

Cloud based Personal Storage Services

Definizione della metodologia di analisi per lo studio e la comparazione di servizi di storage basati su architetture cloud



**POLITECNICO
DI TORINO**

 **Plane**



4° Borsisti Day – 13/09/2013

GARR 



1. Molto diffusi

2. Offrono piani iniziali gratuiti

3. Generano grandi moli di dati da trasferire attraverso la rete

4. Dropbox: 100M utenti, 500M dispositivi, 10¹² upload/giorno



Google Drive

Poche informazioni riguardo a

- Architettura del sistema
- Funzionalità avanzate
- Prestazioni offerte



iCloud

Servizi inclusi e finalità dell'analisi

- Analisi dei cinque servizi più popolari:



Dropbox



Google Drive



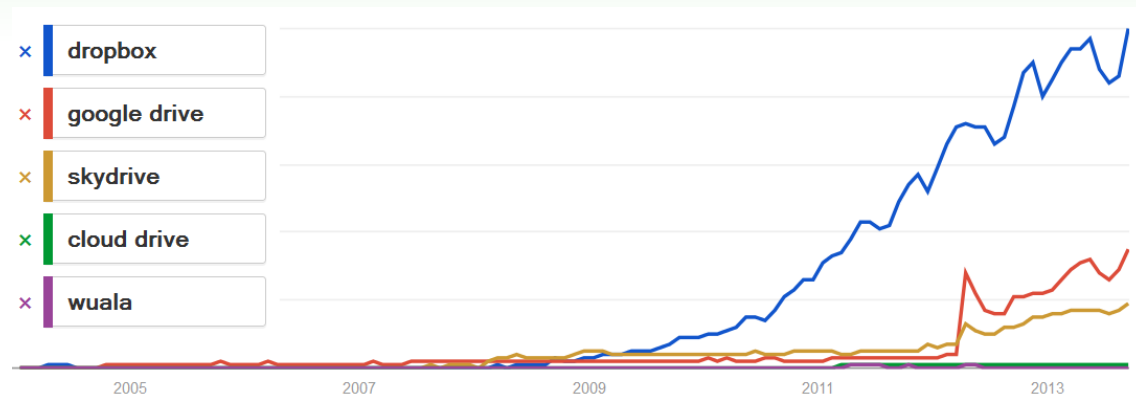
Microsoft SkyDrive



Amazon Cloud Drive



Wuala (by Lacie)



- Finalità della ricerca:

- Misura di prestazioni fornite all'utente finale e qualità del servizio
- Valutazione dell'efficienza dei servizi in analisi
- Valutazione del carico di traffico imposto alla rete
- Punti deboli e tecniche di ottimizzazione per la rete
- Migliorie applicabili a ciascun servizio

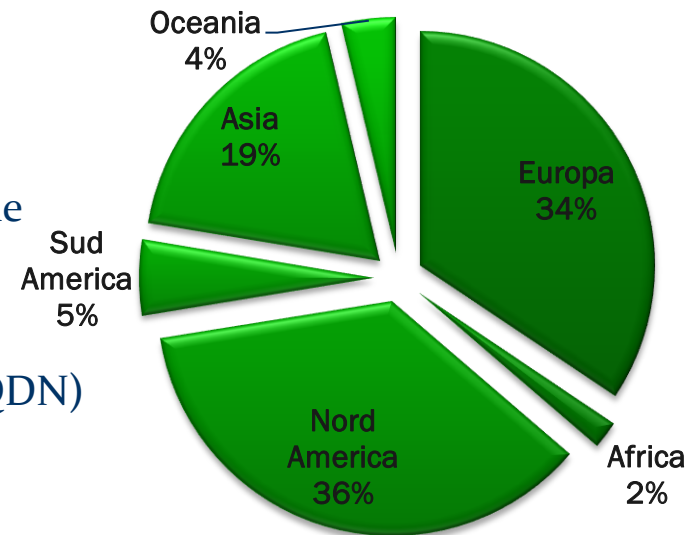
Introduzione all'analisi dei servizi

- I servizi fanno uso di molteplici connessioni a server differenti
- Vengono analizzati rispetto a tre operazioni comuni a tutti:
 - I. Procedura di login
 - II. Notifica di cambiamenti nei file sincronizzati
 - III. Memorizzazione/aggiornamento dei contenuti
- Durante le tre fasi vengono monitorate le connessioni effettuate e raccolti i nomi DNS dei server
- Attraverso i nomi DNS è possibile classificare i flussi

I.	clientX.dropbox.com	Login e identificazione
II.	notifyX.dropbox.com	Notifica cambiamenti
III.	dl-clientX.dropbox.com	Memorizzazione contenuti

Identificazione dell'infrastruttura

- I nomi DNS vengono tradotti in IP da 2000 resolvers in tutto il mondo
 - Assunzione di un punto di vista su scala globale
 - Raccolta degli indirizzi IP corrispondenti ai nomi DNS identificati
 - Identificazione tecniche di load balancing
 - Inclusi più di 100 paesi e 500 ISP
 - Maggiore presenza nelle aree con concentrazione più alta di traffico e datacenter
 - Raccolta dei Fully Qualified Domain Names (FQDN)
 - Utile per localizzazione geografica (airport tag)
 - Verifica dell'ownership dei datacenter



a.resolvers.level3.net (4.2.2.1) resolves dl-web.dropbox.com to [23.23.152.71] (ec2-23-23-152-71.compute-1.amazonaws.com)

- Necessario alle fasi successive

- Basata su quattro fonti d'informazioni indipendenti:

1. Airport tag incluso nel nome

“**mil**o2so6-in-fio.1e100.net” è uno dei FQDN di googleusercontent.com

2. Round Trip Time minimo (ICMP echo e TCP)

- Proporzionale alla distanza fisica dal server
- Triangolazione da diversi punti di misura appartenenti alla rete PlanetLab
- 200 nodi totali di cui si conosce la posizione geografica

Server login di Dropbox:	client-lb.dropbox.com	108.160.161.177
Nodo PlanetLab più vicino:	pli1-pa-6.hpl.hp.com	1.261 1.318 1.356

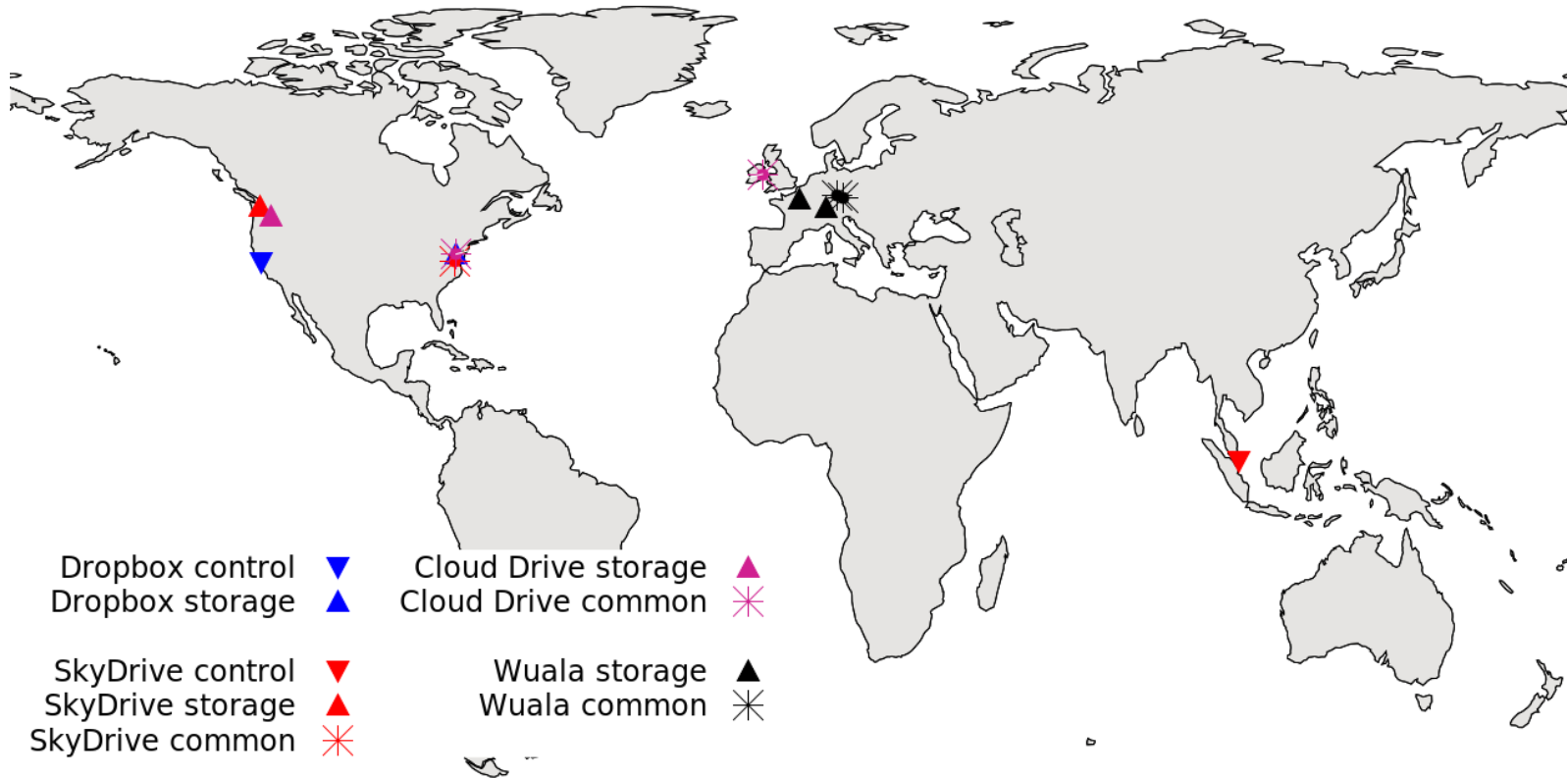
3. Traceroute all'ultimo router e analisi nomi dei router sul percorso

4. Database geografici

- Forniti da terze parti
- Noti per non essere affidabili su infrastrutture cloud

Geolocation, servizi centralizzati

- Dropbox: 2 datacenter, U.S.
- Microsoft SkyDrive: 3 datacenter, U.S. e Singapore
- Amazon Cloud Drive: 3 datacenter, U.S. e Gran Bretagna
- Wuala: 4 datacenter, Europa



Geolocation, Google Drive

- Google Drive: 41 edge-point → sistema distribuito
- L'utente raggiunge il punto più vicino
 - ✓ Offload della rete pubblica
 - ✓ RTT client – edge point ottimo

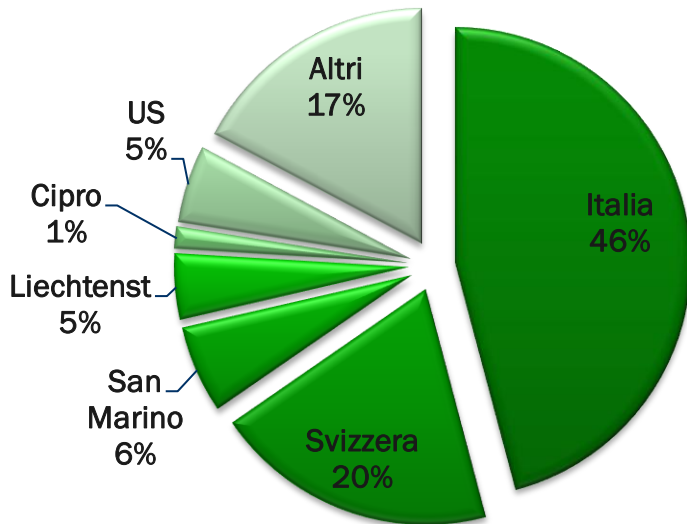


Google Drive, Locality awareness

- Locality awareness: permette di indirizzare l'utente verso la destinazione più vicina alla sua posizione
- I server DNS restituiscono gruppi di indirizzi IP differenti

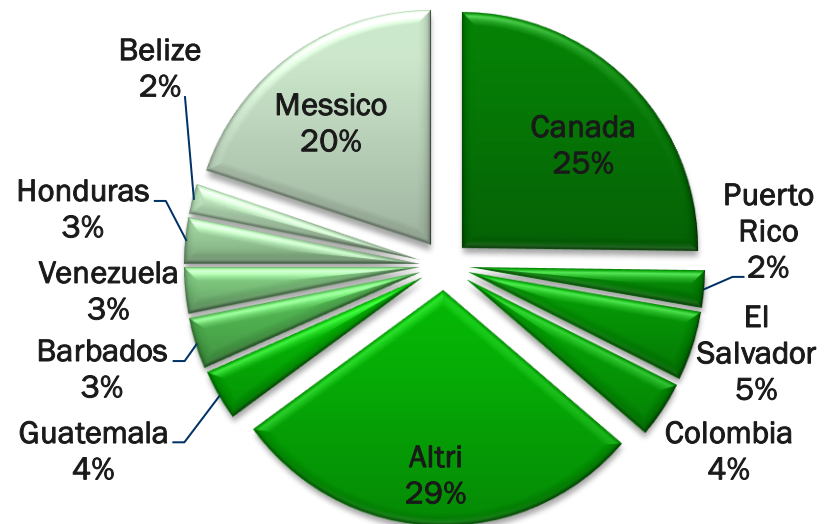
Italia

- ✓ Totalità del traffico nazionale verso datacenter di Milano
- ✓ Traffico estero in maggioranza da paesi confinanti



Stati Uniti

- ✓ I datacenter servono nel 95% dei casi traffico proveniente dagli U.S.
- ✓ Traffico estero da America Centrale e Isole del Pacifico

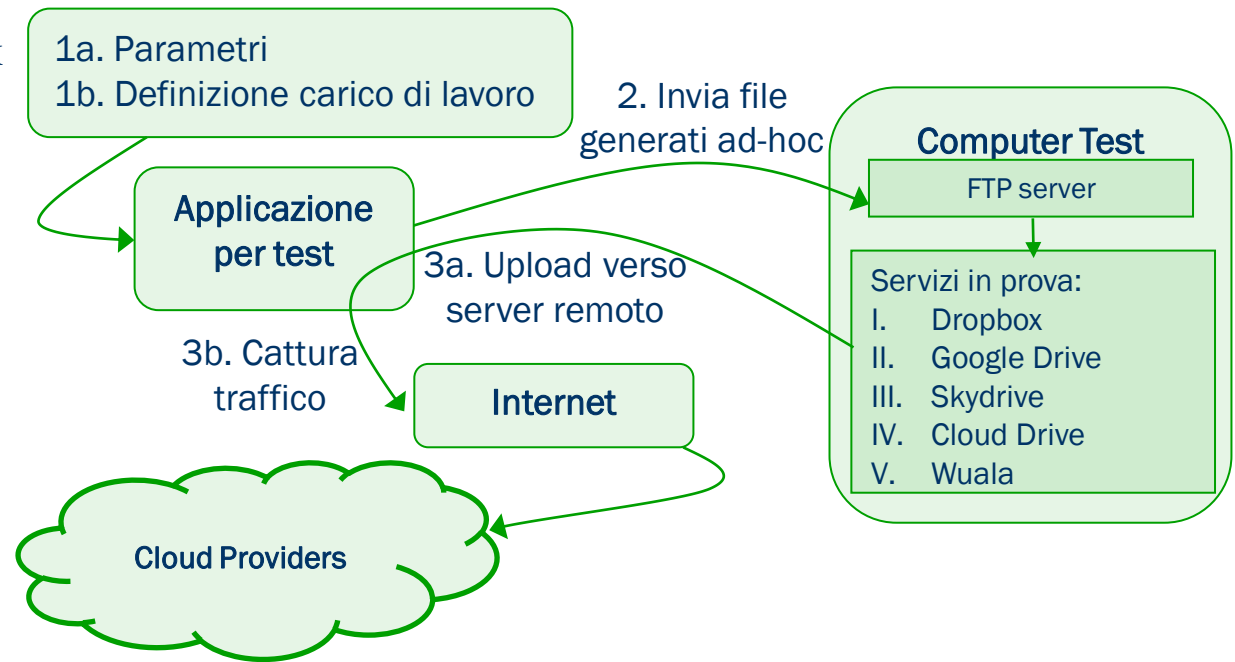


Caratteristiche avanzate - Testbed

- Sviluppo di un'applicazione di controllo per testing dei client software
- Uso estensivo di scripting python
- File generati ad-hoc al momento dell'esecuzione a seconda del test

Testbed composto da:

- ✓ Server basato su OS Linux con ruolo di controllore
- ✓ Virtual machine con OS Windows 7 Pro e client software in analisi
- ✓ Connessione ad internet tramite rete cablata del Politecnico, IP pubblico

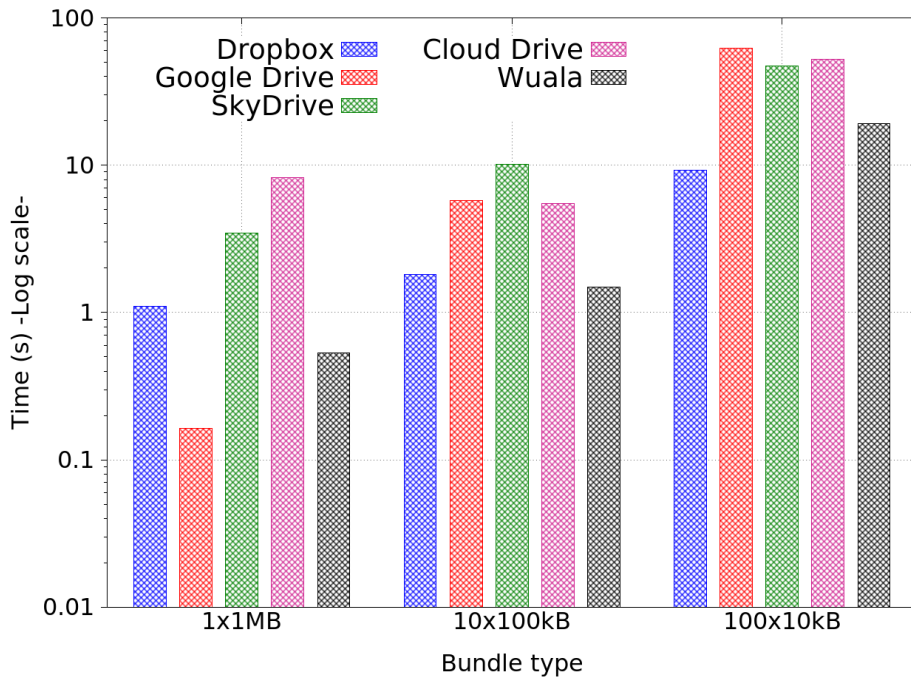


Caratteristiche avanzate - Sommario

- Risiedono nel client software fornito dai servizi di storage
- Necessaria la creazione di file ad hoc per verificarne l'implementazione
- Sei individuate:
 - ✓ Bundling
 - ✓ Compression
 - ✓ Delta encoding
 - ✓ De-duplication
 - ✓ De-deletion
 - ✓ Chunking

	Bundling	Compression	Delta encoding	De-duplication	De-deletion	Chunking
Dropbox	✓	✓	✓	✓	✓	~4MB su connessione TCP unica
Google Drive	X	✓	X	Parziale	X	~8MB su connessioni TCP multiple
Microsoft SkyDrive	Assistito	X	X	Parziale	X	Variabile
Amazon Cloud Drive	X	X	X	Parziale	X	X
Wuala (by Lacie)	Assistito	X	X	✓	✓	Variabile

- Capacità del client software di utilizzare un'unica connessione per il trasferimento in blocco di più file
 - + Riduzione overhead dovuto all'apertura di connessioni verso i server di storage
 - + Aumento throughput complessivo



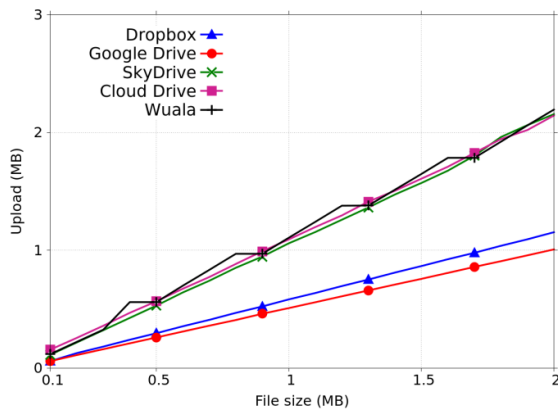
✓ Dropbox: supporto completo, trasmissione continuativa

✓ Wuala e SkyDrive: unica connessione ma attesa di conferma da livello applicazione, trasmissione a burst

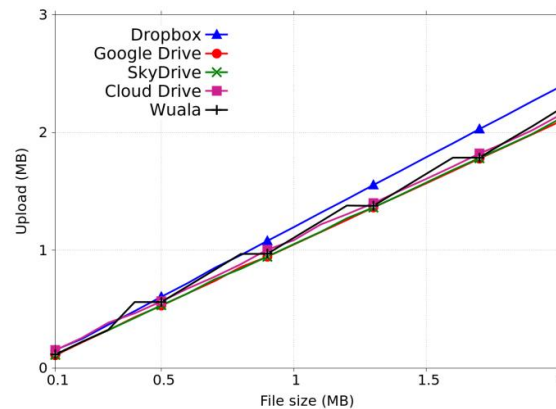
✓ Google Drive e CloudDrive: connessioni differenti per ogni contenuto / chunk

- Capacità di effettuare compressione del contenuto lato utente
 - + Riduzione capacità di upload/download necessaria
 - + Riduzione tempo di completamento sincronizzazione
 - Ritardo necessario alla compressione (solitamente trascurabile rispetto al tempo di download/upload)
 - Vantaggioso solo su alcuni tipi di file
- Smart compression: identificazione di contenuto comprimibile a seconda dell'estensione del file

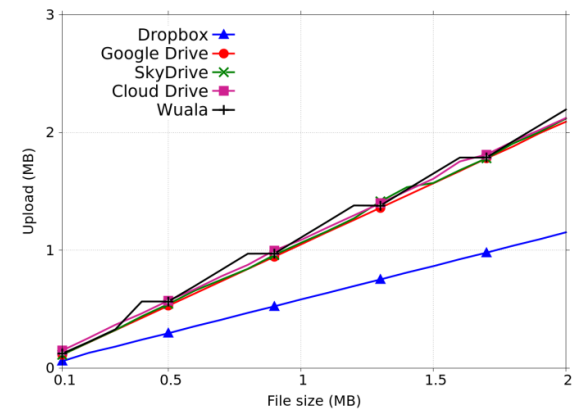
Testo in chiaro



Bytes casuali

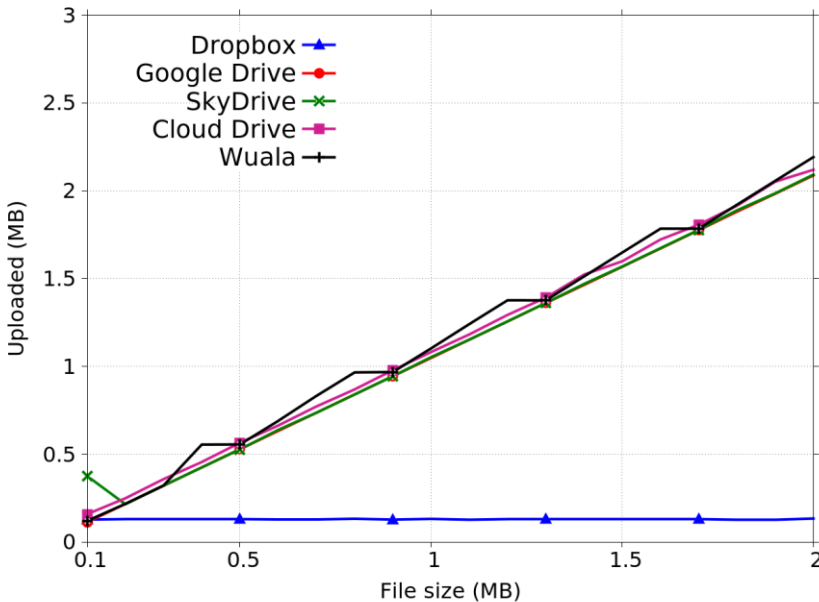


Smart compression

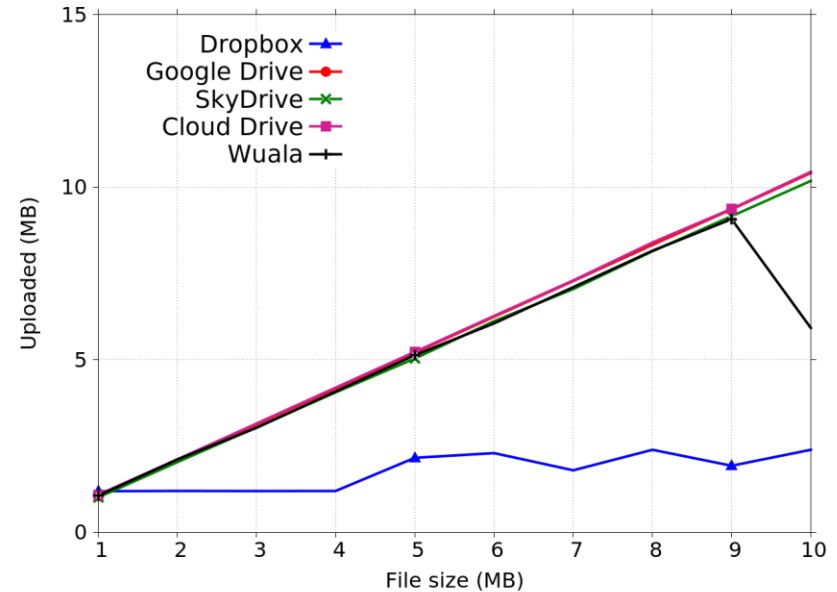


- Capacità di identificazione e aggiornamento delle porzioni di file soggette a modifiche, lasciando inalterato il resto del file
 - + Riduzione traffico di upload/download
 - + Riduzione tempo di completamento sincronizzazione
 - Vantaggioso solo con file frequentemente modificati dall'utente

Aggiunta 100Kb a fine file



Modifica 100KB random



4,5,6. De-duplication, De-deletion, Chunking

De-duplication

- Capacità del client software di identificare doppioni dello stesso file ed evitarne il trasferimento
 - + Riduzione traffico di upload/download necessaria
 - + Riduzione tempo di completamento sincronizzazione

De-deletion

- Capacità del client software di identificare file precedentemente sincronizzati e cancellati localmente
 - + Evita upload dello stesso contenuto al momento del ripristino del file locale
 - + Disponibilità limitata nel tempo

Chunking

- Suddivisione di file grandi in più parti per il trasferimento
 - + Facilita trasferimento in caso di connessioni instabili o propense a errori (wireless)
 - Riduzione throughput complessivo dovuto al tempo di silenzio tra due chunk e all'overhead per l'apertura di nuove connessioni

- Valutazione delle prestazioni offerte:
 - Throughput misurato in upload / download
 - Efficienza e overhead dei servizi in analisi
 - Incidenza delle caratteristiche avanzate sul tempo di completamento
 - Misure ripetute su un periodo di tempo esteso
 - Misure da postazioni differenti (4 servizi su 5 sono centralizzati e situati negli U.S.)
- Verifica evoluzione temporale dei servizi:
 - Nuove release software
 - Aggiunta di caratteristiche avanzate
- Aggiunta di nuovi servizi nel set di analisi