

# CUBIKA BIG INSIGHTS

## DATA SHEET

### HIGHLIGHTS

- **Uncover hidden insights enhanced cognitive search, infused with Natural Language Processing in Thai**
- **Holistic approach to Big Data and Analytics, swift and scalable to meet every business's demands.**
- **Actionable insights in real-time manner**
- **A breakthrough Data Refiner tool, powered by Thai NLP.**
- **Integration of intelligent tools, services and data platform - full suite of tools, connectors and enablers to manage big data and data lake. Ready to scale, cost efficiently.**
- **Built-in PDPA-complied data security tool**
- **Future-proof metadata template for Thai Government metadata.**



Every business is a data business. As the range of digitized devices and services expand, data becomes a gold mine for business. However, by just collecting and storing data is not sufficient. Organizations must leverage power of analytics to tap, uncover and navigate through tsunami of data and transform them into insights that drive business growth with the speed that can stay on top of what customers need and stay ahead in the market.

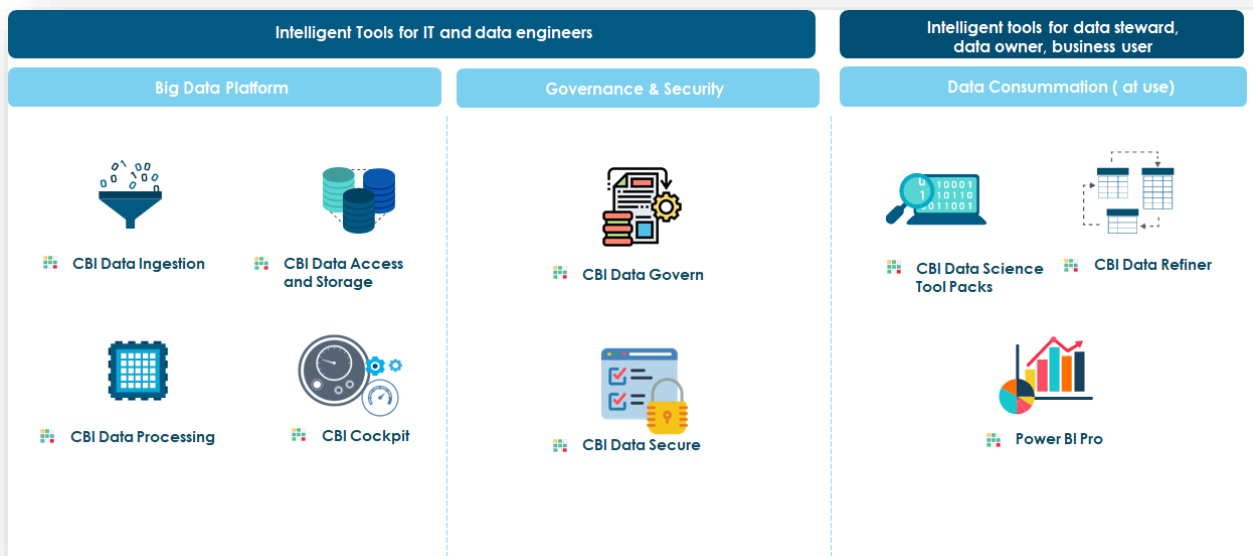
With big data and analytics technologies evolving at the breakneck speed, organizations are facing with the challenge of keeping up with the momentum. The main questions which organizations need to address are which are the right tools and solutions in helping them to become a data-driven organization, swiftly, flexibly, and cost-effectively.

CUBIKA Big Insights, part of CUBIKA AI - Digital Dialogue's suite of intelligent product, helps organizations in democratizing data and analytics, enabling the organizations to harness meaningful and valuable business insights which transforms into opportunities for business growth, powered by Digital Dialogue's proprietary Thai NLP engine and framework.

## KEY FEATURES

CUBIKA Big Insights products applies machine learning, analytics, and Digital Dialogue's Thai Natural language processing in automating tasks, categorizing, and standardizing not only English but Thai data across enterprise's big data environment. CUBIKA Big Insights help everyone from business users, data engineers, analysts to IT in achieving their task, gaining understanding of data, and turning them into actionable insights.

CUBIKA Big Insights enables organizations to capture, curate and consume data with the speed of business, embedded with 5Vs of Big Data (Volume, Variety, Veracity, Velocity and Value) as a cornerstone in order to transform data into meaningful insights for business. CUBIKA Big Insights handles mountains of data (Volume) in various shapes and forms; unstructured, structured and semi-structured data from various sources through API integration (Variety and Velocity), process them through data management, including profiling and deduplicating, leveraging open-source technologies in Hadoop ecosystem, ensuring data's integrity, lineage and accuracy with no anomaly (Veracity). Data is now ready to be analyzed by intelligent analytics model, powered by machine learning and Digital Dialogue's owns NLP engine (TH/ENG) which serve as a brain to understand the context of data and help with data de-duplication efforts to make sure the data is clean even before entering the storage (in-transit context analysis). Users can harness and unlock powerful, actionable insights in data usage stage via industrialized tool such as Power BI (Value) or Tableau or KNIME. With our principle in democratizing data, CUBIKA Big Insights empowers users in accessing data with simple and user-friendly tools and user interface in real-time and on-demand.



**Figure 1: CUBIKA Big Insights Intelligent Data Platform Suite**

## KEY BENEFITS

- **Holistic Approach to Big Data and Analytics**  
CUBIKA Big Insight's visionary solution encompasses data journey from data discovery to data usage by integrating leading open-source technologies and accelerators in big data and analytics. We weave them into a fabric of big data and analytics solutions with the capabilities in natural language understanding and processing, including Thai language.
- **High Flexibility with No Vendor Lock-in**  
Leveraging open-source technologies provides the benefit of shifting technology vendors based on organization's business focus. CUBIKA Big Insights tools and platform are developed with that elasticity in mind, providing organizations a much-needed flexibility in the solution.

- **Thai Government Metadata Compliance**

CUBIKA Big Insights supports off-the-shelf Government Metadata Management, suitable for Thai government Data Framework.

- **Future-Ready Tools to Help Enterprise in Data Management, Ready to Scale, Cost Efficiently**

CUBIKA Big Insights future-ready tools help enterprises in capture and curate raw data, ensuring that they have the high quality and refined data ready for analytics.

The tools in CUBIKA Big Insights (CBI) family are compatible with any data platform and scalable to match any business's needs in data management.

- CBI Data Access and Storage
- CBI Cockpit
- CBI Data Secure
- CBI Data Govern
- CBI Data Ingestion
- CBI Data Processing
- CBI Data Refiner
- CBI Zeppelin
- CBI Phonetics

- **Actionable Insights in Real-time Manner**

Sharpening your competitive edge and staying ahead in the market with real-time insights, enabling business in making insightful real-time decisions with CUBIKA Big Insight's NLP enhanced cognitive search capability where data can be searched and analyzed quicker than ever.

- **Unleash Full Potential of Big Data and Analytics...at a Fraction of Cost**

Becoming a data-driven company, maximizing business growth through insights while letting us take the wheel in managing the system. Digital Dialogue embeds an end-to-end approach from planning, implementation to support in CUBIKA Big Insights, empowering organizations to focus their efforts and resources in leveraging insights to accelerate their business with flexible pricing model and minimal total cost of ownership.

- **Business at the speed of intelligence insights**

Enhancing your business success in Thai and Oversea market with AI-infused analytics machine that understands and fulfills customers and users' intention in Thai language. Integrating NLP capability enables organizations in uncovering more unstructured data in Thai language which in turn creates more valuable insights especially for Thai competitive landscape. Ultimately providing organizations a game-changing boost capture demands, increase customer satisfaction and improve operational efficiencies. CUBIKA Big Insights is powered by an intelligent Thai Natural Language Processing engine, a proprietary of Digital Dialogue.

## TECHNICAL SPECIFICATION

CUBIKA Big Insights augments open-source Hadoop with enterprise-class functionality and integration essential to meet key business requirements. We leverage Apache Hadoop with the main components such as Hadoop Common, HDFS, YARN and MapReduce, which are supporting modules for Apache Hadoop. Powered by the latest and the most stable version of Apache Hadoop, CUBIKA Big Insights incorporates a number of significant enhancements, making it cost-effective such as less disk space required for faster processing<sup>1</sup>, support multiple standby Name Node as much as required. CUBIKA Big Insights also compatibles with all file and big data file systems, including Microsoft Azure Data Lake.

CUBIKA Big Insights (CBI) product portfolio, which is **an integration of intelligent data platform and solutions** that helps enterprise accelerates their shift to become a data-driven enterprise.

CUBIKA Big Insights supports ETL (Extract/Transform/Load), ELT (Extract/Load/Transform), Push Down process from various data source with Hadoop Distribution via Apache NiFi.

CUBIKA Big Insights comes with high availability, resiliency and fault tolerance with nodes and clusters redundancy to ensure the data is intact and always available in case of an unexpected failure<sup>2</sup>.

### Supported Platforms

CUBIKA Big Insights works with operating systems with 64-bit CPU as indicated below:

- CentOS
- Red Hat
- Ubuntu

### Incorporating CUBIKA Big Insights with Hadoop Ecosystem

CUBIKA Big Insights combines Apache Hadoop and its components throughout Hadoop Ecosystem to effectively harness data, facilitating full configuration and server configuration. Recovery is possible as the historical setting log within Hadoop Ecosystem is stored and recorded. Furthermore, no need to reinstall after restart Hadoop Ecosystem as HDFS, YARN, HBase, Kafka and ZooKeeper.

CUBIKA Big insights works with the following Hadoop Ecosystem components below:

<b>HDFS (Hadoop Distributed File System)</b>	A distributed file system for storing and retrieving structured, semi-structured and unstructured data, providing high throughput access data by providing the data access in parallel. CUBIKA Big Insights supports REST API for data accessing in HDFS.
<b>YARN (Yet Another Resource Negotiator)</b>	A job scheduling framework and resource management in the cluster.
<b>MapReduce</b>	A YARN-based system for parallel processing of large data set.

<sup>1</sup> Comparing to other Hadoop 1 and Hadoop 2 distributions, CUBIKA Big Insights can save up to 75% of disk space required on the data nodes to capture the same amount of data while providing the same level of data resiliency and protections.

<sup>2</sup> CUBIKA Big Insights provides system resiliency through multiple Name Nodes deployment which required more physical server or VMs.

<b>Hadoop Management Console</b>	Leveraging CUBIKA Cockpit as a tool for provisioning, managing and monitor Apache Hadoop cluster via easy-to-use web-based dashboard.
<b>HBase</b>	A non-relational database management system, supporting structured data in large table
<b>Hive</b>	A data warehouse software for writing, reading, managing large data sets in HDFS or HBase using SQL.
<b>Kafka</b>	A platform for building real-time data pipelines and streaming apps.
<b>Sqoop</b>	A tool for transferring data between Hadoop and structured relational database, enterprise data warehouses and NoSQL. Sqoop allow import data from various data sources ;Oracle, MYSQL, SQL Server, External File System
<b>ZooKeeper</b>	A centralized service providing distributed coordination, naming and configuration data, provide flexible and synchronization within distributed systems. The services in the cluster are replicated and stored on a set of servers, each of which maintains an in-memory database containing the entire data tree of state as well as a transaction log and snapshots stored frequently, to protect data in case of an unexpected failure).
<b>Spark</b>	A fast and unified analytics engine, capable of processing large sets of data including streaming, combination of batch processing, streaming processing with machine learning.
<b>Azkaban</b>	A batch workflow job scheduler/manager for Hadoop workloads.
<b>Flume</b>	Distributed, reliable, and available service for efficiently collecting, aggregating, and moving large amounts of unstructured from multiple Data source into HDFS or HBase for example. Flume has built-in compatibility with several sources such as Avro, Thrift, Exec (/bin/sh -c for example), Spooling directory, Taildir, JMS, Syslog and many more.
<b>Storm</b>	A lightning-fast and reliable processing engine for real-time/streaming data.

CUBIKA Big Insights utilizes Apache NiFi, a web-based platform software, backed by Web-UI in designing, controlling, feedback, and monitoring data efficiently. NiFi enables users to make a routing decision for data easily from numerous data pipelines and indicating data flow of each data collection to its destination, which can be big data platform or any further data integration process, including but not limited to Data Transformation and Connection.

<b>Data Transformation</b>	Aggregation	LookUp	Sorter
	Data Masking	Normalizer	Transaction control/Two phase commit
	Expression	Rank	Union
	Filter	Router	Unstructured data reader/writer for PDF, Excel, XML, etc.
	Merge/Join	Sequence generator	Strategy Updates
<b>Data Connection</b>	Azure Blob	HDFS (Hadoop Distributed File System)	SQOOP
	AWS S3	RDBMS	Flat File and Binary Flies
	Hadoop	Azure Datawarehouse	Complex File HDFS
	Impala	Spark SQL	SOAP REST
	Avro	CSV	Parquet
	HBase	Hive	ADLS (Azure Data Lake Storage)

CUBIKA Big Insights extracts sound, video, images, newsfeed, and other types of data from source systems via Apache NiFi using FlowFile. Data is processed through advanced analytics model and accessible in visual elements format by leveraging data visualization tools, where business can gain insights and make a data-driven decision from a massive amount of data.



. Ensuring that data can be searched and analyzed for insights in a real time manner, CUBIKA Big Insights uses CBI Phonetics and CBI Data Crawler in diving into your data lake with the enhanced cognitive search capability. It assists the analytics model in curating for the right data for consumption in a much simple and quicker manner

### Internet of Thing (IoT) Data Management

One of the contributions to the explosive growth of big data is IoT where data flows continuously from sensors across millions of devices, containing insights waiting to be uncovered. CUBIKA Big Insights manages and extracts Internet of Thing (IoT) data by leveraging Apache Kafka. The customization code can be added to support the additional transformation logic to make data ingestion process compatible and smooth via Apache NiFi. CUBIKA Big Insights also connect with RDBMS Data Lake (Hadoop) and unstructured data as PDF, Word, Excel, etc.

## CUBIKA BIG INSIGHTS DATA PLATFORM PRODUCTS PORTFOLIO

CUBIKA Big Insight's flexible and scalable data platform solution is ready to enhance your data journey and helps enterprise accelerates their shift to become a data-driven enterprise.

CUBIKA Big Insights for Data Ingestion	CUBIKA Big Insights for Data Governance	CUBIKA Big Insights for Data Security	CUBIKA Big Insights for Data Access and Storage	CUBIKA Big Insights for Data Science Tool Packs
<ul style="list-style-type: none"> <li>• CBI Data Ingestion</li> <li>• CBI Data Processing</li> </ul>	<ul style="list-style-type: none"> <li>• CBI Data Govern</li> </ul>	<ul style="list-style-type: none"> <li>• CBI Data Secure</li> </ul>	<ul style="list-style-type: none"> <li>• CBI Data Access and Storage</li> <li>• CBI Cockpit</li> </ul>	<ul style="list-style-type: none"> <li>• CBI Zeppelin</li> <li>• CBI Refiner</li> <li>• CBI Phonetics</li> </ul>

**Figure 2: CUBIKA Big Insights Suite**

## CUBIKA Big Insights for Data Ingestion

### CBI Data Ingestion

CBI Data Ingestion is an essential part of data on-boarding process. Ingesting data from various sources, including streaming data into big data system for further management and processing.

- Drag and Drop application to upload file via web-based GUI.
- Can create everything from simple to complex through a GUI such as create flows or change data structure by configuring processors.
- Compatible with ETL/ELT/Push Down with Hadoop Distribution. CBI Data ingestion enables ETL workflow management system via web-based GUI.
- All task in the activity log will keep in log system
- Compatible with HDFS (Hadoop Distributed File System) in the following:
  - Supporting multiple data processing pipeline
  - Supporting control over data flow direction for big data platform from design, control, feedback to monitor
- Compatible with connection below:

<b>HDFS (Hadoop Distributed File System)</b>	A distributed file system for storing and retrieving structured, semi-structured and unstructured data, providing high throughput access data by providing the data access in parallel. Edit and create data category management and show in Browser Directory (One of File Explorer)
<b>YARN (Yet Another Resource Negotiator)</b>	A job scheduling framework, resource, and queue management in the cluster.
<b>Hadoop Management Console</b>	Leveraging Apache Ambari as a tool for provisioning, managing and monitor Apache Hadoop cluster via easy-to-use web-based dashboard.

<b>HBase</b>	A non-relational database management system, supporting structured data in large table
<b>Hive</b>	A data warehouse software for writing, reading, managing large data sets in HDFS or HBase using SQL.
<b>Kafka</b>	A platform for building real-time data pipelines and streaming apps.
<b>Sqoop</b>	A tool for transferring data between Hadoop and structured relational database, enterprise data warehouses and NoSQL. Sqoop allow import data from various data sources; Oracle, MYSQL, SQL Server, External File System etc.
<b>ZooKeeper</b>	A centralized service providing distributed coordination, naming and configuration data, provide flexible and synchronization within distributed systems. The services in the cluster are replicated and stored on a set of servers, each of which maintains an in-memory database containing the entire data tree of state as well as a transaction log and snapshots stored frequently, to protect data in case of an unexpected failure).
<b>Spark</b>	A fast and unified analytics engine, capable of processing large sets of data including streaming, combination of batch processing, streaming processing with machine learning.
<b>Azkaban</b>	A batch workflow job scheduler/manager for Hadoop workloads.
<b>Storm</b>	A lightning-fast and reliable processing engine for real-time/streaming data.
<b>Flume</b>	Distributed, reliable, and available service for efficiently collecting, aggregating, and moving large amounts of unstructured from multiple Data source into HDFS or HBase for example. Flume has built-in compatibility with several sources such as Avro, Thrift, Exec (/bin/sh -c for example), Spooling directory, Taildir, JMS, Syslog and many more.
<b>MapReduce</b>	A YARN-based system for parallel processing of large data set.

- Supporting data import and processing from Apache Hadoop with at least ANSI-92
- Applicable with Thai language by importing/exporting data, data processing, data analyze.
- Supporting ODBC and JDBC connection
- Collecting joined data points within the single view of truth
- Compatible with any changes in data structure whether on the data structure itself or user requirement.
- Controlled schema version for the correct reporting

### CBI Data Processing

**CBI Data Processing** is an essential part of data on-boarding process. Ingesting data from various sources, including streaming data into big data system for further management and processing.

- CBI Data processing works with operating systems with 64-bit CPU as indicated below:
  - CentOS
  - Red Hat
  - Ubuntu
- Supporting real-time processing as follow:
  - Compatible with real-time transformation logic for both input and output data
  - Compatible with Real-time Pipeline and Real-time Streaming
  - Compatible with Distributed Messaging Queue
  - Distributed in cluster data management.
  - Resilient architecture with redundancy
  - High fault tolerance
  - High availability (HA)
  - Horizontal-node scalability
  - Support both Java and Scala in transformation logic
  - Supporting publish and subscribe model

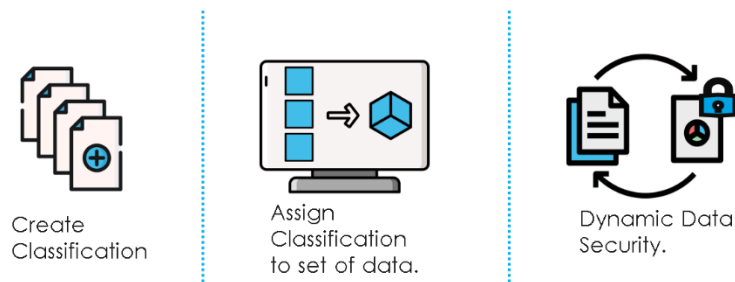
- Compatible with Java Development Kit (JDK)8 64 bit
- Low-code, GUI compatible for drag-and-drop workflow creation. Supporting remote execution and job schedule automation.
- Supporting direct read and write on data source from big data system such as Apache Hive and Impala
- Providing accurate performance record (read/write) such as working parameter and assigned KPIs.
- Supporting parallel executive & computing
- Supporting addition scripts
- Supporting additional libraries or plug-ins
- Supporting job queuing for higher efficiency

## CUBIKA Big Insights for Data Governance

### CBI Data Govern

*CBI Data Govern* a software which empowers enterprises to organize, find and understand data, helping enterprise to quickly discover and use appropriate data they can trust based on rules, policies, and responsibilities.

- **Embedded with Government Metadata Template model**, in compliance with Thai Government Data Catalog and metadata in Government Data Framework in accordance with the announcement from Committee of Digital Government Development on the subject of "Data Governance for Government" on 12<sup>th</sup> March 2020.
- **Smart business glossary and catalog** – easily assign label and description to your data
- **Data classification or tagging for dynamic data security rule setting** – CBI Data Govern equips you with the ability to assign multi label to data set and leverage it as a single point of directory in managing access and security.



- **Data Lineage traceability** - A system of truth with visualized data activities; join, create, and transform
- **Data classification for Dynamic data security rule setting** – CBI Data Govern equips you with the ability to assign multi label to data set and leverage it as a single point of directory in managing access and security.
- **Metadata template management** – a user can create a standardized metadata template for a particular business-unit with file log record along with the capability to track metadata template duplication.
- **Detecting sensitive data and data privacy** - CBI Data Govern can detect sensitive Personally Identifiable Information (PII) in automatic and Semi-Automatic way



# CUBIKA Big Insights for Data Security

## CBI Data Secure

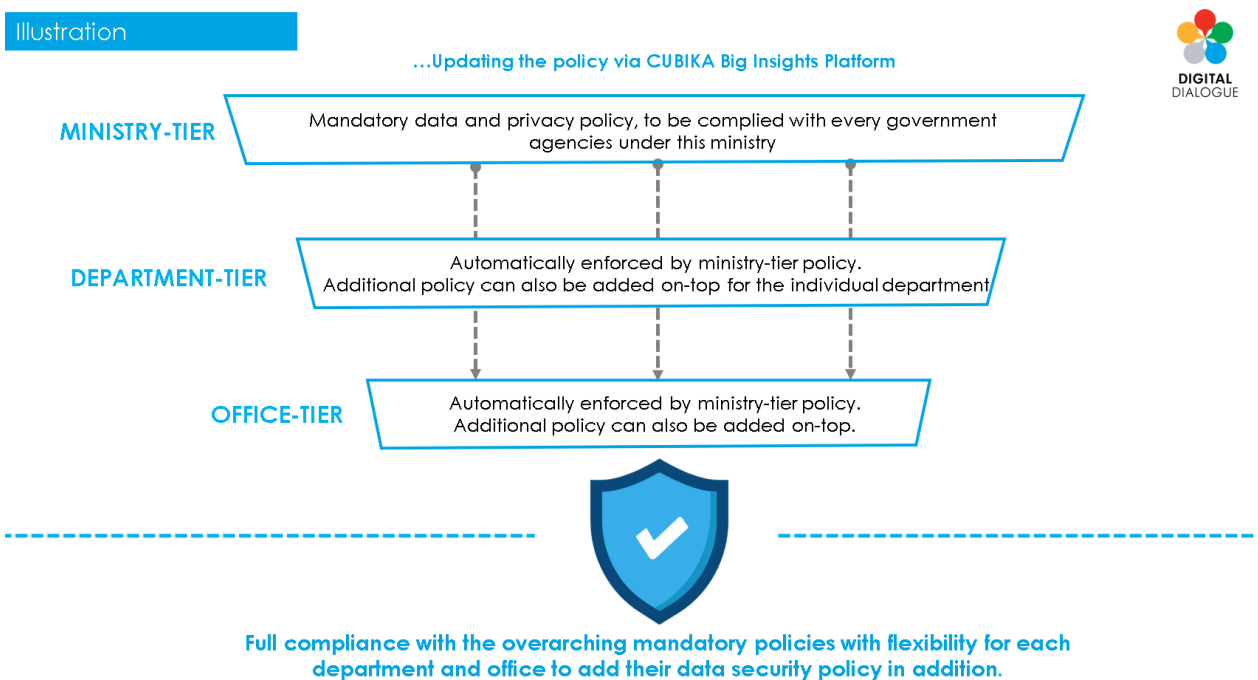
CBI Data Secure is a tool to enhance data security management and make it simple for enterprises to comply with data security, privacy policy, and data governance policy. Easily create access rights management and save data audit trails in one platform.

CBI Data Secure 's main feature is data pseudonymization through 'dynamic data masking' including hash, partial mask and redact, supporting multi-level trickling-down policy setups.

- Supporting user authentication via Native user authentication, LDAP user authentication, OpenID Connect, Kerberos for Single Sign-On (SSO) and SSL/TLS, including managing user authentication mechanisms as below:

<b>Privileges</b>	An authorization level for managing users and policies for preliminary authorizations at startup.
<b>Role Based</b>	An authorization level is based on the assigned role within an organization of a user.
<b>Permission</b>	An authorization level is based on the access policy for a user: 1. View – user with this access level can access and only view details and data within the system. 2. Modify – user with this access level and access and modify any settings and configurations within the system.

- Enabled data warehouse security.
- Providing capability for Role-based Access Control (RBAC) working in tandem with HDFS Access Controls (ACLs)
- Supporting dynamic trickling-down policy enforcement, PDPA & regulatory compliance with automatic policy enforcement from the top to bottom level where top-tier policies are mandatory enforced and cannot be revised by the lower-tier. The lower-tier policies can be added, if needed. The policy flow setup process can be fully configured.



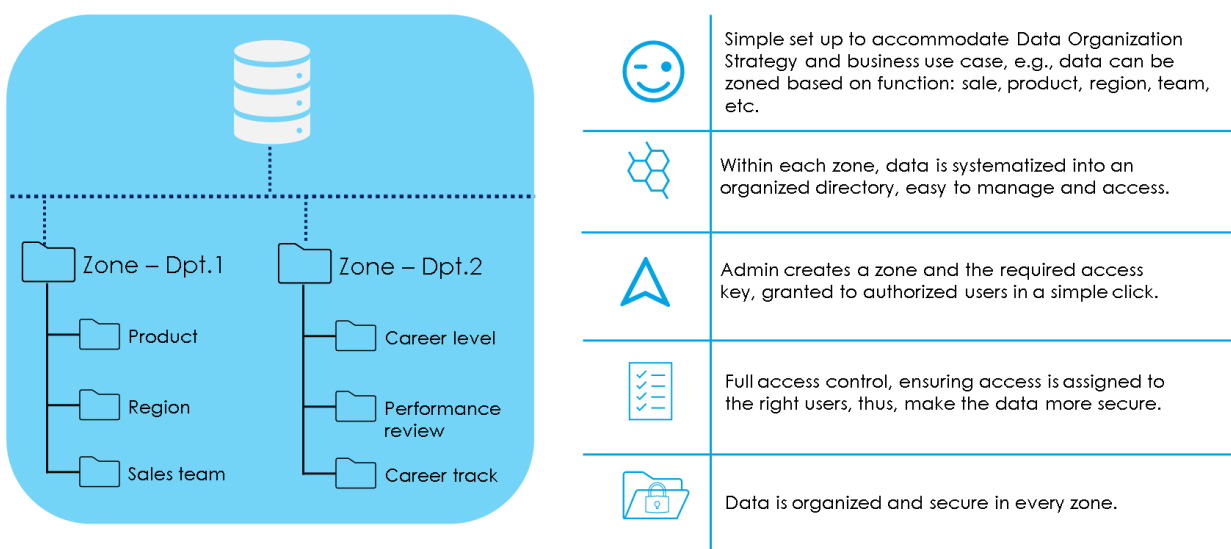
**Figure 3: Illustration of How Dynamic Trickling Down Policy works through CBI Data Secure**

- Dynamic data masking is available for data pseudonymization, enhancing data security.
  - Cryptographic Hash
  - Partial Masking
  - Data Redact
- Data access audit - every previous event and activity in Apache Hadoop can be audited (Back in Time) from securely stored activity logs which can be accessed and searched via GUI by administrators and auditors. The activity log will keep all task below:
  - Hadoop service access
  - User log-in
  - Policy admin access
  - Addition logs as below:

<b>HDFS Audit Logs</b>	2 different types of audit logs for HDFS 1) Audit logs of user activity 2) Audit logs of service activity (Hadoop Service)
<b>MapReduce Audit Logs</b>	2 different types of audit logs for MapReduce 1) Audit logs of user activity 2) Audit logs of service activity (Hadoop Service)
<b>YARN Audit Logs</b>	Audit log events for YARN are logged in daemon log files.
<b>Hive Audit Logs</b>	Using Hive Metastore for logging Hive audit events. Auditors can audit username and IP and Hive service log event by filtering logs with <i>HiveMetaStore, audit</i>
<b>HBase Audit Logs</b>	HBase audit log events and services containing information and event in Column Family, Column and Table.

- Secure data zone setup - Creating the business-focused and encrypted data zone based on data organization strategy. With CBI Data Secure, data zone can be set up easily and it supports ACL and RBAC Role and Authorization setup.

CBI Data Security helps your enterprise in architect and organize data zone with data organization strategy as a backbone, support heavy curation of data for various business use.



**Figure 4: Illustration of Secure Data Zone setup and its benefits**

- Supporting Data-in-Transit security with at least TLS 1.2 encryption

## CUBIKA Big Insights for Access and Storage

### CBI Cockpit

*CBI Cockpit* is a highly intuitive management console for monitoring, provisioning, and managing big data ecosystem. Moreover, the service can help administrators for improve work efficiency.

The dashboard is powered by web-UI and REST API, enable administrators to monitor health, status or key performances of each, Hadoop cluster, get notification of certain events or system alerting with one centralized dashboard for Hadoop Ecosystem - Hadoop Distributed File System (HDFS), Zookeeper, Yet Another Resource Negotiator (Yarn), Spark, Hive, HBase, Hadoop MapReduce.

CBI Cockpit, powered by CBI Cockpit Alert Framework, can send dispatch notifications to alert user of status changes through SNMP or provided email and SMTP. Rest API and Python API are supported to provide historical information of clusters, nodes, applications.

Users can select which module they wish to configure; YARN+MapReduce, Tez, Hive, HBase, Sqoop, ZooKeeper, Kafka, Spark.

CBI Cockpit empowers enterprises in

1. Monitoring resource overview; node, resource and components in Hadoop ecosystem within one dashboard, enabling system administrators to get a quick view on Resource Utilization status right away.
2. Simplifying big data ecosystem configuration - providing system administrator easy, efficient, effortless tool for managing Hadoop ecosystem.
3. Providing capability for **Role-based Access Control (RBAC)** configuration working in tandem with HDFS Access Controls (ACLs) for Apache Hadoop. Also, CBI Cockpit includes tools for installation, migration, uninstallation, and role re-delegation on Hadoop Ecosystem operating in each server independently via web-based GUI without the need to uninstall at the operating system level.
4. Taking proactive measures – Live Notification is available for system errors and warning on each service.
5. Smart Self-Healing - Automatic detection of Data node problems and perform auto-restart remedy process which can also be customizable.

### CBI Data Access and Storage

*CBI Data Access and Storage* enables enterprises to manage and store large data sets in big data in HDFS (Hadoop Distributed File System) and query data in similar format to RDBMS or table format.

All data types can be managed by *CBI Data Access and Storage*.

- *Structured data* - Relational Database Management System such as Oracle, MySQL, MSSQL Server, PostgreSQL and many more.
- *Semi-structured data* - For example, Log file, Text file, csv, xlsx, xls, TXT, JSON and XML)
- *Unstructured data* - Various file types and streaming data, MongoDB, Social Media, IoT)

CBI Data Access and Storage also supports data Access and Storage in HDFS via standard ANSI SQL and Atomic, consistent, isolated, and durable (ACID) transactions.

By leveraging CUBIKA Data Access and Storage, enterprises can tap into full potential of HDFS powerful scalability, supporting thousands of nodes in a single cluster. With the sufficient hardware, it can efficiently and rapidly scale up to over 100 petabytes of raw storage capacity in one cluster.

## CUBIKA Big Insights for Data Science Tool Packs

Powered by natural language processing, data cleansing, preparing and transformation is even easier. Interactively explore, resolve quality issues, classify, standardize, summarize, and join data via web GUI. The outcome is the one true set of data with high quality.

### CBI Zeppelin

A "lab" for data science professionals to test their hypothesis and discover new data analytics model.

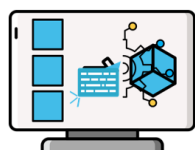
- Web-based GUI Zeppelin notebook for Python code interactive testing
- Machine learning engine for Classification, Regression, Clustering, Dimension Reduction, Model Selection, Preprocessing through R, Python, and Spark.

### CBI Data Refiner

A breakthrough and interactive approach to data wrangling, from exploration, refinery to preparation powered by Thai NLP.



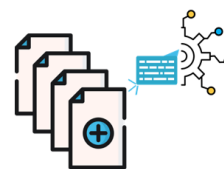
Clean Data



Transform Data

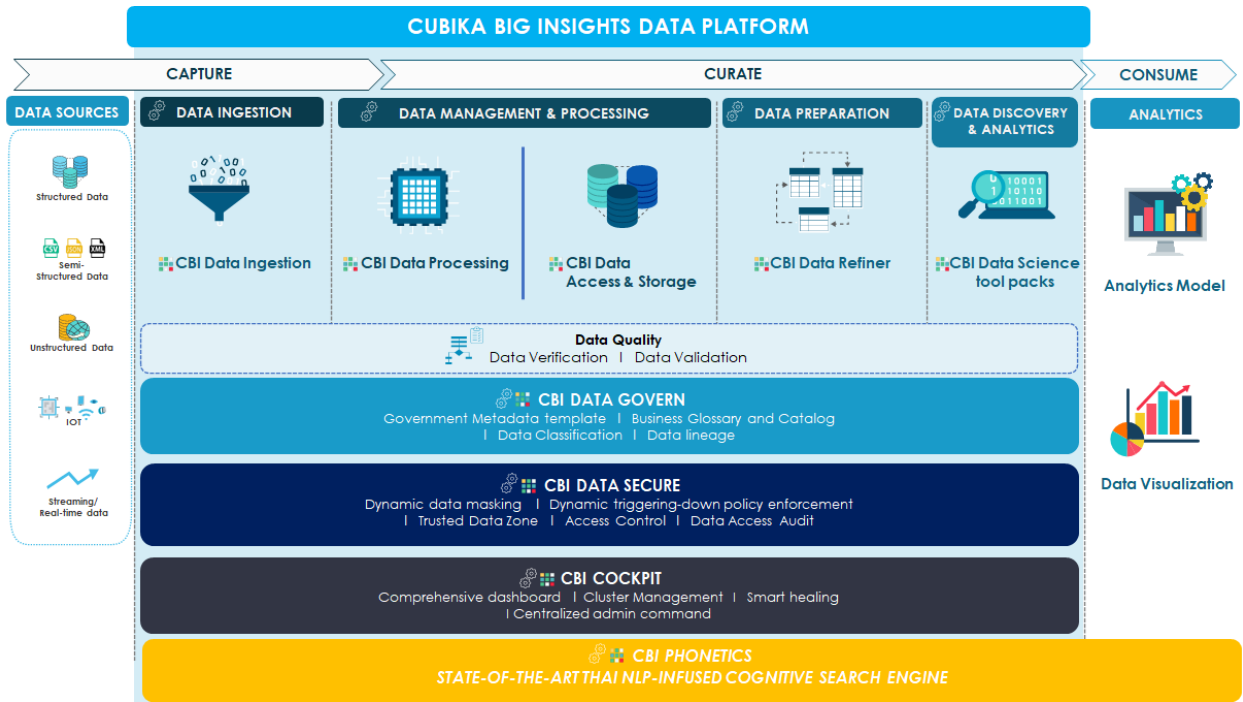


Convert Data Type



Detect redundant data

- Web-based GUI, providing user with Excel-like data view in column
- Powered by Thai NLP (Natural Language Processing), CBI Data Refiner can automatically detect data duplication and redundancy based on Thai phonetics. For example, the detection in the likelihood of ไกรสอน ไกรสร and ไกรศร as the same entity with different spelling.
- Compatible with column merging and splitting.
- Data error detection system during data cleansing process
- Supporting data transformation as below:
  - Text (string) to integer or number
  - Integer or number to text (string)
  - Text (string) to date format as below:
    - dd/mm/year
    - year/mm/day
    - dd-mm-year
    - year-mm-dd
    - year/mm/dd
    - year/dd/mm
    - year/dd/mm
    - year-dd-mm
    - Supporting Thai format date such 31 ตุลาคม 2563 and 31 ต.ค. 2563 as the same entity.
- Supporting changing date to text(string)
- Supporting automatic letter case format, for example, uppercase, lowercase, capitalized word.
- Address format detection and formatting. Embedded with machine learning, the system can automatically detect if the zip code doesn't match the associated district.
- Activity logs are recorded, and redo process is possible.
- Exportable activity logs in a template format for data cleansing and maintenance
- Data cleansing can be completed via
  - Parsing
  - Data correction
  - Data standardization
  - Data deduplication



**Figure 5: CUBIKA Big Insights Data Management Platform Suite and Data journey from managing data to revealing insights**



## WHY DIGITAL DIALOGUE?

Digital Dialogue is one of the pioneers in Thai Natural Language Processing (NLP) engine and framework with unparalleled processing speed and capability for Thai language. Faster than market by 30 microseconds in processing Thai complex-intent query (Based on 1,000 complex-intent queries).

We have extensive experiences and long-standing credentials in helping many prestigious organizations across industries in Thailand harnessing meaningful insights through CUBIKA Big Insights, resulting in timesaving, tangible increase in profit and revenue.

Digital Dialogue is committed in partnering with enterprises to accelerate their business growth through meaningful insights, create long-lasting value. We extend technology and business capabilities through a powerful alliance ecosystem across digital and industry landscape with the objective of helping our clients in becoming a data-driven organization where data is democratized. Everyone has access to the right and meaningful data at the right time, making an insightful decision for business.

Innovation is the key part of our core value. We are proud to have been recognized by the industry for our innovative solutions which make a positive impact with clients and marketplace. Our accolades include "Solution Innovator for AI and Data of the Year", "ISV Partner of the Year" and "Partner of the Year" in several years by Microsoft, "Top Ten Government Tech Solutions Provider in APAC" by CIO Outlook Magazine.

CUBIKA Big Insights can be integrated with other CUBIKA intelligent products; Conver, Voice Analytics, and Engage, enable organizations to unlock full capabilities of AI and machine learning and re-imagine the way business can thrive.



### DIGITAL DIALOGUE

© Copyright 2020 Digital Dialogue. This document is provided for information purposes only and the contents hereof are subject to change without notice. This document is not warranted to be error-free, nor subject to any other warranties or conditions, whether expressed orally or implied in law, including implied warranties and conditions of merchantability or fitness for a particular purpose. Digital Dialogue shall not be liable for technical or editorial errors or omissions contained herein. CUBIKA is a registered trademark of Digital Dialogue in Thailand. All third-party marks are property of their respective owners.

**Contact our specialist to get started  
with CUBIKA Big Insights and harness  
valuable business insights.**



**[contact@ddlghq.com](mailto:contact@ddlghq.com)**



**02-088-0795**