

Inside Dropbox: Understanding Personal Cloud Storage Services

- *Idilio Drago*
- Marco Mellia
- Maurizio M. Munafò
- Anna Sperotto
- Ramin Sadre
- Aiko Pras

- **Personal cloud storage** services are gaining popularity
- Dropbox:
 - *“the largest deployed networked file system in history”*
 - *“over 50 million users – one billion files every 48 hours”*

- **Personal cloud storage** services are gaining popularity
- Dropbox:
 - *“the largest deployed networked file system in history”*
 - *“over 50 million users – one billion files every 48 hours”*
- **Little public information about the system**
 - How does Dropbox work?
 - How is the system performing?
 - Are there typical usage scenarios?

How does Dropbox work?

- Some public information about Dropbox
 - **Native client**, Web interface, LAN-Sync, etc
 - Storage in **Amazon S3**
 - Files are split in **chunks of up to 4MB**
 - Delta encoding, **encrypted** communication

How does Dropbox work?

- Some public information about Dropbox
 - **Native client**, Web interface, LAN-Sync, etc
 - Storage in **Amazon S3**
 - Files are split in **chunks of up to 4MB**
 - Delta encoding, **encrypted** communication
- To understand the client protocol
 - **MITM** against our own client
 - Squid proxy, SSL-bump and a self-signed CA certificate
 - **Replace a trusted CA certificate in the heap at run-time**
- Proxy logs and decrypted packet traces

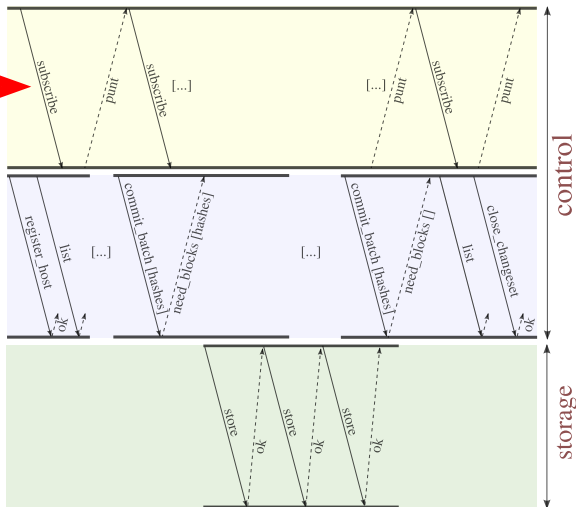
How does Dropbox (v1.2.52) work?

- Clear separation between **storage** and meta-data/client **control**
- Sub-domains identifying parts of the service
 - `notifyX` → client notification
 - `client-lb/clientX` → meta-data control
 - `dl-clientX` → Client storage
 - etc
- **HTTP/HTTPs** in all functionalities

How does Dropbox (v1.2.52) work?

■ Notification

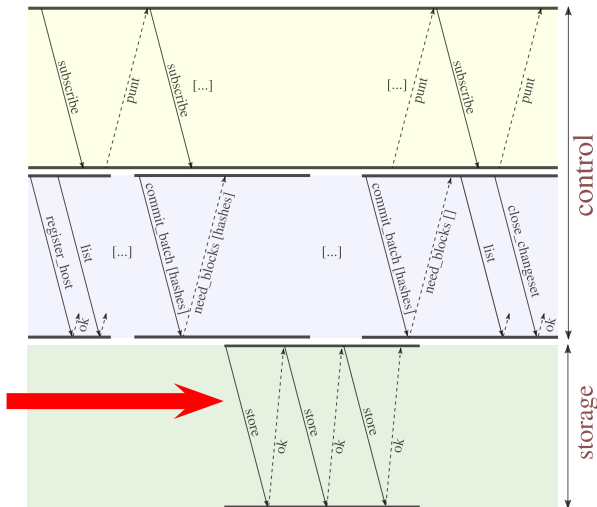
- Kept open
- Not encrypted**
- Device ID**
- Namespaces



How does Dropbox (v1.2.52) work?

■ Storage

- Retrieve vs. Store
- Sequential acks
- Amazon EC2



- Rely on **Tstat**¹ to export **layer-4 flows**
- Isolate Dropbox flows
 - **DN-Hunter**², TSL/SSL certificates, IP addresses

1 – <http://tstat.polito.it/>

2 – DNS to the rescue: Discerning Content and Services in a Tangled Web

- Rely on **Tstat**¹ to export **layer-4 flows**
- Isolate Dropbox flows
 - **DN-Hunter**², TSL/SSL certificates, IP addresses
- **Device IDs** and namespace IDs
- Use the knowledge from our own decrypted flows to
 - **Tag Dropbox flows** – *e.g.* as *store* and *retrieve* flows
 - Estimate the **number of chunks** in a flow

1 – <http://tstat.polito.it/>

2 – DNS to the rescue: Discerning Content and Services in a Tangled Web

	Type	IP Addr.	Dropbox		
			Flows	Vol. (GB)	Devices
Campus 1	Wired	400	167,189	146	283
Campus 2	Wired/Wireless	2,528	1,902,824	1,814	6,609
Home 1	FTTH/ADSL	18,785	1,438,369	1,153	3,350
Home 2	ADSL	13,723	693,086	506	1,313
Total			4,204,666	3,624	11,561

- 42 consecutive days in March and April 2012
 - 4 vantage points in Europe

	Type	IP Addr.	Dropbox		
			Flows	Vol. (GB)	Devices
Campus 1	Wired	400	167,189	146	283
Campus 2	Wired/Wireless	2,528	1,902,824	1,814	6,609
Home 1	FTTH/ADSL	18,785	1,438,369	1,153	3,350
Home 2	ADSL	13,723	693,086	506	1,313
Total			4,204,666	3,624	11,561

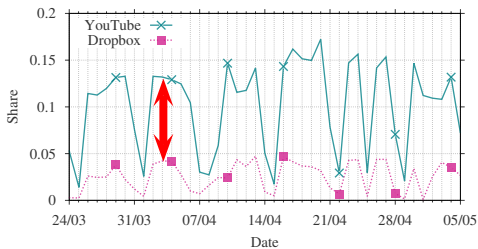
- 42 consecutive days in March and April 2012
 - 4 vantage points in Europe
 - Number of IP addresses in home probes \approx installations

	Type	IP Addr.	Dropbox		
			Flows	Vol. (GB)	Devices
Campus 1	Wired	400	167,189	146	283
Campus 2	Wired/Wireless	2,528	1,902,824	1,814	6,609
Home 1	FTTH/ADSL	18,785	1,438,369	1,153	3,350
Home 2	ADSL	13,723	693,086	506	1,313
Total			4,204,666	3,624	11,561

- 42 consecutive days in March and April 2012
 - 4 vantage points in Europe
 - Number of IP addresses in home probes \approx installations
 - 11,561 unique devices
- 2nd capture in Campus 1 in June 2012

How much traffic to personal cloud storage?

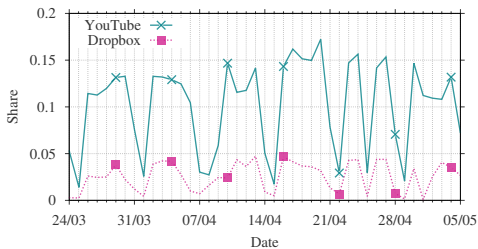
7/14



- Equivalent to **1/3 of YouTube volume at campus**

How much traffic to personal cloud storage?

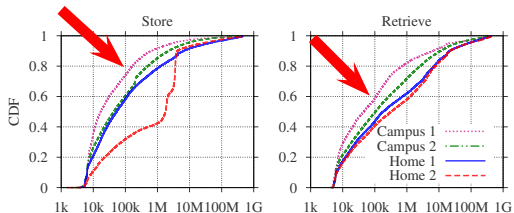
7/14



- Equivalent to 1/3 of YouTube volume at campus
- **Popularity: 6–12% adoption in home networks**
- **90% of the Dropbox traffic is from the native client**

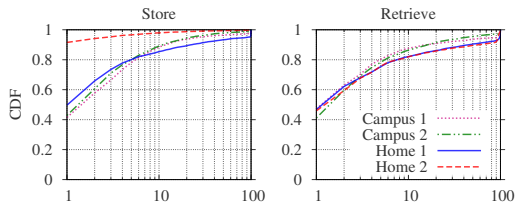
How does the storage traffic look like?

8/14



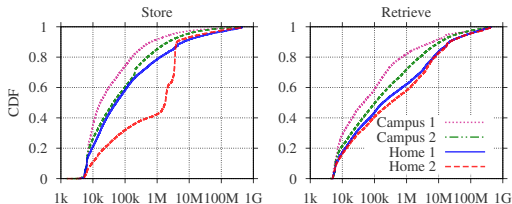
■ Flow size

- **Store: 40%–80% < 100kB**
→ Deltas, small files
- **Larger retrieve flows**



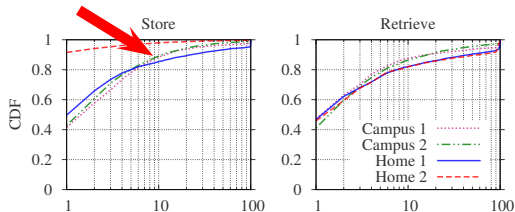
How does the storage traffic look like?

8/14



■ Flow size

- **Store: 40%–80% < 100kB**
→ Deltas, small files
- **Larger retrieve flows**

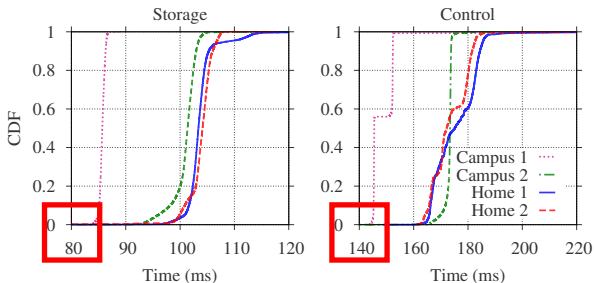


■ Chunks per flow

- **80% ≤ 10 chunks**
- **Remaining: up to 100**
→ Limited by the client

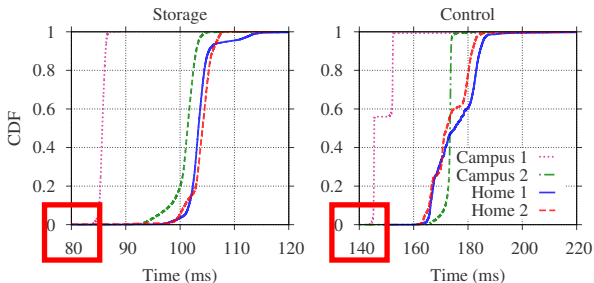
Where are the servers located?

9/14



- Minimum RTT per flow → stable over 42 days
- PlanetLab experiments → **the same U.S. data-centers worldwide**

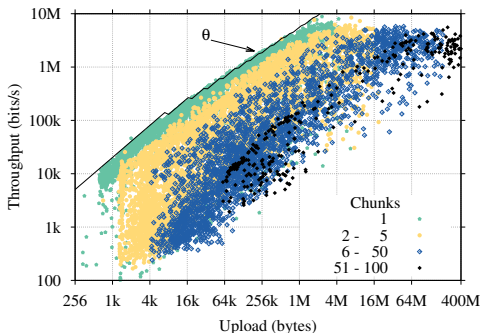
Where are the servers located?



- Minimum RTT per flow → stable over 42 days
- PlanetLab experiments → **the same U.S. data-centers worldwide**
- *"less than 35% of our users are from the USA"*

How is the performance far from the data-centers?

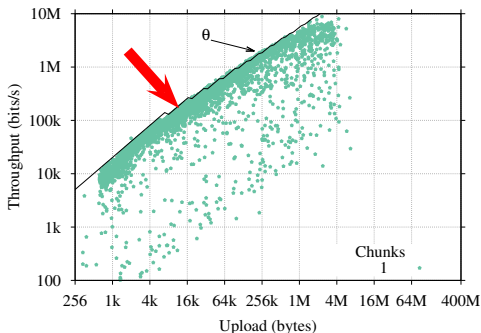
10/14



- Storage throughput in campuses
- **Most flows experience a low throughput**

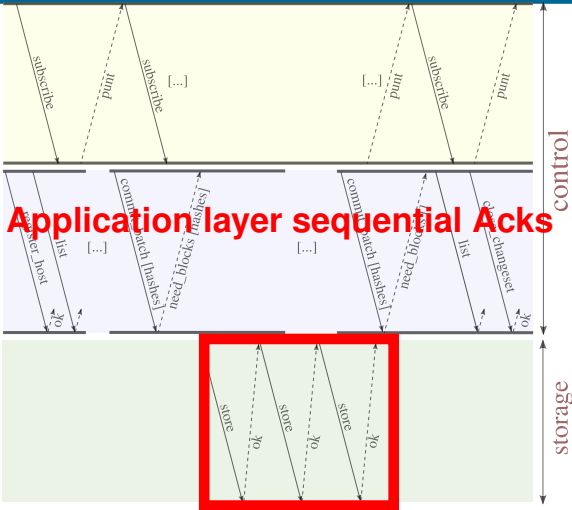
How is the performance far from the data-centers?

10/14



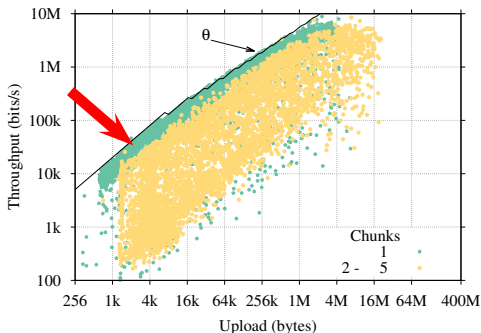
- Flows carrying 1 chunk
 - Size \leq 4MB, RTT \approx 100ms
 - Most of them finish in **TCP slow-start**

How is the performance far from the data-centers?



How is the performance far from the data-centers?

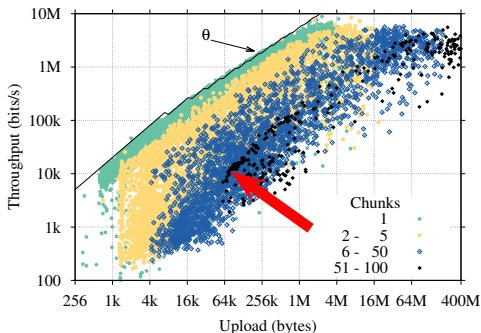
10/14



- Flows carrying several chunks
 - **Pause between chunks** → RTT and client/server reaction

How is the performance far from the data-centers?

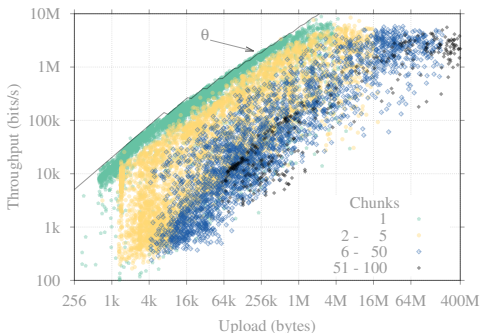
10/14



- Flows carrying several chunks
 - **Pause between chunks** → RTT and client/server reaction
 - Transferring **100 chunks takes more than 30s**
 - RTTs → 10s of inactivity

How is the performance far from the data-centers?

10/14



- **Delaying acknowledgments**
- **Bundling chunk** → recently deployed
- **Distributing servers** → storage traffic is heavy!

How much improvement from chunk bundling?

- New version released on Apr 26th (v1.4.0)
 - Small chunks are bundled together

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

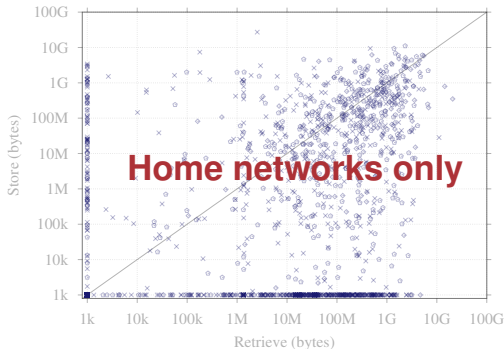
How much improvement from chunk bundling?

- New version released on Apr 26th (v1.4.0)
 - Small chunks are bundled together

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

- **Less small flows** → less TCP slow-start effects
- Average **throughput** is up to **65% higher**

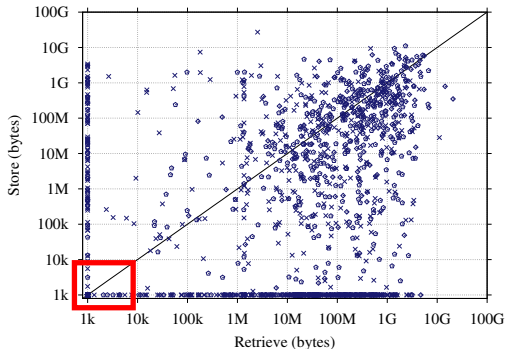
Are there typical usage scenarios?



- **More downloads** → download/upload ratio up to **2.4**
- What about download/upload per user?

Are there typical usage scenarios?

12/14



Occasional:

Users: 31%

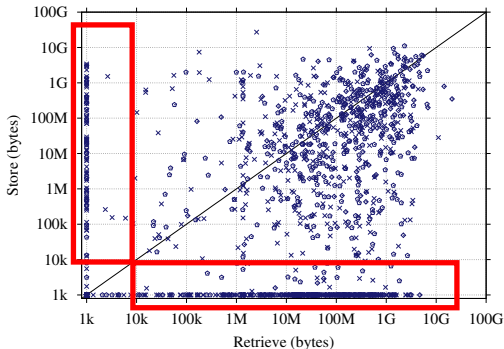
Devices per user: 1.22

Abandoned Dropbox clients

No storage activity for 42 days

Are there typical usage scenarios?

12/14



■ Upload-only:

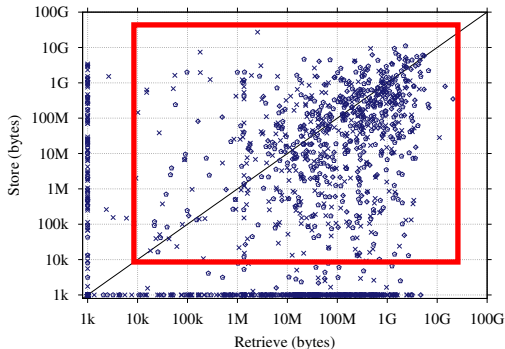
- Users: 6%
- Uploads: 11–21%
- Devices per user: **1.36**

■ Download-only:

- Users: 26%
- Downloads: 25–28%
- Devices per user: **1.69**

- Backup and content sharing
- Geographically dispersed devices

Are there typical usage scenarios?



■ Heavy:

- Users: **37%**
- Uploads: 79–89%
- Downloads: 72–75%
- Devices per user: **2.65**

■ Synchronization of content in a household

- **1st to analyze Dropbox usage on the Internet**

- **1st to analyze Dropbox usage on the Internet**
- Adoption above 6% in our datasets, data hungry application!

- **1st to analyze Dropbox usage on the Internet**
- Adoption above 6% in our datasets, data hungry application!
- Architecture and performance
 - **Bottlenecks from system design choices**

- **1st to analyze Dropbox usage on the Internet**
- Adoption above 6% in our datasets, data hungry application!
- Architecture and performance
 - **Bottlenecks from system design choices**
- **Extensive characterization of workload and usage**
 - User groups, number of devices, daily activity etc.

- **Thank You.**

- Anonymized traces and scripts

 - `http://traces.simpleweb.org/dropbox/`