



What do we need if we want secure, durable, consistent, and efficient access to mutable files?



The slide is taken from "Maintenance-Free Global Data Storage," by S. Rhea, C. Wells, P. Eaton, D. Geels, B. Zhao, H. Weatherspoon, and J. Kubiatowicz, published in IEEE Internet Computing, September-October 2001. Other OceanStore papers can be found at http://oceanstore.cs.berkeley.edu/publications/index.html. It is also covered in the textbook starting on page 422.











Pond: the OceanStore Prototype, Sean Rhea, Patrick Eaton, Dennis Geels, Hakim Weatherspoon, Ben Zhao, and John Kubiatowicz. Appears in Proceedings of the 2nd USENIX Conference on File and Storage Technologies (FAST '03), March 2003



This and the next two slides are from the first lecture. They illustrate (roughly) the logical layout of files in OceanStore. The actual file layout is a form of B-tree, but such details aren't important here.





A file is a sequence of versions. One modifies a file by creating a new version. Blocks not modified in creating a new version are shared with the old version. In principle, versions last forever.







GUID = globally unique identifier. Each link is a GUID. The intent is that GUIDs are published in Tapestry.





Note that there are two copies of each block.

i is the number of unavailable copies

This is, for all possible numbers of unavailable copies that still allow reading,

the product of (# ways to arrange unavailable copies in unavailable servers) x (# of ways to arrange available copies in available servers) / # of ways to arrange all copies on all servers

For this particular case, the formula would be simpler if you consider the complement case: what is the probability that NO blocks are available? However, we use this formula as it is the same as the one in the next slide.

This and the next two slides are based on material from "Erasure Coding vs. Replication: A Quantitative Comparison," by H. Weatherspoon and J. Kubiatowicz, published in *Proceedings of the First International Workshop on Peer-to-Peer Systems (IPTPS 2002)*, March 2002. A copy may be found at http://oceanstore.cs.berkeley.edu/publications/papers/pdf/erasure_iptps.pdf.















Here we pause to look at options for replicating a mutable file



With group replication, updates may be sent to any member of the group (where the group here is all computers). This is convenient, except that getting all computers to agree on the order of updates could be a problem. We will study approaches for doing this, but they don't scale to hundreds of thousands of computers.



With master replication, updates are sent just to one computer, which then propagates them to all the others, determining the effective order. However, now the system has a bottleneck.





With this approach, the order of updates is determined by the master group, which then propagates them to the others.





The primary replica's function is too important to be trusted to a single server. Thus all its actions must be agreed upon by a set of servers known as the *inner ring*.



Copies of the current-version information are found via the file's AGUID, which serves to identify the file.



The servers in a file's inner ring must agree on all updates to the file, and hence must agree on its current-version information. Getting agreement among a group of servers, some of which may be (maliciously) faulty, is known as the Byzantine Generals Problem. We take this up in a later lecture. However, we will see that it is solvable if more than two-thirds of the participants do what they're supposed to do.



The current-version information agreed upon by the inner ring is called the *heartbeat* and is represented as a security certificate, signed by at least one more than two-thirds of the inner-ring servers. As we will discuss in a later lecture, clients, upon getting a copy of such a certificate, can verify that it really was signed by such super majority. In addition to the latest VGUID, the signed heartbeat contains the AGUID, a timestamp, and a version sequence number.









This slide taken from the original presentation of "Pond: the OceanStore Prototype", Sean Rhea, Patrick Eaton, Dennis Geels, Hakim Weatherspoon, Ben Zhao, and John Kubiatowicz. Appears in Proceedings of the 2nd USENIX Conference on File and Storage Technologies (FAST '03), March 2003,





Performance					
	LAN		WAN		Predominant
Phase	Linux NFS	Pond	Linux NFS	Pond	benchmark
1	0.0	1.9	0.9	2.8	Read and write
2	0.3	11.0	9.4	16.8	Read and write
3	1.1	1.8	8.3	1.8	Read [status only]
4	0.5	1.5	6.9	1.5	Read
5	2.6	21.0	21.5	32.0	Read and write
Total	4.5	37.2	47.0	54.9	
8			IV-39	Copyright © 2015 Thomas W. Doeppner and Rodrigo Fonseca.	

Performance vs NFS

Writes dominated by cryptographic operations (group signing), erasure coding

This is a killer in the local area

Latency in the wide area hides most of this overhead

In conclusion, it is feasible!

This slide is taken from the textbook (Coulouris, Dollimore, and Kindberg, 4th edition) and shows the performance of Pond vs. NFS. The text of the figure, also taken from the textbook, is: "The figures show times in seconds to run different phases of the Andrew benchmark. It has five phases: (1) creates subdirectories recursively; (2) copies a source tree; (3) examines the status of all the files in the tree without examining their data; (4) examines every byte of data in all the files; and (5) compiles and links the files."