Security requirement

# Hadoop

Deutsche Telekom Group

Version     1.1
Date        Jul 1, 2021
Status      Released

# Publication Details

Published by
Deutsche Telekom AG
Vorstandsbereich Technology & Innovation
Chief Security Officer

Reuterstrasse 65, 53315 Bonn
Germany

| File name | Document number | Document type |
|---|---|---|
| | 3.98 | Security requirement |

| Version | State | Status |
|---|---|---|
| 1.1 | Jul 1, 2021 | Released |

| Contact | Validity | Released by |
|---|---|---|
| Telekom Security | Jul 1, 2021 - Jun 30, 2026 | Stefan Pütz, Leiter SEC-T-TST |
| psa.telekom.de | | |

Summary

# Table of Contents

# 1. Introduction

This document addresses Hadoop installations storing and using not only pseudo- or anonymized data but also raw data. Installations of the former type require pseudo-/anonymization at the data sources or at the points of ingress to the Data Lake and according to predefined use cases. Installations of the latter type require the implementation of above-baseline measures to ensure adequate protection of sensitive data. The challenge being is to protect both large amounts of data and most sensitive personal data.

Further, within this document an on-premises installation of Hadoop is assumed, in contrast to managed or cloud installations. In the latter cases further requirements, e.g. in regard to commissioned data processing or processes in regard to interfaces to third parties, and strict audit requirements might be essential to define and comply with.

The requirements defined in this document are addressing the aspects of data privacy and data security for big data of "Design Guideline for Big Data projects based on Hadoop" by Dr. Lars Hofmann et al.

While the author of this document on security requirements strives to define distribution-independent requirements, they might be biased towards Hortonworks, as that is the distribution in place at the authors company.

# 2. Hadoop

## 2.1. Choice of Distribution

| Req 1 | Hadoop distributions supporting central management and central security management must be used. |
|---|---|

In practice, it will not be possible to ensure rollout and enforcement of consistent security policies if non-central management solutions are in place.
At the time of writing, Apache Ranger and Apache Sentry were available for central security management and Apache Ambari for central operational management.

ID: 3.98-1/1.1

| Req 2 | Hadoop distributions supporting authentication and authorization of all components of a Hadoop installation must be used. |
|---|---|

To ensure complete and consistent implementation of security policies as well as to prevent loopholes, authentication of users and services in regard to data ingest, data access, and for execution of jobs must be possible.

ID: 3.98-2/1.1

| Req 3 | Hadoop distributions must support the forwarding of user and service rights. |
|---|---|

A user or service must access data directly, e.g. data stored on a data node should not accessed directly by a user, but through some other software component from the Hadoop set of components. In such typical case, the access rights of the original user must be used to grant data access, not the access rights of the respective other software component

ID: 3.98-3/1.1

## 2.2. Architecture

| Req 4 | Hadoop-specific network segments and systems must be separated from other segments and systems by stateful packet filters and by using separate VLANs. Ingress nodes can share a dedicated network segment with high volume data sources without a separate stateful packet filter, if additional security measures are implemented. |
|---|---|

For Hadoop's security framework to be effective, it is crucial to use all and only components from the Hadoop set of components. To prevent circumvention of security measures, such components must be segregated from non-Hadoop components at least at the data link layer, for example by using separated and dedicated VLANs, and by using stateful packet filters for connection of different network segments.
"Separation" regarding this document implies allowing only necessary and explicitly approved traffic from machines in one segment to machines in other segments while blocking all other traffic, which in turn implies the necessity to create a maintain a communication matrix.
Due to large amounts of data potentially being ingested by specific data sources, such high volume data sources and ingress node can share a common network segment which is used solely for the purpose of transferring data from specific high volume data sources to ingress edge nodes. In such cases, host based packet filters must be used on the ingress edge nodes, e.g. iptables on Linux systems.

ID: 3.98-4/1.1

| Req 5 | At least four network segments must be used to separate traffic; regarding data ingest, data use and data access, monitoring and utility, and name and data nodes; separated by stateful packet fil- |
|---|---|

ters.

Edge nodes are the nodes through which data is ingested in the Data Lake and which control access to data stored in the Data Lake, and all systems of a Hadoop installation require management and monitoring functions. To prevent by-pass of access controls those four different functions (data ingest, data use and data access, monitoring and utility, and master and data nodes from the Data Lake) require segregation from each other as well as segregation from other systems.

To note: Some Hadoop-components might not be clearly assignable to one of the four network segments, e.g. Hue, which might be used as a user access component and therefore is to be placed on an edge node or as a management component, to be placed on a monitoring and utility node.

ID: 3.98-5/1.1

| Req 6 | Edge nodes used to ingest data into the data lake must be dual homed, one leg homed in the data ingest segment and the other leg homed in the name and data node segment if access to such dual homed hosts is secured with an extra layer of security. |
|---|---|

This requirement attributes to the large amounts of ingested data flowing directly from edge nodes to data nodes. For performance reasons another packet filter in this flow might not be feasible. To ensure security in such cases, access to such dual homed edge nodes require extra protection, e.g. a shell control box through which all shell access is granted only through a proxy and which whitelists only allowed commands permitted to be executed on such edge nodes.

Independent from the architectural design decision, i.e. using dual homed or not dual homed edge nodes, access from edge nodes to the data lake must be possible to allow the edge nodes' HDFS clients access to the data lake.

ID: 3.98-6/1.1

| Req 7 | Non-system-management access to data in a Hadoop installation must be performed trough egress edge nodes only. |
|---|---|

For accessing data in the Date Lake, users must connect to edge nodes, they must not access data in the data lake directly. This requirement implies that not all services available in principle in fully accessible installation can be used or that such services require special drivers being aware of edge nodes. For more detailed requirements on egress edge nodes, see Knox configuration.

ID: 3.98-7/1.1

| Req 8 | Systems used for Hadoop components must not be used for non-Hadoop related applications, i.e. Hadoop components must not share a system with application not related to Hadoop. |
|---|---|

To avoid bypass of security policies, resulting in risk of unauthorized data access, such sharing of systems must not take place, neither by having Hadoop components as well as non-Hadoop related applications installed under the same operating systems nor by using different virtual machines on the same platform.

ID: 3.98-8/1.1

| Req 9 | Master services must run on master nodes separated from worker services, which must run on worker nodes, separated from management services, which must run on management nodes, and user-facing services, which must run on edge nodes. |
|---|---|

The Hadoop architecture often distinguishes between master services/nodes and worker services/nodes, for example name nodes and data nodes. While name nodes manage access to data, data nodes actually store blocks of data. Because of the importance of master services they should not run on the same machines as worker services, as such separation allows for:

- Resource contention: E.g., user jobs exhausting resources on a worker node will not degrade essential master nodes services or even cause a denial-of-service to master node services.
- Reduction of the effects of privilege escalation: E.g., master services require different authorizations than worker nodes; in case of mixed master/worker services on one node access rights for one type of service might be used to access data of the other type of service.
- Containment of security vulnerabilities: E.g., vulnerabilities specific to worker nodes will not influence essential master node services.
- One might like to consider installing different master services on different machines for added reliability and added protection of services, resulting again in resource contention and containment of security vulnerabilities.

Non-exhaustive list of master services:
- HDFS name node, secondary name node (standby name node), failover controller, journal node
- MapReduce job tracker and failover controller
- YARN resource manager and job history server
- Hive metastore server
- ZooKeeper server
- HBase master

Accumulo master, tracer, and garbage collector
Non-exhaustive list of worker services:
- HDFS data node
- MapReduce task tracker
- YARN node manager
- HBase region server
- Accumulo tablet server
- Solr server

Non-exhaustive list of management services:
- configuration management
- monitoring
- alerting
- software repositories
- backend databases
- DHCP servers

Non-exhaustive list of edge node services:
- HDFS HttpFS and NFS gateway
- KNOX
- client CLI tools

ID: 3.98-9/1.1

## 2.3. Choice of Hadoop Software Components, Use of Plugins

| Req 10 | Components from the Hadoop set of components must be used to cover the security framework blocks of administration (central management and consistent security), authentication (authentication of users and systems), authorization (provision access to data), audit (maintain records of data access), data protection (protection of data at rest an in motion), and perimeter security & security |

baseline (protection of network boundaries and systems).

All software components must supplement and extend each other to ensure security, which is only ensured for components from the Hadoop set of components. Non-Hadoop components (i.e. components not part of the Hadoop set of components) might or might not fit, but the risk of incompatibilities and lack of long-term maintenance is quite high. Exception: As currently, i.e. at the time of writing, no Hadoop component exists covering pseudo-/anonymization, such component must be developed separately and integrated into the Hadoop installation.

ID: 3.98-10/1.1

| Req 11 | All non-security Hadoop components (i.e. all components allowing access to data) must use plug-ins enforcing the security policy. |

Some components can be used without enforcement of security policies. Such approach might be sensible in cases where only public, generally accessible data is used. Therefore, all non-security components must use according plug-ins, allowing for policy enforcement.

ID: 3.98-11/1.1

## 2.4. Encryption of Data in Rest and in Motion

| Req 12 | Credentials must be encrypted during transmission. |

For example, HTTPS, SSH, or SFTP must be used in favour of HTTP, Telnet, or FTP.

ID: 3.98-12/1.1

| Req 13 | Transmission of data between components in regard to system management or security management must be encrypted. |

All such data, e.g. data transmitted by Apache Ambari to its clients, must be protected in regard to confidentiality and integrity. Therefore, for such communications, TLS must be configured and used.

ID: 3.98-13/1.1

| Req 14 | Data transmitted to or accessed from the outside of the Hadoop installation must be encrypted. |

To ensure confidentiality and integrity of sensitive data, all such data must be encrypted.

ID: 3.98-14/1.1

| Req 15 | Administrative access to any Hadoop component must be encrypted. |

To ensure confidentiality and integrity of sensitive data, all such data must be encrypted.

ID: 3.98-15/1.1

| Req 16 | Data transmitted within the Hadoop installation must be encrypted if classified "confidential" or "strictly confidential". Such encryption can be omitted in case the underlying cables, physical layer infrastructure and data link layer infrastructure is used exclusively by Hadoop systems within one Hadoop network segment and such infrastructure does not cross site boundaries. |

To ensure confidentiality and integrity of sensitive data, all such data must be encrypted.

The preferred option to do so is to use end-to-end encryption, i.e. encryption should be performed on the endpoints of the Hadoop components. Nevertheless, such data can be transmitted using encrypted site-to-site VPNs, if the VPN endpoints terminate within the Hadoop installation or the firewall separating the Hadoop installation's network segments. The exception stated addresses cases where data is transferred from one node to another and both respective nodes are housed within the same rack, racks in the same data centre room or adjacent data centre rooms, and a private networks infrastructure with own cables and switches is installed.

ID: 3.98-16/1.1

| Req 17 | Data classified "confidential" or "strictly confidential" must be stored encrypted on all parts of the Hadoop installation. |
|--------|--------------------------------------------------------------------------------------------------------------------------------|

To ensure confidentiality and integrity of sensitive data, all such data must be encrypted. Compliance with this requirement will also ensure compliance with standards requiring to keep data inaccessible even if some hat physical access to a server, e.g. PCI-DSS or HIPPA.

ID: 3.98-17/1.1

## 2.5. Data, Data Access, and Roles

| Req 18 | During transformation of data from the raw data stage into cleansed data stage, all data must be labelled. |
|--------|------------------------------------------------------------------------------------------------------------|

As authorization of data access depends on role, stage, and label, all non-raw data must be labelled.

ID: 3.98-18/1.1

| Req 19 | Labelling of data should reflect information classification principles as defined by corporate policy. |
|--------|--------------------------------------------------------------------------------------------------------|

Each type of label must be part of one class according to corporate policy.

ID: 3.98-19/1.1

| Req 20 | Data not pseudo- or anonymized must be stored in the Hadoop file system separately from other data. |
|--------|-----------------------------------------------------------------------------------------------------|

Such separation can be ensured by implementing stringent staging, which stores data of different classification in different directories, combined with according policies or ACLs.

ID: 3.98-20/1.1

| Req 21 | The Hadoop file system must be configured to allow access only if such access is controlled via the central security management system. |
|--------|-----------------------------------------------------------------------------------------------------------------------------------------|

E.g., data nodes must not accept requests access to data from users directly logged in a data node. It is understood that administrative users might have more extensive access rights, possibly allowing for circumvention of the central management system.

ID: 3.98-21/1.1

| Req 22 | Assuming roles for production, testing, and development is mutually exclusive, e.g. one person cannot assume roles for both production and testing or for all three group of roles, within one Hadoop-installation. |
|--------|--------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|

Implementation of this requirement hinders circumvention of the role concept. One person can be developer on a development installation and assuming a production role on a production system.

ID: 3.98-22/1.1

| Req 23 | User must not access data nodes directly, unless for system administrative purposes. Such access must take place via edge nodes only. |
|---|---|

For the defined Hadoop security framework being effective it is essential to strictly gate any non-administrative access through the defined gates, i.e. the edge nodes, only.

ID: 3.98-23/1.1

## 2.6. Basic System Configuration

| Req 24 | Hadoop installation must be configured to run in "secure mode". |
|---|---|

By default Hadoop runs in non-secure mode in which no actual authentication is required. By configuring Hadoop to run in secure mode, each user and service needs to authenticate using Kerberos in order to use Hadoop services.

ID: 3.98-24/1.1

| Req 25 | Big Data installations is the ability to quickly add and remove nodes. Assignment of configuration data like IP addresses must be centrally controlled e.g. by DHCP. |
|---|---|

Key advantage of Big Data installations is the ability to quickly add and remove nodes to adapt to current processing and storage requirements. To do so in a controlled and secure way and to ensure consistent and central system management, manual assignment of configuration data like IP addresses is not feasible. Therefore, DHCP should be used for data nodes.

ID: 3.98-25/1.1

| Req 26 | For remote administration of nodes of a Hadoop installation Secure Shell (SSH) must be used via Admin-LAN, but no other protocols. Exception: In case the Hadoop distribution used provides integrated means for administration to ensure central administration, those means can be used, even if not SSH, as long as they use encrypted protocols. |
|---|---|

To ensure effective security of Big Data installations central administration and management is essential. Therefore, other protocols than SSH are acceptable, as long as the protocols used transmit user data encrypted.

ID: 3.98-26/1.1

| Req 27 | Each Hadoop component, i.e. each daemon, must run as a different Unix user. |
|---|---|

HDFS-, YARN-, MapRedue JobHistory- and other daemons must all run as different Unix users.

ID: 3.98-27/1.1

| Req 28 | Unix users of Hadoop components, i.e. each daemon, must not be the root user. |
|---|---|

As root is the privileged system user it must not be used for Hadoop components as vulnerabilities in a component can result in much greater damage and in circumvention of authorization policies if the root user is used.

ID: 3.98-28/1.1

| Req 29 | Hadoop components, i.e. each daemon, must share a common Unix group. |
|---|---|

Recommended is the use of "hadoop" for the name of such group.

ID: 3.98-29/1.1

| Req 30 | All hosts of a Hadoop installation must be time-synchronized with trusted timeservers. |
|---|---|

As Kerberos is essential for Hadoop security and as Kerberos is time-sensitive, all hosts must be time-synchronized with trusted timeservers. Such trust will usually be accomplished only by using timeservers controlled by a Deutsche Telekom group company. If clocks differ by five or more minutes, authentication will fail.

ID: 3.98-30/1.1

| Req 31 | Certificates used to secure access to components from outside a production Hadoop installation must be issued by a standard certification authority. |
|---|---|

While the use of self-signed certificates or other certificates not signed by a recognized certification authority is technically possible, for production systems only certificates issued by standard certification authorities must be used, if used for securing access from outside of the respective Hadoop installation.
As such standard certification authority a globally operating one should be used, e.g. TeleSec of Deutsche Telekom.

ID: 3.98-31/1.1

| Req 32 | Authentication of users of management tools must be based on existing central directory services. |
|---|---|

In most cases, such central directory service will be an existing Active Directory server used for authentication of users accessing office systems like workstations or of users accessing production systems.
It is understood that administrators must login as local administrators on some systems for some tasks. Nevertheless, login on a management system using existing central directory services must occur before.

ID: 3.98-32/1.1

| Req 33 | The central security management of Hadoop (e.g. Apache Ranger, used by Hortonworks) must synchronise its user base via TLS. |
|---|---|

Typically, many users will directly or indirectly have access to parts of the Data Lake. Therefore, implementing and maintaining another user directory in addition to the one used for authentication to use office infrastructure components will not be feasible; the solution being is to synchronise with an existing directory.

ID: 3.98-33/1.1

| Req 34 | The central security management of Hadoop should synchronise its user base continually; such synchronisation must take place at the very least once per day. Synchronisation of administrative user can be performed manually as required. |
|---|---|

Such synchronisation is necessary to have changes in the user base, e.g. by termination of work contracts and new hires.

ID: 3.98-34/1.1

| Req 35 | Auditing must be enabled. |
| --- | --- |

To allow for investigation of misuse and to detect loopholes in the use concept, enabling of auditing is required. Auditing must be enabled at least on the central security management component (which is Apache Ranger on Hortonworks installations).

ID: 3.98-35/1.1

## 2.7. HDFS Configuration

| Req 36 | For HDFS, a umask of 077 must be set. |
| --- | --- |

While any access to data in the HDFS should be controlled by a central security solution anyway, it is good practice to implement the security principle of "defense in depth" by setting according file permissions. In case Apache Ambari is used, such setting can be performed by setting "Service" -> "HDFS" -> "Advanced hdfs-site" and setting "fs.permissions.umask-mode = 077".

ID: 3.98-36/1.1

| Req 37 | The system user used to run the HDFS daemon must not be used for any other service. |
| --- | --- |

The system user used to run the HDFS daemon must not be used for any other service.

ID: 3.98-37/1.1

| Req 38 | Fallback for authentication to local authentication mechanisms must be disabled. |
| --- | --- |

Fallback for authentication to local authentication mechanisms must be disabled. In case Apache Ambari is used, such setting can be performed by setting "Service" -> "HDFS" -> "Advanced ranger-hdf-security" and setting "xasecure.add-hadoop-authorization=false".

ID: 3.98-38/1.1

## 2.8. Kerberos Configuration

| Req 39 | Kerberos must be used as a basis for identification, authentication, and authorization; i.e., Hadoop "kerberos" authentication mechanisms must be used. |
| --- | --- |

Identification, authentication, and authorization are essential parts of granting access to data. As Kerberos is the supported service to do so by Hadoop, it must be used.
In Hadoop available is also the "simple" mechanism, which is insufficient in regard to security for production environments. It might be an option on small test or development systems if only synthetic data is used; in such case, even pseudo- or anonymized data must not be used.

ID: 3.98-39/1.1

| Req 40 | The KDC must be installed on two or more servers, i.e. in a master/slave configuration. |
| --- | --- |

As Kerberos is essential for security functions within a Hadoop installation, it must be configured on two or more servers to allow for redundancy.

ID: 3.98-40/1.1

| Req 41 | /etc/krb5.conf and /var/kerberos/krb5kdc/kdc.conf must contain installation-specific values allow- |
| --- | --- |

ing for secure operation.

Those files can be updated on the KDC host, krb5.conf can then be copied to all other servers. The Deutsche Telekom Security Requirement Cryptographic Algorithms and IP Security (IPSec), must be honored. Care must be taken to remove potentially insecure cryptographic algorithms, cp. "supported enctypes" below. For example, one might need to delete "RC4" from this list.

Implementation example: Example of a krb5.conf:

```
[logging]
 default = FILE:/var/log/krb5libs.log
 kdc = FILE:/var/log/krb5kdc.log
 admin_server = FILE:/var/log/kadmind.log[libdefaults]
 dns_lookup_realm = false
 ticket_lifetime = 24h
 renew_lifetime = 7d
 forwardable = true
 rdns = false
 default_realm = DATA.LAKE
 default_ccache_name = KEYRING:persistent:%{uid}[realms]
 DATA.LAKE = {
  kdc = kdc.data.lake
  admin_server = kdc.data.lake
 }[domain_realm]
 .data.lake = DATA.LAKE
 data.lake = DATA.LAKEExample of a kdc.conf:
[kdcdefaults]
 kdc_ports = 88
 kdc_tcp_ports = 88[realms]
 DATA.LAKE = {
  master_key_type = aes256-cts
  acl_file = /var/kerberos/krb5kdc/kadm5.acl
  dict_file = /usr/share/dict/words
  admin_keytab = /var/kerberos/krb5kdc/kadm5.keytab
  supported_enctypes = aes256-cts:normal aes128-cts:normal des3-hmac-sha1:normal arcfour-hmac:normal camel-
lia256-cts:normal camellia128-cts:normal
  default_principal_flags = +preauth
  max_renewable_life = 7d 0h 0m 0s
 }
```

ID: 3.98-41/1.1

| Req 42 | Every user of a Hadoop-service and every service on each host used by a Hadoop user must have a unique Kerberos user principal name (UPN) and service principal name (SPN), resp. |

The same services running on different hosts must use different SPNs, because if credentials on one node are compromised, other nodes are not affected, and because Kerberos falsely detect reply attacks if by change different services with the same SPN request and use tickets at the same time.

ID: 3.98-42/1.1

| Req 43 | The FQDN must be appended to a service principal's name. |

Adding the FQDN to a service principal's name ensures uniqueness.

ID: 3.98-43/1.1

| Req 44 | For each Hadoop installation, a different Kerberos realm must be used. |

While it is technically possible to use the same Kerberos realm for production, test, and development installations, administrators and user will sooner or later assume they are working in e.g. a test installation, while actually they are

working in the production setup. Such mistakes can mitigate by using different realm names for different installations.

ID: 3.98-44/1.1

| Req 45 | All UPNs and SPNs must be registered in a network-based directory service such as Active Directory or OpenLDAP. |
|---|---|

All users of a cluster must be provisioned on servers of a cluster. Two options are available to do so: Provisioning all users locally, e.g. in the /etc/passwd, or using a network-based directory service. In practice, consistency of accounts can only be ensured if using network-based directory services.

ID: 3.98-45/1.1

| Req 46 | Non-administrative users must not be enabled for shell access. |
|---|---|

Regular users must not utilize a system's shell access. Best practice on Unix-like systems is to assign each such the default shell /sbin/nologin and to use AllowUsers, DenyUsers, AllowGroups, and DenyGroups settings in /etc/ssh_/sshd_config.

ID: 3.98-46/1.1

| Req 47 | Keytab files must be accessible only by the system user running the respective service. Such access must only be read access. |
|---|---|

File permissions of keytab files must be set to 400, the owner of the keytab file must be the system owner running the respective service.
As the keytab files contain credentials unencrypted, those files must be secured.

ID: 3.98-47/1.1

| Req 48 | „Java Cryptographic Extension" must be installed on all nodes. |
|---|---|

This requirement is a hint for administrators setting up a Hadoop installation. This package is required for cluster operation.

ID: 3.98-48/1.1

## 2.9. Knox Configuration (Hortonworks)

| Req 49 | For all but administrative access, Knox must be used for edge nodes allowing access to data in the Data Lake. |
|---|---|

The Apache Knox Gateway („Knox") is a system to extend Hadoop services to users outside of a Hadoop installation while ensuring enforcement of security policies. Further, Knox simplifies security for users who access data from the Data Lake and who execute jobs.
Knox simplifies access by extending Hadoop's REST/HTTP services by encapsulating Kerberos.
Knox enhances security by exposing Hadoop's REST/HTTP services without revealing Hadoop installation network details and by providing TLS out of the box.
Knox centralizes controls by central enforcement of REST API security.
Knox integrates with Hadoop's Kerberos installation, which in turn integrations with a companies identity management solutions like Active Directory.
As of the time of writing, Knox supports the following Hadoop services, subject to specific versions:

- YARN
- WebHDFS

- WebHCat/Templeton
- Oozie
- HBase/Stargate
- Hive (via WebHCat)
- Hive (via JDBC)

ID: 3.98-49/1.1

---

| Req 50 | The security directory, $gateway/data/security, and its contents must be readable and writeable only by the operating system user running the Knox component. |

In this directory the master secret and other confidential information is stored. Therefore, care must be taken to ensure security of this data.

ID: 3.98-50/1.1

---

| Req 51 | Web application security filters available for Knox must be installed and used. |

As of the time of writing, one filter addressing cross site request forgery (CSRF) exists.

ID: 3.98-51/1.1

---

| Req 52 | All requests to Knox must be documented to allow for both developers knowing how to access Knox and administrators of possibly used web application firewalls to configure their systems. |

As of the time of writing, one filter addressing cross site request forgery (CSRF) exists.

ID: 3.98-52/1.1