

# ***i*SHAKE: Incremental Hashing With SHAKE128 and SHAKE256 for the Zettabyte Era**

Danilo Gligoroski and Simona Samardjiska

Department of Telematics,

Faculty of Information Technology, Mathematics and Electrical Engineering

Norwegian University of Science and Technology - NTNU, NORWAY

# Incremental cryptography

- Introduced 20 years ago in Crypto '94 paper by Bellare, Goldreich and Goldwasser, *“Incremental Cryptography: The Case of Hashing and Signing”*
- Bellare, Goldreich and Goldwasser, STOC '95, *“Incremental Cryptography and Applications to Virus Protection”*
- Bellare and Micciancio, *“A New Paradigm for Collision-Free Hashing: Incrementality at Reduced Cost”*, Eurocrypt '97

# Incremental cryptography

- Introduced 20 years ago in Crypto '94 paper by Bellare, Goldreich and Goldwasser, *“Incremental Cryptography: The Case of Hashing and Signing”*
- Bellare, Goldreich and Goldwasser, STOC '95, *“Incremental Cryptography and Applications to Virus Protection”*
- Bellare and Micciancio, *“A New Paradigm for Collision-Free Hashing: Incrementality at Reduced Cost”*, Eurocrypt '97

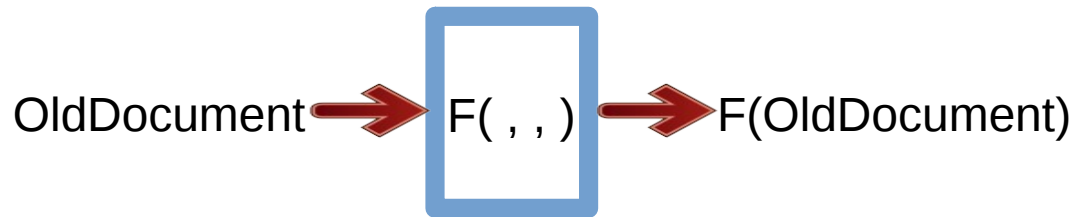
**WARNING:**

# Basic idea of Incrementality

- If we have already computed the function on some document, and this document is modified, then we update the function value based on the old value, and introduced changes rather than re-computing it from scratch.

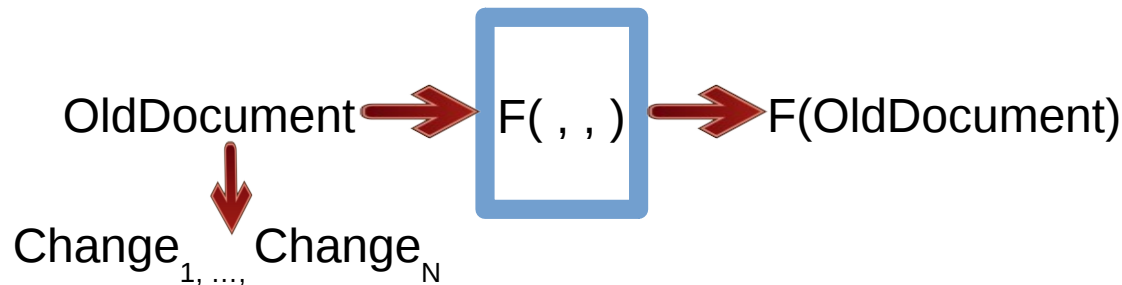
# Basic idea of Incrementality

- If we have already computed the function on some document, and this document is modified, then we update the function value based on the old value, and introduced changes rather than re-computing it from scratch.



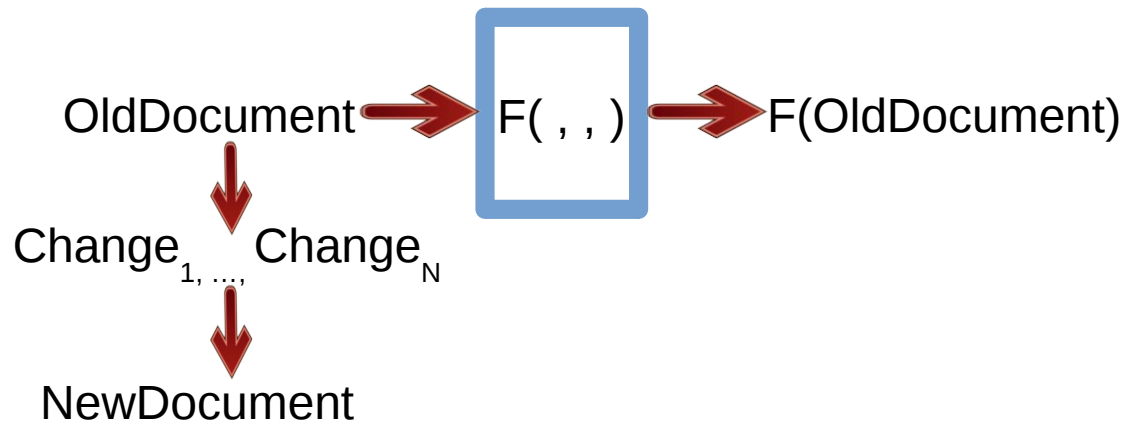
# Basic idea of Incrementality

- If we have already computed the function on some document, and this document is modified, then we update the function value based on the old value, and introduced changes rather than re-computing it from scratch.



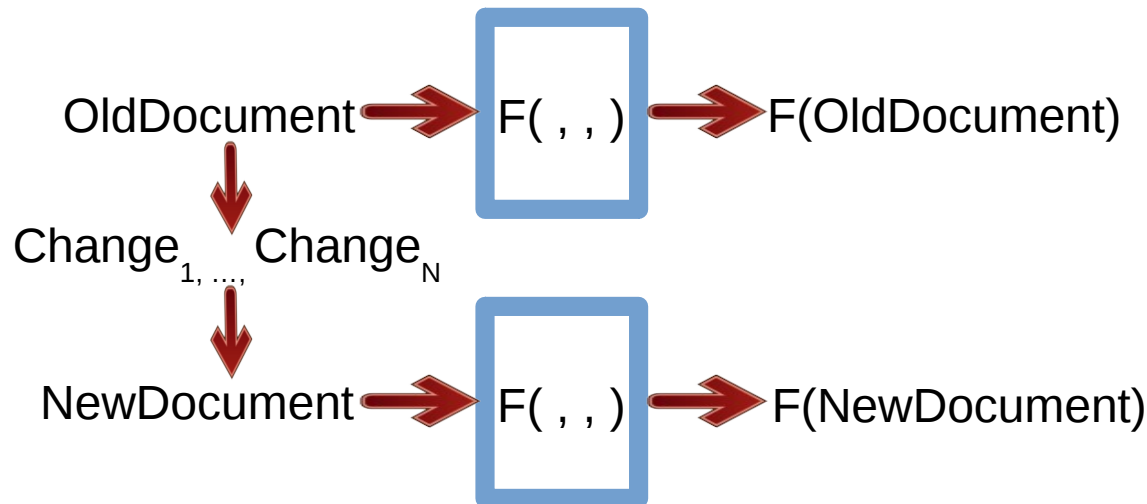
# Basic idea of Incrementality

- If we have already computed the function on some document, and this document is modified, then we update the function value based on the old value, and introduced changes rather than re-computing it from scratch.



# Basic idea of Incrementality

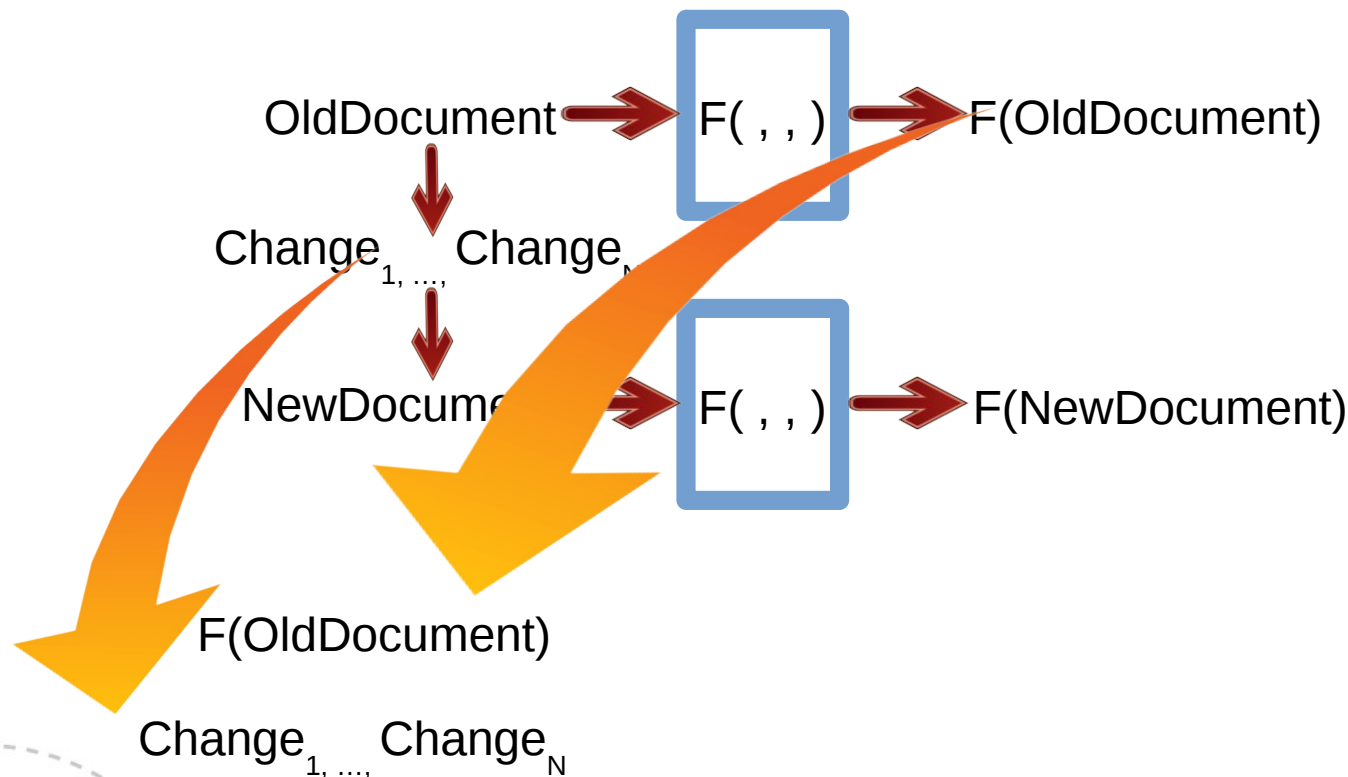
- If we have already computed the function on some document, and this document is modified, then we update the function value based on the old value, and introduced changes rather than re-computing it from scratch.





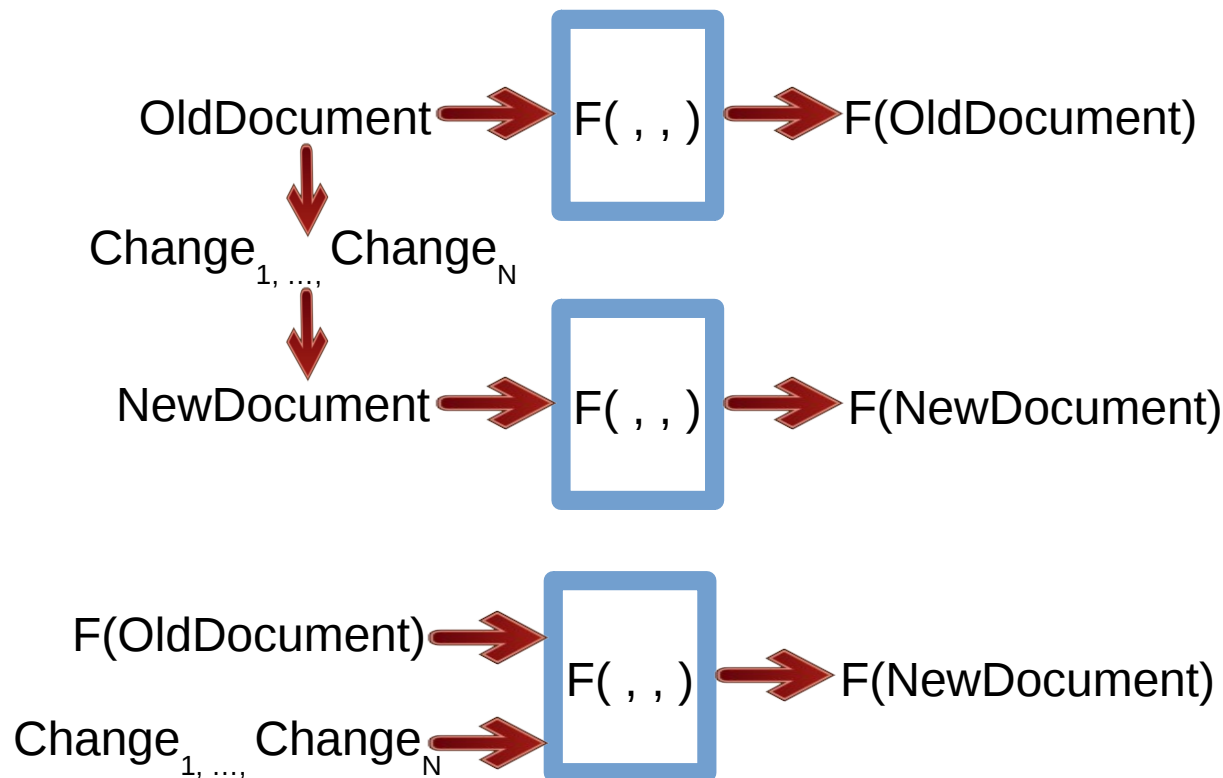
# Basic idea of Incrementality

- If we have already computed the function on some document, and this document is modified, then we update the function value based on the old value, and introduced changes rather than re-computing it from scratch.



# Basic idea of Incrementality

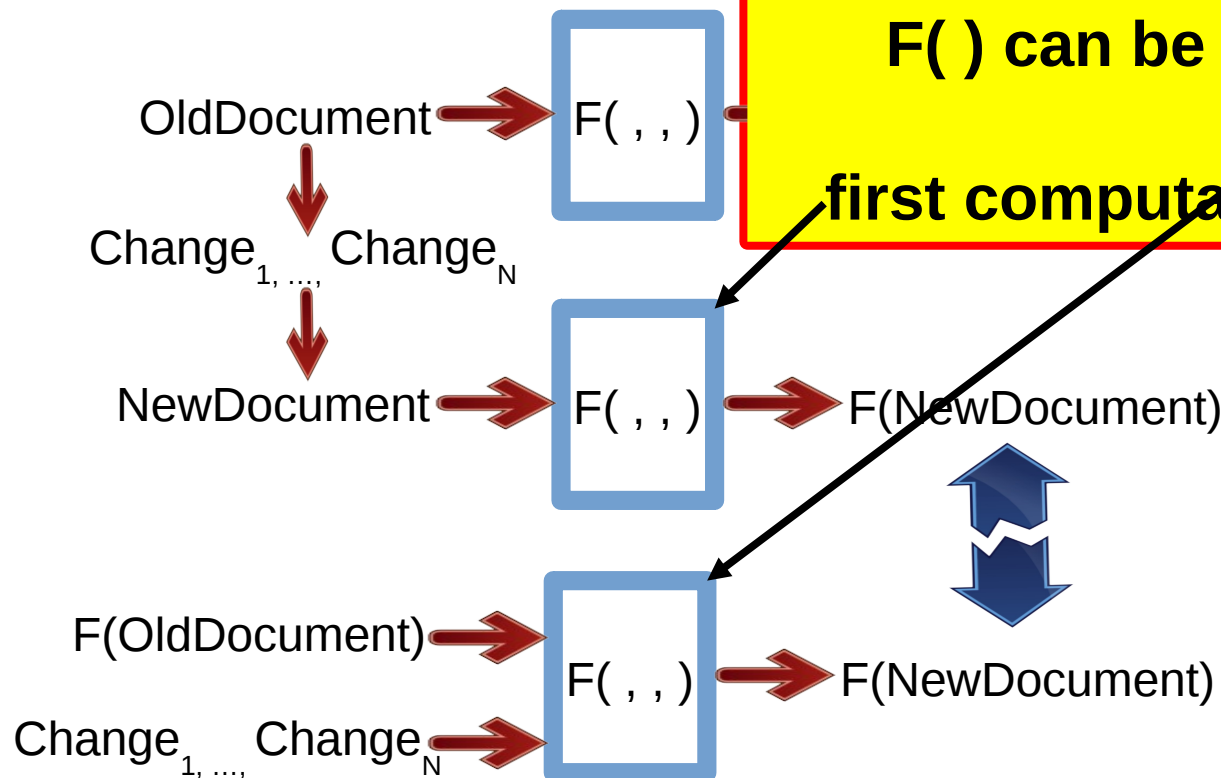
- If we have already computed the function on some document, and this document is modified, then we update the function value based on the old value, and introduced changes rather than re-computing it from scratch.



# Basic idea of Incrementality

- If we have already computed the function on some document, and this document is modified, then we update the function value based on the changes rather than re-computing the entire function value.

**The main motivation is that the second computation of  $F()$  can be orders of magnitude faster than the first computation of  $F()$ .**



# Bellare-Micciancio concrete proposals for construction of Collision-free Hashing

M

$M_n$   
↓  
 $M_n$

- Group  $G$  with combining operation  $\odot$
- Randomizer or compression function  $h$  mapping fixed size strings to element of  $G$
- Original message

$$M = M_1 M_2 \dots M_n$$

that we transform into the message

$$M' = \langle 1 \rangle M_1 \langle 2 \rangle M_2 \dots \langle n \rangle M_n$$

(prepend the index of each block to the block)

- Apply  $h$  to each block  $M'_i = \langle i \rangle M_i$  to obtain  $y_i = h(M'_i)$
- Combine  $y_i$  in  $G$  by the operation  $\odot$ :  
$$y(M) = y_1 \odot y_2 \odot \dots y_n$$

$\odot$   
↓  
 $y(M)$

# Operation $\odot$ , the size of the hash and properties of $h$

- Bellare and Micciancio proposed three variants for the operation  $\odot$ :
  - MuHASH: Multiplication in a group
    - special cases:
      - multiplication in groups of prime order
      - integer multiplication modulo  $p$
  - AdHASH: Addition modulo  $M$
  - LtHASH: Vector addition
- They concluded that in order the incremental hash function  $y(\ )$  to be collision free the following conditions must be satisfied:
  - In the group  $(G, \odot)$  the balance problem should be hard,
  - The compression function  $h(\ )$  should be collision free
  - The size  $k$  of the hash output should be around 1024 bits

# Operation $\odot$ , the size of the hash and properties of $h$

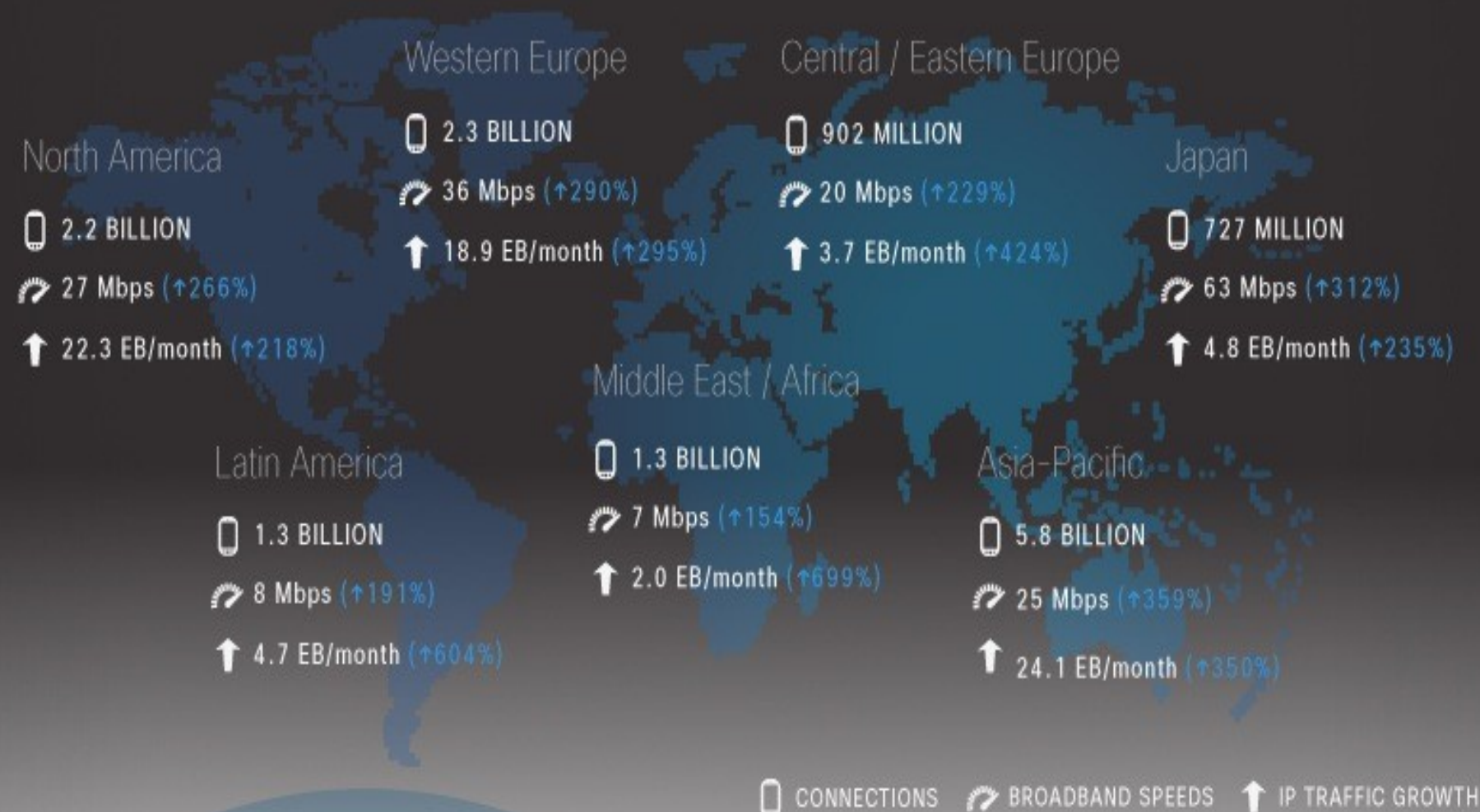
- However Wagner in “*A generalized birthday problem*” CRYPTO 2002, showed that the size  $k$  of the hash should be much bigger (for standard security levels, even up to tens of thousands of bits).
- Basically those findings killed the attractiveness of the concept of incremental hashing.
- Even more: None of the SHA-3 candidates explicitly explained can they offer some form of incremental hashing

# HOWEVER

# Welcome to the Zettabyte Era

Cisco Visual Networking Index

Global IP Traffic Forecast 2010-2015

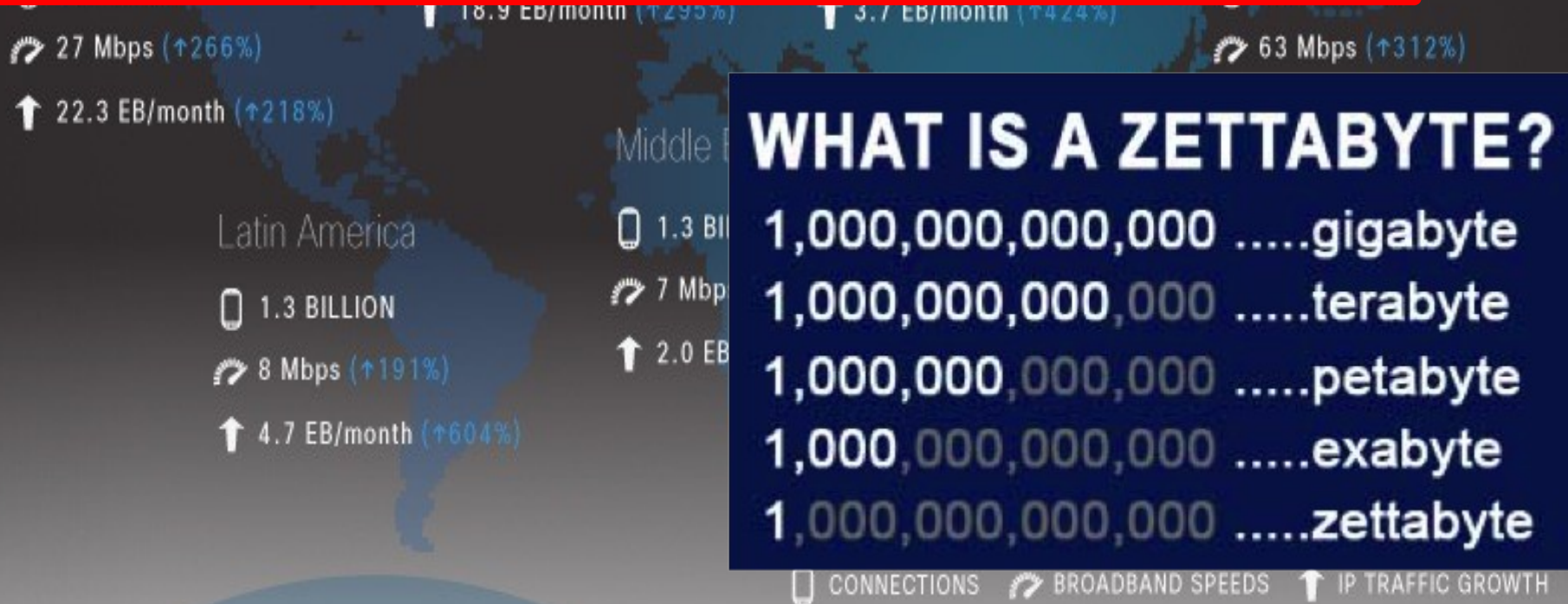




# Welcome to the Zettabyte Era

Cisco Visual Networking Index

**Annual global IP traffic will pass the zettabyte threshold by the end of 2015, and will reach 1.4 zettabytes per year by 2017. In 2015, global IP traffic will reach 1.0 zettabytes per year or 83.8 exabytes per month, and by 2017, global IP traffic will reach 1.4 zettabytes per year or 120.6 exabytes per month.**

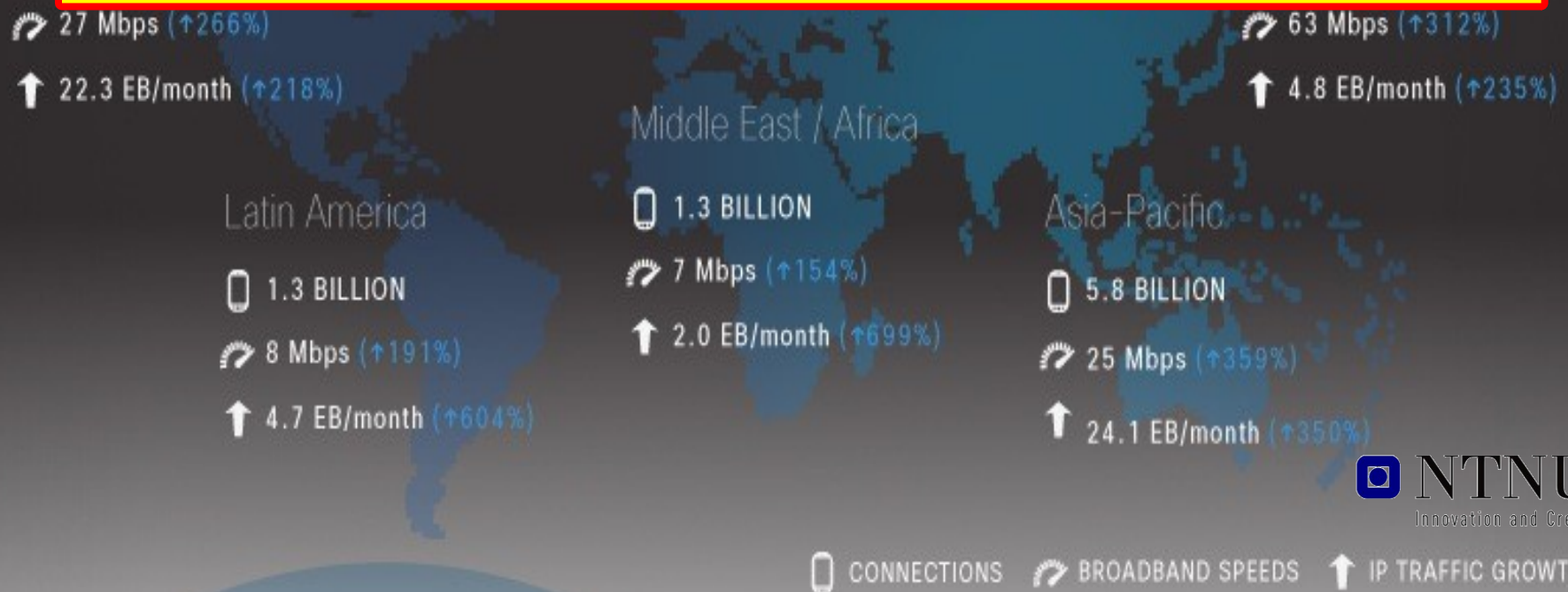


# Welcome to the Zettabyte Era

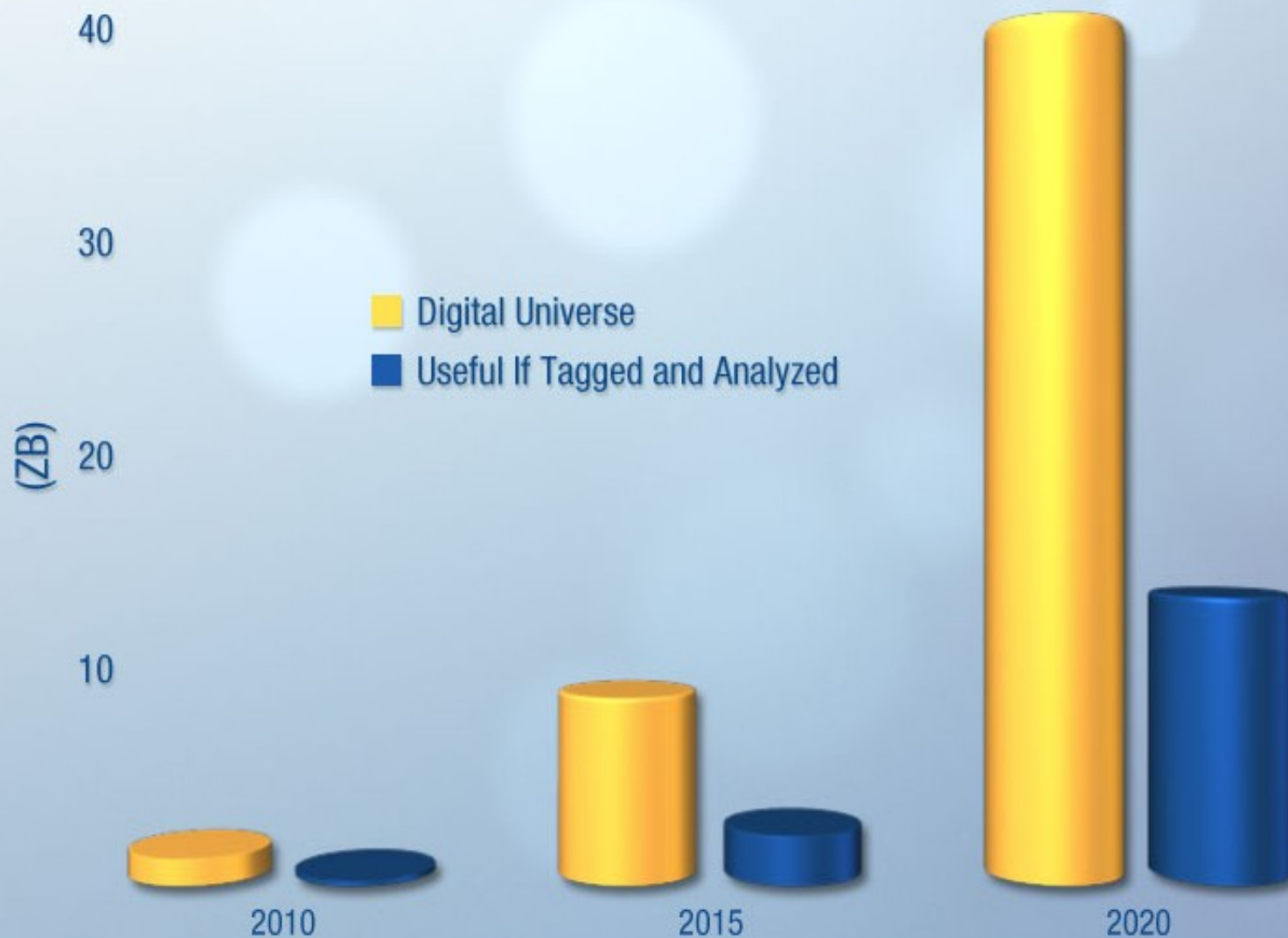
Cisco Visual Networking Index

Global IP Traffic Forecast 2010-2015

**In 2017, the gigabyte equivalent of all movies ever made will cross the global Internet every 3 minutes. The global Internet networks will deliver 13.8 petabytes every 5 minutes in 2017.**

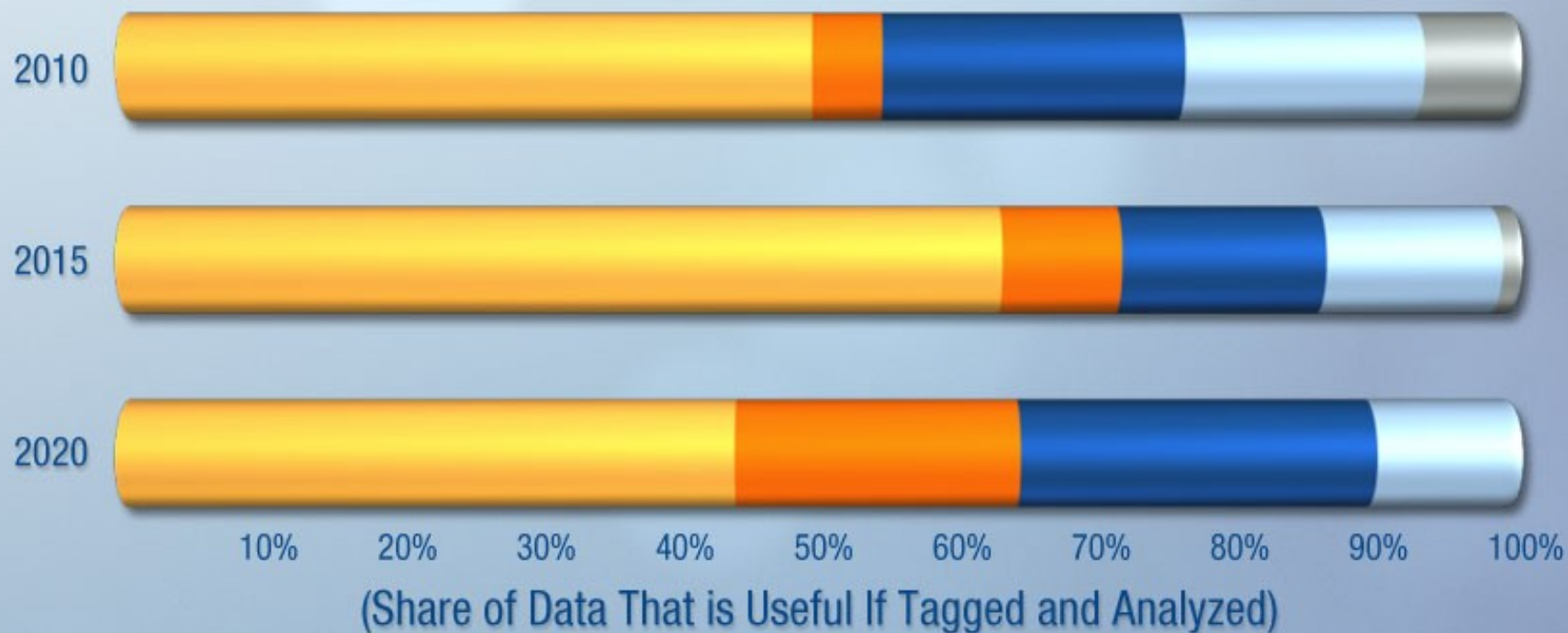


# Opportunity for Big Data



# Candidates for Big Data

- Surveillance
- Embedded and Medical
- Data Processing
- Entertainment and Social Media
- Consumer Images and Voice





# Are the practical needs for incremental hashing present today?


I would like to use MurmurHash3 to uniquely identify large pieces of data. This implementation:

0

<http://code.google.com/p/smhasher/source/browse/trunk/MurmurHash3.h>

Doesn't seem to provide a way to update the hash incrementally, though--it seems to compute one separate hash per block of data given. For example, if I were hashing 512MB of data from disk I might not want to load it all in memory at once, or if I were hashing an unknown amount of data from the network. Has anyone ever used MurmurHash3 in such a context before (hashing a large amount of

# Are the practical needs for incremental hashing present today?

 **crypto-js**  
JavaScript implementations of standard and secure cryptographic algorithms

[Project Home](#) [Downloads](#) [Wiki](#) **Issues** [Source](#)

[New issue](#)   for   [Advanced search](#) [Search tips](#) [Subscriptions](#)

★ **Issue 24: Support for computing hashes incrementally**  
2 people starred this issue and may be notified of changes.

**Status:** Fixed  
**Owner:** ---  
**Closed:** May 2012  
**Type:** Enhancement  
**Priority:** Medium

[Add a comment below](#)

Reported by [jhaber...@gmail.com](#), Jan 26, 2012

I have an application where I want to incrementally compute hashes of data, eg.

```
SHA256(A)
SHA256(A+B)
SHA256(A+B+C)
etc.
```

I don't want this to be an  $O(n^2)$  algorithm, so I was wondering if the API could support a lower-level object with an API like:

```
var sha256 = new Crypto.SHA256()
sha256.append(A)
var digestA = sha256.digest();
sha256.append(B)
var digestB = sha256.digest();
// etc.
```

# Are the practical needs for incremental hashing present today?

## Duplicity (software)

---

From Wikipedia, the free encyclopedia

**Duplicity** is a [software suite](#) that provides [encrypted](#), [digitally signed](#), [versioned](#), [remote backup](#) of files requiring little of the remote server.<sup>[1]</sup> Released under the terms of the [GNU General Public License](#) (GPL), Duplicity is [free software](#).

Duplicity devises a scheme where the first archive is a complete (full) backup, and subsequent (incremental) backups only add differences from the latest full or incremental backup.<sup>[2]</sup> Chains consisting of a full backup and a series of [incremental backups](#) can be recovered to the point in time that any of the incremental steps were taken. If any of the incremental backups are missing then the incremental backups following it cannot be reconstructed. It does this using [GnuPG](#), [libsync](#), [tar](#), and [rdiff](#).<sup>[1]</sup> To transmit data to the backup repository it can use [SSH/SCP/SFTP](#), local file access, [rsync](#), [FTP](#), [Amazon S3](#),<sup>[3]</sup> [Google Cloud Storage](#),<sup>[4]</sup> [Rackspace Cloud Files](#),<sup>[5]</sup> and others. Refer to its [man page](#) [↗](#) for the constantly growing list of back-ends.

# Are the practical needs for incremental hashing present today?



## Duplicity - Bandwidth Efficient Encrypted Backup

[Overview](#) [Code](#) **[Bugs](#)** [Blueprints](#) [Translations](#) [Answers](#)

### Very slow incremental backup of modified large files

Duplicity » Bugs » **Bug #1301892**

Reported by Andreï on 2014-04-03

This bug affects 1 person

Affects	Status	Importance	Assigned to	Milestone
Duplicity	New	Undecided	Unassigned	

Also affects project Also affects distribution/package Nominate for series

#### Bug Description

I'm using duplicity to backup to a LOCAL drive my virtual machine files.

The initial full backup goes swimmingly. Subsequent incremental backups also work very well as long as NONE of those large files (> 30 GB) have changes made to them. IF however those large >30 GB files have changes to them, it may take duplicity up to 10 hours to process the changes made to a single one of these large files.

ity



# Are the practical needs for incremental hashing present today?



## Duplicity - Bandwidth Efficient Encrypted Backup

[Overview](#) [Code](#) **[Bugs](#)** [Blueprints](#) [Translations](#) [Answers](#)

### Very slow incremental backup of modified large files

Duplicity » Bugs » **Bug #1301892**

Reported by Andrei

This bug affects 1 person

#### Affects



Duplicity

Milestone

Also affects project Also affects distribution/package Nominate for series

#### Bug Description

I'm using duplicity to backup to a LOCAL drive my virtual machine files.

The initial full backup goes swimmingly. Subsequent incremental backups also work very well as long as NONE of those large files (> 30 GB) have changes made to them. IF however those large >30 GB files have changes to them, it may take duplicity up to 10 hours to process the changes made to a single one of these large files.

In the latest version of Duplicity this issue is resolved (by using hash trees and storing intermediate hash values and signatures).

ity

# Are the practical needs for incremental hashing present today?

Signature files are not required to restore a backup set, but without an up-to-date signature, duplicity cannot append an incremental backup to an existing archive.

To save bandwidth, duplicity generates full signature sets and incremental signature sets. A full signature set is generated for each full backup, and an incremental one for each incremental backup. These start with **duplicity-full-signatures** and **duplicity-new-signatures** respectively. These signatures will be stored both locally and remotely. The remote signatures will be encrypted if encryption is enabled. The local signatures will not be encrypted and stored in the archive dir (see **--archive-dir** ).

# Are the practical needs for incremental hashing present today?

Signature files are not required to restore a backup set, but without an up-to-date signature, duplicity cannot append an incremental backup to an existing archive.

To save ba  
set is gene  
with dupl  
stored bot  
local signa

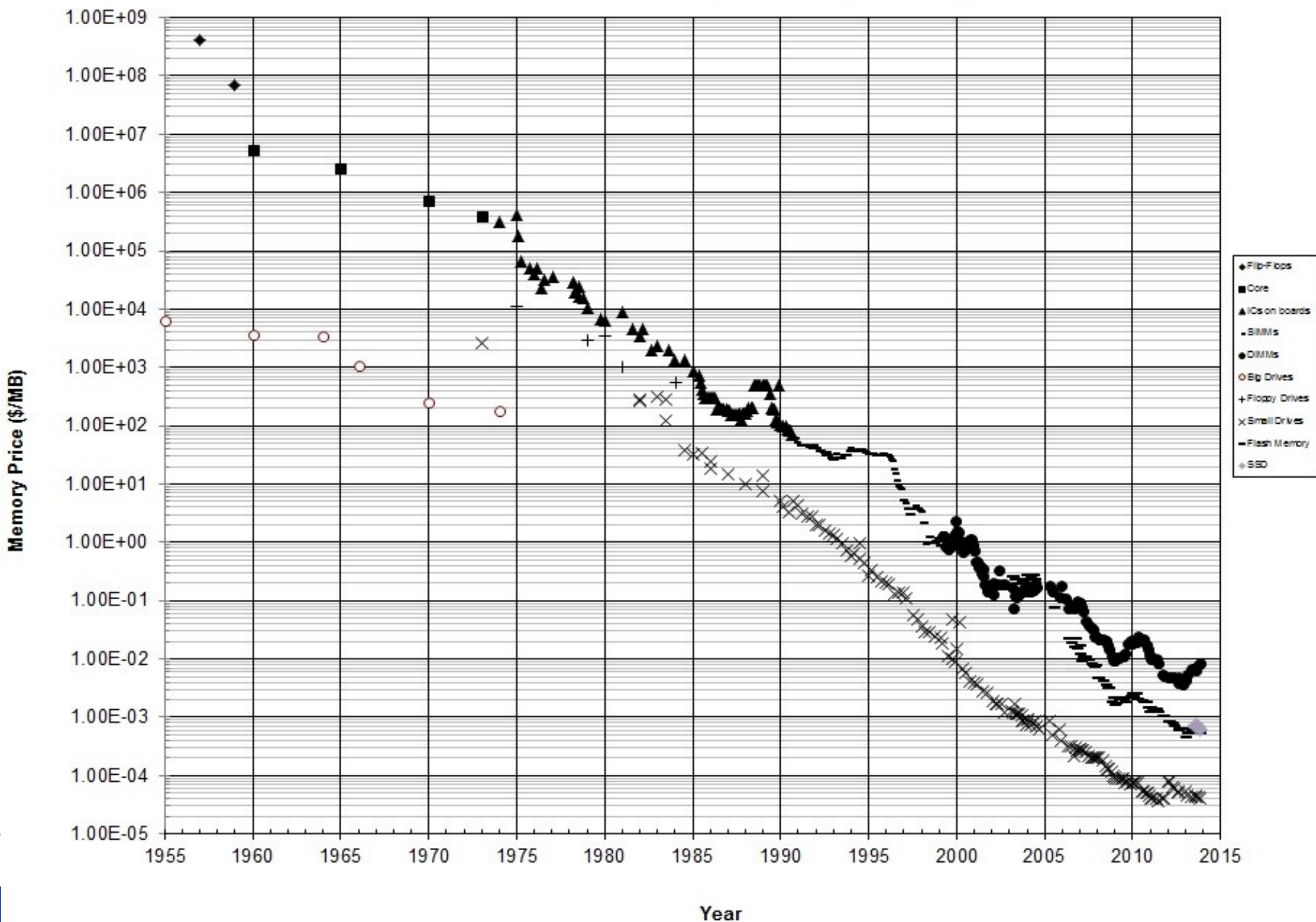
Having an incremental hash function would increase the speed of production of signed backups by order of magnitude, while decreasing the extra storage for intermediate signature sets.

. A full signature  
p. These start  
tures will be  
is enabled. The

# Observations

**Observation 1.** In the Zettabyte era the trend for reducing the cost of data storage will diminish the importance of the fact that incremental hashing needs longer digests than traditional hash functions.

# Historical Cost of Computer Memory and Storage



# Observations

**Observation 2.** In the Zettabyte era there will be an increased need for an efficient and secure cryptographic primitive that will perform incremental collision-free hashing.

# Observations

**Observation 3.** For the compression functions  $h : \{0, 1\}^b \rightarrow \{0, 1\}^k$  where  $k$  is multiple of 64 bits i.e.  $k = 64 \cdot L$ , on modern 64-bit CPUs, instead of modular operations with  $k$  bit prime numbers in the group  $(G, \odot)$ , much more efficient operations would be word wise operations of addition in the group  $((\mathbb{Z}_{64})^L, \boxplus_{64})$ .

# Formalizing previous notations

## Definition 1.

1. Let  $h : \{0,1\}^b \rightarrow \{0,1\}^k$  be a compression function that maps  $b$  bits into  $k$  bits.
2. Let the message  $M$  be represented as a concatenation of  $n$  blocks, where  $n < N$  for some predefined number  $N$  which is larger than the number of blocks in any message we plan to hash, i.e.,  $M = M_1 || M_2 || \dots || M_n$ .
3. The size of each block  $M_i$  is determined by the following relation:  $|M_i| = b - \lg(N)$ .
4. For each block  $M_i$ ,  $i = 1, \dots, n$ , prepend a  $\lg(N)$ -bit binary encoding  $\langle i \rangle$  of the block index  $i$  to the block content  $M_i$  to get an augmented block  $\overline{M}_i = \langle i \rangle || M_i$ .
5. For each  $i = 1, \dots, n$ , apply  $h$  to  $\overline{M}_i$  to get a hash value  $y_i = h(\overline{M}_i)$ .
6. Let  $(G, \odot)$  be a commutative group with operation  $\odot$  where  $G \subseteq \{0,1\}^k$ .
7. Combine  $y_1, \dots, y_n$  via a combining group operation  $\odot$  to get the final hash value  $y = y_1 \odot y_2 \odot \dots \odot y_n$ .

Denote the incremental hash function as:

$$y(M) = \text{HASH}_{\langle G \rangle}^h(M_1 || M_2 || \dots || M_n) = \bigodot_{i=1}^n h(\langle i \rangle || M_i) \quad (1)$$



# Formalizing previous notations

**Proposition 1.** *Let  $\text{HASH}_{\langle G \rangle}^h$  be an incremental hash function defined by Definition 1. For any  $Y \in \{0,1\}^k$  the complexity of finding a preimage message  $M = M_1 || M_2 || \dots || M_K$  of length  $K \leq N$  blocks such that  $Y = \text{HASH}_{\langle G \rangle}^h(M)$  is:*

$$\min_{K \leq N} O(K \cdot 2^{\frac{k}{1+\lceil \lg K \rceil}}) \quad (2)$$

*If the length of the messages is not restricted, then the minimum in equation (2) is achieved for messages of  $K = 2^{\sqrt{k}-1}$  blocks.*

# Formalizing previous notations

**Proposition 1.** *Let  $\text{HASH}_{\langle G \rangle}^h$  be an incremental hash function defined by Definition 1. For any  $Y \in \{0,1\}^k$  the complexity of finding a preimage message  $M = M_1 || M_2 || \dots || M_K$  of length  $K \leq N$  blocks such that  $Y = \text{HASH}_{\langle G \rangle}^h(M)$  is:*

$$\min_{K \leq N} O(K \cdot 2^{\frac{k}{1+\lceil \lg K \rceil}}) \quad (2)$$

*If the length of the messages is not restricted, then the minimum in equation (2) is achieved for messages of  $K = 2^{\sqrt{k}-1}$  blocks.*

**Proof.** Just adopt the notation from Wagner's Crypto 2002 paper [Sec. 2, Summary] to match the notation of variables in Definition 1.

# Extendable-Output Functions SHAKE128 and SHAKE256

- DRAFT FIPS 202 “*SHA-3 Standard: Permutation-Based Hash and Extendable-Output Functions*”

# Extendable-Output Functions SHAKE128 and SHAKE256

$$\text{SHAKE128}(M, d) = \text{RawSHAKE128}(M || 11, d),$$

where

$$\text{RawSHAKE128}(M, d) = \text{KECCAK}[256](M || 11, d),$$

and

$$\text{SHAKE256}(M, d) = \text{RawSHAKE256}(M || 11, d),$$

where

$$\text{RawSHAKE256}(M, d) = \text{KECCAK}[512](M || 11, d).$$

# iSHAKE128

## Definition 2.

1. Let  $h : \{0, 1\}^{1344} \rightarrow \{0, 1\}^{2688}$  be defined as the function  $h(m) = \text{SHAKE128}(m, 2688)$  where  $|m| = 1344$ .
2. Let the message  $M = M_1 || M_2 || \dots || M_n$  be represented as a concatenation of  $n$  blocks, where  $n < N$ , and  $N = 2^{25}$  is the largest number of blocks in any message we plan to hash.
3. The size of each block  $M_i$  in bits is determined by the following relation:  $|M_i| = 1344 - 64 = 1280$ .
4. For each block  $M_i$ ,  $i = 1, \dots, n$ , prepend a 64-bit binary encoding  $\langle i \rangle$  of the block index  $i$  to the block content  $M_i$  to get an augmented block  $\overline{M}_i = \langle i \rangle || M_i$ .
5. For each  $i = 1, \dots, n$ , apply  $h$  to  $\overline{M}_i$  to get a hash value  $y_i = h(\overline{M}_i) = \text{SHAKE128}(\overline{M}_i, 2688)$ .
6. Let  $((\mathbb{Z}_{64})^{42}, \boxplus_{64})$  be a commutative group with the operation  $\boxplus_{64}$  that represents a 64-bit word wise addition of 42 words.
7. Combine  $y_1, \dots, y_n$  via a combining group operation  $\boxplus_{64}$  to get the final hash value

$$y = y_1 \boxplus_{64} y_2 \boxplus_{64} \dots \boxplus_{64} y_n.$$

Denote the incremental hash function iSHAKE128 as:

$$i\text{SHAKE128}(M) = \boxplus_{64}^N \text{SHAKE128}(\overline{M}_i, 2688). \quad (3)$$

# *i*SHAKE128

*Corollary 1:* Let  $b = 1280$ ,  $k = 2688$ . and let the maximal allowed number of blocks be  $N = 2^{25}$ . Then

$$\min_{K \leq N} O(K \cdot 2^{\frac{k}{1+\lg\lfloor K \rfloor}}) = 2^{128.385} \quad (4)$$

# *i*SHAKE128

- For a  $2^{128}$  level of security, the maximal size of the files that can be hashed with iSHAKE128 is 5 GB.
- For small file sizes such as 160 KB the complexity of finding collisions with Wagner's generalized birthday attack is  $2^{254}$

and

- For files long 1.25 TB the complexity of finding collisions drops down to  $2^{112}$ .

# *i*SHAKE256

## Definition 3.

1. Let  $h : \{0, 1\}^{1088} \rightarrow \{0, 1\}^{6528}$  be defined as the function  $h(m) = \text{SHAKE256}(m, 6528)$  where  $|m| = 1088$ .
2. Let the message  $M = M_1 || M_2 || \dots || M_n$  be represented as a concatenation of  $n$  blocks, where  $n < N$ , and  $N = 2^{28}$  is the largest number of blocks in any message we plan to hash.
3. The size of each block  $M_i$  in bits is determined by the following relation:  $|M_i| = 1088 - 64 = 1024$ .
4. For each block  $M_i$ ,  $i = 1, \dots, n$ , prepend a 64-bit binary encoding  $\langle i \rangle$  of the block index  $i$  to the block content  $M_i$  to get an augmented block  $\overline{M}_i = \langle i \rangle || M_i$ .
5. For each  $i = 1, \dots, n$ , apply  $h$  to  $\overline{M}_i$  to get a hash value  $y_i = h(\overline{M}_i) = \text{SHAKE256}(\overline{M}_i, 6528)$ .
6. Let  $((\mathbb{Z}_{64})^{102}, \boxplus_{64})$  be a commutative group with the operation  $\boxplus_{64}$  that represents a 64-bit word wise addition of 102 words.
7. Combine  $y_1, \dots, y_n$  via a combining group operation  $\boxplus_{64}$  to get the final hash value

$$y = y_1 \boxplus_{64} y_2 \boxplus_{64} \dots \boxplus_{64} y_n.$$

Denote the incremental hash function *i*SHAKE256 as:

$$i\text{SHAKE256}(M) = \bigoplus_{i=1}^N \text{SHAKE256}(\overline{M}_i, 6528). \quad (5)$$



# *i*SHAKE256

*Corollary 2:* Let  $b = 1024$ ,  $k = 6528$ . and let the maximal allowed number of blocks be  $N = 2^{28}$ . Then

$$\min_{K \leq N} O(K \cdot 2^{\frac{k}{1+\lg\lfloor K \rfloor}}) = 2^{253.103} \quad (6)$$

# *i*SHAKE256

- For a  $2^{253}$  level of security, the maximal size of the files that can be hashed with iSHAKE256 is 32 GB.
- For small file sizes such of 1 MB the complexity of finding collisions with Wagner's generalized birthday attack is  $2^{479}$

and

- For files long 8 TB the complexity of finding collisions drops down to  $2^{212}$ .

# Properties of iSHAKE parallel operations

- Trivial parallelization without a need for special scheduling
  - (except that every message bloc  $M_i$  has to be prepended with its index  $\langle i \rangle$  in the message)
- The obtained hash is the same, regardless of the level of parallelism used in the computation
- No need to keep extra intermediate values in order to achieve incrementality
  - What is bigger: the amount of stored intermediate chain values in tree-hash modes or the size of hash values in iSHAKE (2688 / 6528 bits)?

Thank you for your attention!