

# A Security Pattern for Data Integrity in P2P Systems

Benjamin Schleinzer

Email: schleinzer (at) comsys.rwth-aachen.de <sup>\*†</sup>

Nobukazu Yoshioka

Email: nobukazu (at) nii.ac.jp <sup>‡</sup>

September 26, 2010

Peer-To-Peer-systems (P2P) brought new methods to distribute large amounts of data to end users. The basic concept is to split larger data into chunks and distribute these chunks among peers so that each chunk can be requested from different peers instead of just one server. This allows increased distribution speed as resources from all participating peers are used, and therefore the workload on own resources is decreased. However, new security risks in terms of confidentiality, integrity, and availability need to be addressed. Security patterns targeted specifically for P2P-systems help designers to identify threats and show appropriate countermeasures.

In P2P-systems a chunk's hash identifies the peer storing the chunk, and each peer is responsible for a specific range of the hash's keyspace. This mechanism is well established in P2P-software and solves routing problems. We introduce a new security pattern for P2P-systems that uses these hashes to guarantee the integrity of the distributed data even though chunks are sent to, received from, and stored by multiple, possible untrustworthy, peers.

## 1 Introduction

P2P-networks can distribute large quantities of data to multiple users without relying on classic client server based scenarios [1]. To use the resources of all participating peers in the network and to distribute the workload among them the data is split into multiple even-sized data blocks called chunks [2]. All peers participating in the network act, while they download content, as servers for the other peers. Therefore chunks are disseminated among the downloading peers, and a missing chunk can be downloaded from either another peer that already downloaded the chunk or, if no such peer exists, directly from the initial service provider. This process is shown in figure 1. The initial service provider (blue laptop) disseminates chunks of his data (blue folder) to peers of the P2P-system (green laptops) which in turn also disseminate the chunks (blue folder trees) to other requesting peers (white laptops) while malicious peers (red laptop) try to spread manipulated chunks (red folder trees).

One goal of P2P-systems is to increase the distribution speed of the data even when only limited resources are available. Downloading content is done by requesting chunks simultaneously from many sources which decreases the use of own resources [3]. A negative effect of this approach is the data storage on untrustworthy peers. Therefore the security, defined as confidentiality, integrity, and availability by [4], in P2P-systems needs to be assessed.

Data availability in P2P-systems is achieved as downloading peers redundantly store identical chunks of data. However, providing confidentiality and integrity proofs to be difficult. Chunks downloaded from and thus stored on different peers are vulnerable to attacks [5]. Possible attacks include: information disclosure (confidentiality) or tampering (integrity). Individual encryption keys for each user provide confidentiality but negate the increased distribution speed which results from spreading identical chunks to all peers. Signatures can be used to provide

---

\*Chair of Communication and Distributed Systems, RWTH Aachen University, Ahornstr. 55, 52074 Aachen, Germany

†This work was partially funded by the DFG research group 733

‡GRACE Center, National Institute of Informatics, 2-1-2 Hitotsubashi, Chiyoda-ku, Tokyo

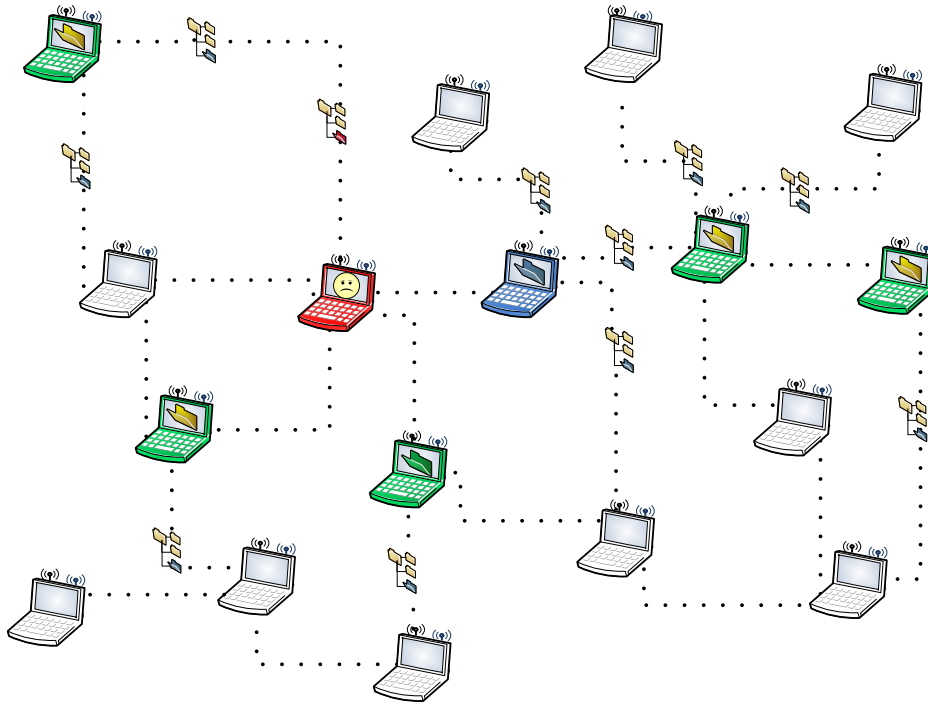


Figure 1: Basic P2P-system

integrity, for individual chunks or the data, but do not provide routing information for finding peers that store a specific chunk. Thus, to address confidentiality and integrity in P2P-system we need different approaches.

We concentrate primarily on integrity and show how it can be ensured using secure hash algorithms. Confidentiality is considered but depends on the intended application and therefore providing it can not be generalized. References for different approaches are given in section 2.10. In contrast to signatures hashes are not only used to identify valid data blocks but also identify peers storing a specific chunk [6]. Additionally the validity of a chunk is computed by the downloading peer alone. Therefore the overall system performance is increased, and communication among peers is minimized.

Section 2 will discuss the pattern in detail motivated by a short example followed by the context in which to apply the pattern and forces to consider. Afterwards in section 2.6 we will go into different details of the pattern. First, for a high-level of abstraction a class diagram is used, second a sequence diagrams shows the flow of information, and third implementation details are given. We consider the reader to have a basic understanding of how P2P-systems work and how hashes are computed. More information on P2P-systems can be found in [1]. For secure hash algorithms refer to [7].

## 2 Data integrity in P2P-systems through hash addressing

### 2.1 Intent

Secure hashes can be used to ensure data integrity. In a P2P-network data is split into smaller chunks, hashed, and spread over multiple peers. Depending on the strength of the hash's algorithm tampering attacks are prevented even though the chunks are stored and handled by untrusted entities. Additionally the storage location of a chunk and the route to it is identified by the hash.

### 2.2 Example

Alice wants to retrieve a large document from Bob that he stored using a P2P-enabled software as Bob has only very limited resources. To store the document he split it into several small chunks and distributed them over the peers of the network. Alice does not know which peers store these chunks. To find them Bob sends her a message

containing a list of the chunks for his document and how to address the peers storing the chunks. Alice now asks peers in the network to send the chunks that belong to Bob's document. Eve wants to spread her worm and instead of sending Alice the requested chunk she will send her worm. Alice has no way of detecting this misuse and to remove or quarantine the worm.

## 2.3 Context

This pattern can be applied by service providers with limited resources, who want to distribute large quantities of data, and want the distribution process to be resilient against tampering attacks and implemented by using a P2P-based application.

## 2.4 Problem

In a P2P-network all peers disseminate the data to different peers to increase the overall available bandwidth. This means that, among others, the data is vulnerable to tampering attacks. However, service providers need to utilize peers to decrease the use of their own resources, but still want to ensure the integrity of the data. An approach is needed that provides security while at the same time negative effects on the performance of the system should be minimal.

## 2.5 Forces

For the proposed solution we have to consider the following forces:

**Minimize own resources for data distribution** Large quantities of data need to be distributed with limited own resources. Instead of relying only on own resources, downloading peers help in disseminating the data by acting as servers to others.

**Fast resource identification** Peers store different chunks of data. If downloads happen from multiple peers simultaneously, identifying a peer that stores a specific chunk has to be done efficiently.

**Data integrity** Data must be secured against tampering, as chunks stored on foreign peers can be altered for different misuse scenarios.

**Data Confidentiality** As with tampering, the confidentiality of the data needs to be ensured.

**Minimal communication with data provider** The original owner of the data has only limited resources, thus the validity of a downloaded chunk has to be identified with minimal communication with the owner of the data.

## 2.6 Solution

It has been shown that strong cryptographic hashes ensure data integrity, and at the same time identify chunks in a P2P-network [6]. These characteristics can be used to prevent tampering attacks in P2P-networks and efficiently distribute data. To do so the initial data provider splits the data into even-sized chunks of smaller data. For each chunk the hash is computed and stored, and users requesting the data get a list with these hashes over a secure connection. Each entry in the list identifies a different chunk that needs to be retrieved. The following sections present the structure of the solution, its dynamics with involved parties and common pitfalls during implementation.

### 2.6.1 Structure

The structure of the secure pattern is given by the class diagram in figure 2. The service provider wants to distribute some data to a large user base. The ChunkGenerator divides the data into chunks and generates a list of hashes. The list, together with the chunks, is sent to the peers participating in the system. Depending on the implementation different peers receive different chunks during this initial dissemination. Other peers use the list to request needed chunks and verify the validity of each chunk using the ChunkVerifier after a successful download. A valid chunk is then stored to either assemble the original data when all chunks are downloaded or send it to requesting peers.

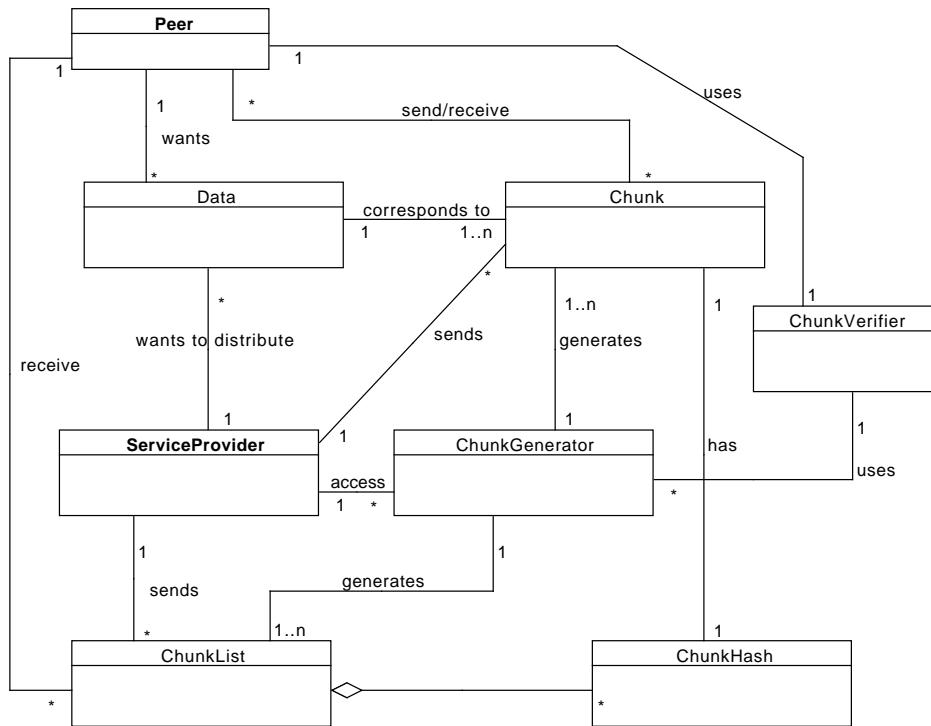


Figure 2: Class diagram for Secure Hash Addressing Pattern

### 2.6.2 Dynamics

The following two sequence diagrams, given in figure 3 and figure 4, show the dynamics for using hashes to download chunks and achieve integrity. The first diagram shows how the list of hashes is generated:

**Summary:** A service provider wants to distribute a large amount of data. The data is divided into equal-sized chunks that are disseminated to peers of the P2P-system.

**Actors:** The initial service provider

**Precondition:** None

**Description:**

1. The data is passed to the ChunkGenerator to be divided into equal-sized chunks.
2. The ChunkGenerator divides the data and computes the hash for each chunk.
3. The ChunkGenerator returns a list of hashes and the generated chunks.
4. The ChunkList is published using a secure method that allows peers to validate the integrity of the list and the identity of the service provider.
5. For each chunk the responsible peer is searched and the chunk is sent to the peer for storage.

**Postcondition:** The chunks are ready for distribution over the network. Clients can request the list of needed chunks.

After the chunks and the list are generated they can be downloaded by all peers. This process and the following verification is illustrated in the sequence diagram shown in Figure 4:

**Summary:** A peer wants to download the data offered by the service provider.

**Actors:** The downloading peer and other peers of the P2P-system.

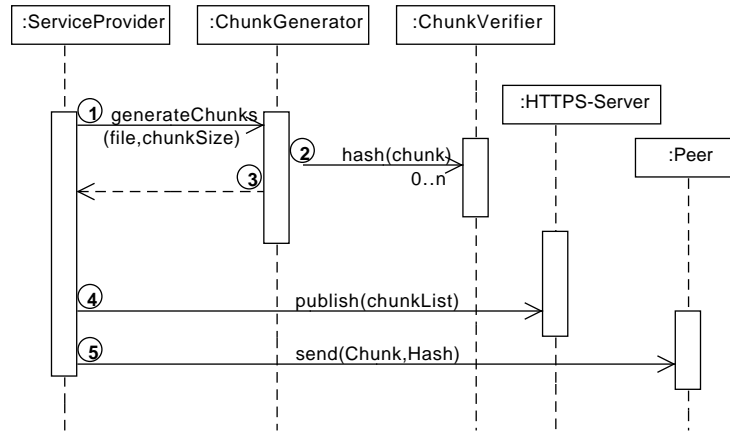


Figure 3: Sequence diagram for computing the list of hashes

**Precondition:** The peers know how to retrieve the list of valid hashes (ChunkList) and use a secure connection that validates the integrity of the list and the identity of the service provider.

**Description:**

1. The peer asks the initial service provider for a ChunkList.
2. The service provider sends the ChunkList via a secure connection.
3. For each chunk the peer nearest to the chunk's hash is requested to send it. The nearest peer has an ID that is smaller than the hash of the chunk and there is no other peer known whose ID is smaller than the chunk's hash and bigger than the ID of the contacted peer.
4. If the requested chunk is not stored by the peer contacted in step three the send request is forwarded to the peer storing it. This process is repeated until the correct peer is found or if no such peer exists the chunk is requested from the initial service provider. For details cf. section 2.7.
5. The chunk is downloaded from the peer storing it.
6. The retrieved chunk is hashed.
  - a) If the computed hash and the hash stored in the list are equal the chunk is verified and gets accepted.
  - b) If the computed hash and the hash stored in the list are not equal the chunk is discarded and downloaded again from a different peer.
7. The process is repeated until all chunks from the list are retrieved and verified.
8. The original data is reassembled.

**Postcondition** The data is downloaded and verified.

**2.6.3 Implementation**

The following problems need to be considered when implementing this pattern. A common mistake is to distribute the list of hashes without securing the transmission. This leaves the list open to attacks which could be used to add malicious chunks or change a chunk's hash. To prevent this attack the download of the list has to be secured against tampering attacks, and the initial data provider's identity has to be verified using the secure pipe pattern, e.g. using SSL encrypted connections [8].

Another problem arises when the algorithm used to compute the hashes is not resilient against collision and preimage attacks [9, 10]. Attackers could perform these attacks to generate valid looking chunks that contain malicious data.

Apart from these problems the peers and the initial service provider have to use the same algorithm to generate the hashes. Only then can the peer verify the downloaded chunk.

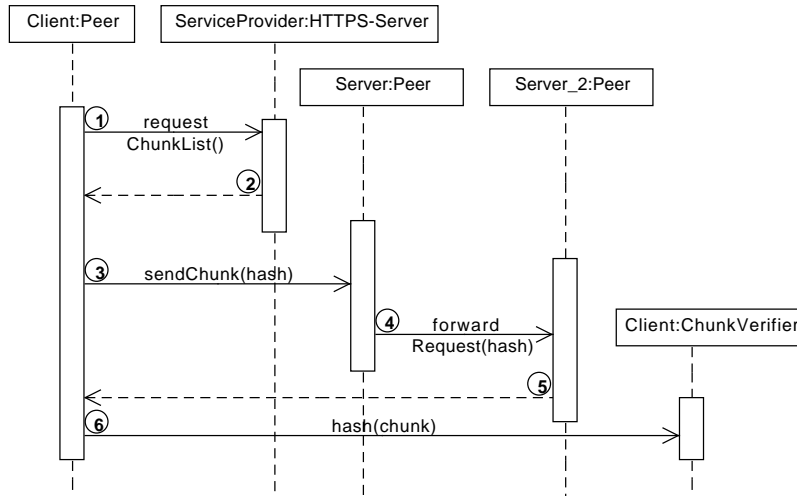


Figure 4: Sequence diagram to download a single chunk of data

## 2.7 Known Uses

Some known implementations for this pattern are, e.g. Freenet [11], CFS [12], Kadmila [13] or Chord [6]. Freenet and CFS are P2P-filesystems that use hashes to identify individual files. Kademila implementations of a Distributed Hash Table (DHT) that was adapted for and used by BitTorrent [14] to find and verify downloadable chunks. Chord also implements a DHT, and provides fast and efficient routing in the P2P-system. More variants of DHTs exist with Can [15], Pastry [16] or Tapestry [17] that have different properties when it comes to route length, number of know peers or exchanged message when a new peer joins.

All DHTs use a ring like communication structure and distribute the hash's keyspace over available peers. Peers are responsible for storing chunks that belong to their part of the keyspace and forward request to all other chunks to other peers. Therefore they always know their direct predecessor and successor and, depending on the implementation, various other peers.

## 2.8 Consequences

Using secure hashes in P2P-systems to verify the integrity of chunks and to address chunks has the following advantages:

- The initial service provider is only contacted when the list of valid hashes is requested or when a chunk is not stored by any peer.
- The integrity of the chunk is validated when the hash of a retrieved chunk is equal to the one stored in the list. No communication to the initial service provider is needed.
- To increase the download speed data is downloaded by requesting multiple chunks from different peers.
- The hash of a chunk is used to locate the peer storing it and to forward routing requests.
- Integrity checks and data assembly can be done by the peers without communicating with other peers.

The pattern presents some liabilities, which in part can be negated employing other patterns (see Section 2.10):

- The list of valid hashes is susceptible to different misuses, e.g. tampering, man in the middle attacks. An appropriate pattern has to be applied to ensure the list's security.
- Some algorithms for computing hashes are vulnerable to collision attacks. Using these algorithms threatens the security of the whole pattern.

- The confidentiality of the data is not protected using this pattern. If confidentiality is a concern this needs to be addressed by employing other patterns.
- If only a limited number of peers exist, the overhead of finding available peers to download the chunk from negates the advantage that arises from using multiple peers.

## 2.9 Example resolved

Alice wants to retrieve a large document from Bob that he stored using a P2P-enabled software as Bob has only very limited resources. To store the document he split it into several small chunks and distributed them over the peers of the network. Alice does not know which peers store these chunks. To find them Bob uses a secure connection to send her a message containing a list of the chunks for his document and how to address the peers storing the chunks. Alice now asks peers in the network to send the chunks that belong to Bob's document. Eve wants to spread her worm and instead of sending Alice the requested chunk she will send her worm. After downloading a chunk Alice computes the chunk's hash, detects Eve tampering attack and discards Eve's chunk. Alice repeats the downloading and verification steps until all chunks on the list have been retrieved and verified. She then reassembles the chunks and has the original document that Bob wanted to send her without using his resources, but ensuring the document's integrity.

## 2.10 Related patterns

Other patterns should be used together with this pattern. To transfer the list of chunks the secure pipe pattern [18] is needed, otherwise tampering attacks are still possible. Hashing and appropriate algorithms are described in [7]. However, recent developments in attacks against these algorithms needs to be taken into account.

## 2.11 Conclusion

We showed how even with limited own resources large amounts of data can be distributed to a wide user base while insuring the integrity of the disseminated data. The data is split into chunks and for each chunk a hash is computed. The initial service provider only communicates the computed hashes and individual chunks are downloaded directly from other peers. Through the hash the integrity of each downloaded chunk and thus the integrity of the original data can be verified. Additionally hashes are used for routing requests to find peers that store specific chunks, which is important when high-performing P2P-systems are designed.

## References

- [1] S. Ratnasamy, P. Francis, S. Shenker, R. Karp, and M. Handley, "A scalable content-addressable network," in *In Proceedings of ACM SIGCOMM*, 2001, pp. 161–172.
- [2] M. Rabin, "Efficient dispersal of information for security, load balancing, and fault tolerance," *Journal of the ACM*, vol. 36, pp. 335–348, 1989.
- [3] S. Saroiu, K. P. Gummadi, and S. D. Gribble, "A measurement study of peer-to-peer file sharing systems," in *Multimedia Computing and Networking (MMCN)*, January 2002.
- [4] M. Bishop, *Introduction to Computer Security*. Addison-Wesley Professional, 2004.
- [5] F. Cornelli, E. Damiani, S. D. Capitani, S. Paraboschi, and P. Samarati, "Choosing reputable servants in a p2p network," in *In Proceedings of the 11th World Wide Web Conference*, 2002, pp. 376–386.
- [6] I. Stoica, R. Morris, D. Karger, F. Kaashoek, and H. Balakrishnan, "Chord: A Scalable Peer-To-Peer Lookup Service for Internet Applications," *Computer Communication Review*, vol. 31, no. 4, pp. 149–160, Oct. 2001.
- [7] B. Schneier, *Applied cryptography (2nd ed.): protocols, algorithms, and source code in C*. New York, USA: John Wiley & Sons, Inc., 1995.

- [8] E. Rescorla, *SSL and TLS: designing and building secure systems*. Boston, MA, USA: Addison-Wesley Longman Publishing Co., Inc., 2001.
- [9] X. Wang, X. Lai, D. Feng, H. Chen, and X. Yu, “Cryptanalysis of the hash functions md4 and ripemd,” in *In Proceedings of Eurocrypt ’05, volume 3494 of LNCS*. Springer-Verlag, 2005, pp. 1–18.
- [10] B. S. P. Hoffman, “Attacks on cryptographic hashes in internet protocols,” RFC 4270, November 2005. [Online]. Available: <http://www.ietf.org/rfc/rfc4270.txt>
- [11] I. Clarke, O. Sandberg, B. Wiley, and T. W. Hong, “Freenet: A distributed anonymous information storage and retrieval system,” in *Designing Privacy Enhancing Technologies, volume 2009 of LNCS*. Springer-Verlag, 2001, pp. 46–66.
- [12] F. Dabek, M. F. Kaashoek, D. Karger, R. Morris, and I. Stoica, “Wide-area cooperative storage with CFS,” in *SOSP ’01: Proceedings of the eighteenth ACM symposium on Operating systems principles*. New York, NY, USA: ACM, 2001, pp. 202–215.
- [13] P. Maymounkov and D. Mazières, “Kademlia: A peer-to-peer information system based on the xor metric,” in *IPTPS ’01: Revised Papers from the First International Workshop on Peer-to-Peer Systems*. London, UK: Springer-Verlag, 2002, pp. 53–65.
- [14] J. Pouwelse, P. Garbacki, D. Epema, and H. Sips, “The bittorrent p2p file-sharing system: Measurements and analysis,” in *In Proceedings of Peer-to-Peer Systems IV, volume 3640 of LNCS*. Springer-Verlag, 2005, pp. 205–216.
- [15] S. Ratnasamy, P. Francis, M. Handley, R. Karp, and S. Shenker, “A scalable content-addressable network,” in *SIGCOMM ’01: Proceedings of the 2001 conference on Applications, technologies, architectures, and protocols for computer communications*, New York, NY, USA, October 2001, pp. 161–172.
- [16] A. Rowstron and P. Druschel, “Pastry: Scalable, decentralized object location and routing for large-scale peer-to-peer systems,” in *IFIP/ACM International Conference on Distributed Systems Platforms (Middleware)*, Nov. 2001, pp. 329–350.
- [17] B. Y. Zhao, J. Kubiatowicz, A. D. Joseph, B. Y. Zhao, J. Kubiatowicz, and A. D. Joseph, “Tapestry: An infrastructure for fault-tolerant wide-area location and routing,” University of California at Berkeley, Tech. Rep., 2001.
- [18] C. Steel, *Applied J2ee Security Patterns: Architectural Patterns & Best Practices*. Upper Saddle River, NJ, USA: Prentice Hall PTR, 2005.