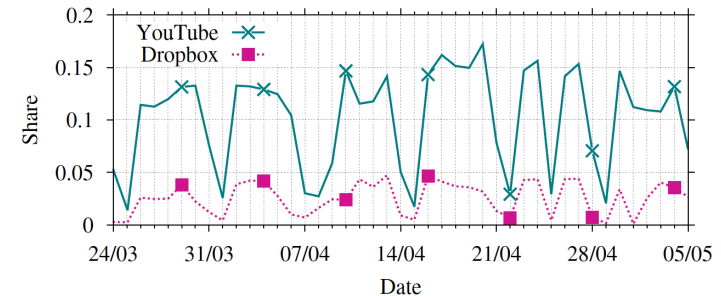


[D+13] Drago et al., "Benchmarking Personal Cloud Storage," *Proc. of the 13th ACM SIGCOMM Conf. on Internet Measurement (IMC '13)*, 2013

[D+12] Drago et al., "Inside DropBox: Understanding Personal Cloud Storage Services," *Proc. of the 12th ACM SIGCOMM Conf. on Internet Measurement (IMC '12)*, 2012

## Motivation

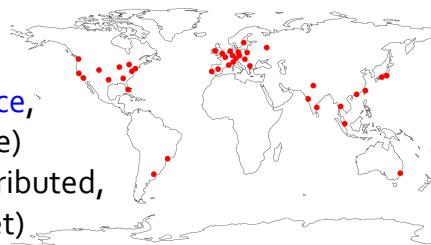
Personal cloud storage is gaining share of Internet traffic, e.g., at one European university campus, DropBox traffic accounted for 4% of traffic or about  $\frac{1}{3}$  of YouTube traffic [D+12]



## Goals of Study

Investigate the **performance improvements** employed by various personal cloud service providers to **synchronize clients' files** and their **effectiveness**

Biggest determining factor in performance is **client-storage distance**, some providers (Google) are geographically distributed, most others are not (yet)



## Providers Studied

**DropBox**: most popular service, established 1997

- control traffic goes to DropBox's data centers
- storage provided by Amazon EC2 and S3

**Google Drive**: public launch in April 2012

**Microsoft SkyDrive**: public launch in April 2012

**LaCie Wuala**: does client-side encryption

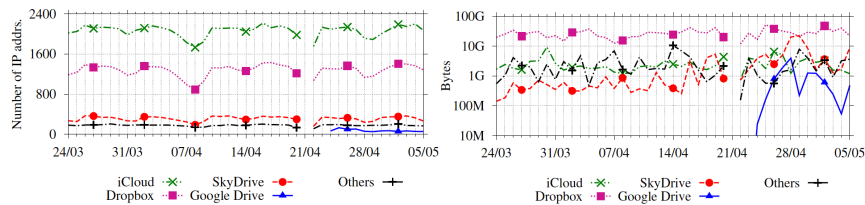
**Amazon Cloud Drive**: included because like DropBox, it relies on Amazon Web Services (AWS)

## Providers Not Studied [D+12]

**Apple iCloud:** could have more clients than DropBox

- but doesn't carry as much traffic
- and doesn't allow clients to store arbitrary files

Others: **SugarSync**, **Box**, **Ubuntu One**: not as popular



## Performance Improvements

**Chunking:** split large content into fixed size data units

- unit of deduplication, delta encoding, and compression

**Bundling:** send multiple small files as one chunk

**Deduplication:** avoid sending chunks already stored on servers

- eliminates duplication in transmission and storage (across users?)
- only DropBox and LaCie implement deduplication, **even for previously deleted files**

**Delta encoding:** send only the diff of old and new chunks

**Data compression,** per chunk

## System Traffic

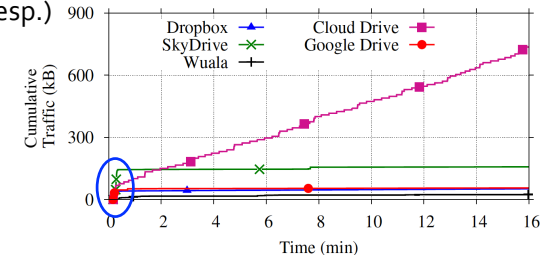
All providers require **client authentication**

- Microsoft uses **4x more traffic** than others, to contact 13 Microsoft Live servers (Why?)

All providers periodically poll server for update:

- LaCie: every 5 mins (generating 60 bps)
- Google: 40 secs interval (42 bps)
- DropBox and Microsoft: 1 min interval

- **Amazon:** once every 15 secs (6 kbps)



## Chunking

Simplifies fault recovery: allows for partial retransmissions, but each chunk is delimited by a pause, introducing delay

- Amazon doesn't do chunking
- Google uses 8 MB chunks
- DropBox uses 4 MB chunks [D+12]:
  - each treated as an independent object
  - identified by a SHA256 value, part of a file's meta data
  - each device keeps a database of meta-data info
  - > 40% of flows have at least 2 chunks
- Microsoft and LaCie uses variable-size chunks

# Bundling

Only DropBox implements bundling, starting April 2012, improving throughput dramatically (by 65%), but each chunk is still sent sequentially [D+12]

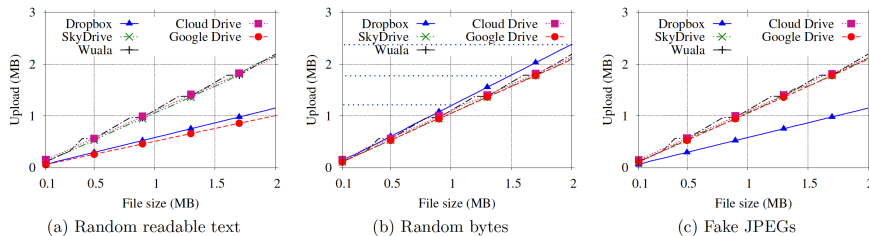
	Mar/Apr		Jun/Jul	
	Median	Average	Median	Average
Flow size				
Store	16.28kB	3.91MB	42.36kB	4.35MB
Retrieve	42.20kB	8.57MB	70.69kB	9.36MB
Throughput (kbits/s)				
Store	31.59	358.17	81.82	552.92
Retrieve	57.72	782.99	109.92	1293.72

Google and Amazon open a separate TCP/SSL connection for each file (as did HTTP 1.0)

Microsoft and LaCie reuse TCP connections, but files are sent sequentially, waiting for application-layer ACK for each file

# Data Compression

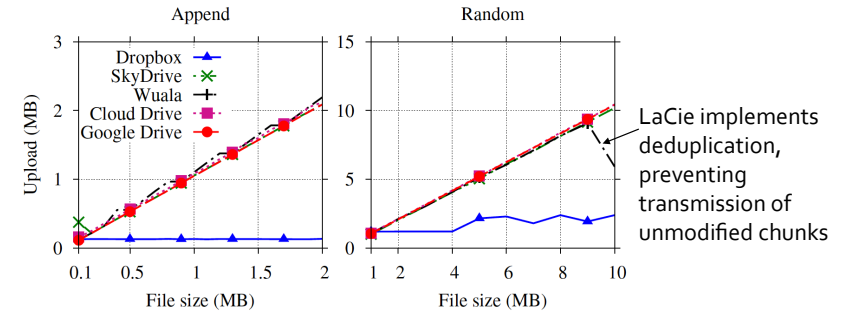
- a) Only DropBox and Google implement compression
- b) "Compression" of already compressed file could result in **larger file**
- c) Google checks file extension and magic number (in file header) before compression



# Delta Encoding

Only DropBox implements delta encoding

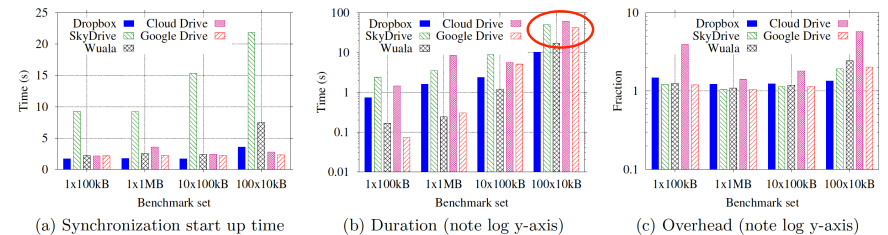
- appending up to file size of 2 MB results in only the addition being sent
- random addition in the middle of large files causes data to shift across chunks, resulting in more data to be sent (delta-encoding is done at chunk granularity)



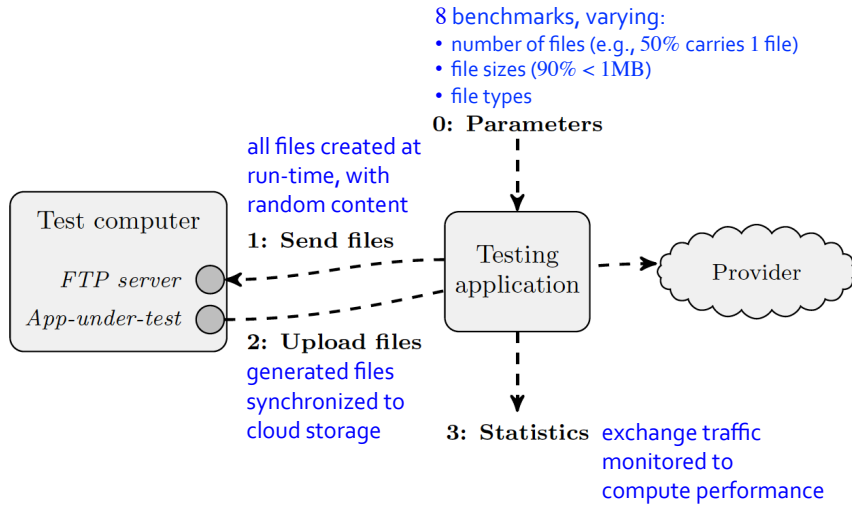
# Overhead and Completion Time

Sending 1 MB of data as one 1 MB file, 10 100 KB files, and 100 10 KB files

- a) Microsoft has the highest overhead, for no clear reason
- b) bundling reduces completion time for small files, **encryption** doesn't seem to affect it
- c) in all cases overhead is higher than data size!



# Benchmarking Methodology



# DropBox Usage [D+12]

Datasets overview 3/24/12-5/5/12

Name	Type	IP Addrs.	Vol. (GB)
Campus 1	Wired	400	5,320
Campus 2	Wired/Wireless	2,528	55,054
Home 1	FTTH/ADSL	18,785	509,909
Home 2	ADSL	13,723	301,448

DropBox users tend to download more than upload, with download/upload ratio:

- Campus 2: 2.4
- Campus 1: 1.6
- Home (Residential ISP) 1: 1.4
- Home (Residential ISP) 2: 0.9

# DropBox Usage [D+12]

Group IP addresses according to behavior:

- **occasional** users: upload and download < 10KB
- **upload/download only**: upload/download 3 orders of magnitude > in the other direction (1 GB vs. 1 MB)
- **heavy** users: all other active users
- **idle** users: client running, no file exchanged (30%)

Upload-only:

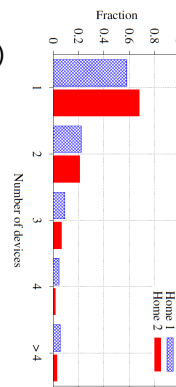
- 7% of IP addresses
- 21% of Home1 transfer volume, 11% of Home2

Download-only:

- 26% of Home1 IPs, 28% of Home2
- 25% of Home1 transfer volume, 28% of Home2

Heavy users, households have multiple devices

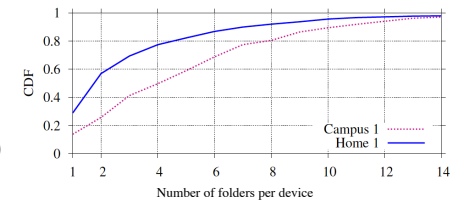
- 37% of Home1 IPs, 33% of Home2
- 50% of DropBox sessions are from heavy users



# DropBox Usage [D+12]

Shared folders: to what extent DropBox is used for content sharing

- Campus1: 13% has 1 folder (Home1: 28%)
- 50% has more than 5 folders (Home1: 23%)



Usage follows the usual daily and weekly patterns

Sessions can last up to 4 hours

Only a small percentage of direct link downloads is bigger than 10 MB, i.e., not one-click hosting movies or archives