



ΠΑΝΕΠΙΣΤΗΜΙΟ ΚΡΗΤΗΣ
UNIVERSITY OF CRETE

HY590.45

Modern Topics in Scalable Storage Systems

Kostas Magoutis

magoutis@csd.uoc.gr

<http://www.csd.uoc.gr/~hy590-45>

Απαιτήσεις του Μαθήματος

- Project (50%)
 - Μελέτη και πειραματική αποτίμηση συστήματος/ων
 - Ατομική ανάθεση
 - Περιοδικές αναφορές προόδου
- Δύο παρουσιάσεις σχετικής δουλειάς (15%)
 - Επιλογή από πρόσφατα συνέδρια (όπως FAST, SOSF, κλπ.)
- Ένα γραπτό κουίζ, Απρίλιος 2021 (5%)
 - Θέματα από ένα paper, θα ανακοινωθεί πριν
- Τελική εξέταση (30%)
 - Επιλογή τριών papers, θα ανακοινωθούν πριν

Syllabus (tentative)

Select project topic (by 3/3, prepare brief proposal, see instructions on site)

Select related work papers (by 10/3), recent publications @ [FAST](#), [SOSP](#), [OSDI](#), etc.

Syllabus

Date	Topic	Readings
Mon 15/2	Course overview	-
Wed 17/2	Background	See recommended readings
Mon 22/2	Extending file systems over the network	Sandberg: Design and Implementation of the Sun Network Filesystem
Wed 24/2	NFS (contd.)	Macklem: Not Quite NFS, Soft Cache Consistency for NFS
Mon 1/3	Distributed coordination	Lampert: Paxos made simple
Wed 3/3	Paxos (contd.)	-
Mon 8/3	Distributed virtual disks	Lee: Petal: Distributed Virtual Disks
Wed 10/3	Petal (contd.)	-
Mon 15/3	Distributed file systems I	Thekkath: Frangipani: A Scalable Distributed File System
Wed 18/3	Frangipani (contd.)	-
Mon 22/3	Distributed file systems II	Ghemawat: The Google File System
Wed 24/3	Google file system (contd.)	-
Mon 29/3	Related work presentations I	-
Wed 31/3	Related work presentations I	-
Mon 5/4	Application-specific storage systems	Saito: Manageability, Availability and Performance in Porcupine: A Highly Scalable, Cluster-based Mail Service
Wed 7/4	Porcupine (contd.)	-
Mon 12/4	Structured data	Chang: A Distributed Storage System for Structured Data
Wed 14/4	BigTable (contd.)	-
Mon 19/4	In-class quiz	-
Wed 21/4	Project status updates	-
Mon 26/4 - Fri 7/5	Easter recess	-
Mon 10/5	Related work presentations II	-
Wed 12/5	Related work presentations II	-
Mon 17/5	Distributed transactions	Corbett: Spanner: Google's Globally-Distributed Database
Wed 19/5	Spanner (contd.)	-

Final report & presentation dates in June, exact date TBA

Course themes

- Fundamentals: Organization, metadata, journaling, Paxos
- Distributed file systems: NFS
- Scalable virtual disks: Petal
- Shared-disk distributed file systems: Frangipani
- Separate data from metadata: Google file system
- Application-specific scalable storage: Porcupine
- Scalable storage of structured data: BigTable
- Scalable transactions: Spanner

Scalable Storage Systems: Goals, Requirements

- Expandability
 - Increase system size (capacity) as needed
- Performance
 - Increase linearly with system size
- Availability
 - Survive failures gracefully
- Manageability
 - React to changes automatically

Concepts

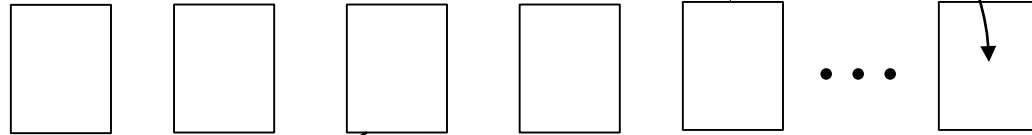
API, semantics

Data model

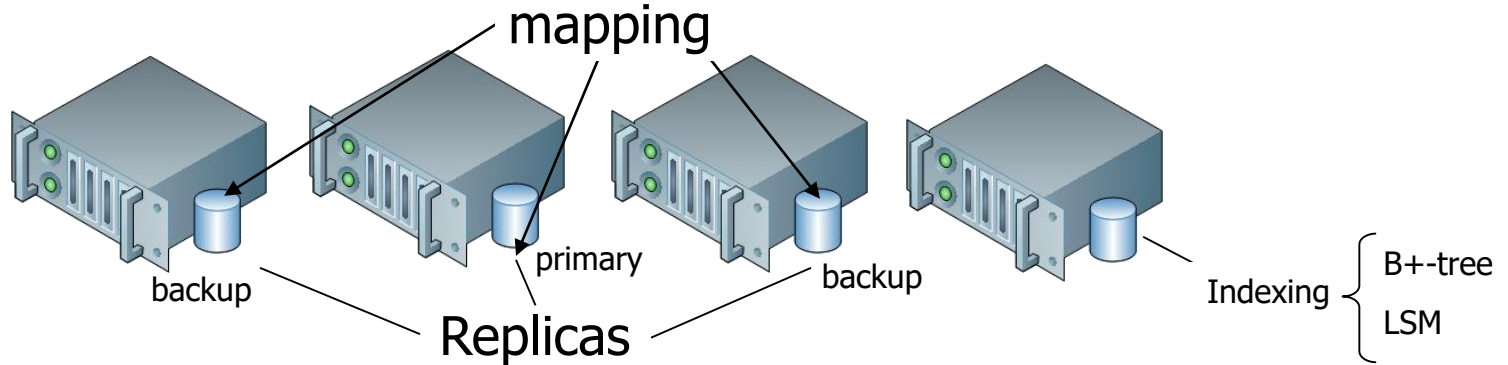
	col-A	col-B	col-Foo	col-XYZ	foobar
row-1					
row-10					
row-18	A18 - v1	B18 - v3	Foo18 - v1	XYZ18 - v2	foobar18 - v1
row-2					
row-5					
row-6					
row-7					

mapping

Horizontal partitions
(shards)



Servers



Storage device & networking technologies

Leverage large-scale infrastructures,
address challenges in doing so

