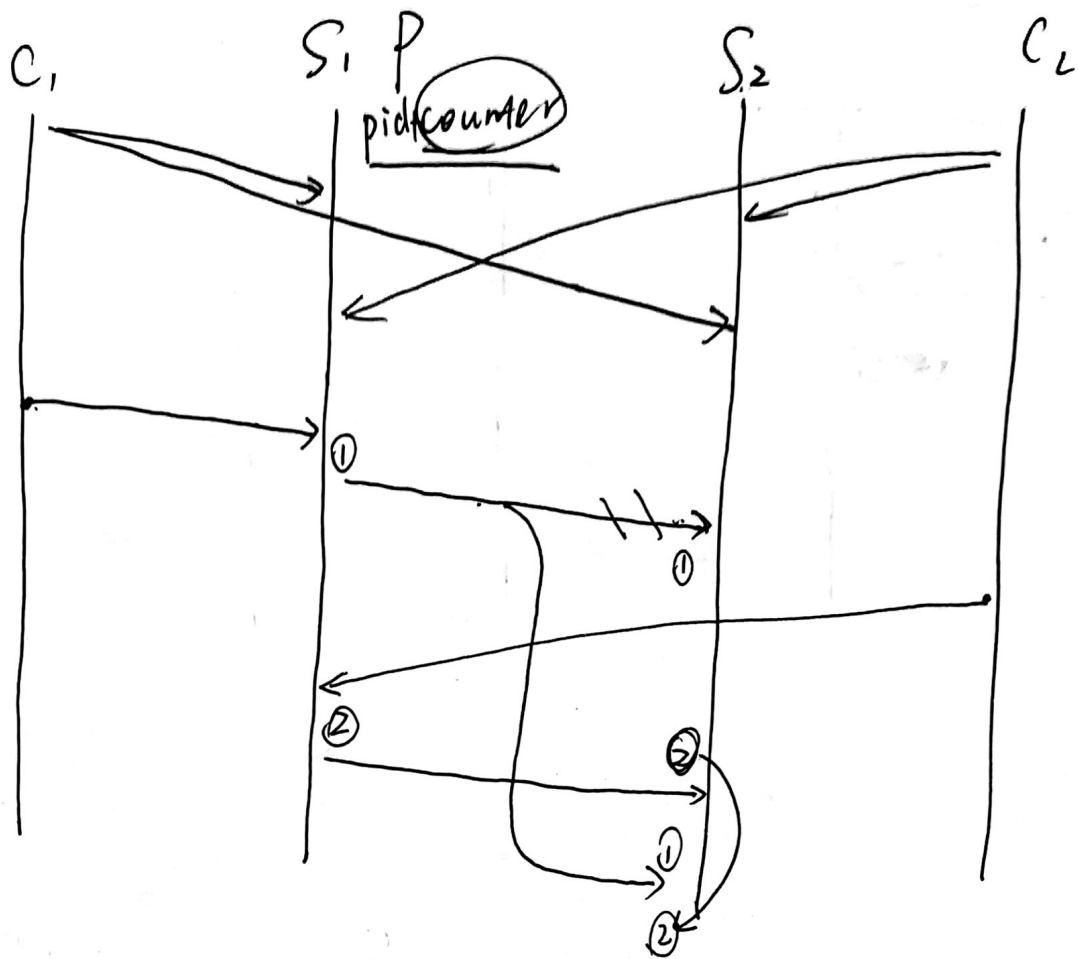
random read/write      git  
   compiler

# Consistency.

correctness : outcome = expectation

concurrent failures.





W1 W2.  $R^3(1)$   $R^3(2)$   $R^3(1)$

