

第一章

大数据技术全栈

大数据技术全栈

大数据技术栈



```
graph LR; A[大数据技术栈] --- B[基础能力]; A --- C[数据采集]; A --- D[数据存储]; A --- E[数据查询]; A --- F[数据计算]; A --- G[其它]; B --- B1[java, python, scala]; B --- B2[linux]; C --- C1[flume, kafka, logstash, filebeat...]; D --- D1[hdfs, hbase, redis...]; E --- E1[hive, spark sql, presto, kylin, impala, durid, clickhouse, greeplum...]; F --- F1[storm, spark stream, flink...]; G --- G1[分布式协调器 --- zookeeper]; G --- G2[资源管理器 --- yarn, mesos]; G --- G3[调度管理器 --- oozie, azkaban, airflow, dalphine scheduler];
```

基础能力

java, python, scala

linux

数据采集

flume、kafka、logstash、filebeat...

数据存储

hdfs、hbase、redis...

数据查询

hive、spark sql、presto、kylin、impala、durid、clickhouse、greeplum...

数据计算

storm、spark stream、flink...

其它

分布式协调器 — zookeeper

资源管理器 — yarn、mesos

调度管理器 — oozie、azkaban、airflow、dalphine scheduler

大数据技术知识体系

基
础
技
术

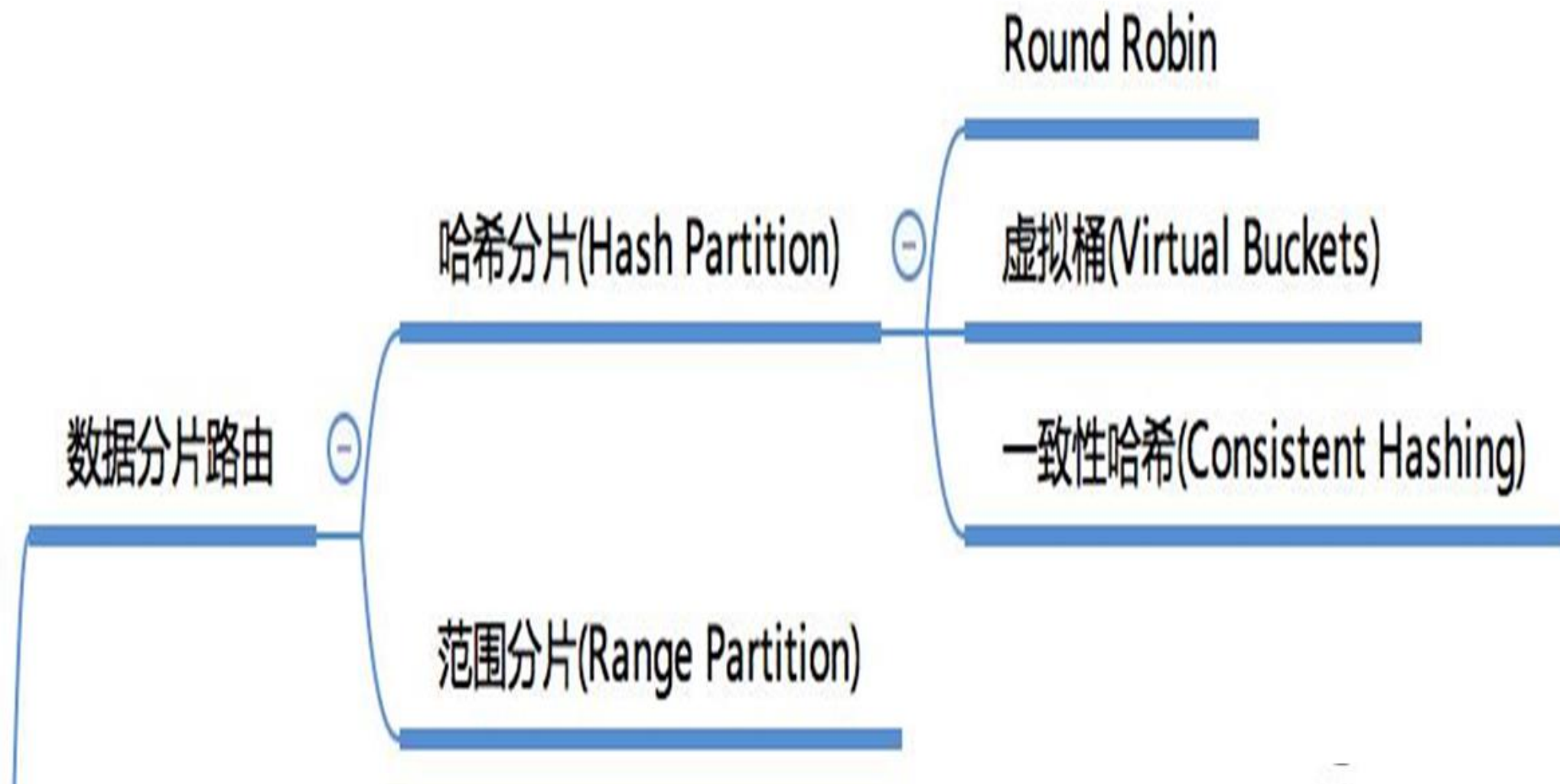
数据应用

数据组织集成

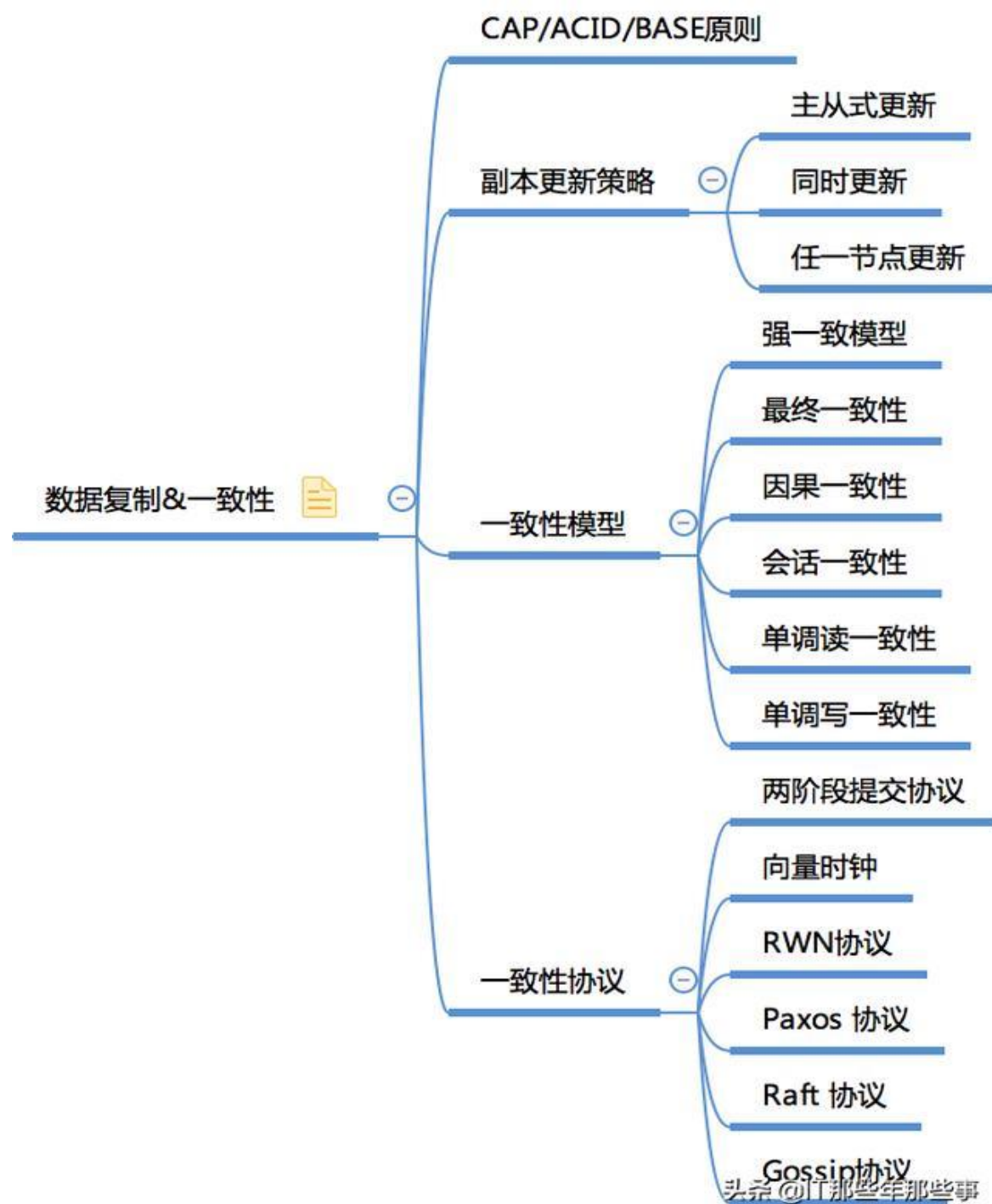
数据传输

数据采集

数
据
治
理



大数据基础技术



大数据基础技术

大数据常用算法与数据结构

SkipList

Bitmap

LSM树

Snappy与LZSS算法

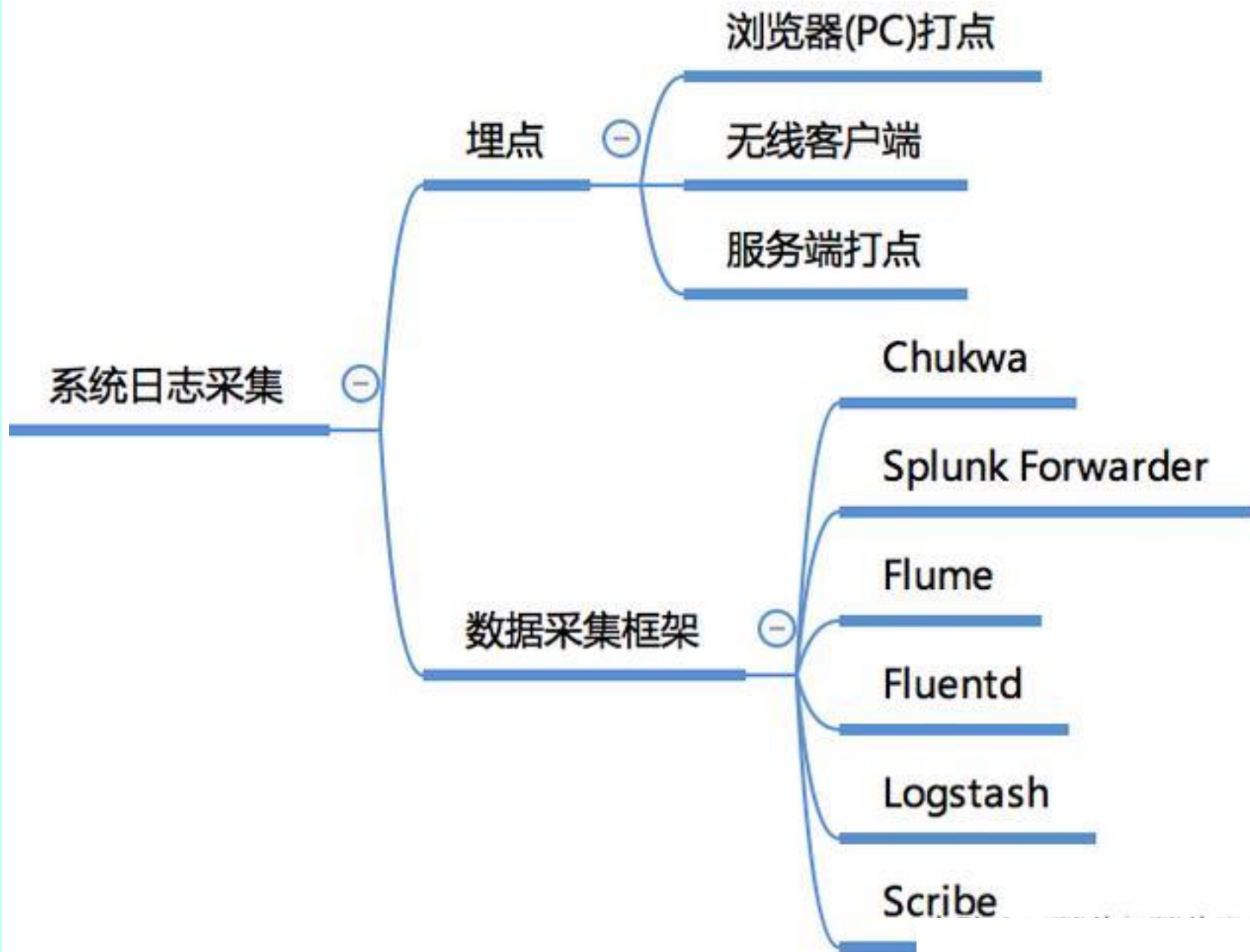
Cuckoo哈希

Mekle哈希树

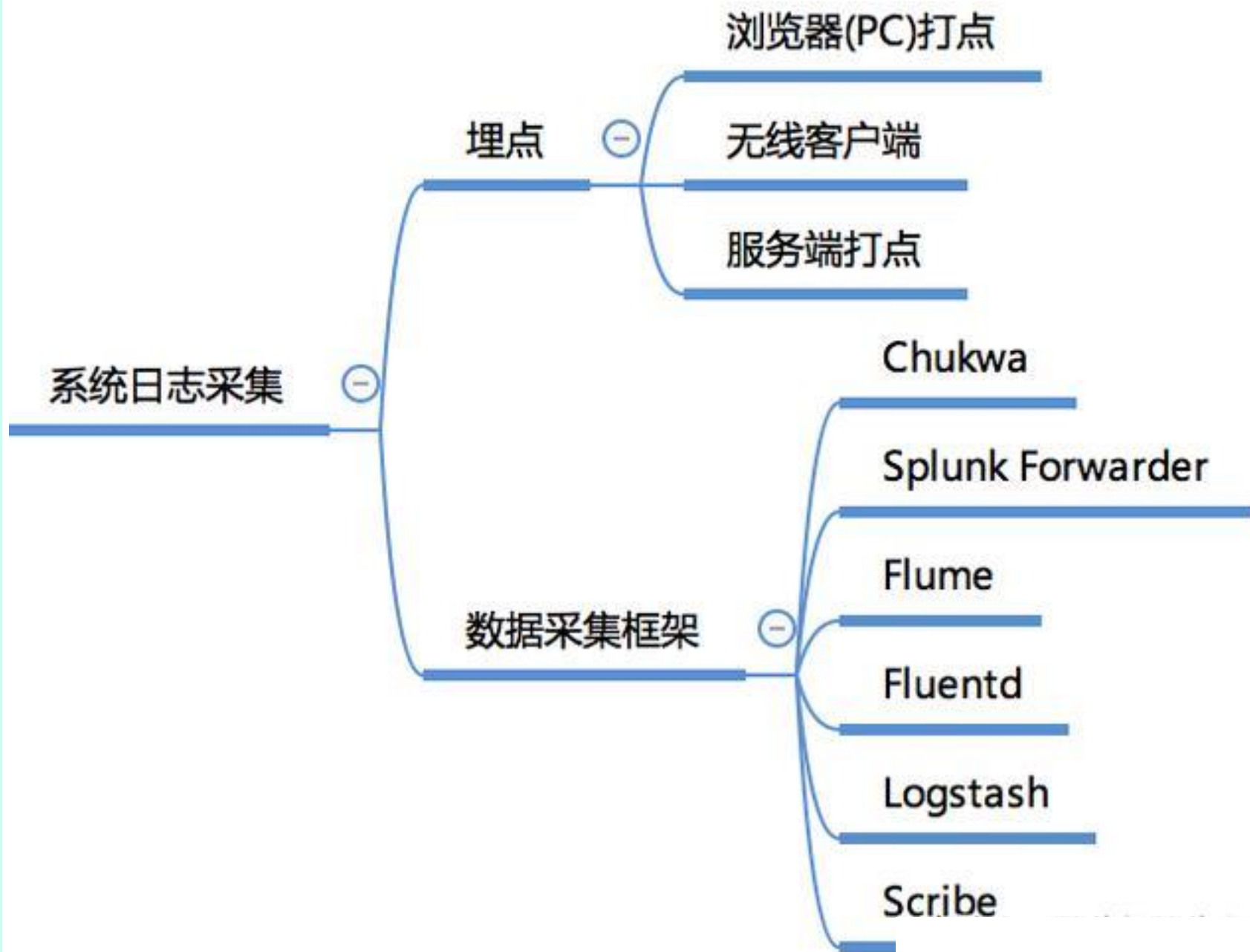
布隆过滤器(Bloom Filter)

Trie树

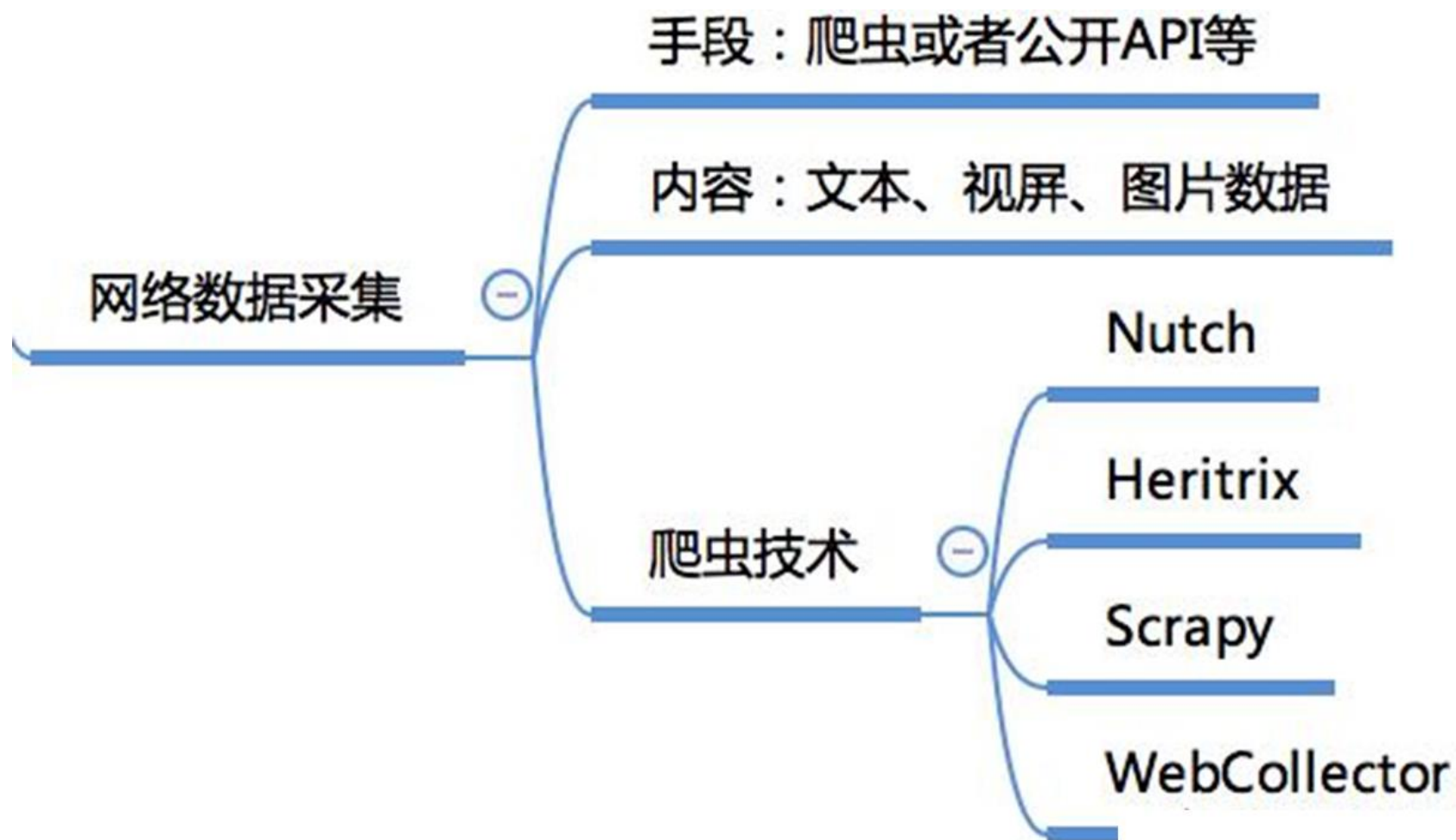
数 据 采 集



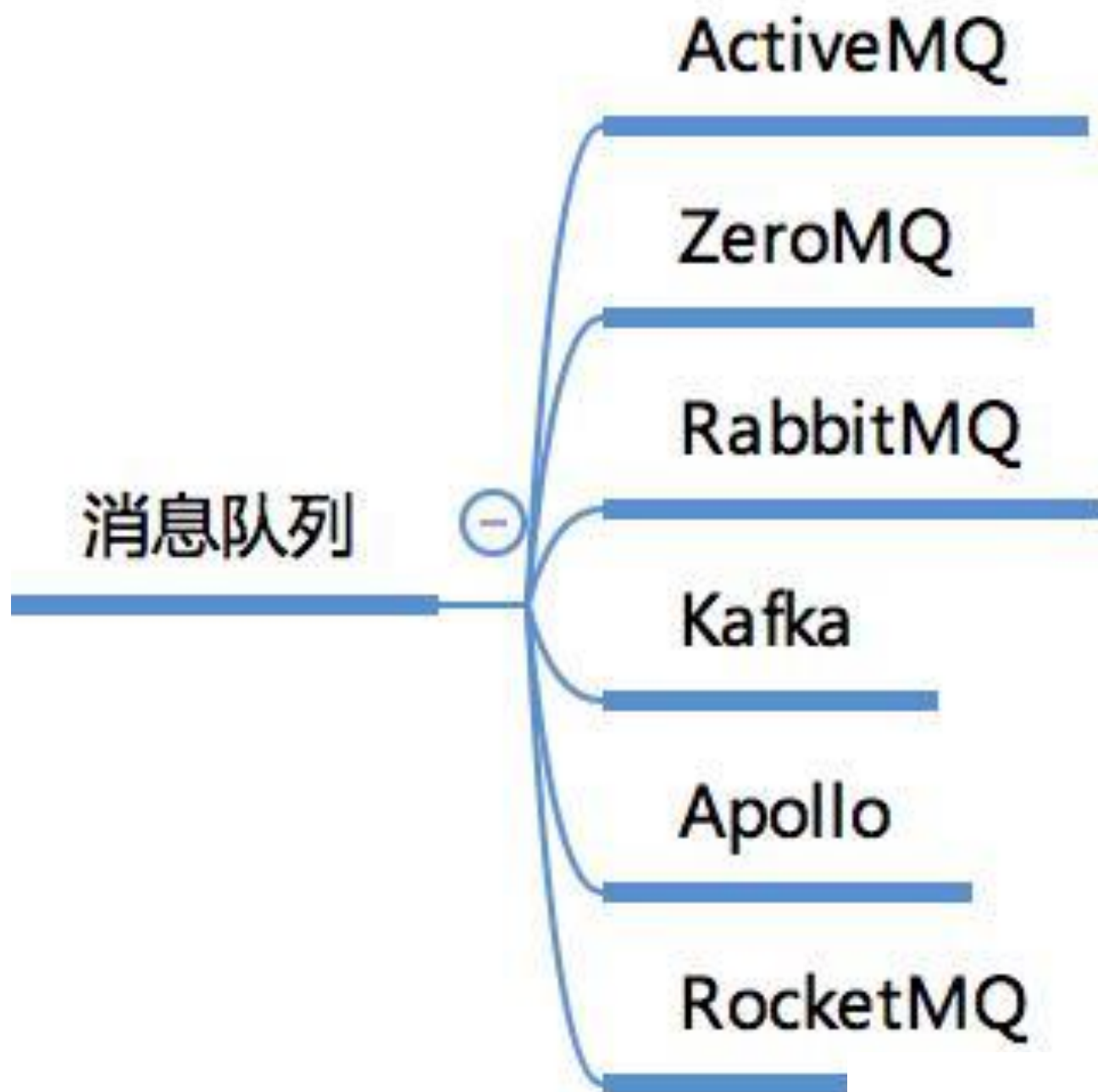
数 据 采 集



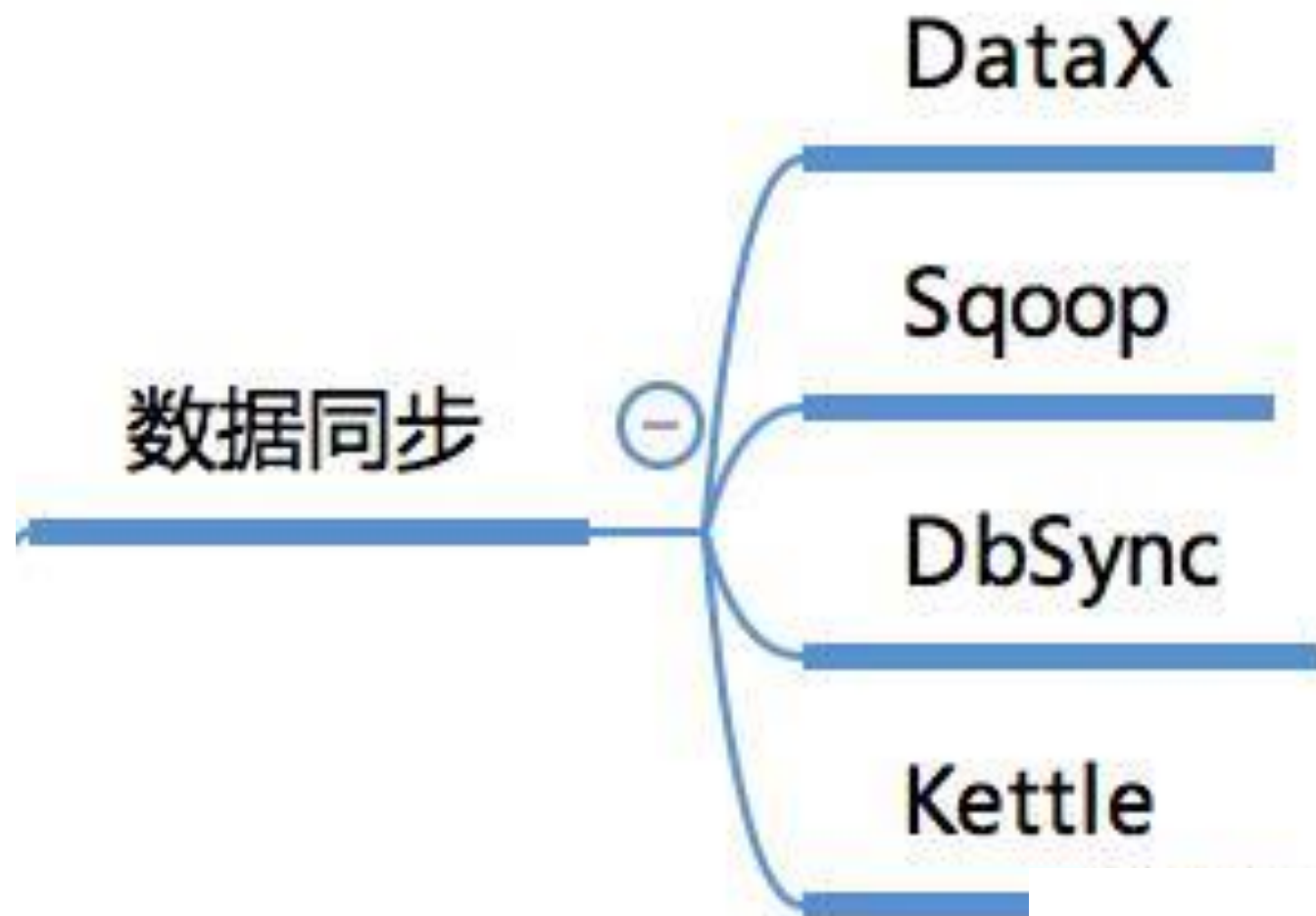
数 据 采 集



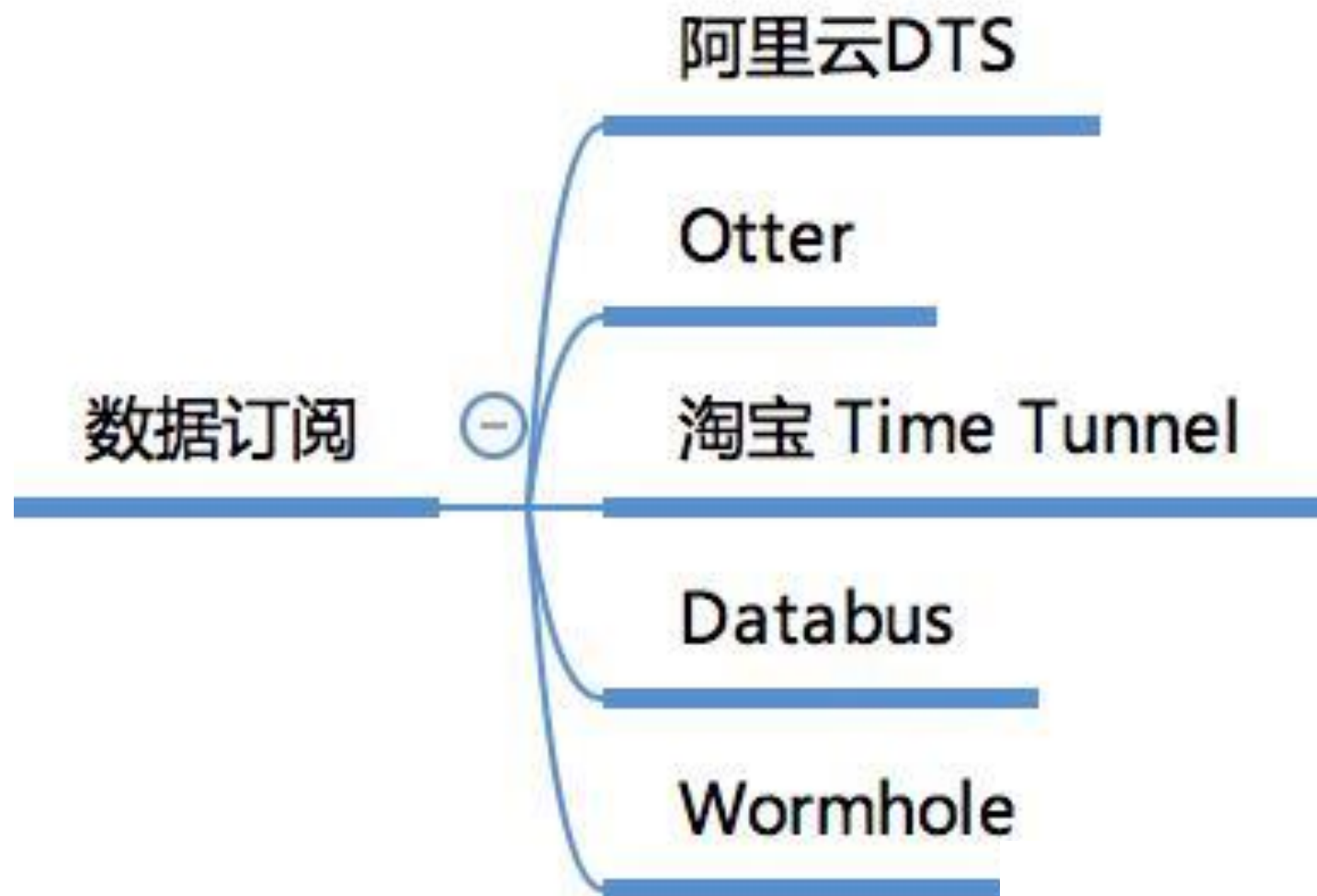
数
据
传
输



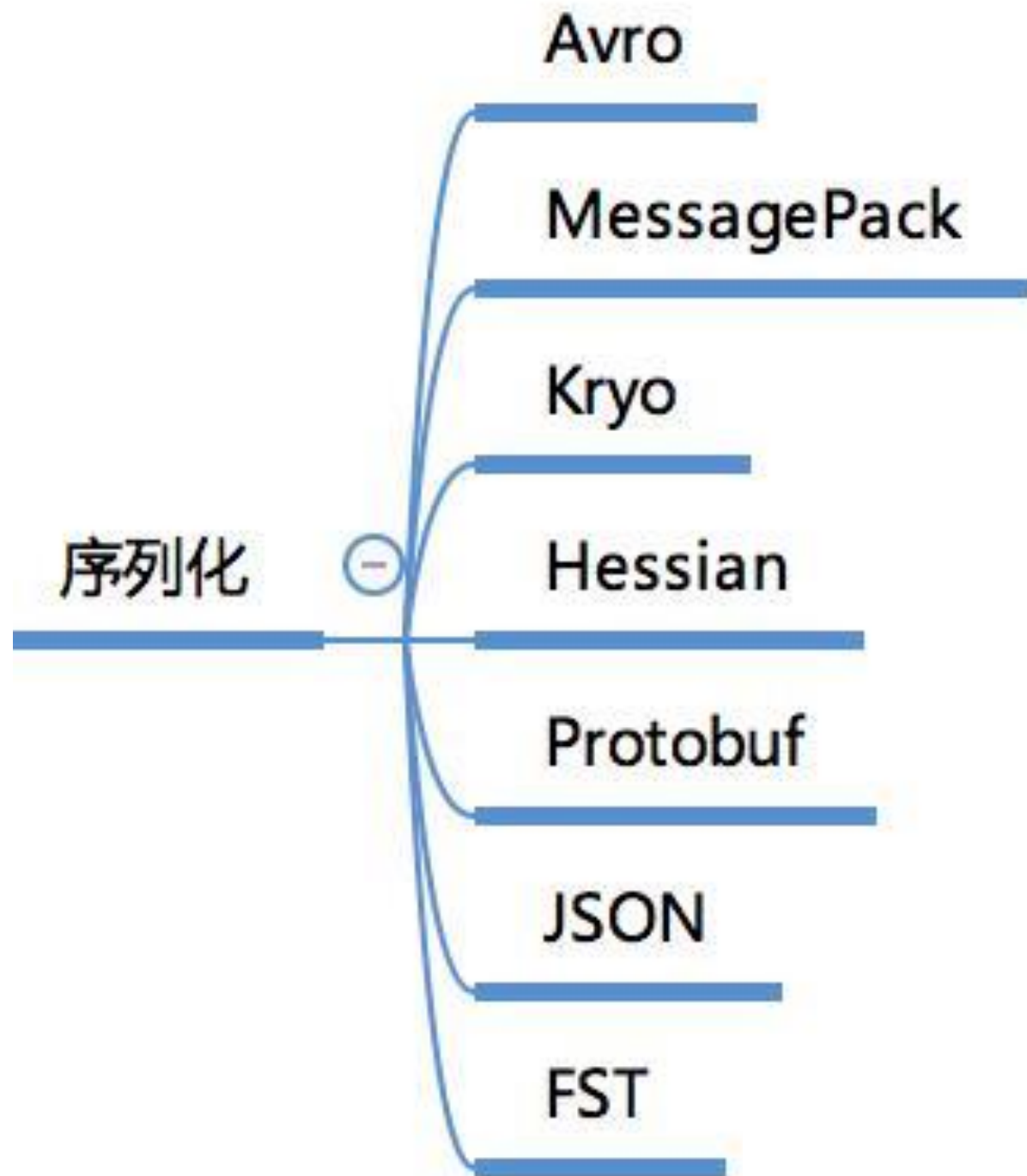
数
据
传
输



数
据
传
输



数
据
传
输

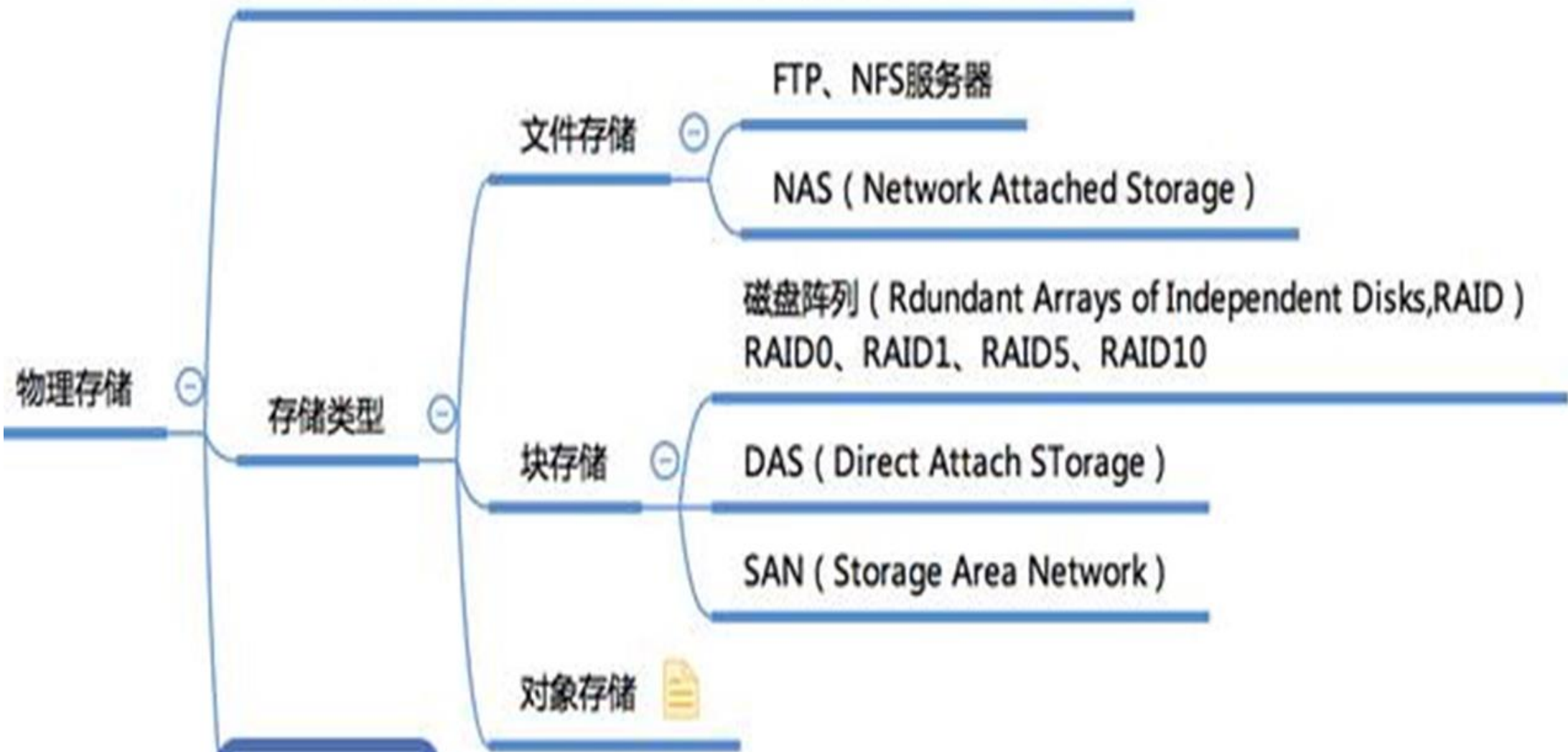


数据组织集成



数据存储

主流的存储系统网络架构有DAS、NAS、SAN三种网络架构



数

据

存

储

分布式文件/对象存储系统



OSS

HDFS

OpenStack Swift

Ceph

GlusterFS

Facebook Hasystack

Lustre

AFS

数

据

存

储

分布式关系型数据库



DRDS

TiDB

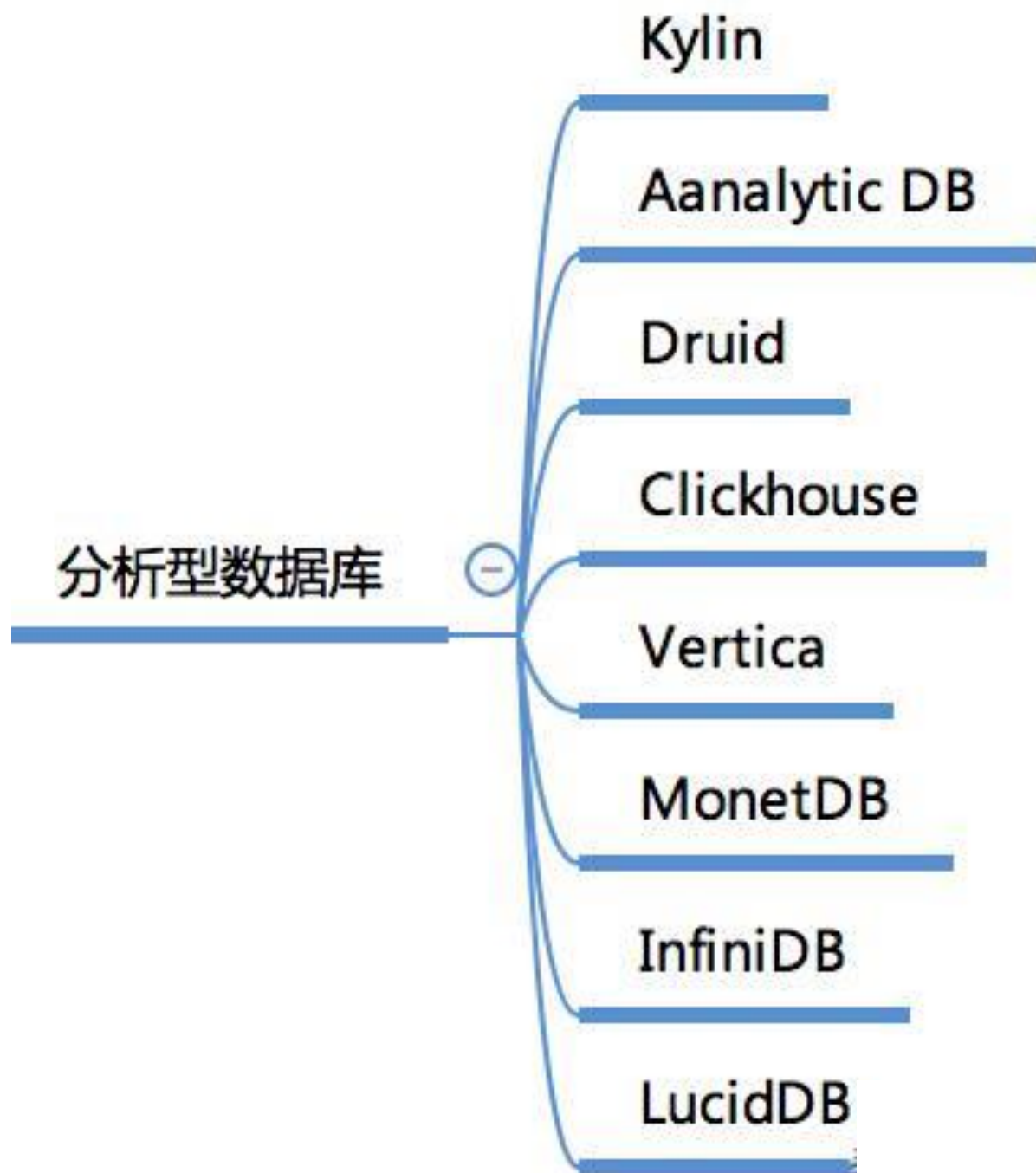
GreenPlum

Mycat

Cobar

Aurora

数
据
存
储

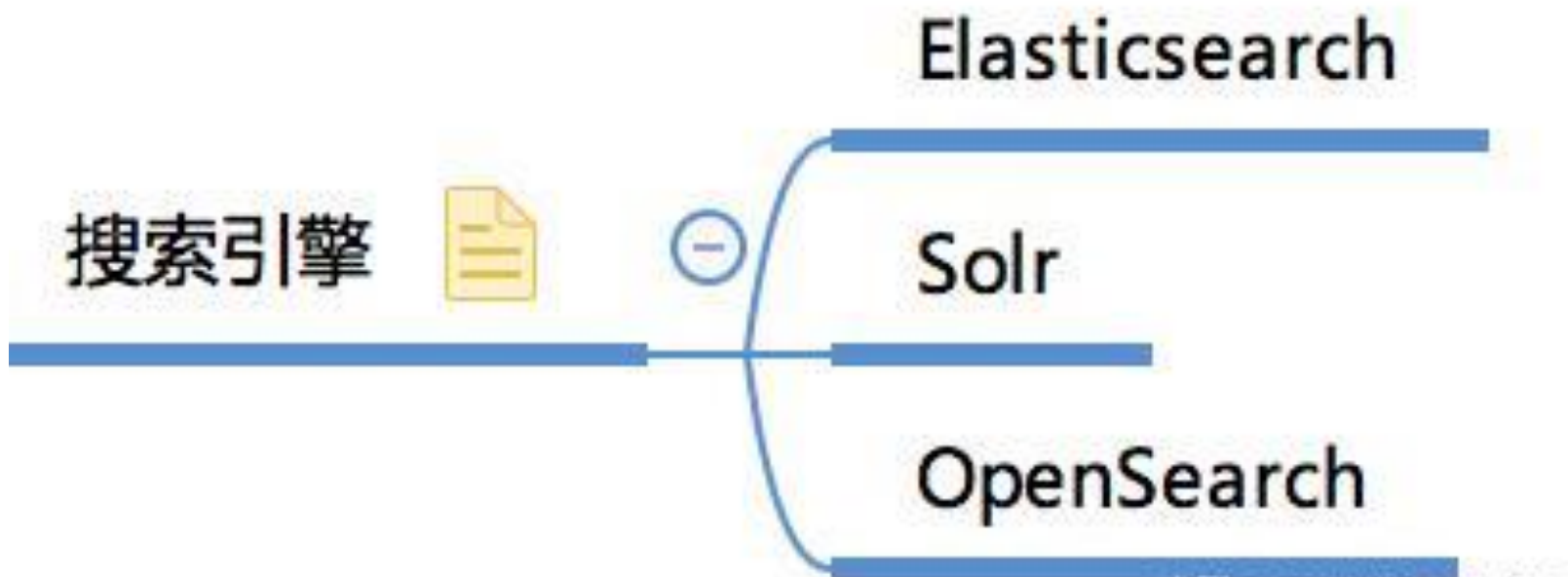


数

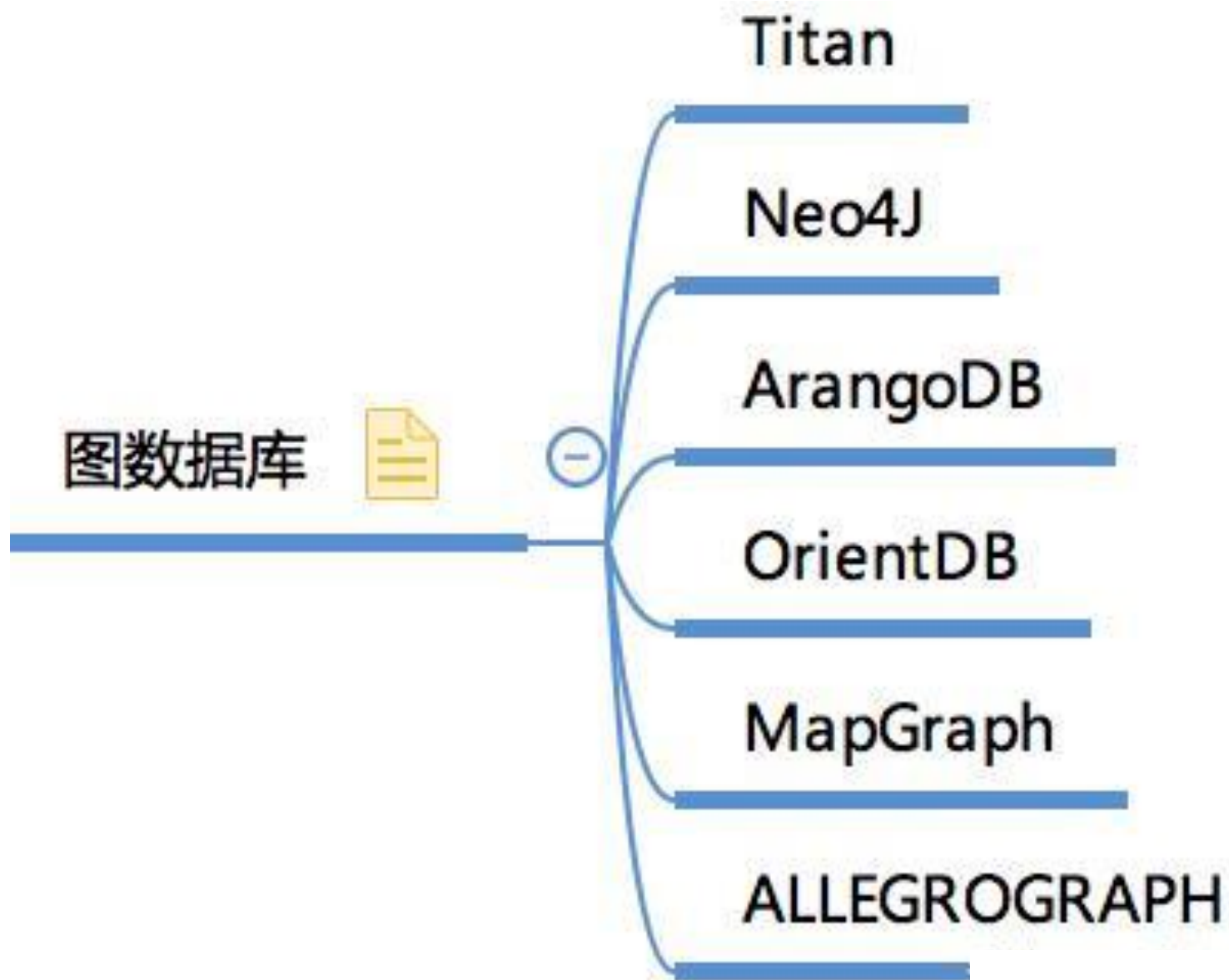
据

存

储



数
据
存
储



数
据
存
储

列存储数据库



Cassandra

Kudu

Hbase

Hypertable

Apache Accumulo

数

据

存

储

文档数据库



MongoDb

CouchDB

MarkLogic

OrientDB

数

据

存

储

键值存储数据库



Redis

Memcached

Tair

数据组织集成

数据计算

流式计算 (stream computing)

大规模批量计算 (batch computing)

即席查询分析(ad-hoc computing)

全量计算&增量计算

图计算

分布式协调系统

集群资源管理和调度

工作流管理引擎

数 据 计 算

流式计算 (stream computing)



Storm

Flink

Yahoo S4

Kafka Stream

Twitter Heron

Apache Samza

Spark Streaming

数 据 计 算

大规模批量计算 (batch computing)



Tez

✓ MapReduce

Hive

Spark

Pig

大数据的编程模型：Apache Beam

数 据 计 算

即席查询分析(ad-hoc computing)



Impala

Hawq

Dremel

Drill

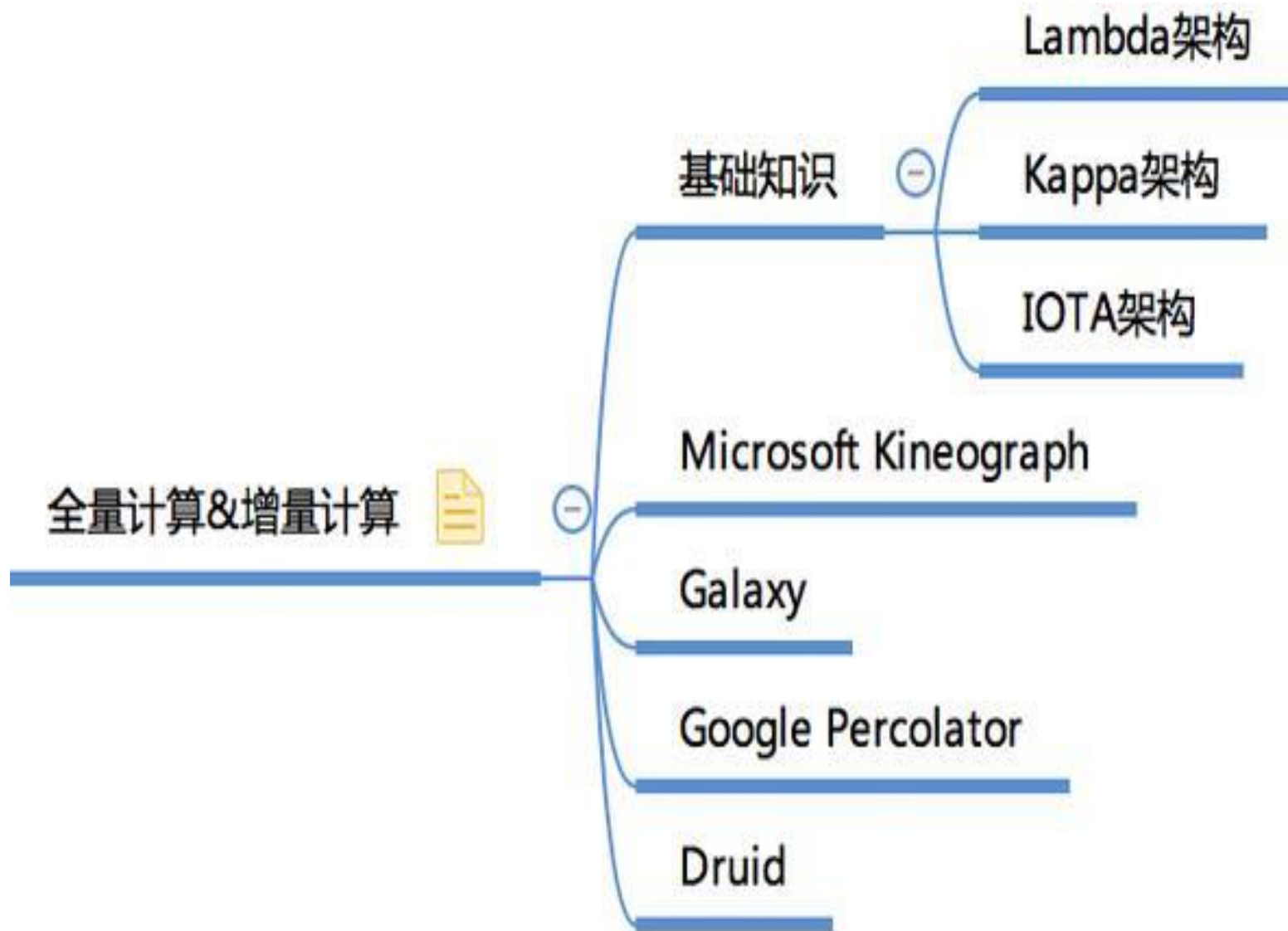
Phoenix

Tajo

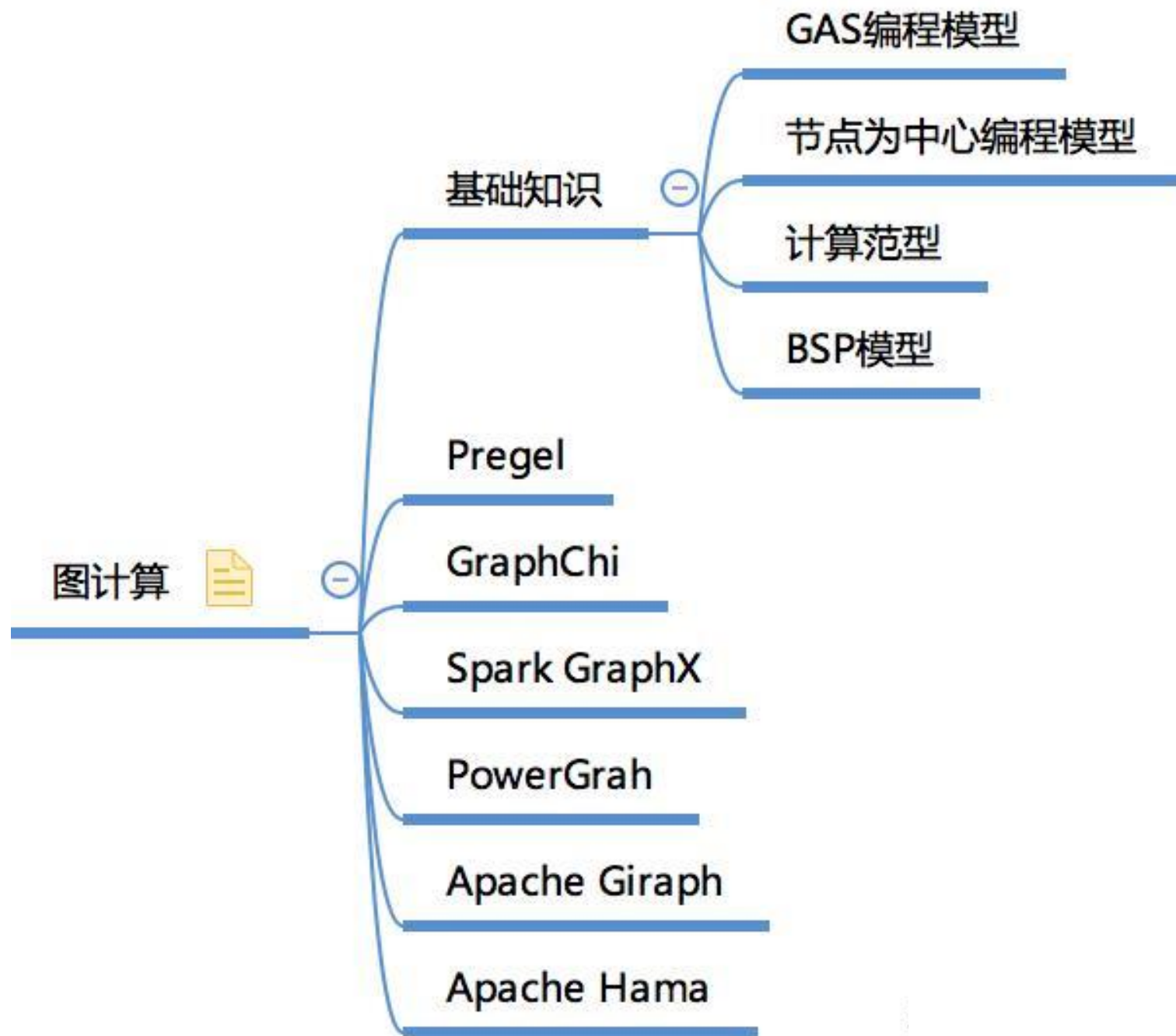
Presto

Hortonworks Stinger

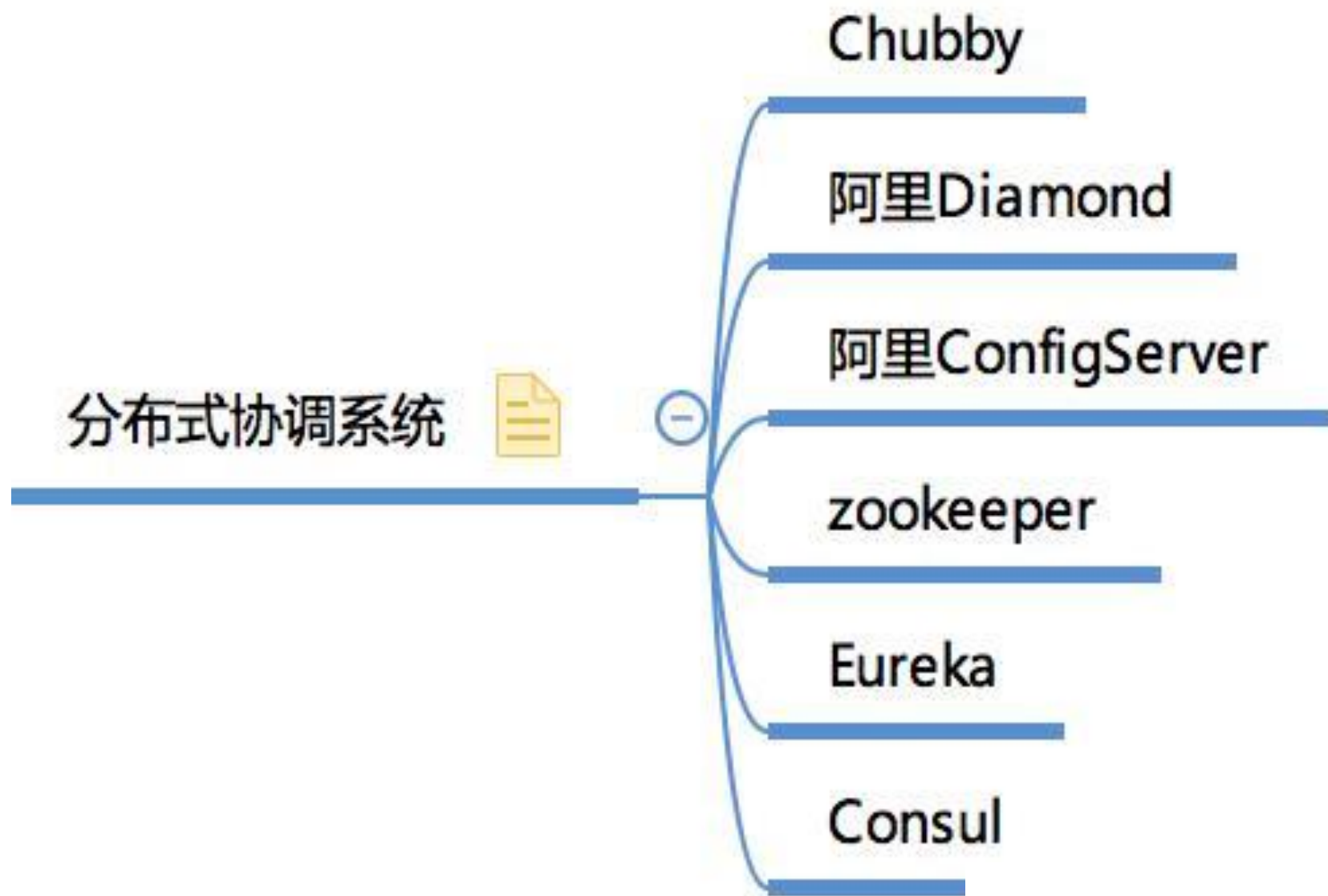
数 据 计 算



数 据 计 算



数据组织集成



数据组织集成



数
据
组
织
集
成

workflows management engine



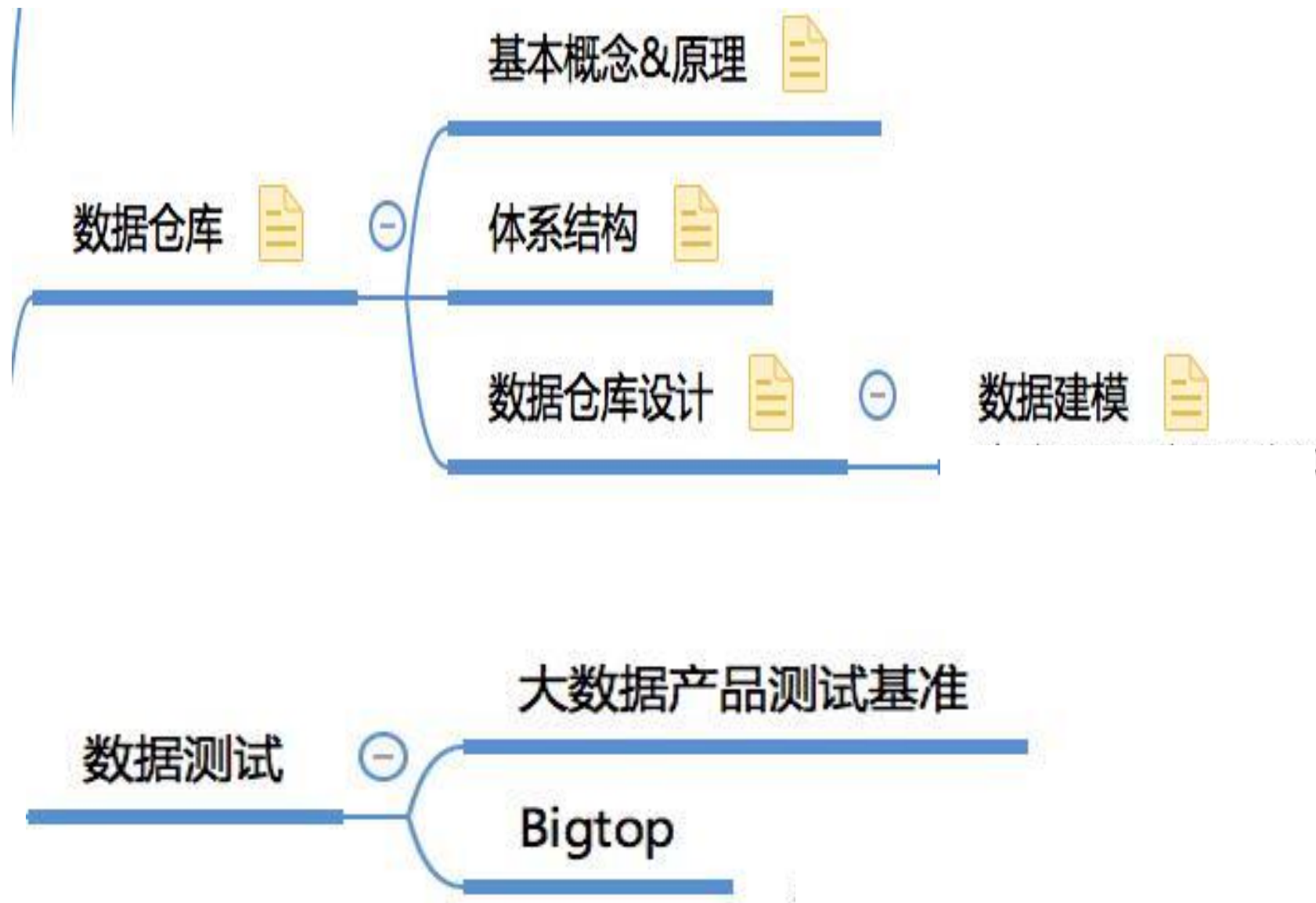
Oozie

Azkaban

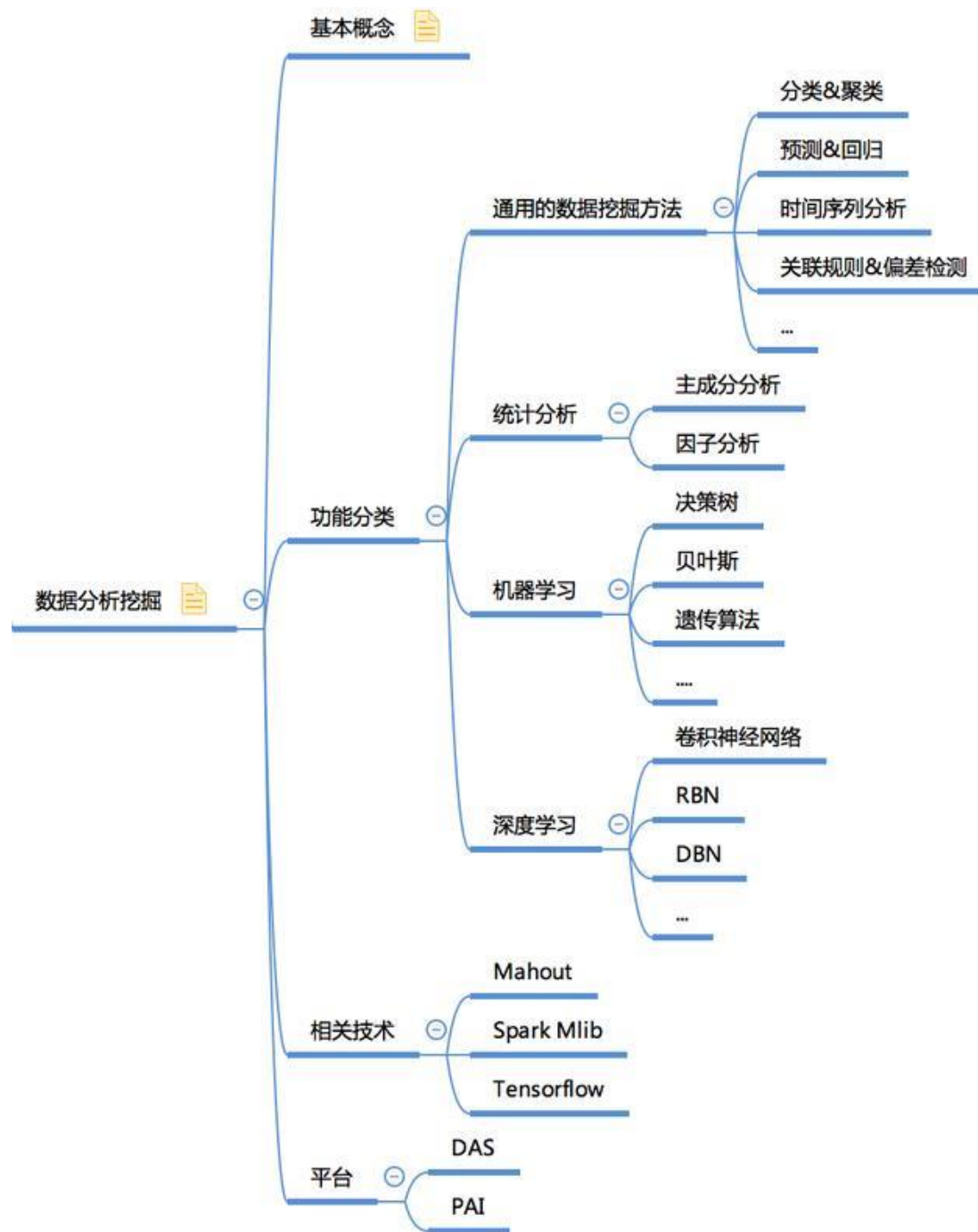
Luigi

Airflow

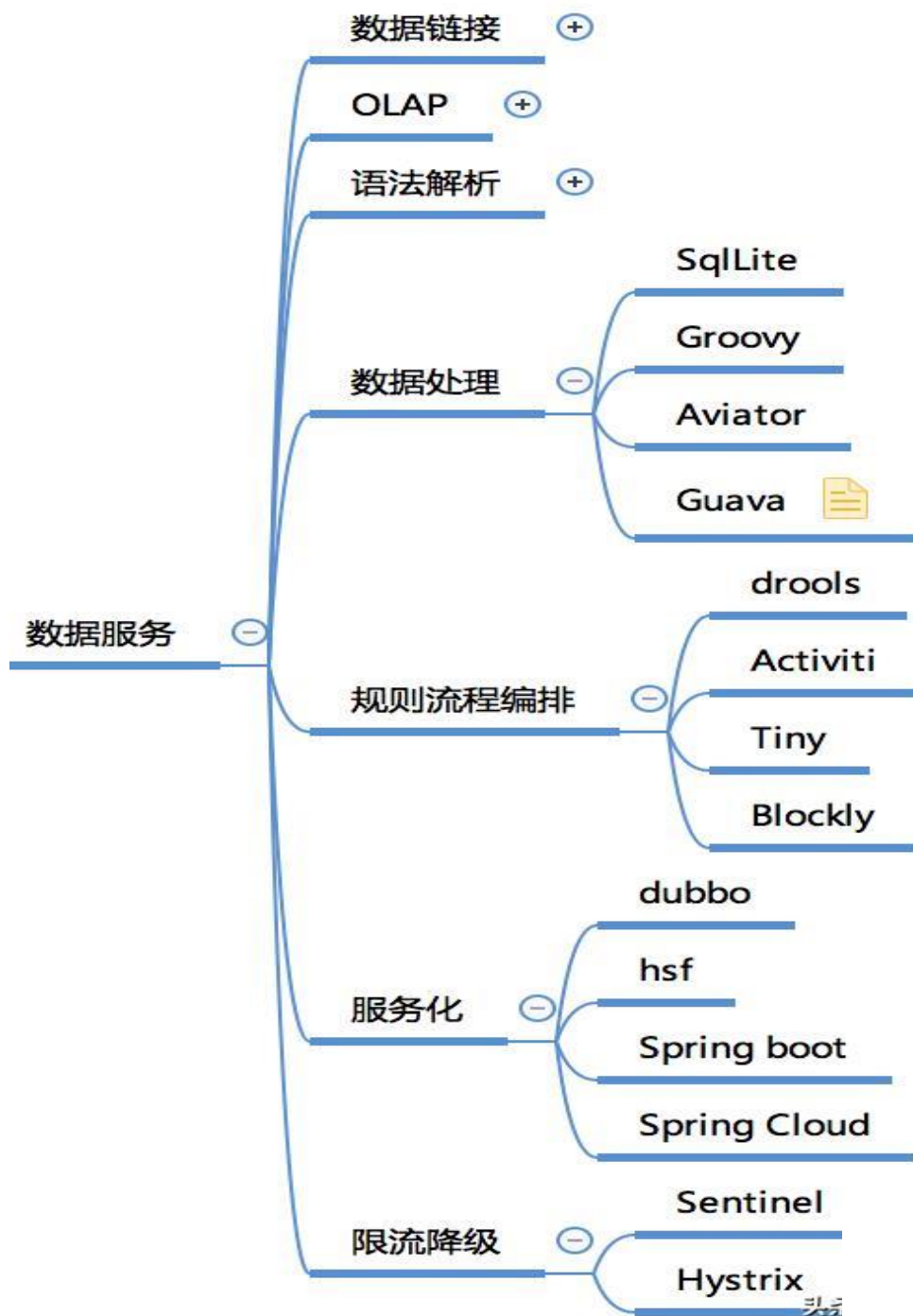
数据组织集成



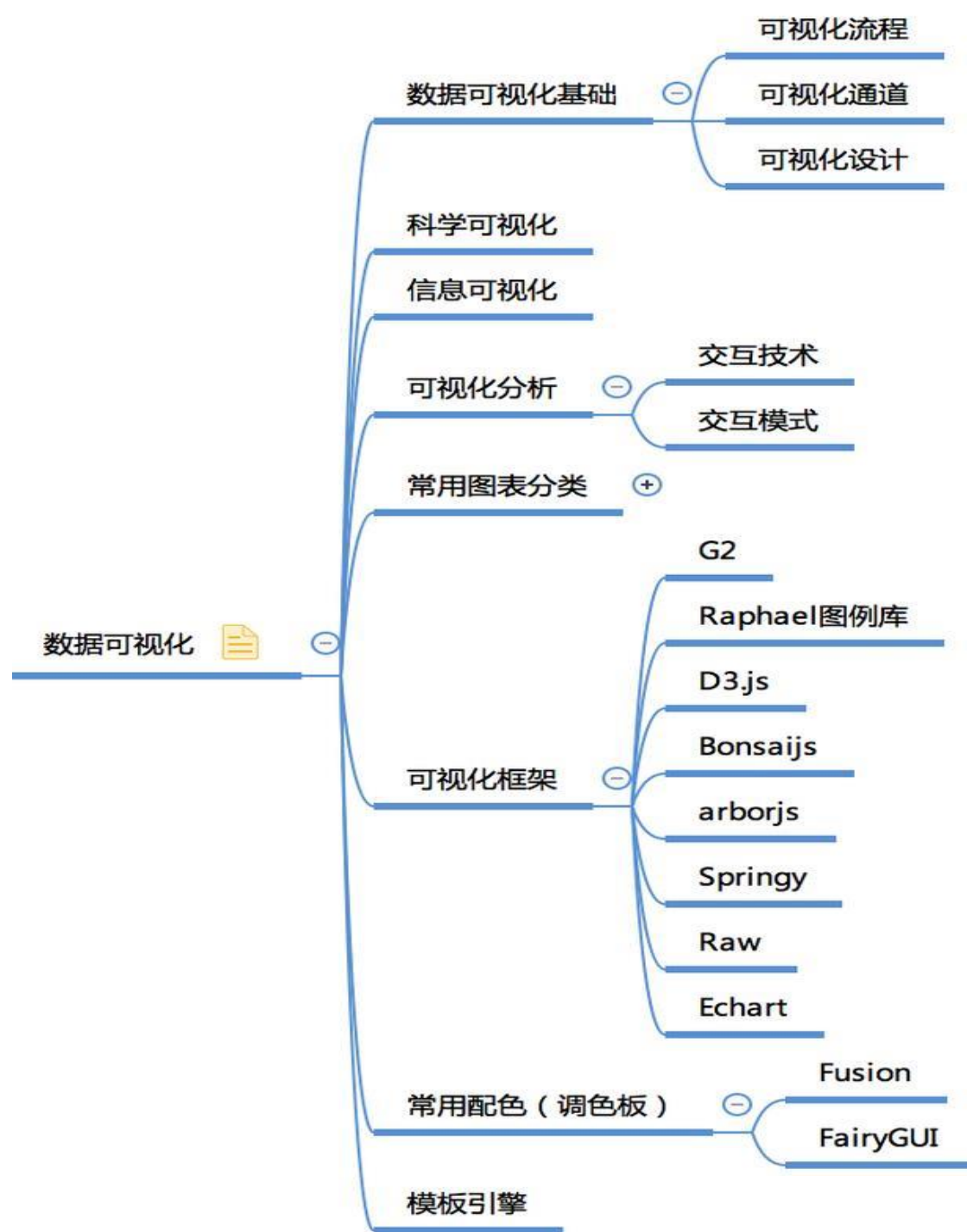
数 据 应 用



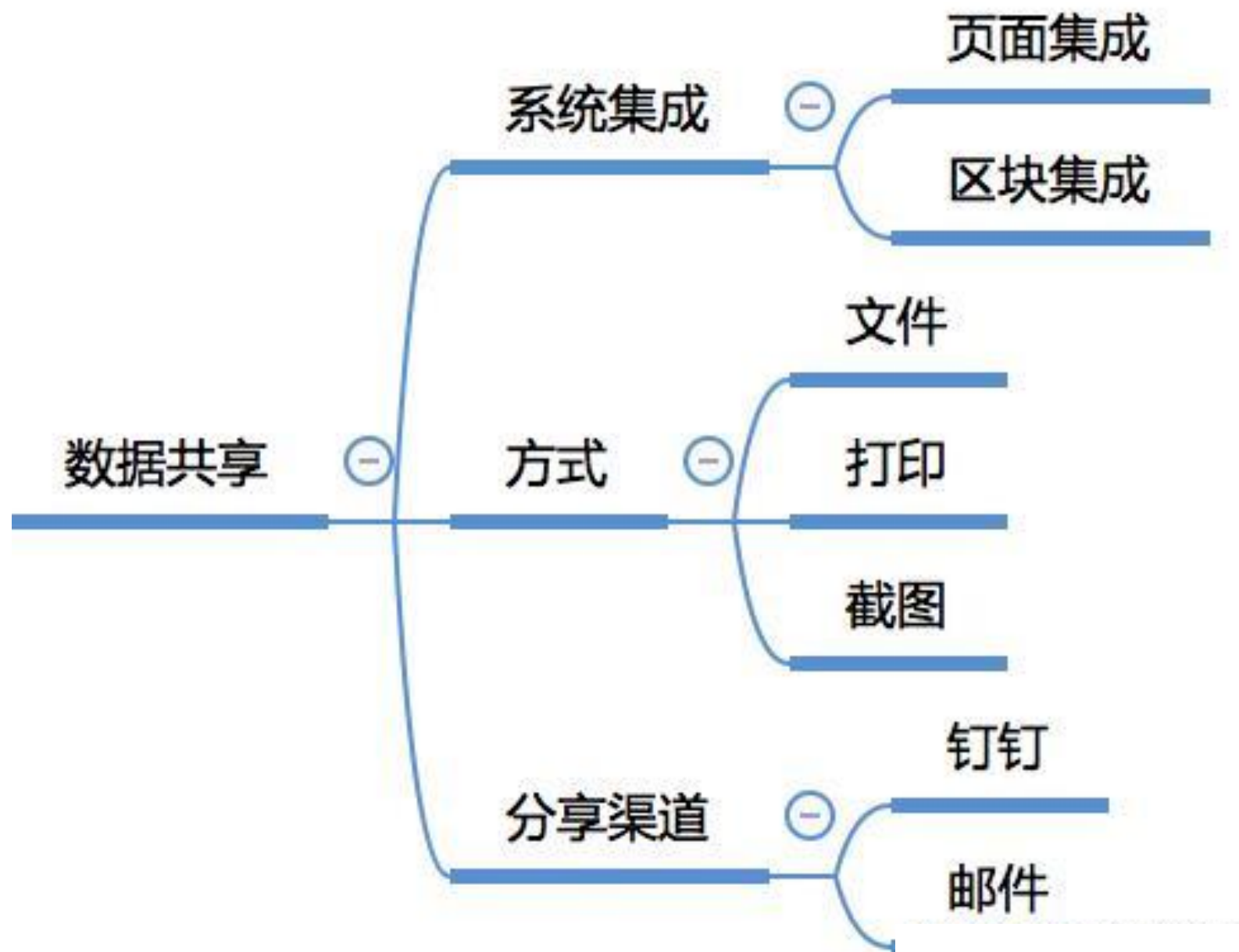
数据应用



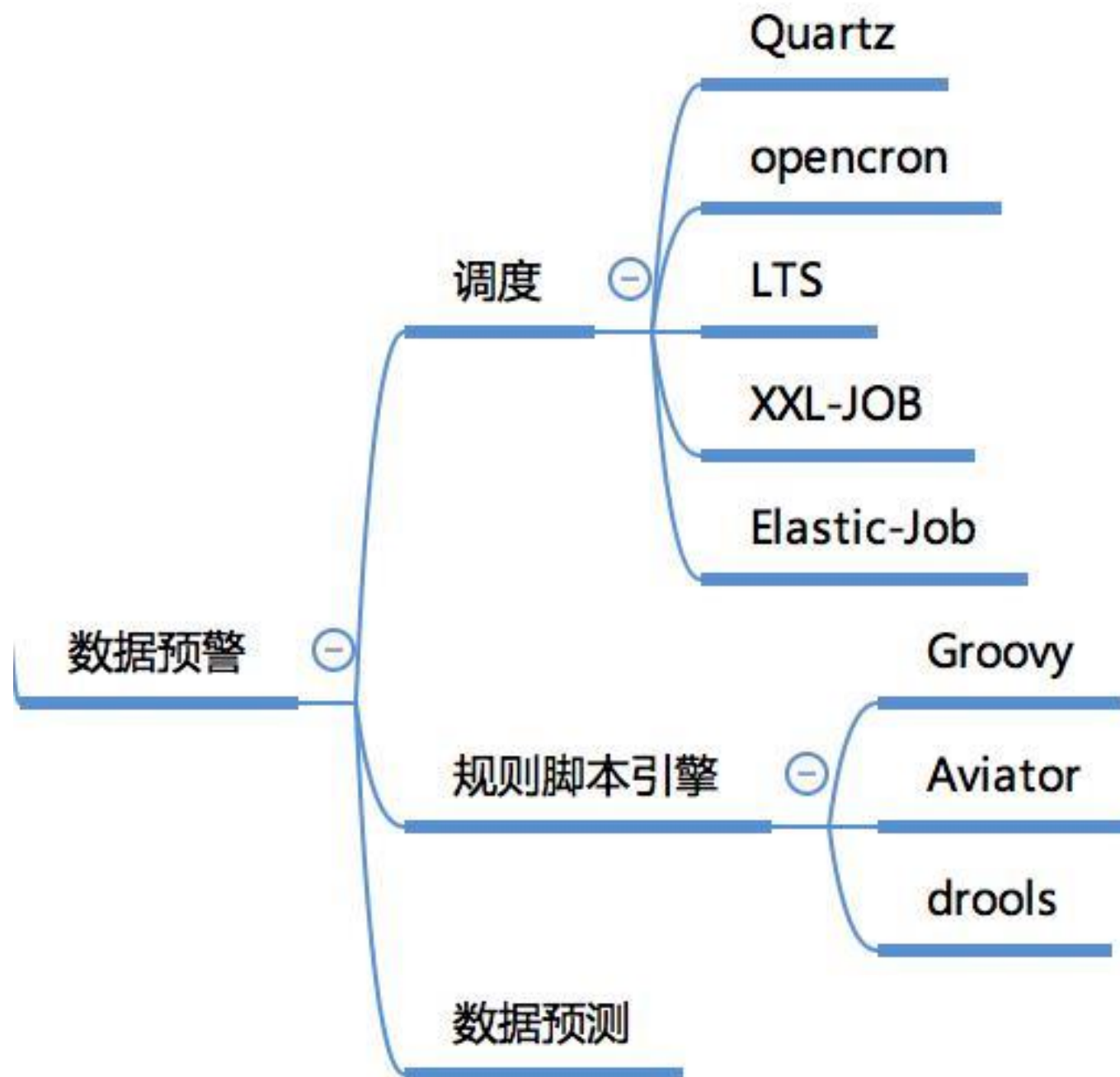
数据应用



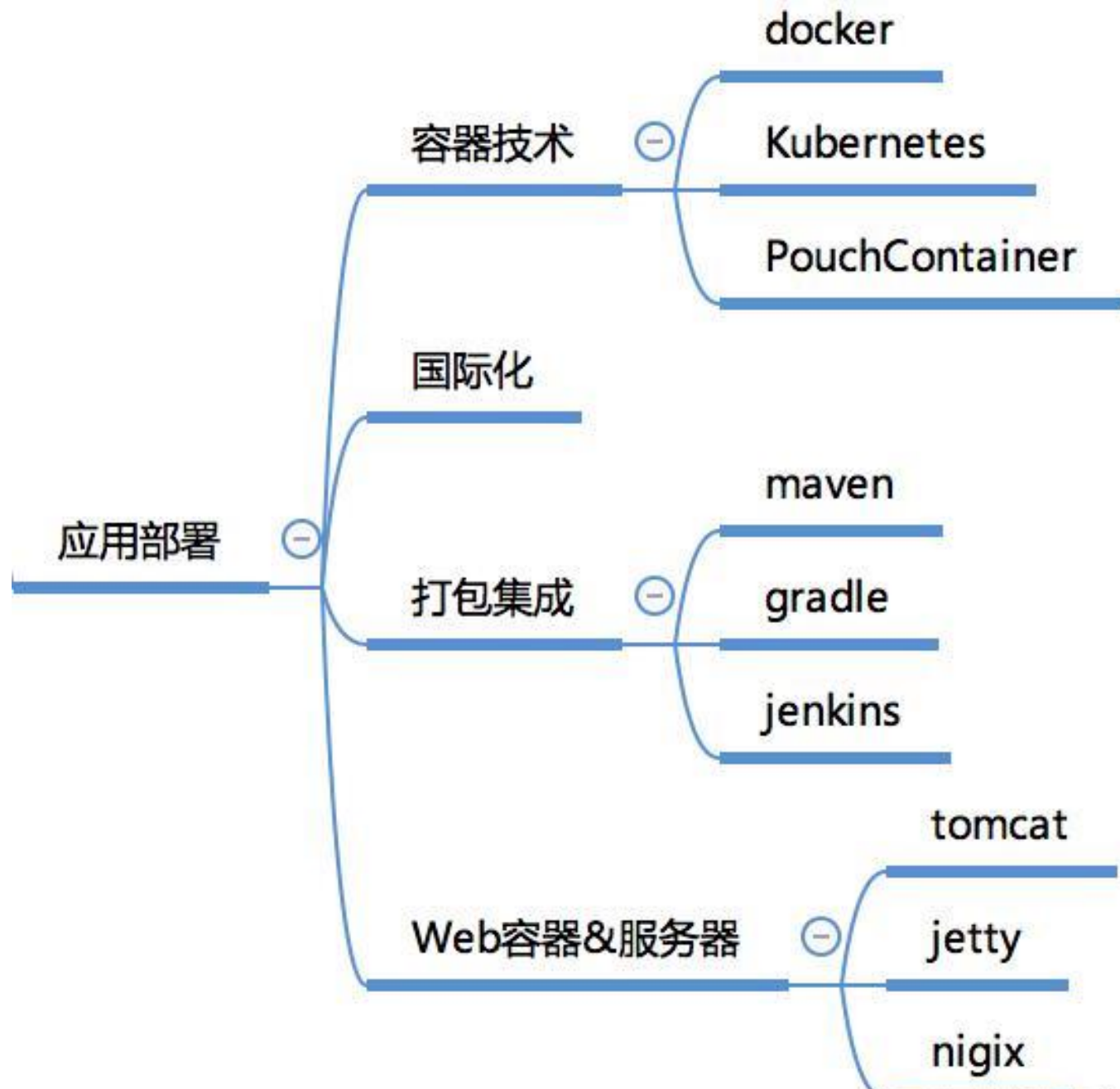
数据应用



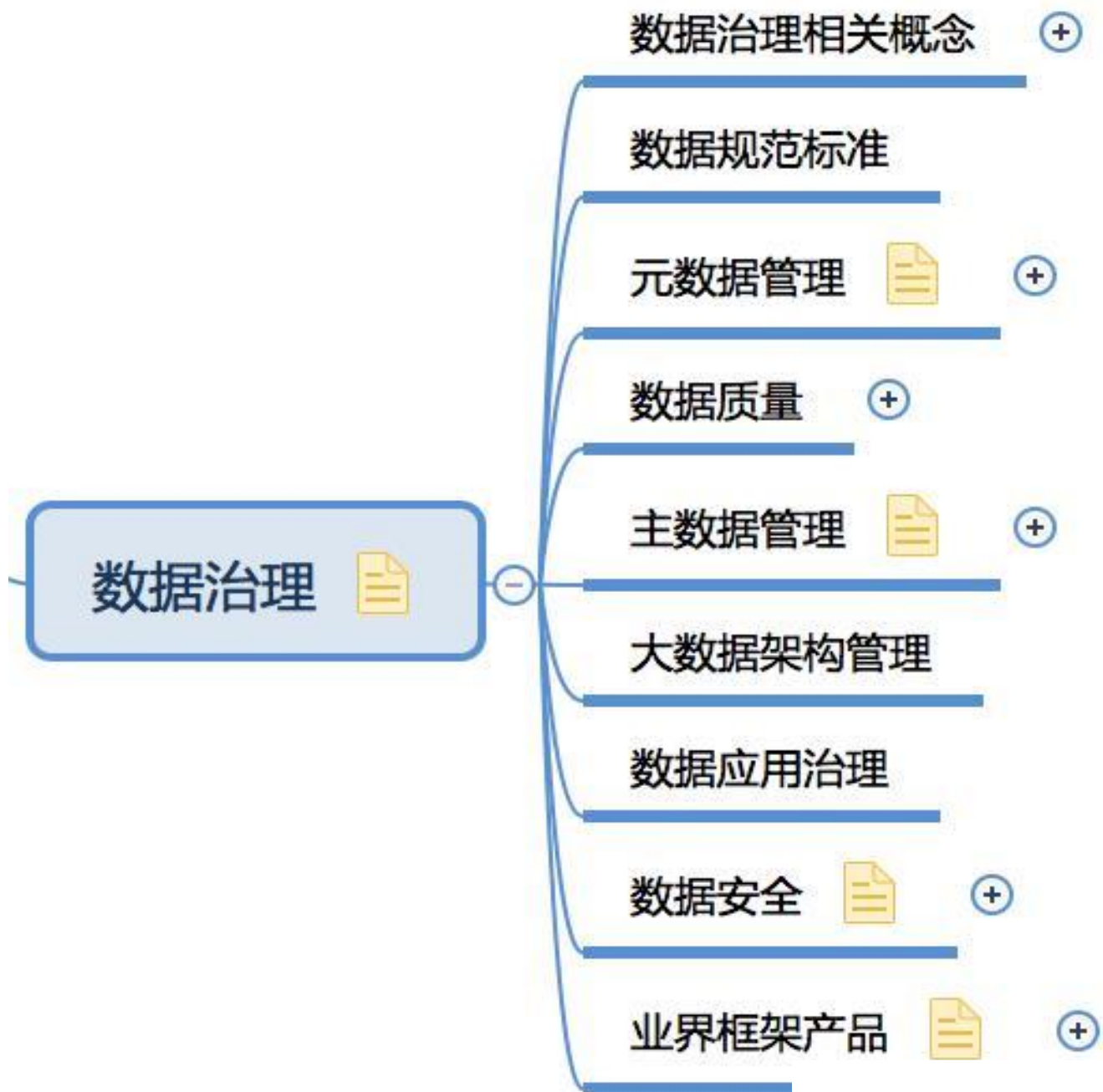
数据应用



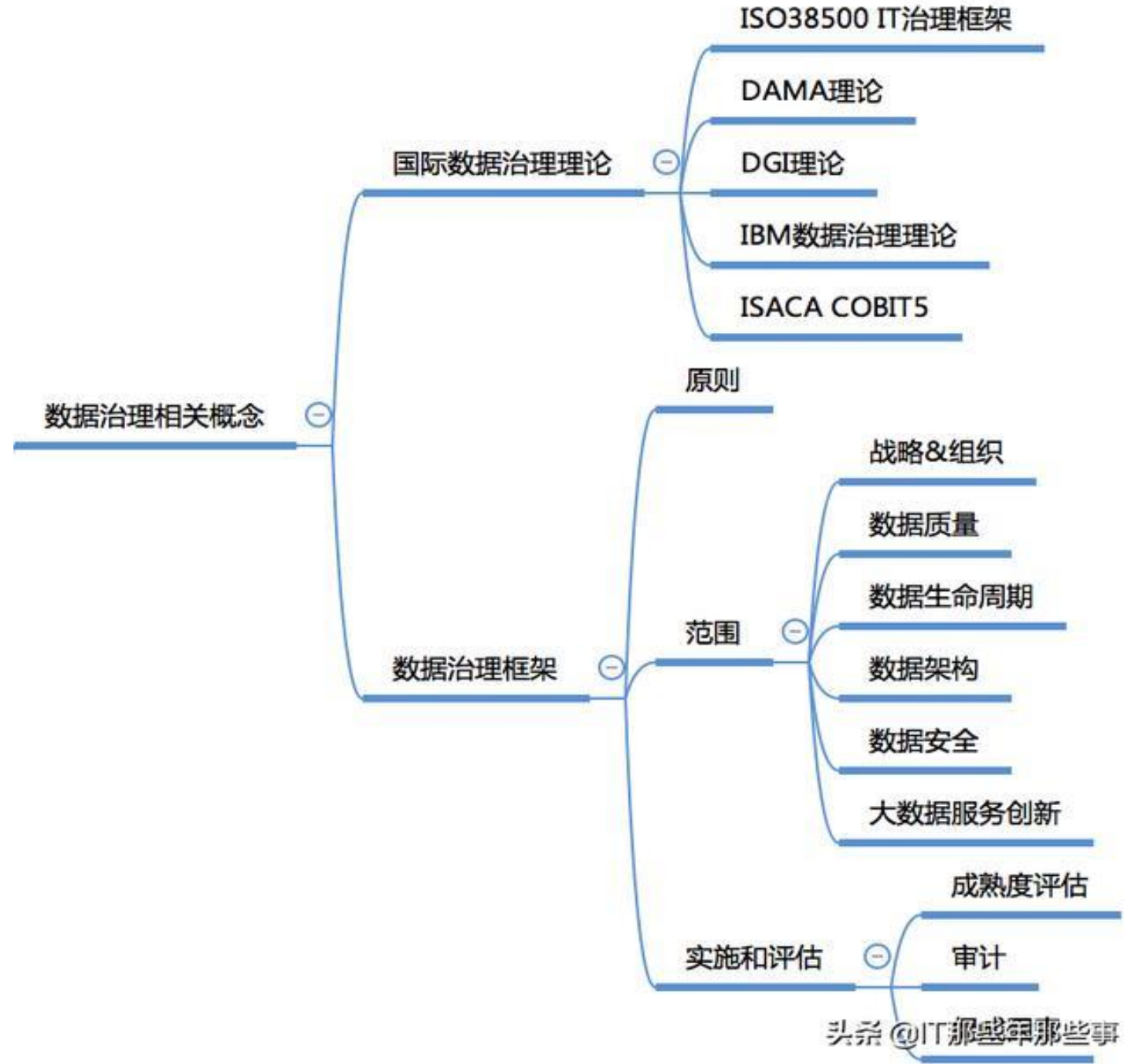
数据应用



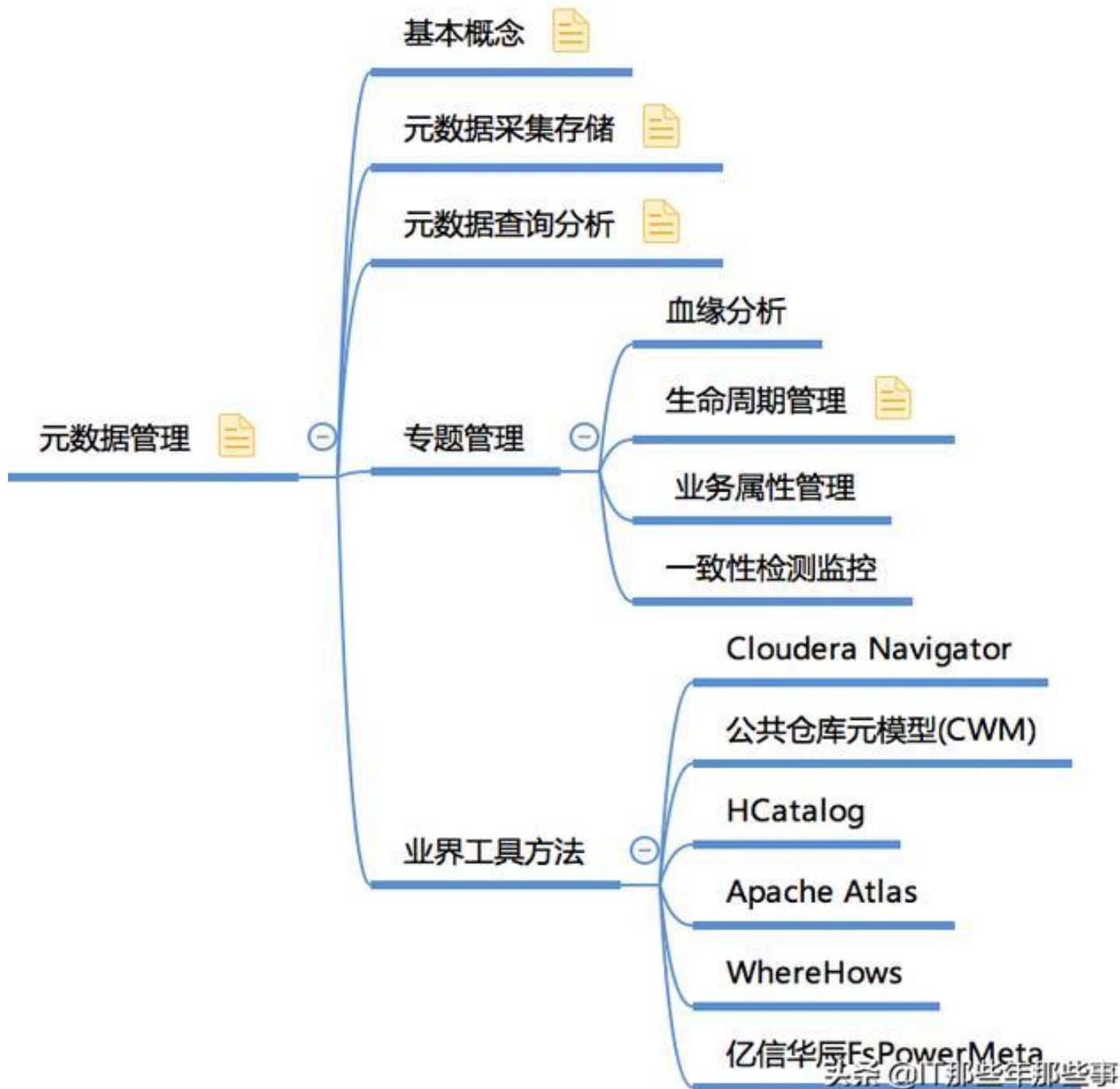
数据治理



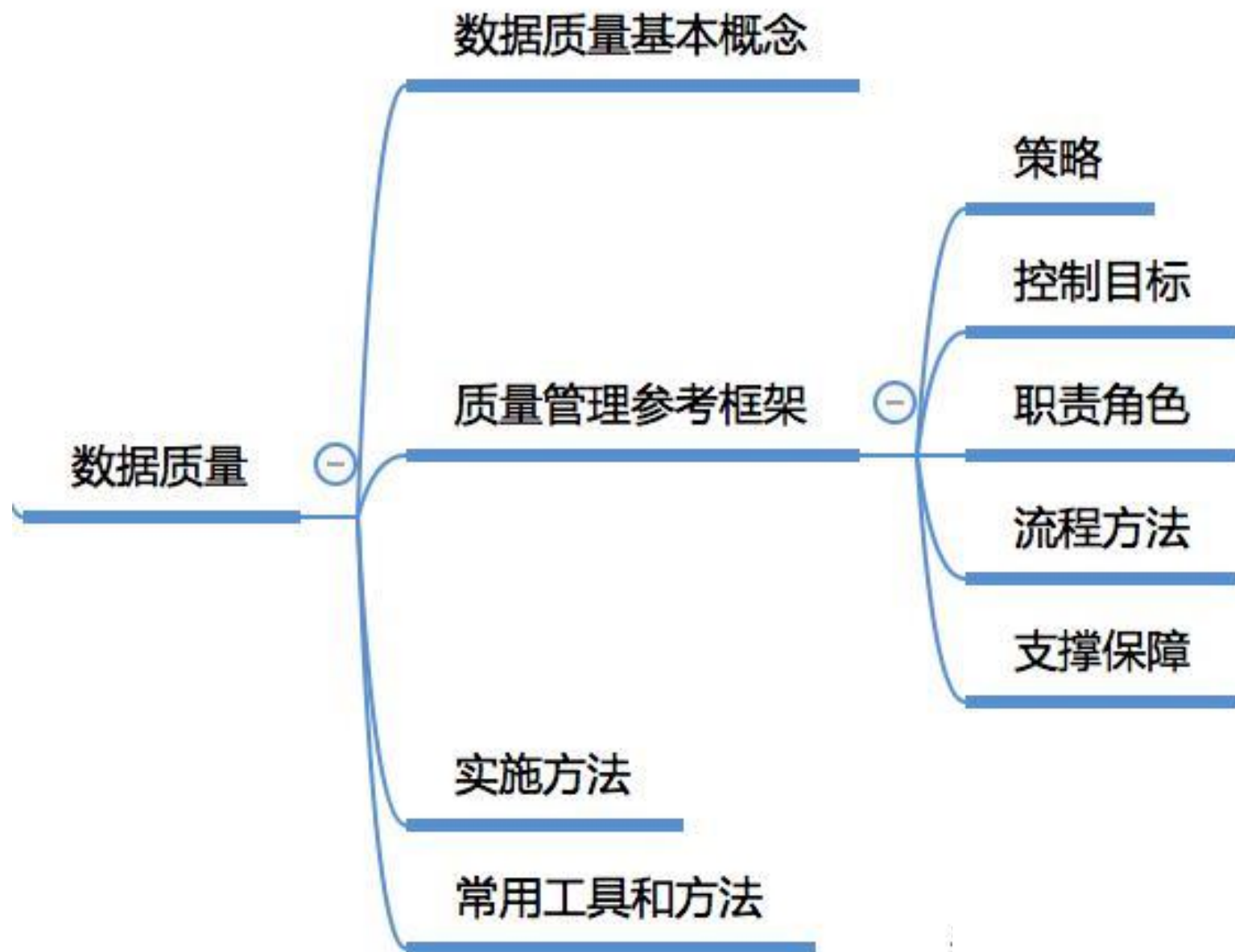
数据治理



数 据 治 理



数据治理

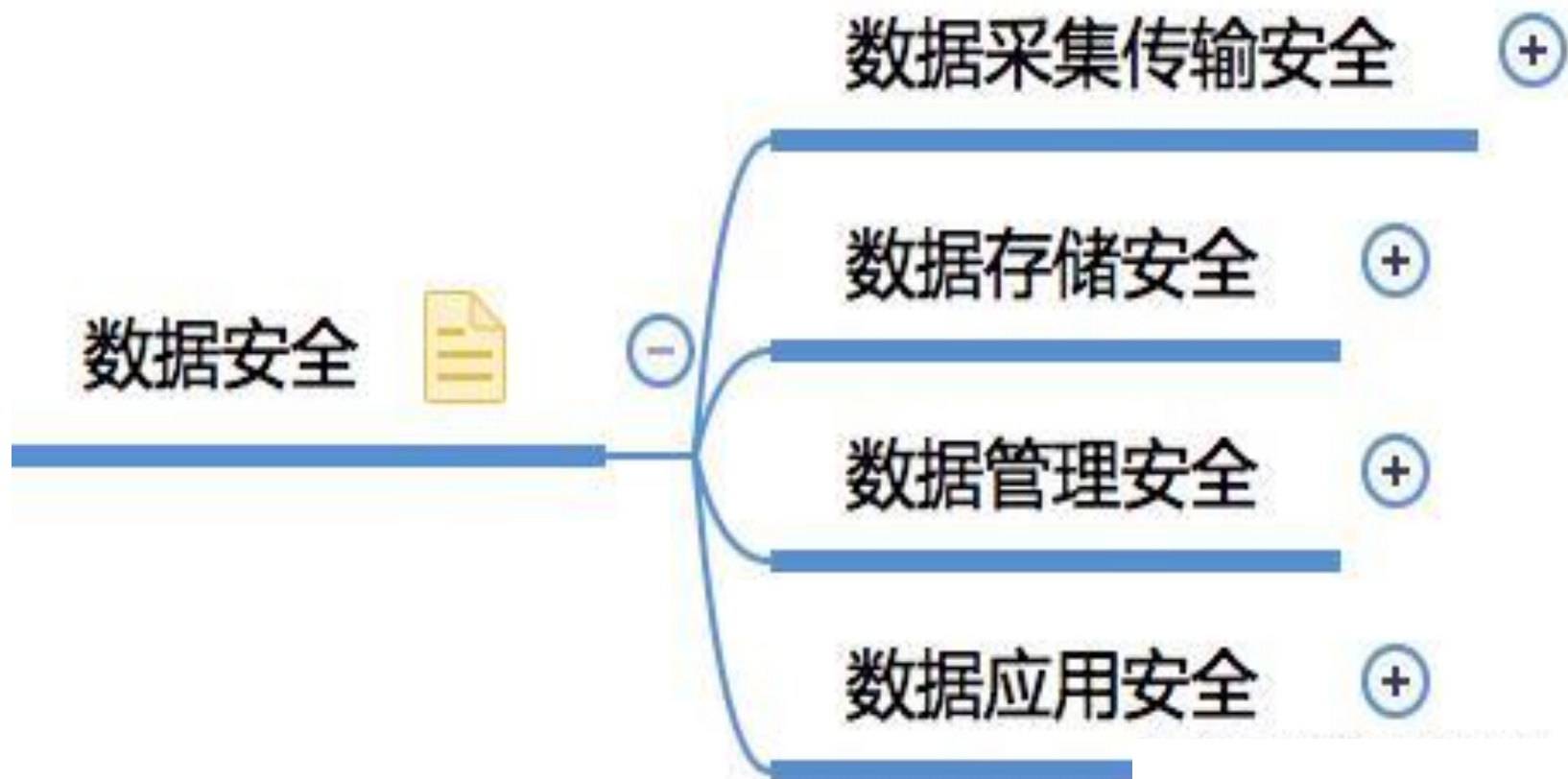


数

据

治

理



数

据

安

全

数据采集传输安全

VPN

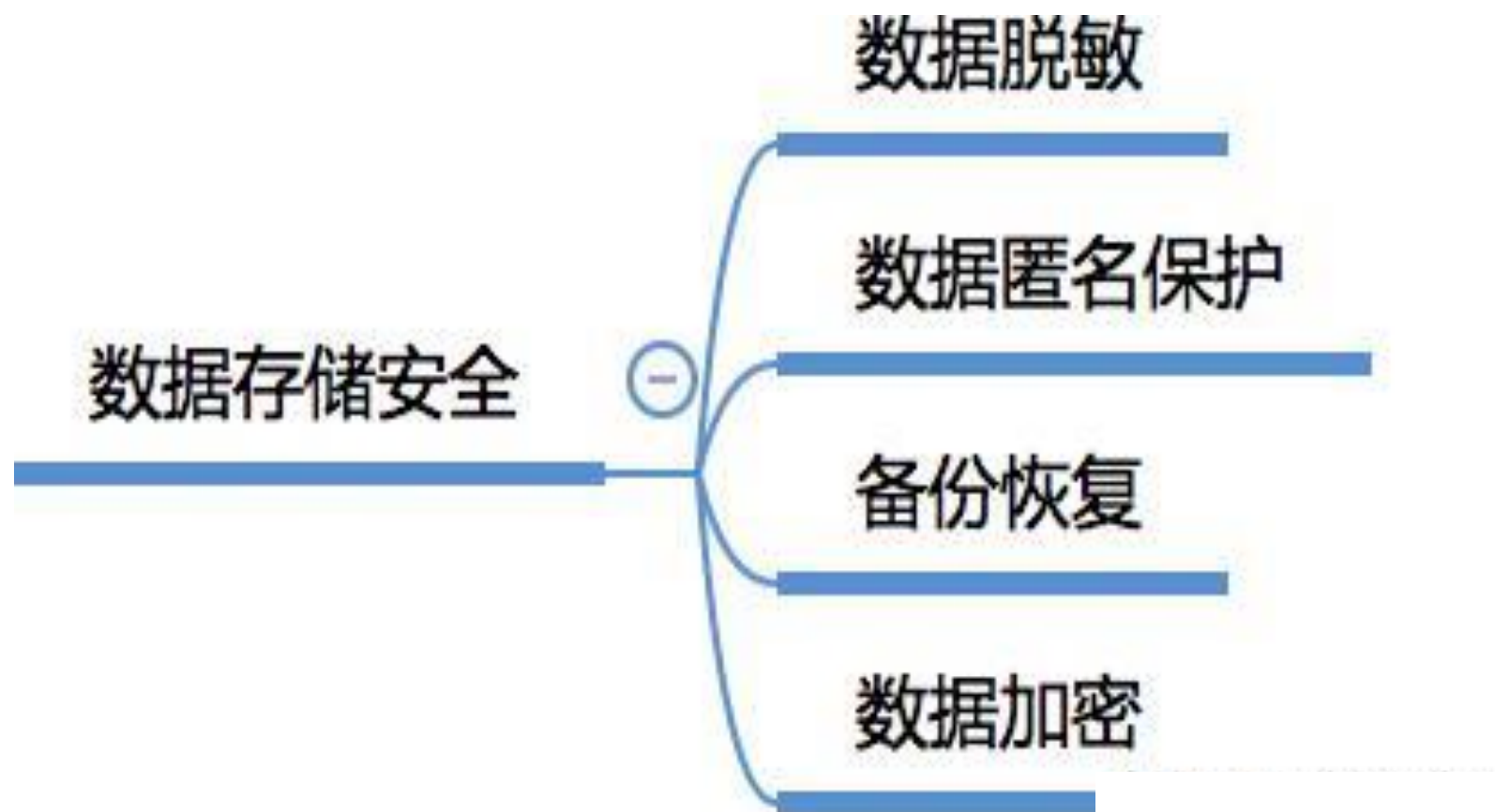
SSL& HTTPS

数字签名

数据加密



数据安全



数据安全



数据安全

