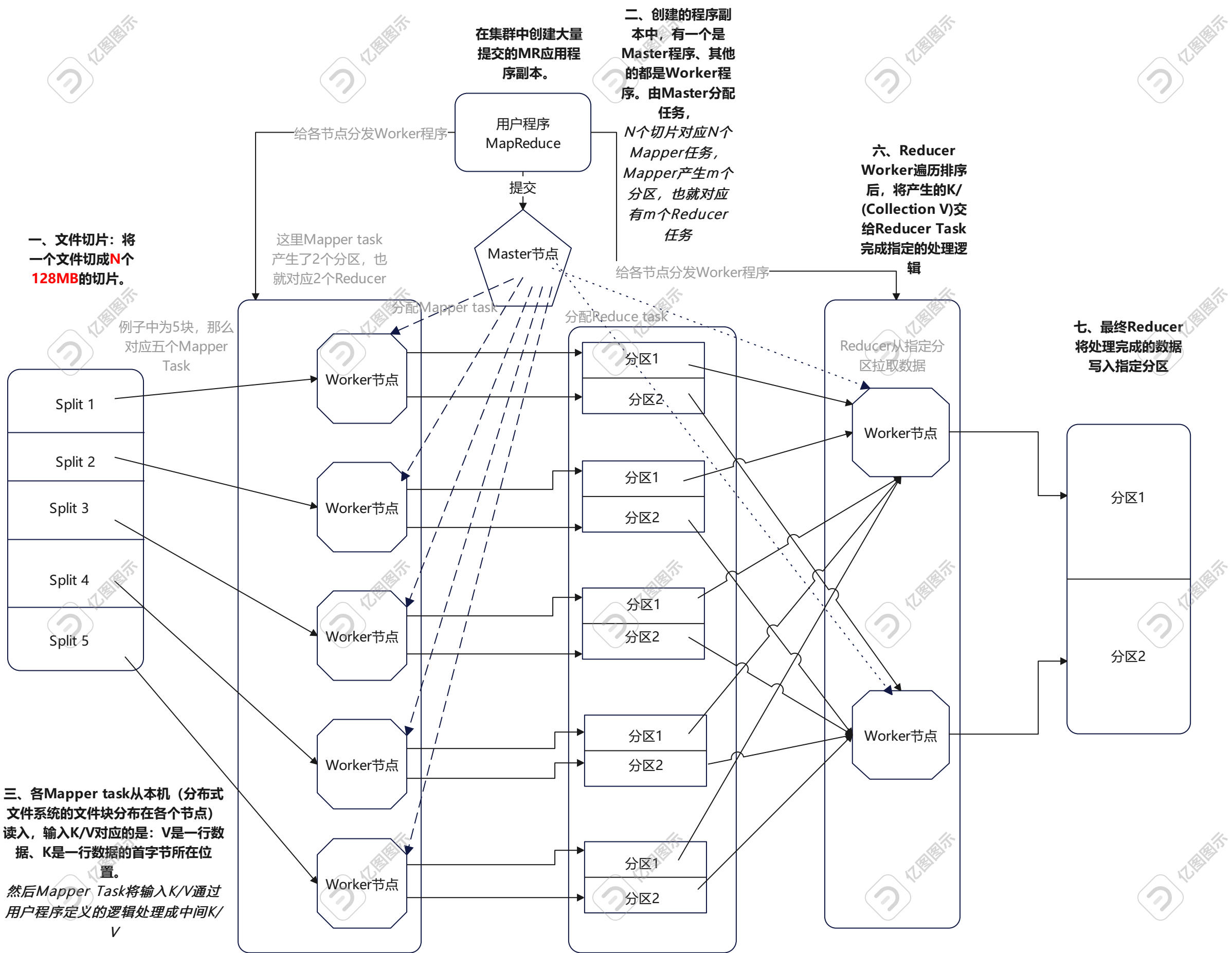
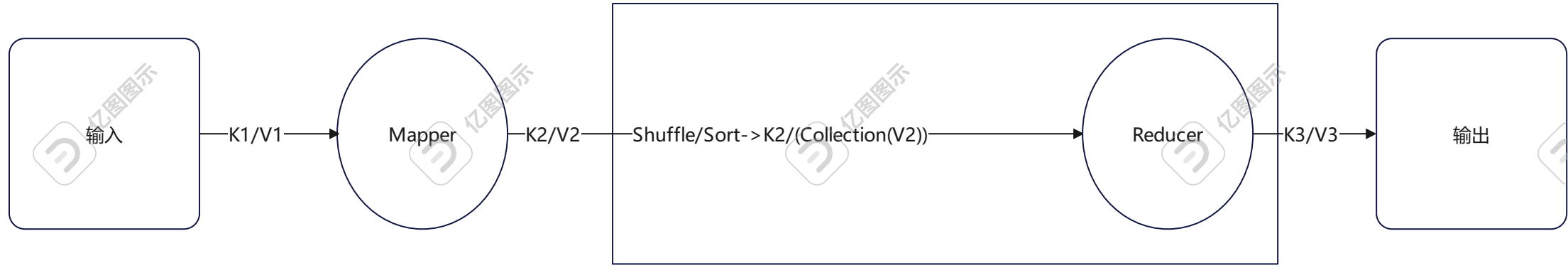


MR数据流



输入数据不一定是文件，可以实现RecordReader接口来自定义输入数据

Map执行完成后，数据可能出现大量相同的K，为了减少Reducer阶段远程拉取数据产生的网络流量和等待时间，可以指定一个Combiner，提前聚合这些K相同的值

分区的行为可以自定义指定，比如让什么样的Key进入哪个分区。需要继承Partitioner类

远程数据传输涉及到序列化，可以把自定义的类进行序列化，当然要实现Writable接口

二次排序：将Reducer的Sort完成后的K/(Collection V)中的一系列的V进行排序

输出不一定是文件，可以实现RecordWriter接口自定义输出到哪里