# Deduplication: Overview & Case Studies

CSCI 333 – Spring 2020

Williams College

# Lecture Outline

## Background

Content Addressable Storage (CAS)

Deduplication

Chunking

The Index

Other CAS applications

# Lecture Outline

Background

## Content Addressable Storage (CAS)

Deduplication

Chunking

The Index

Other CAS applications

# Content Addressable Storage (CAS)

Deduplication systems often rely on Content Addressable Storage (CAS)

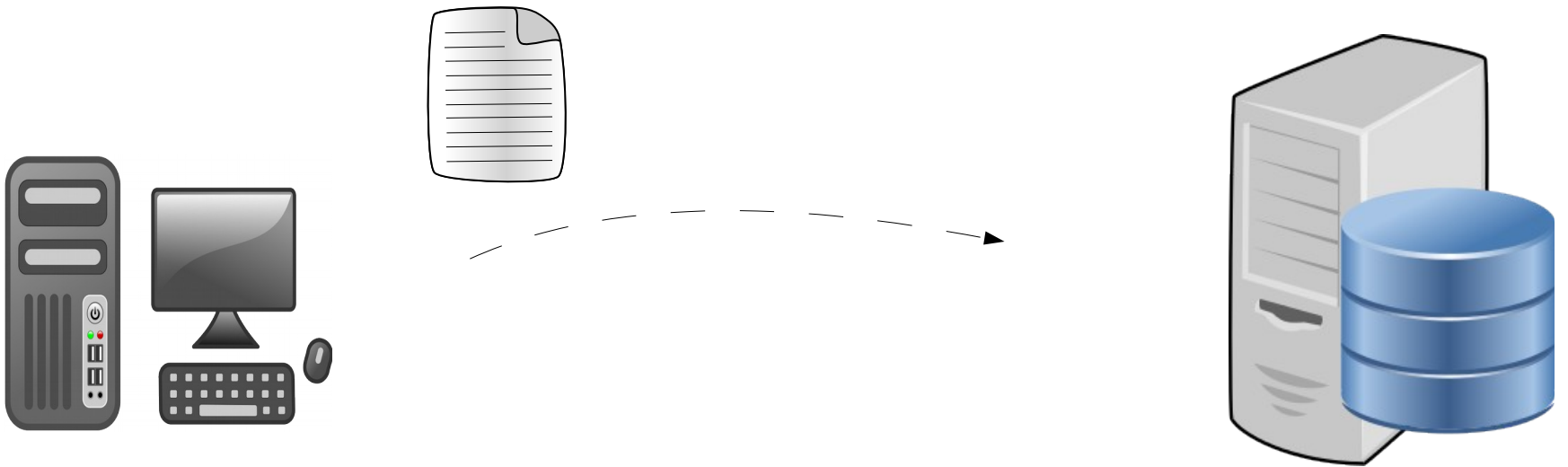Data is indexed by some **content identifier**

The **content identifier** is determined by some function over the data itself
  - often a cryptographically strong hash function

# CAS

Example:

I send a document to be stored remotely
on some content addressable storage

# CAS

Example:

The server receives the document, and calculates a unique identifier called the data's **fingerprint**

# CAS

The **fingerprint** should be:

        unique to the data
           - *NO* collisions

        one-way
           - hard to invert

# CAS

The **fingerprint** should be:

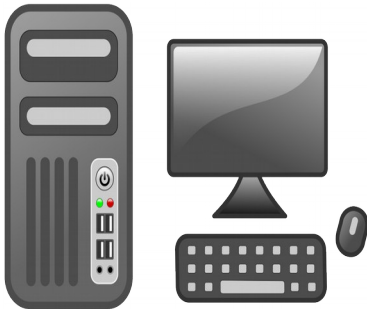unique to the data
- *NO* collisions

one-way
- hard to invert

$10^{24}$ objects before it is more likely than not that a collision has occurred

SHA-1:

20 bytes (160 bits)

$P(collision(a,b)) = (½)^{160}$
$coll(N, 2^{160}) = (_N C_2)(½)^{160}$

# CAS

Example:

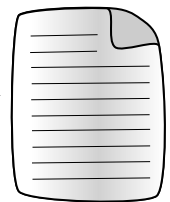$$\text{SHA-1(} \text{📄} \text{)} = \text{de9f2c7fd25e1b3a...}$$
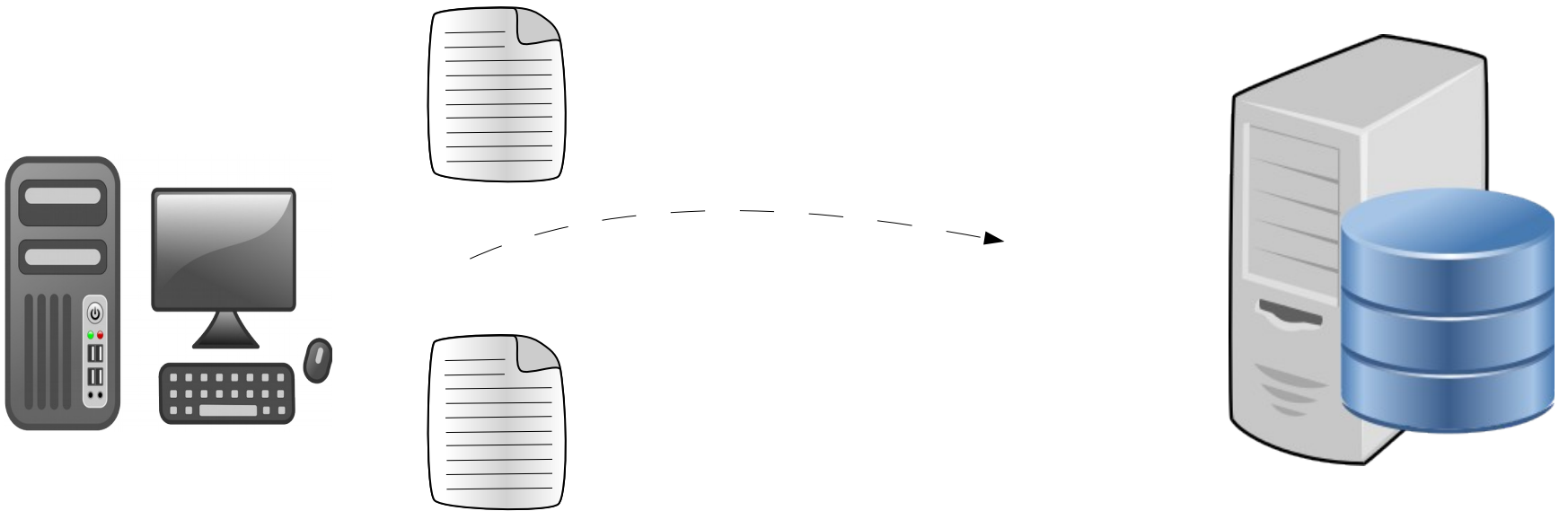


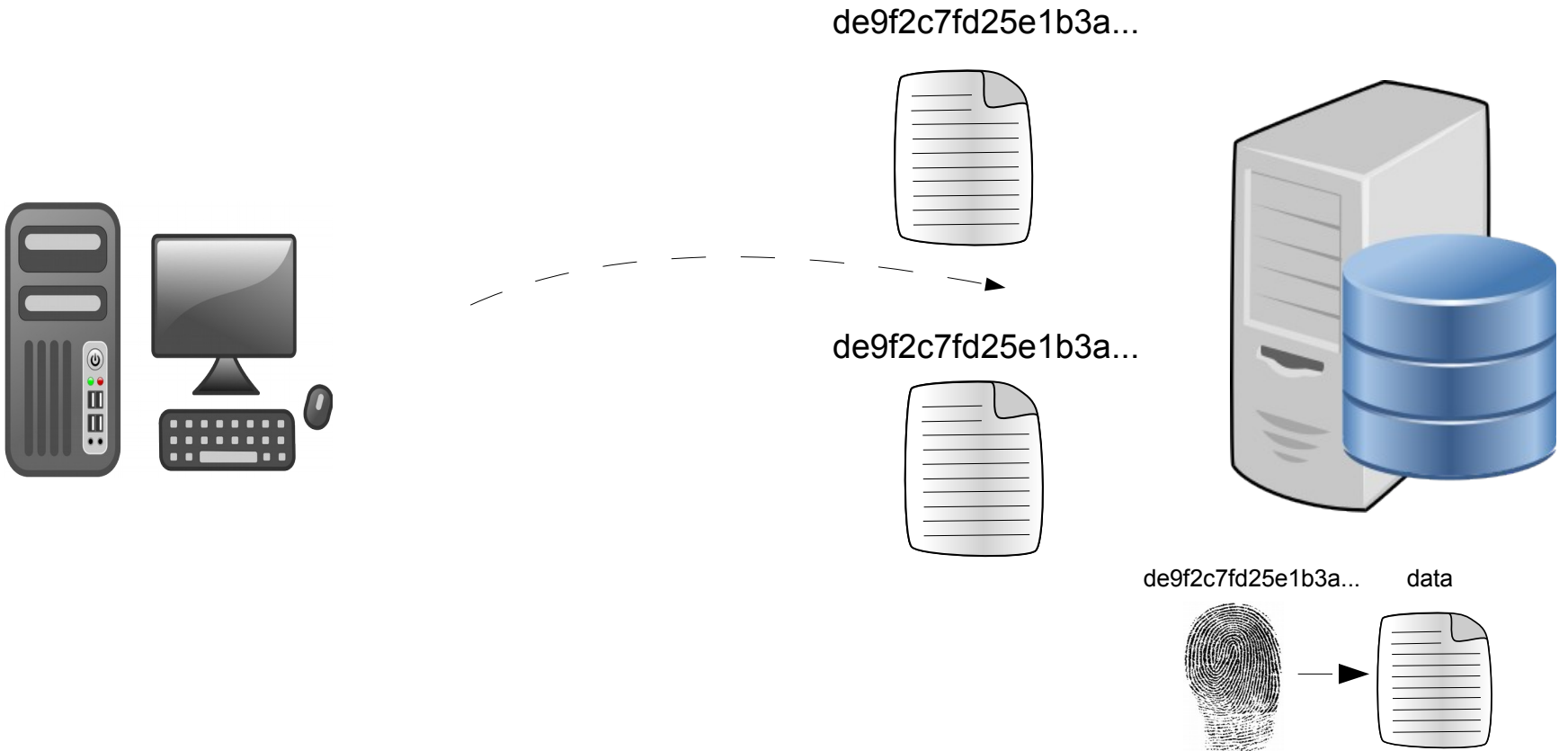Name      de9f2c7fd25e1b3a...

homework.txt

de9f2c7fd25e1b3a...  data

# CAS

Example:

I submit my homework, and my "buddy"
Harold also submits my homework...
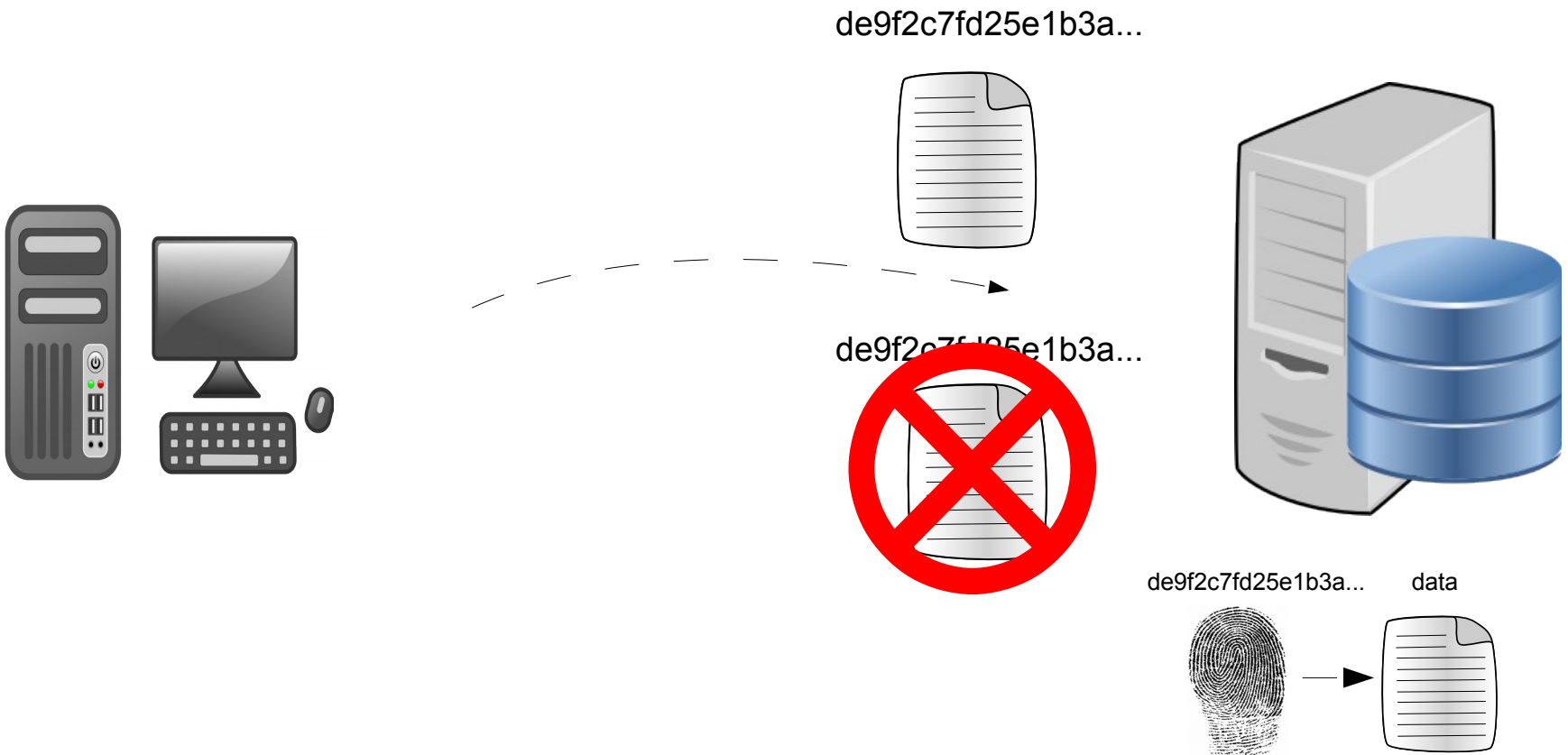
# CAS

Example:

Same contents, same fingerprint.

de9f2c7fd25e1b3a...

de9f2c7fd25e1b3a...

de9f2c7fd25e1b3a...    data

# CAS

Example:

Same contents, same fingerprint.

The data is only stored once!

de9f2c7fd25e1b3a...
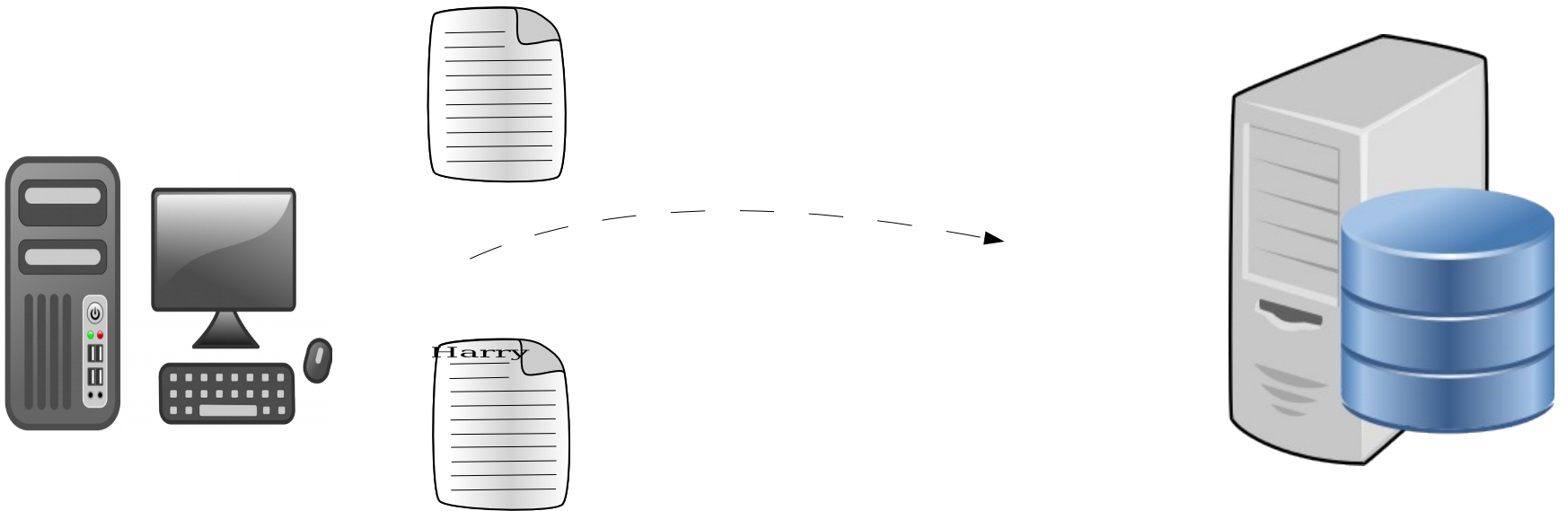
de9f2c7fd25e1b3a...

de9f2c7fd25e1b3a...    data
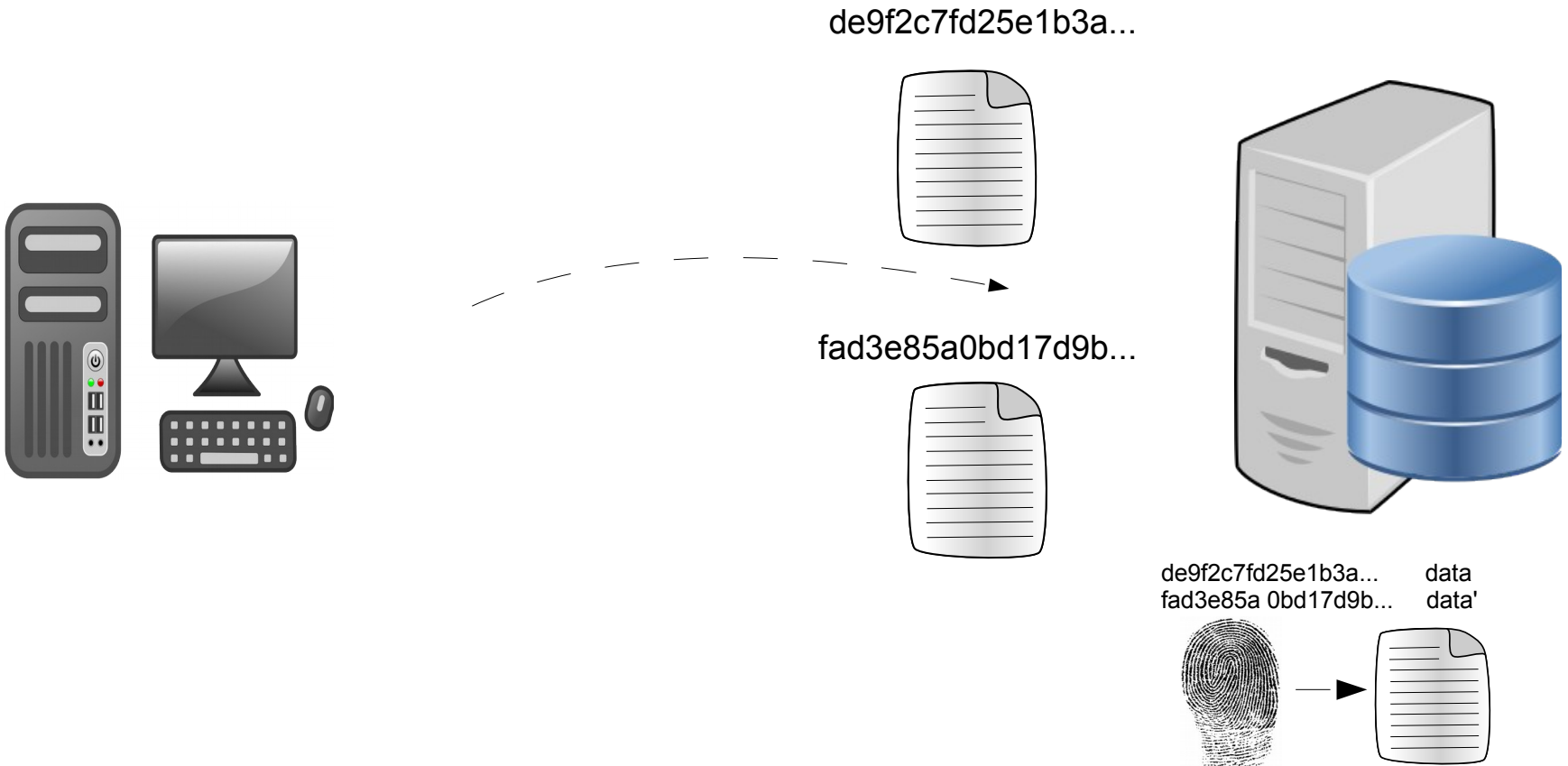
# Background

Background

# CAS

Example:

Now suppose Harry writes his name at the top of my document.

# CAS

Example:

The fingerprints are completely different, despite the (mostly) identical contents.

de9f2c7fd25e1b3a...

fad3e85a0bd17d9b...

de9f2c7fd25e1b3a...    data
fad3e85a 0bd17d9b...    data'

# CAS

**Problem Statement**:

What is the appropriate granularity to address our data?

What are the tradeoffs associated with this choice?

# Background

# Deduplication

**Chunking** breaks a data stream into segments

SHA1( DATA )  becomes

SHA1( CK1 ) + SHA1( CK2 ) + SHA1( CK3 )

How do we divide a data stream?

How do we reassemble a data stream?

# Deduplication

**Division.**

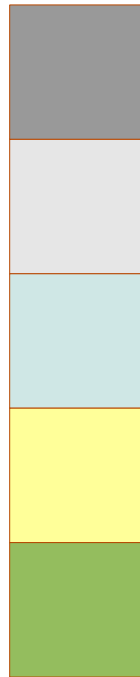Option 1: fixed-size blocks

- Every (?)KB, start a new chunk

Option 2: variable-size chunks

- Chunk boundaries dependent on chunk contents

# Deduplication

**Division:** fixed-size blocks

# Deduplication

**Division:** fixed-size blocks

hw-bill.txt          hw-harold.txt

=|=

=|=

=|=

=|=

=|=

=|=

Harold

Suppose Harold adds his name
to the top of my homework

This is called the
**boundary shifting
problem**.

# Deduplication

**Division.**

Option 1: fixed-size blocks

- Every 4KB, start a new chunk

Option 2: variable-size chunks

- Chunk boundaries dependent on chunk contents
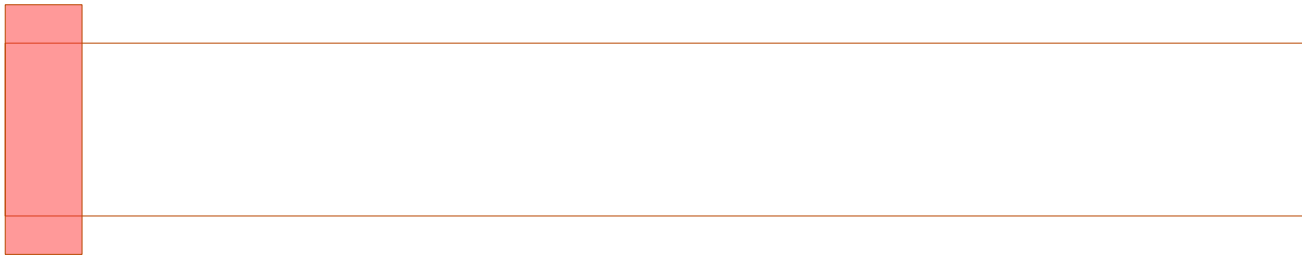
# Deduplication

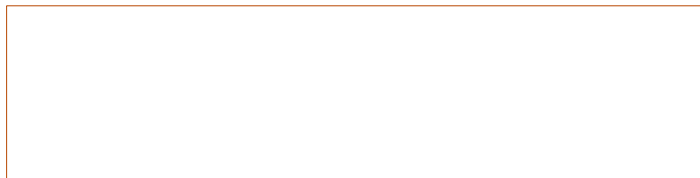**Division:** variable-size chunks

parameters:

Window of width $w$
Target pattern $t$

- Slide the window byte by byte across the data, and compute a window fingerprint at each position.

- If the fingerprint matches the target, $t$, then we have a **fingerprint match** at that position

# Deduplication

**Division:** variable-size chunks
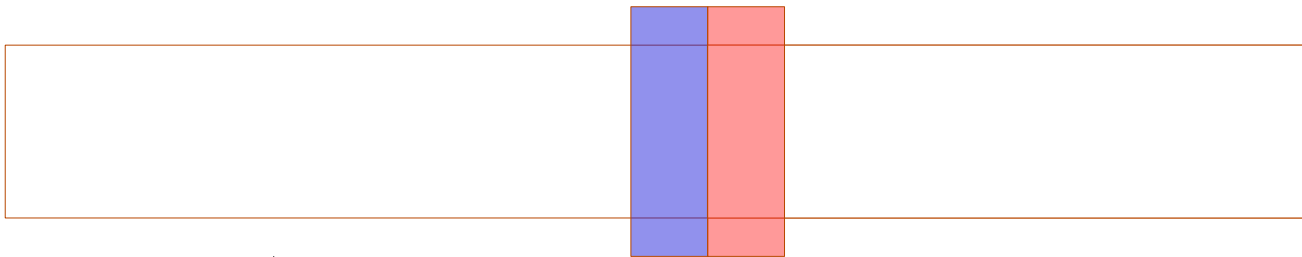
- Slide the window byte by byte across the data, and
  compute a window fingerprint at each position.

- If the fingerprint matches the target, $t$, then we
  have a **fingerprint match** at that position

# Deduplication

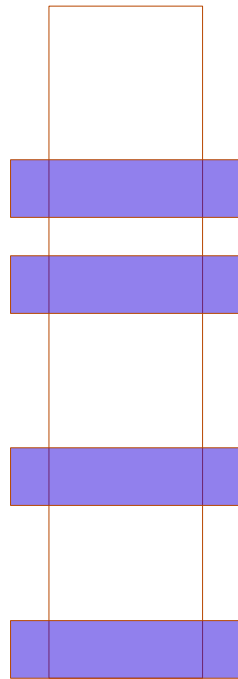**Division:** variable-size chunks

# Deduplication

**Division:** variable-size chunks

hw-wkj.txt          hw-harold.txt

=|=                 Harold ← Suppose Harold adds his name to the top of my homework

Only introduce one new chunk to storage.

# Deduplication

**Division:** variable-size chunks

Sliding window properties:

- collisions are OK, but
  - average chunk size should be configurable

- reuse overlapping window calculations

➡ Rabin fingerprints

# Deduplication

**Division:** variable-size chunks

Rabin fingerprint: preselect divisor $D$, and an irreducible polynomial

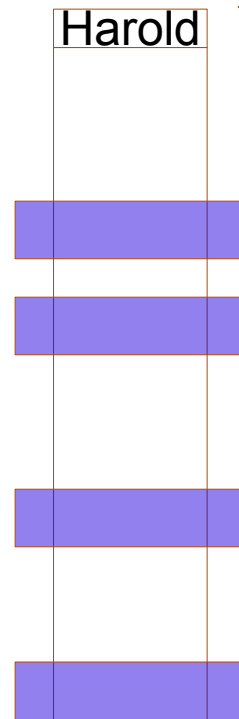$$\mathbf{R}(b_1, b_2, ..., b_w) = (b_1 p^{w-1} + b_2 p^{w-2} + \ldots + b_w) \textbf{ mod } D$$

$$\mathbf{R}(b_i, ..., b_{i+w-1}) = ((\mathbf{R}(b_{i-1}, \ldots, b_{i+w-2}) - b_{i-1} p^{w-1})p + b_{i+w-1}) \textbf{ mod } D$$

| Arbitrary window of width $w$ | previous window calculation | previous first term |
|---|---|---|

# Deduplication

**Recap:**

Chunking breaks a data stream into smaller segments

→ What do we gain from chunking?

→ What are the tradeoffs?

+ Finer granularity of sharing

+ Finer granularity of addressing

- Fingerprinting is an expensive operation

- Not suitable for all data patterns

- Index overhead

# Deduplication

**Reassembling chunks:**

Recipes provide directions for reconstructing files from chunks

# Deduplication

**Reassembling chunks:**

Recipes provide directions for reconstructing files from chunks

| Metadata |
| --- |
| <SHA1> |
| <SHA1> |
| <SHA1> |
| ... |

| DATA BLOCK | | DATA BLOCK | | DATA BLOCK |

# CAS

Example:



| Name | de9f2c7fd25e1b3a... | | de9f2c7fd25e1b3a... | recipe/data |

Name

de9f2c7fd25e1b3a...

homework.txt

| Metadata |
| --- |
| <SHA1> |
| <SHA1> |
| <SHA1> |
| ... |

de9f2c7fd25e1b3a...

recipe/data

???

# Deduplication

# Deduplication

**The Index:**

SHA-1 fingerprint uniquely identifies data, but
the index translates fingerprints to chunks.

$$<sha\text{-}1_1> \longrightarrow <chunk_1>$$
$$<sha\text{-}1_2> \longrightarrow <chunk_2>$$
$$<sha\text{-}1_3> \longrightarrow <chunk_3>$$
$$\ldots \qquad\qquad \ldots$$
$$<sha\text{-}1_n> \longrightarrow <chunk_n>$$

$<chunk_i>$ = {location, size?, refcount?, compressed?, ...}

# Deduplication

**The Index:**

For small chunk stores:
- database, hash table, tree

For a large index, legacy data structures won't fit in main memory
- each index query requires a disk seek

- why?
SHA-1 fingerprints independent and randomly distributed
- no locality

Known as the **index disk bottleneck**

# Deduplication

**The Index:**

Back of the envelope:

Average chunk size: 4KB
Fingerprint: 20B

20TB unique data = 100GB SHA-1 fingerprints

# Deduplication

**Disk bottleneck:**

Data Domain strategy:
- filter unnecessary lookups
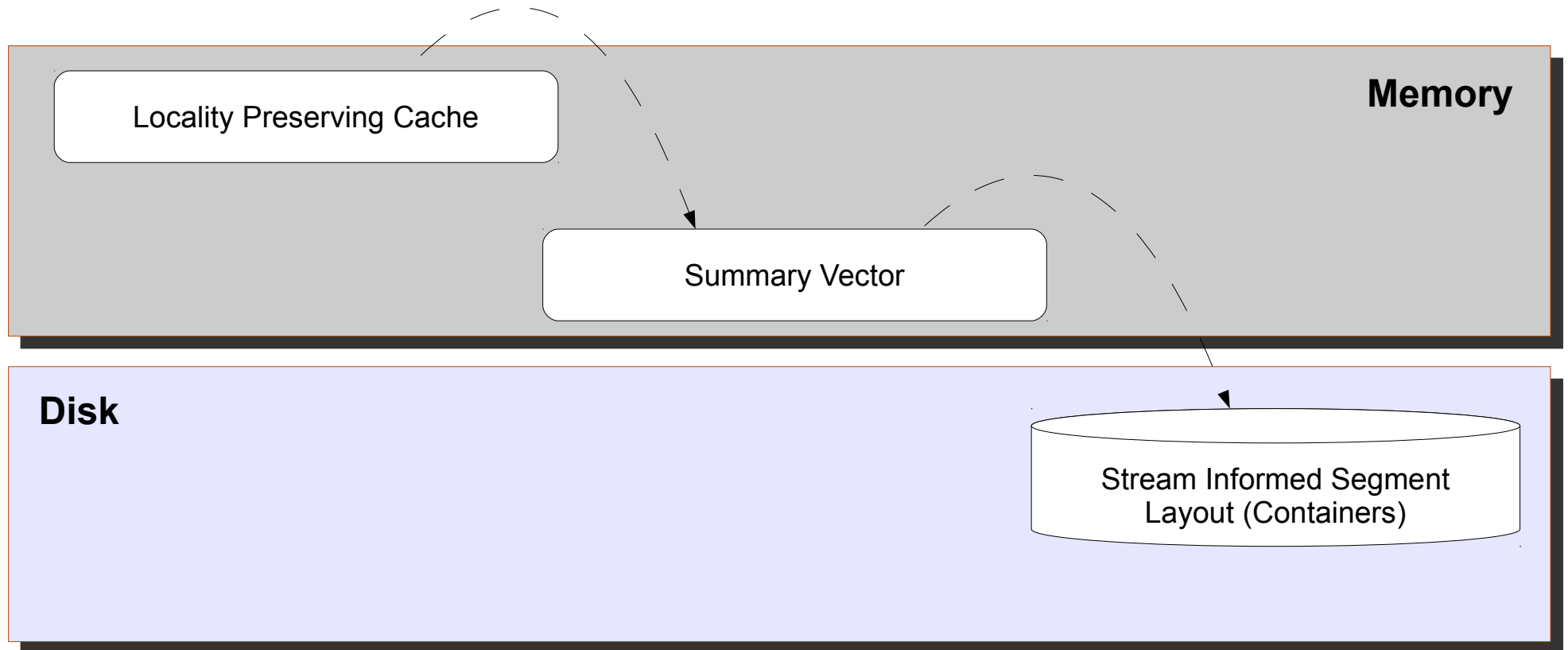- piggyback useful work onto the disk lookups that *are* necessary

**Memory**

Locality Preserving Cache

Summary Vector

**Disk**

Stream Informed Segment Layout (Containers)

# Deduplication

**Disk bottleneck:**

Summary vector

- Bloom filter (any AMQ data structure works)

| | ... | 1 | 0 | 1 | 1 | 1 | 0 | 0 | 1 | 0 | 1 | 0 | 1 | 1 | 1 | 0 | 1 | ... |

$h_1$

$h_2$

$h_3$

**Filter properties:**
- No false negatives
  - if an FP is in the index, it is in summary vector
- Tuneable false positive rate
  - We can trade memory for accuracy

Note: on a false positive, we are no worse off
  - We just do the disk seek we would have done anyway

# Deduplication

**Disk bottleneck:**

Data Domain strategy:
- filter unnecessary lookups
- piggyback useful work onto the disk lookups that *are* necessary

# Deduplication

**Disk bottleneck:**

Stream informed segment layout (SISL)

- variable sized chunks written to fixed size containers
- chunk descriptors are stored in a list at the head
    →"temporal locality" for hashes within a container



Principle:

- backup workloads exhibit **chunk locality**

# Deduplication

**Disk bottleneck:**

Data Domain strategy:
- filter unnecessary lookups
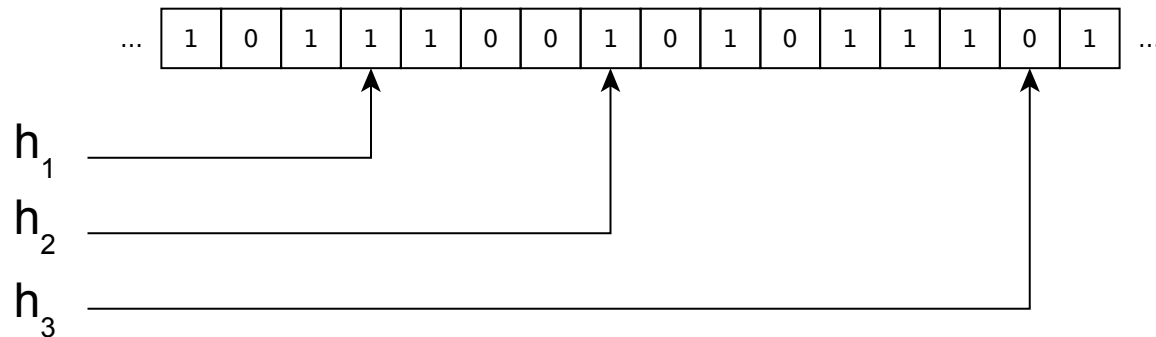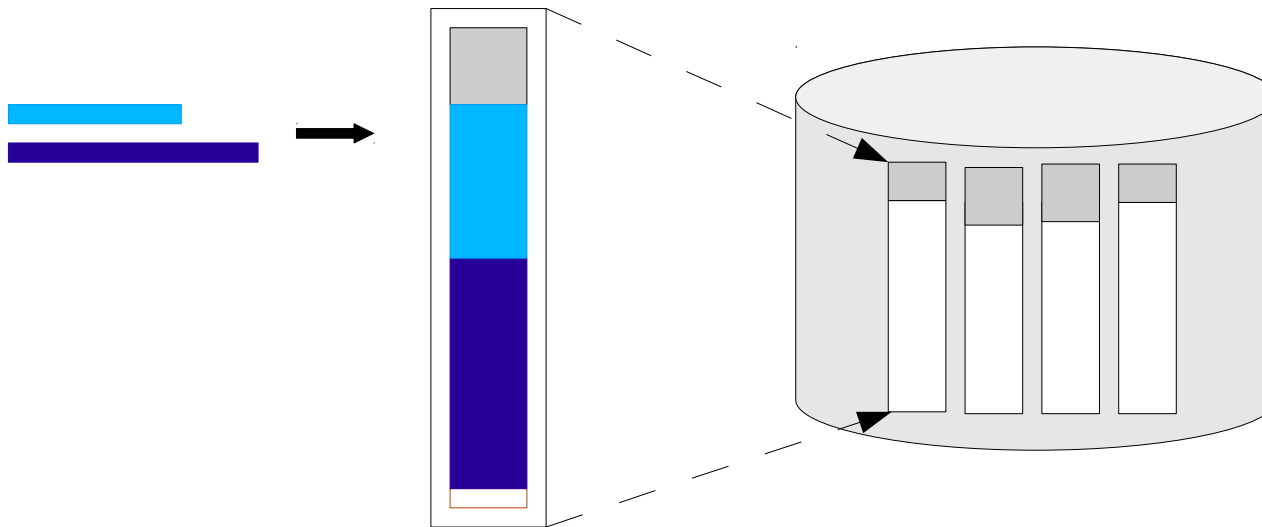- piggyback useful work onto the disk lookups that *are* necessary

**Memory**

Locality Preserving Cache

Bloom Filter

Summary Vector

Group Fingerprints: Temporal Locality

**Disk**

Stream Informed Segment Layout (Containers)

# Deduplication

**Disk bottleneck:**

Locality Preserving Cache (LPC)

- LRU cache of candidate fingerprint groups



On-disk container

Principle:

- if you must go to disk, make it worth your while

# Deduplication

**Disk bottleneck:**



START

Read request for chunk fingerprint

Fingerprint in Bloom filter?

No → No Lookup Necessary → END

Yes

Fingerprint in LPC?

No → On-disk fingerprint index lookup: get container location

Yes → Read data from target container. → END

Prefetch fingerprints from head of target data container. → Read data from target container.

# Deduplication

**Summary:  Dedup and  the 4 W's**

Dedup Goal: eliminate repeat instances of identical data

*What* (granularity) to dedup?

*Where* to dedup?

*When* to dedup?

*Why* dedup?

# Deduplication

**Summary:  Dedup and  the 4 W's**

*What* (granularity) to dedup?

> Hybrid?
> Context-aware.

|  | Whole-file | Fixed-size | Content-defined |
|---|---|---|---|
| **Chunking overheads** | N/A | offsets | Sliding window fingerprinting |
| **Dedup Ratio** | All-or-nothing | Boundary shifting problem | Best |
| **Other notes** | Low index overhead, compressed/ encrypted/ media | (Whole-file)**+**<br><br>Ease of implementation, selective caching, synchronization | Latency, CPU intensive |

# Deduplication

**Summary:  Dedup and  the 4 W's**

*Where* to dedup?

source

destination

Dedup before sending
data over the network
   + save bandwidth
   - client complexity
   - trust clients?

Dedup at storage server
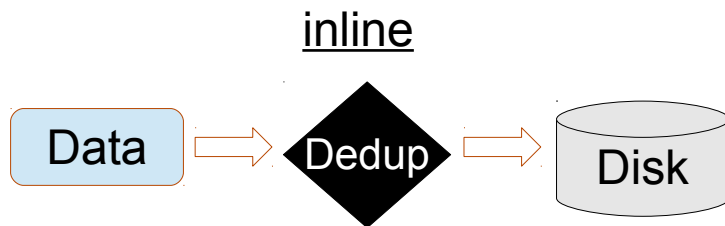   + server more powerful
   - centralized data structures

hybrid

Client index checks membership,
Server index stores location
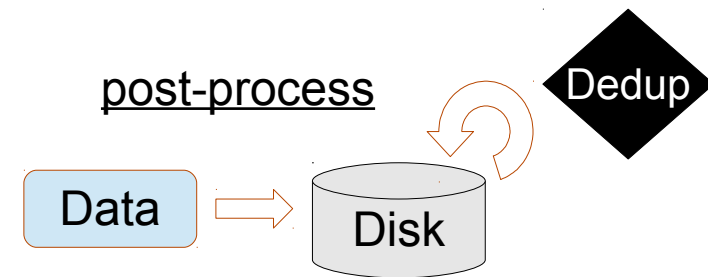
# Deduplication

**Summary:  Dedup and  the 4 W's**

*When* to dedup?



inline

Data → Dedup → Disk

+ never store duplicate data
- slower → index lookup per chunk
+ faster → save I/O for duplicate data

post-process

Data → Disk    Dedup

- temporarily wasted storage
+ faster → stream long writes, reclaim in
                the background
- may create (even more) fragmentation

hybrid

→ post-processing faster for initial commits
→ switch to inline to take advantage of I/O savings

# Deduplication

**Why dedup?**

Perhaps you have a loooooot of data...
- enterprise backups

Or data that is particularly amenable to deduplication...
- small or incremental changes
- data that is not encrypted or compressed

Or that changes infrequently.
- blocks are immutable → no such thing as a "block modify"
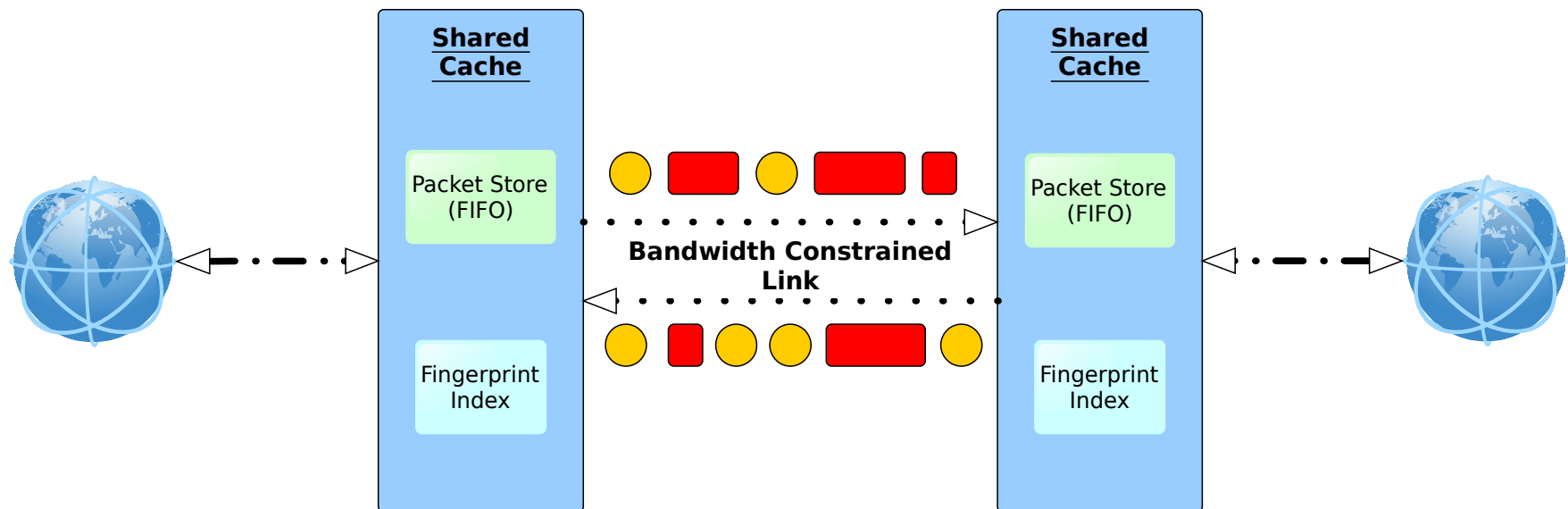- rate of change determines container chunk locality

**Ideal use case: "Cold Storage"**

# Deduplication

**Why dedup?**

Perhaps your bottleneck isn't the CPU

- Use dedup if you can favorably trade other resources



**Example:** Protocol Independent Technique for Eliminating Redundant Network Traffic

# Background

# Other CAS Applications

**Data verification**

CAS can be used to build tamper evident storage. Suppose that:

- you can't fix a compromised server,

- but you never want be fooled by one

**Insight**: Fingerprints uniquely identify data

- hash before storing data, and save the fp locally
- rehash data and compare fps upon receipt