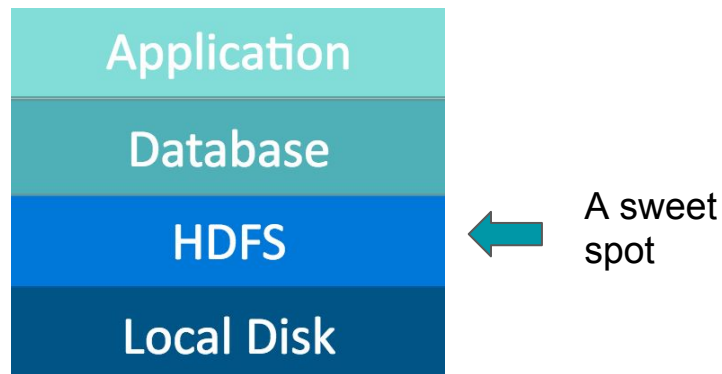# Apache HAWQ
## TDE security integration

hma@pivotal.io 2017-07

**Pivotal**

# Agenda

- TDE introduction
- HAWQ-TDE integration
- Demo
- Future
- Reference
- Q & A

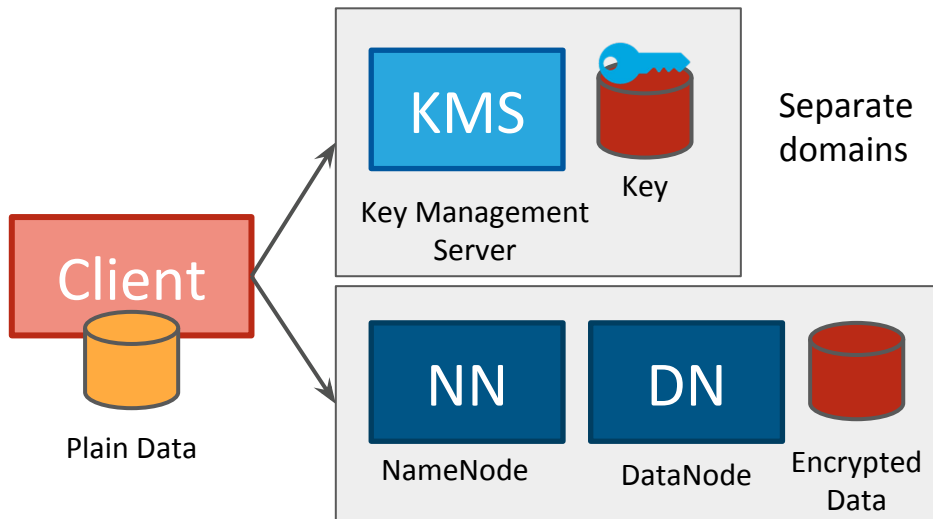# TDE introduction - Attack Vectors

- Physical access
  - Data at-rest on hard drives
- Network sniffing
  - Data in-transit on the network
- Rogue user or admin
  - Insider access to encrypted data, encryption keys



Physical Access

Network Sniffing

Rogue User/Admin

# TDE introduction - Background

- Transparent Data Encryption in HDFS
- https://issues.apache.org/jira/browse/HDFS-6134
- *HDFS implements **transparent**, **end-to-end** encryption. Once configured, data read from and written to special HDFS directories is transparently encrypted and decrypted without requiring changes to user application code.*
- Why Encrypt in HDFS?
  - Excellent performance
  - Safe from OS attacks
  - App transparency
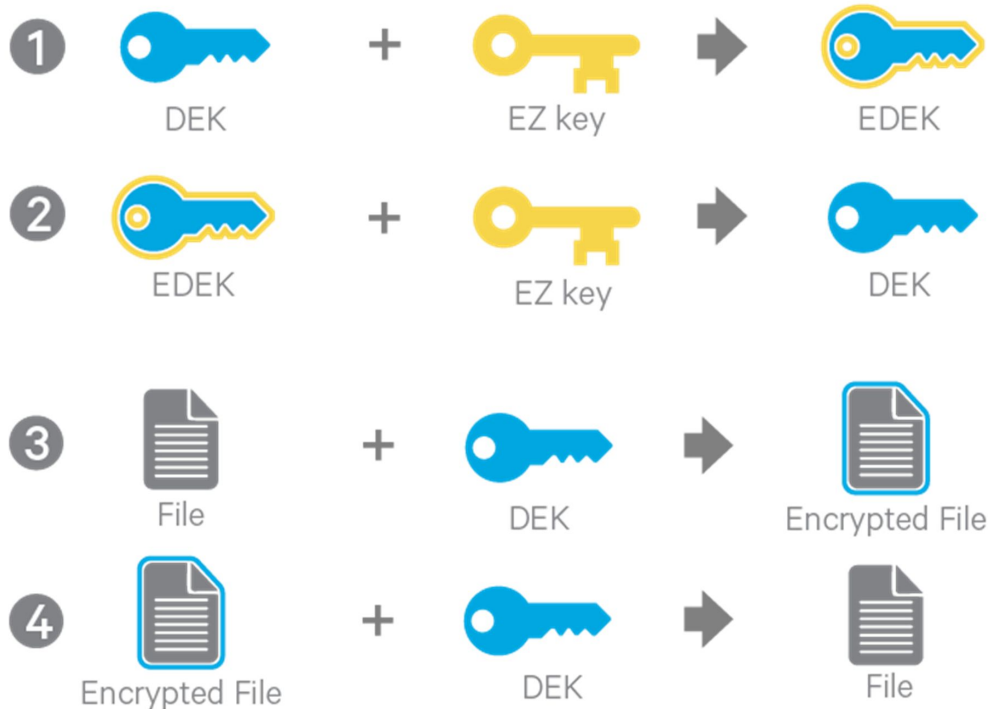  - Easy to deploy



A sweet spot

# TDE introduction - Components

- Hadoop KMS
  - The key server manage encryption keys.
  - Granular per-file encryption keys and per-key ACLs
    - Rogue user can only leak their own data.
  - Encryption algorithm
    - AES-CTR 128/256
    - High-performance
- HDFS
  - Only stores the encrypted data.
  - Encryption Zone, a HDFS directory
- Separate administrative domains for KMS and HDFS
  - Rogue admins do not have access to both keys and data
- API transparency for existing applications



Client

Plain Data

KMS

Key

Key Management Server

Separate domains

NN

DN

NameNode

DataNode

Encrypted Data

# TDE introduction - Keys

- Encryption Zone Key (**EZK**)
  - Unique key per encryption zone
  - key material is stored in the KMS
- Data Encryption Key (**DEK**)
  - **Unique DEK per file**
  - Used by client to encrypt and decrypt data in HDFS
  - Only handled by client and KMS, never HDFS
  - DataNodes never handle plaintext, only ciphertext
- Encrypted Data Encryption Key (**EDEK**)
  - DEK encrypted with corresponding EZK
  - Generated by KMS for NameNode
  - Stored in HDFS metadata as a file xattr

# TDE introduction - Back to Attack Vectors

- Compromised EZK
    - Also requires EDEK and ciphertext
    - Only leaks a single zone
- Compromised EDEK
    - Also requires EZK and ciphertext
    - Only leaks a single file
- Compromised ciphertext
    - Also requires corresponding EZ key and EDEK
    - Only leaks a single file
- Attack vectors
    - Physical access: ciphertext, EDEK, etc.
    - Network sniffing: ciphertext, EDEK, etc.
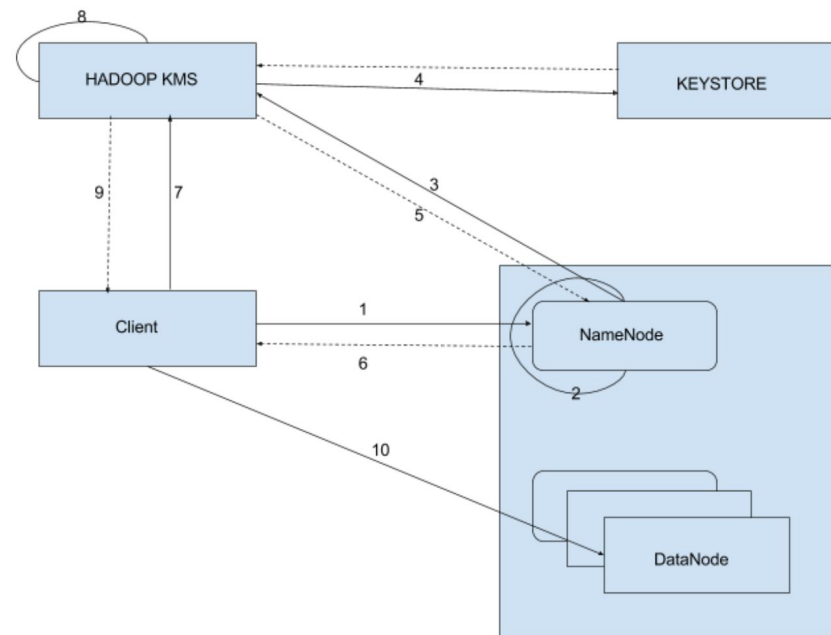    - Rogue admin: only one domain.

# HAWQ-TDE integration - Introduction

- https://issues.apache.org/jira/browse/HAWQ-1193
- HAWQ data is stored in HDFS
  - `hdfs://namenode:9000/hawq_default`
  - The whole directory(filespace) as an Encryption Zone
- libhdfs3
  - **A C/C++ native HDFS client**, written by Pivotal
  - HAWQ visit HDFS by it, need to upgrade to support TDE
  - https://github.com/apache/incubator-hawq/tree/master/depends/libhdfs3
- Other HAWQ components don't need to be modified
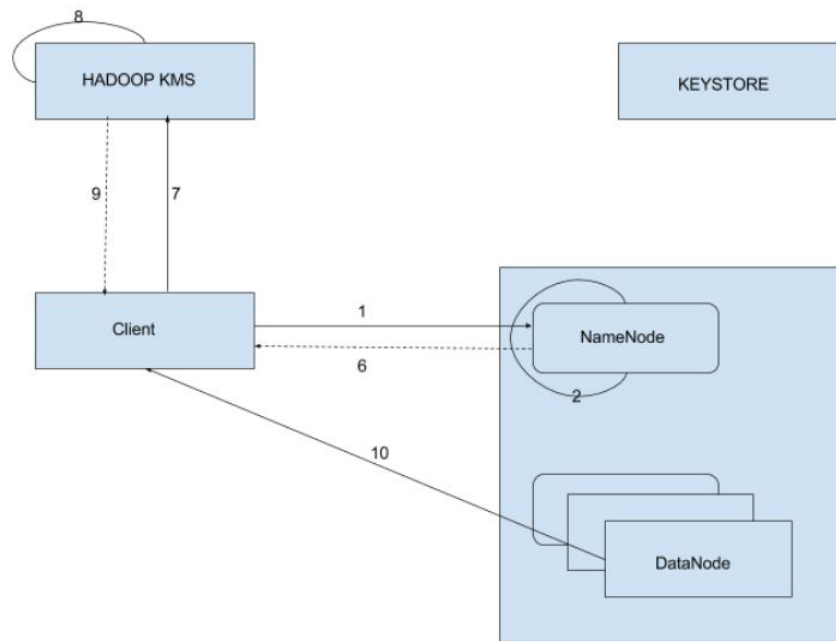  - API transparency

# HAWQ-TDE integration - Write process

1. The client issues a command to store a new file.
2. The NameNode checks if the file is under the Encryption Zone.
3. The NameNode passes the EZK name to KMS and requests KMS to create a new DEK.
4. The KMS retrieves the EZK from the key store and encrypts the DEK using EZK to generate the EDEK.
5. The KMS provides the EDEK to the NameNode and NameNode persists the EDEK as an extended attribute for the file metadata.
6. The NameNode provides the EDEK to the HDFS client.
7. The HDFS client sends the EDEK to the KMS, requesting the DEK.
8. The KMS checks if the user running the HDFS client has access to the EZK.
9. The KMS decrypts the EDEK using the EZK and provides the DEK to the HDFS client.
10. The HDFS client encrypts data using the DEK and writes the encrypted data blocks to HDFS.

# HAWQ-TDE integration - Read process

1. The client issues a command to read a new file.
2. The NameNode checks if the file is under the Encryption Zone.
3. No needed
4. No needed
5. No needed
6. The NameNode provides the EDEK to the HDFS client.
7. The HDFS client sends the EDEK to the KMS, requesting the DEK.
8. The KMS checks if the user running the HDFS client has access to the EZK.
9. The KMS decrypts the EDEK using the EZK and provides the DEK to the HDFS client.
10. The HDFS client encrypts data using the DEK and writes the encrypted data blocks to HDFS.

# HAWQ-TDE integration - Write/Read summary

- Writing a file
  - Client asks NameNode to create a file in an Encryption Zone
  - NameNode gets a new EDEK from KMS
  - NameNode stores EDEK in file's xattr and returns EDEK to client
  - Client asks KMS to decrypt EDEK
  - KMS decrypts EDEK and returns DEK to client
  - Client uses DEK to encrypt file data
- Reading a file
  - Client asks NameNode for file's EDEK
  - Client asks KMS to decrypt EDEK
  - KMS decrypts EDEK and returns DEK to client
  - Client uses DEK to decrypt data

# HAWQ-TDE integration - Conclusion

- Hadoop TDE bring end-to-end, client-side encryption to HAWQ
  - Data protected at-rest and in-transit
- The whole HAWQ data directory as an Encryption Zone
  - It cannot transfer between Encryption Zone and Normal Zone
- Transparent
  - API transparent for existing applications
  - Performance transparent
    - Less 10% performance impact
- Separate administrative domains for key management and HDFS

# Demo

```
// as the key administrator user

$ hadoop key create mykey

// as the HDFS administrator user

$ hadoop fs —mkdir /users/hongxu/my-ez

$ hadoop fs —chown hongxu /users/hongxu/my-ez

$ hadoop fs —chmod 700 /users/hongxu/my-ez

$ hadoop crypto —createZone —keyName mykey —path /users/hongxu/my-ez
```

# Demo (cont.)

```
// as the encryption user

$ echo "hello world" > hello-world.txt

$ hadoop fs –put hello-world.txt /users/hongxu/my-ez/

$ hadoop fs –cat /users/hongxu/my-ez/hello-world.txt

    hello world

// as the HDFS admin

$ hadoop fs –cat /.reserved/raw/users/hongxu/my-ez/hello-world.txt

    ��嵤]□N^���+���
```

# Future & Reference

- TDE with HDP 2.6 support will be released in HAWQ 2.3.0.0
- Future
  - KMS security: Kerberos, HTTPS
  - Finer granularity (tablespace, table)
- Reference
  - http://hadoop.apache.org/docs/current/hadoop-project-dist/hadoop-hdfs/TransparentEncryption.html
  - https://www.slideshare.net/Hadoop_Summit/transparent-encryption-in-hdfs
  - https://issues.apache.org/jira/browse/HDFS-6134
  - https://issues.apache.org/jira/browse/HAWQ-1193

# Thanks

# Q & A

Pivotal is hiring: GPDB, HAWQ, **pivotalrnd_china_jobs@pivotal.io**

Welcome to join!

Pivotal