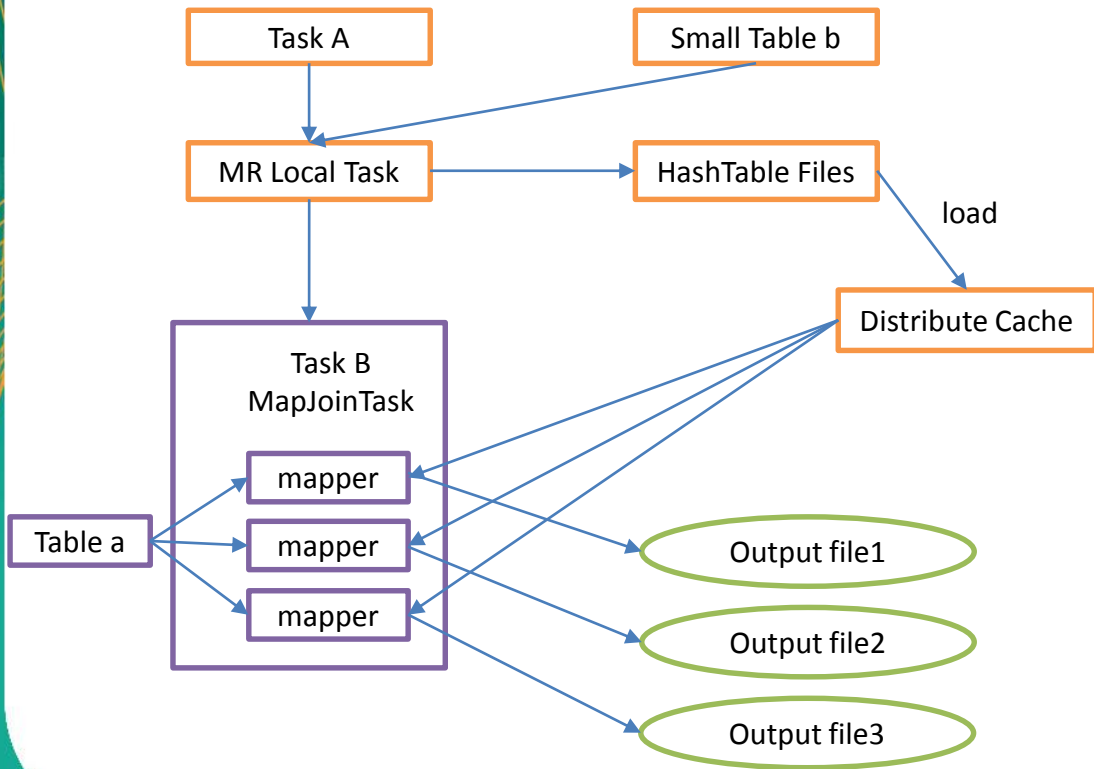




大数据技术图解

签名：汪辉

MapJoin

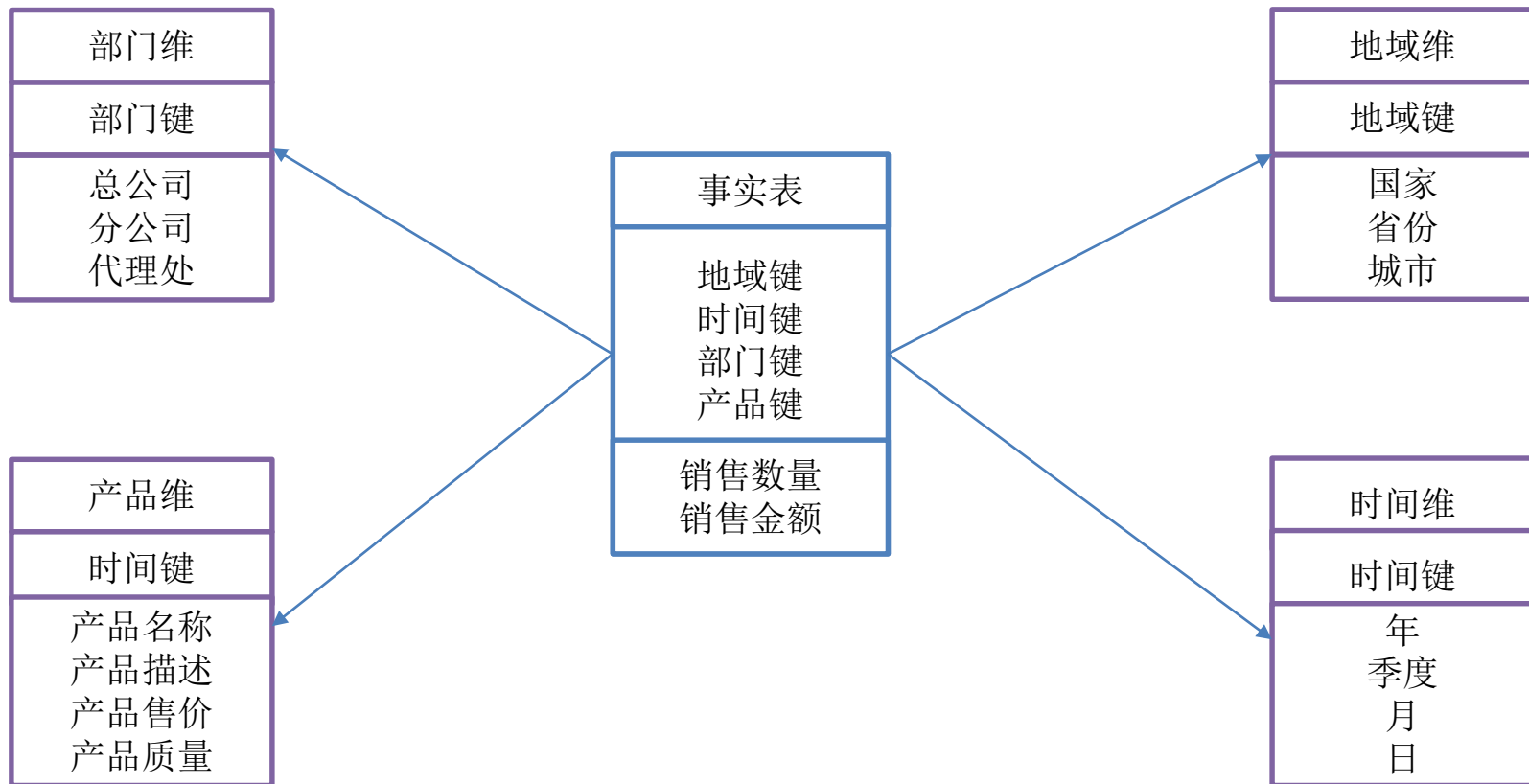


1) Task A, 它是一个Local Task（在客户端本地执行的Task），负责扫描小表b的数据，将其转换成一个HashTable的数据结构，并写入本地的文件中，之后将该文件加载到DistributeCache中。

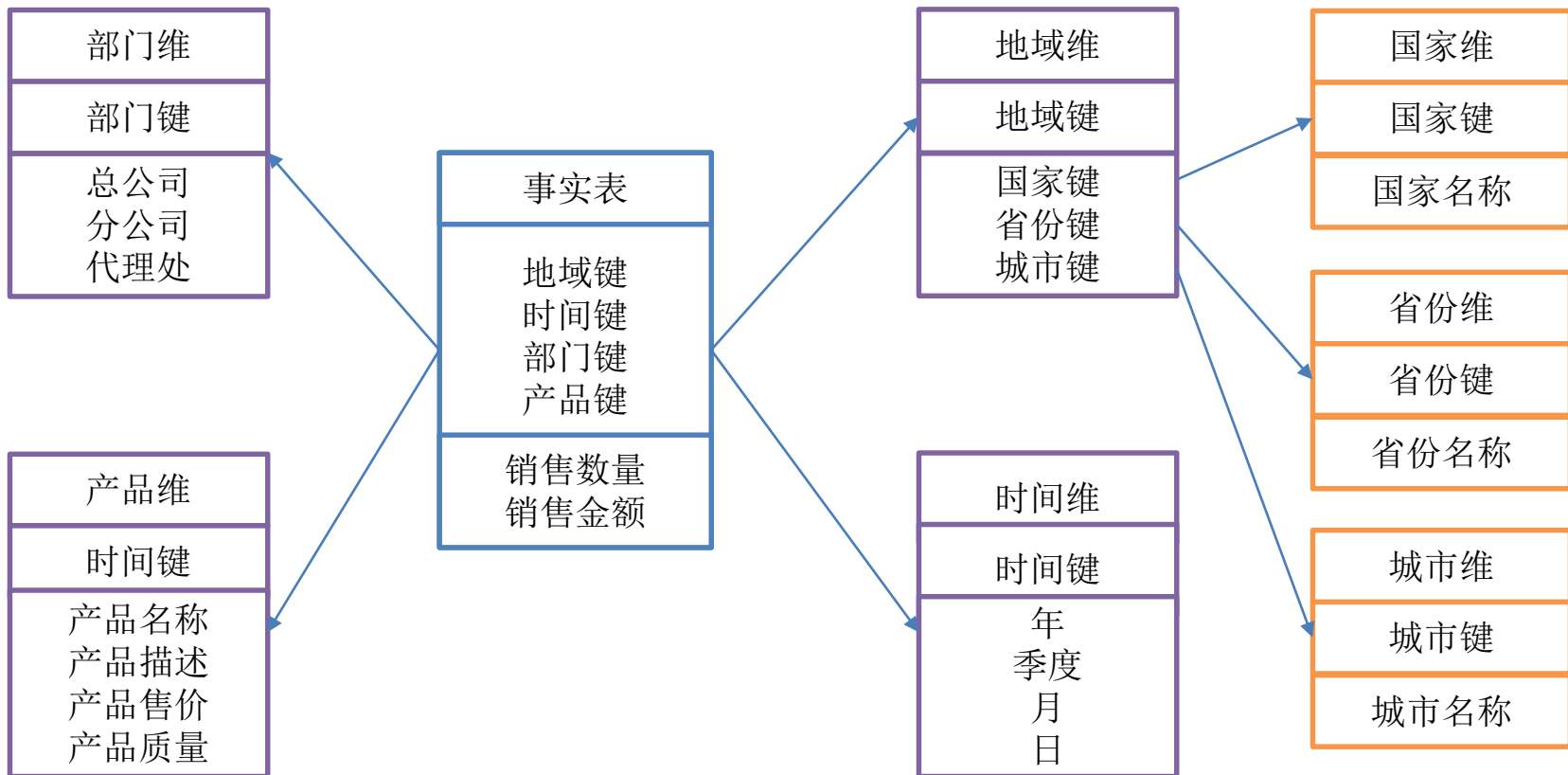
2) Task B, 该任务是一个没有Reduce的MR，启动MapTasks扫描大表a，在Map阶段，根据a的每一条记录去和DistributeCache中b表对应的HashTable关联，并直接输出结果。

3) 由于MapJoin没有Reduce，所以由Map直接输出结果文件，有多少个Map Task，就有多少个结果文件。

星型模型



雪花模型



Hive编译器组成

