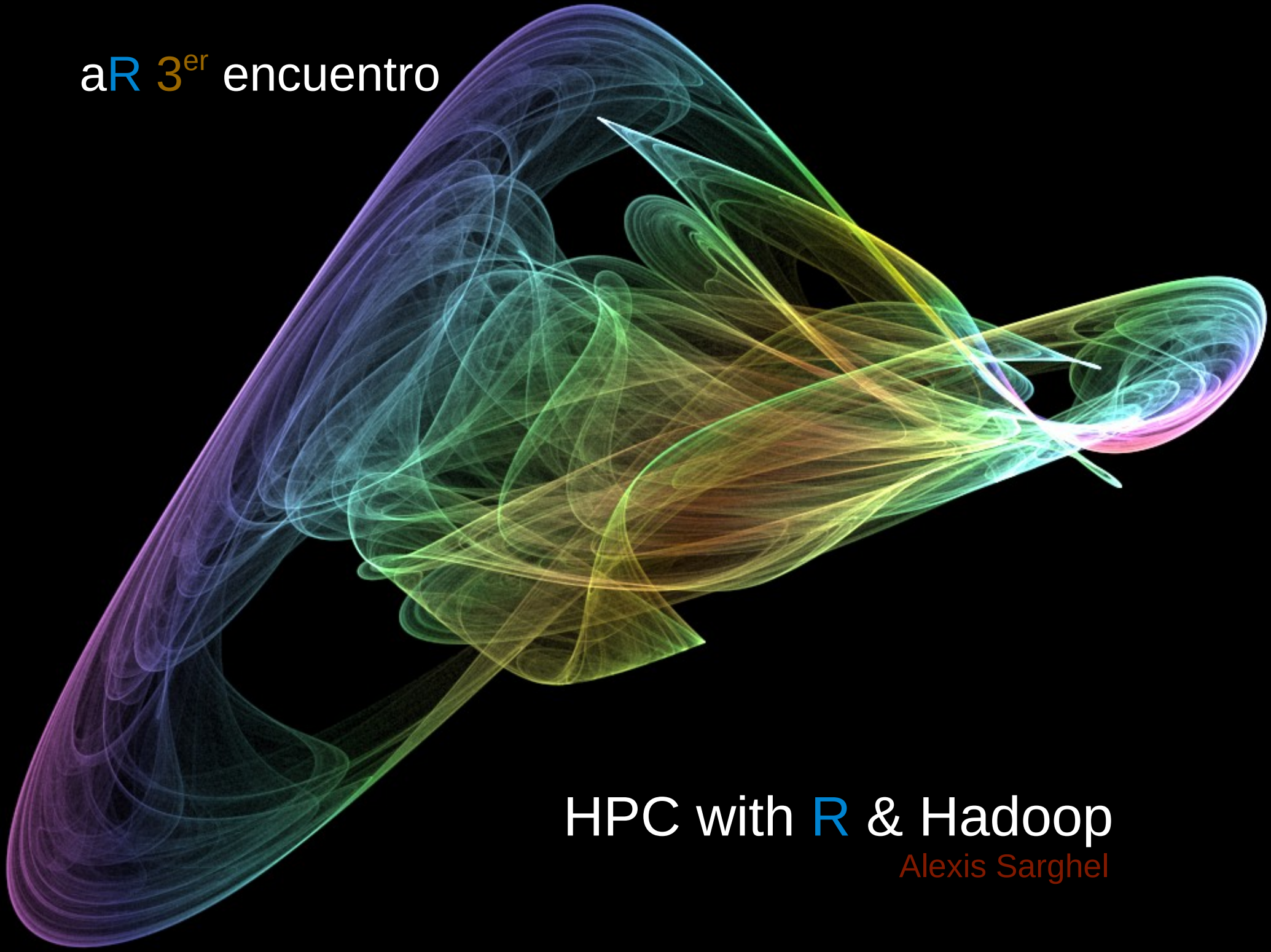


aR 3^{er} encuentro



HPC with R & Hadoop

Alexis Sarghel



HPC ?



Hight-performance computing

A supercomputer is a computer
with a high-level computational capacity
compared to a general-purpose computer

Performance of a **supercomputer** is
measured in floating point operations per second
(**FLOPS**)
instead of million instructions per second (**MIPS**)

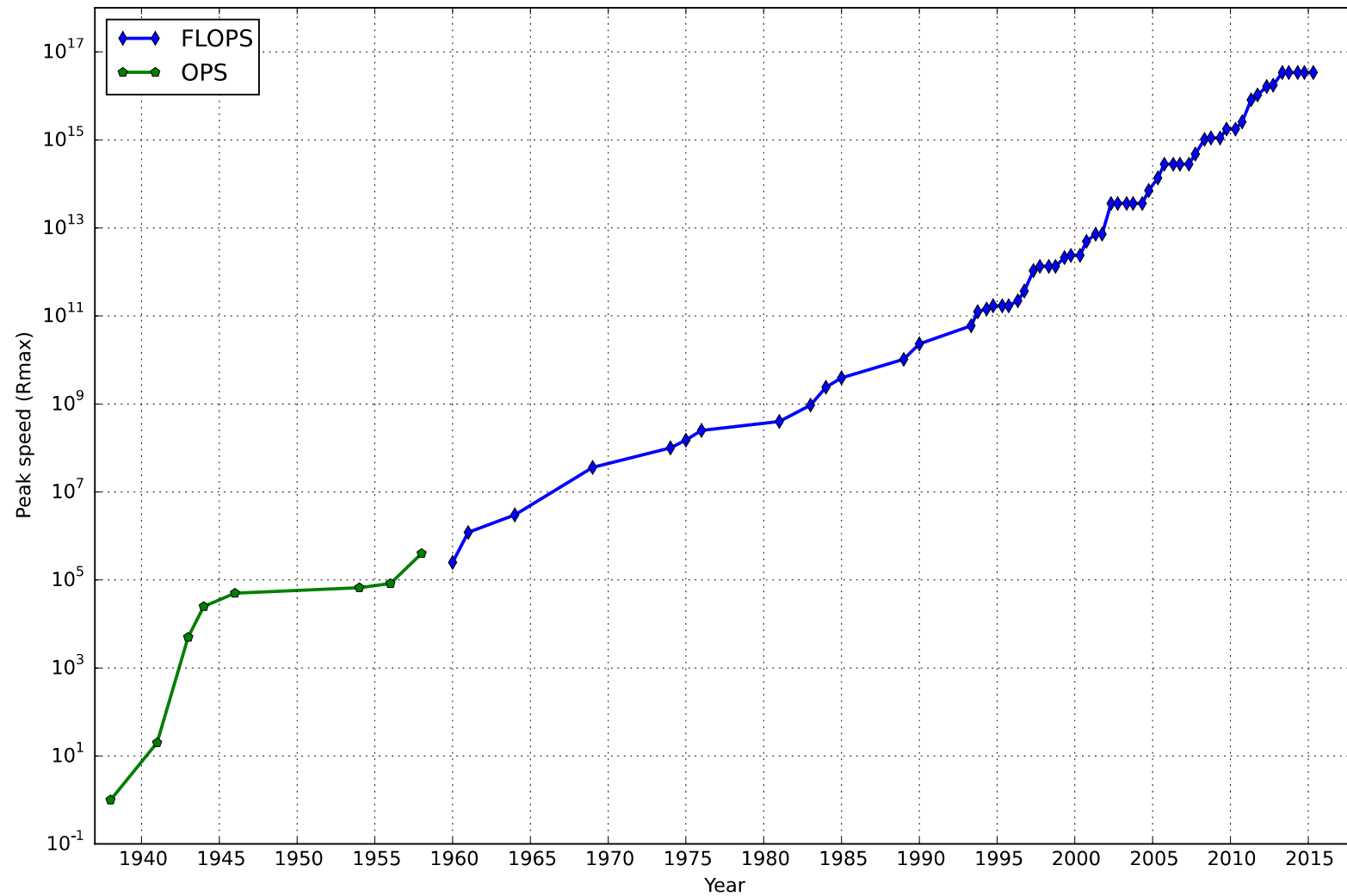
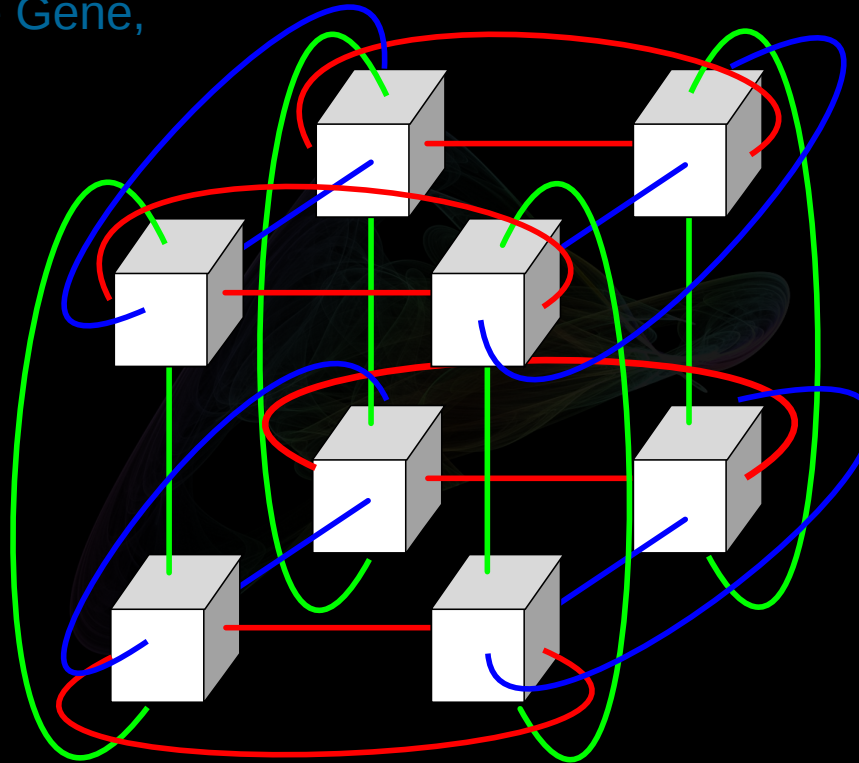


Diagram of a 3-dimensional torus interconnect used by systems such as Blue Gene, Cray XT3, etc.



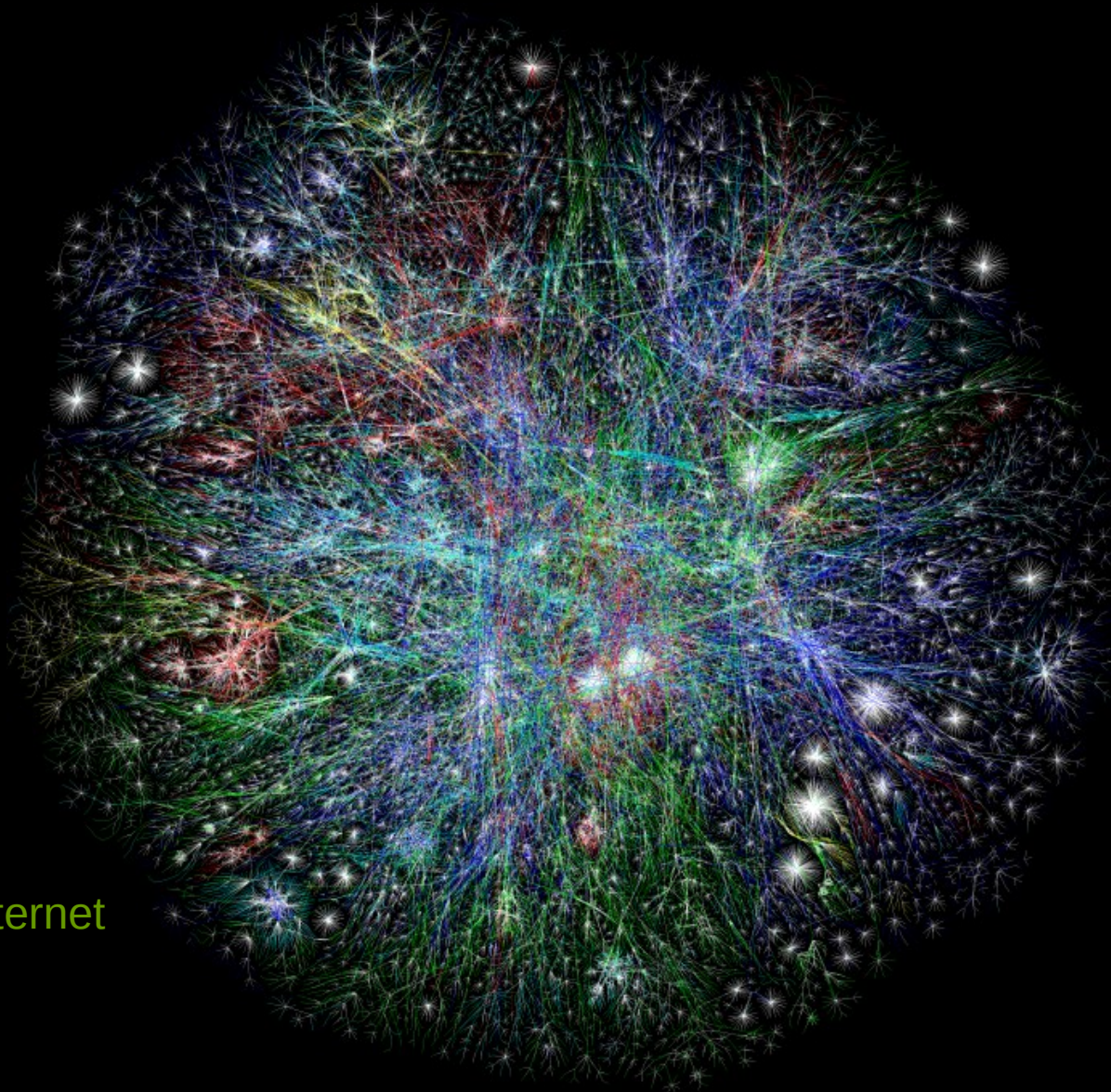
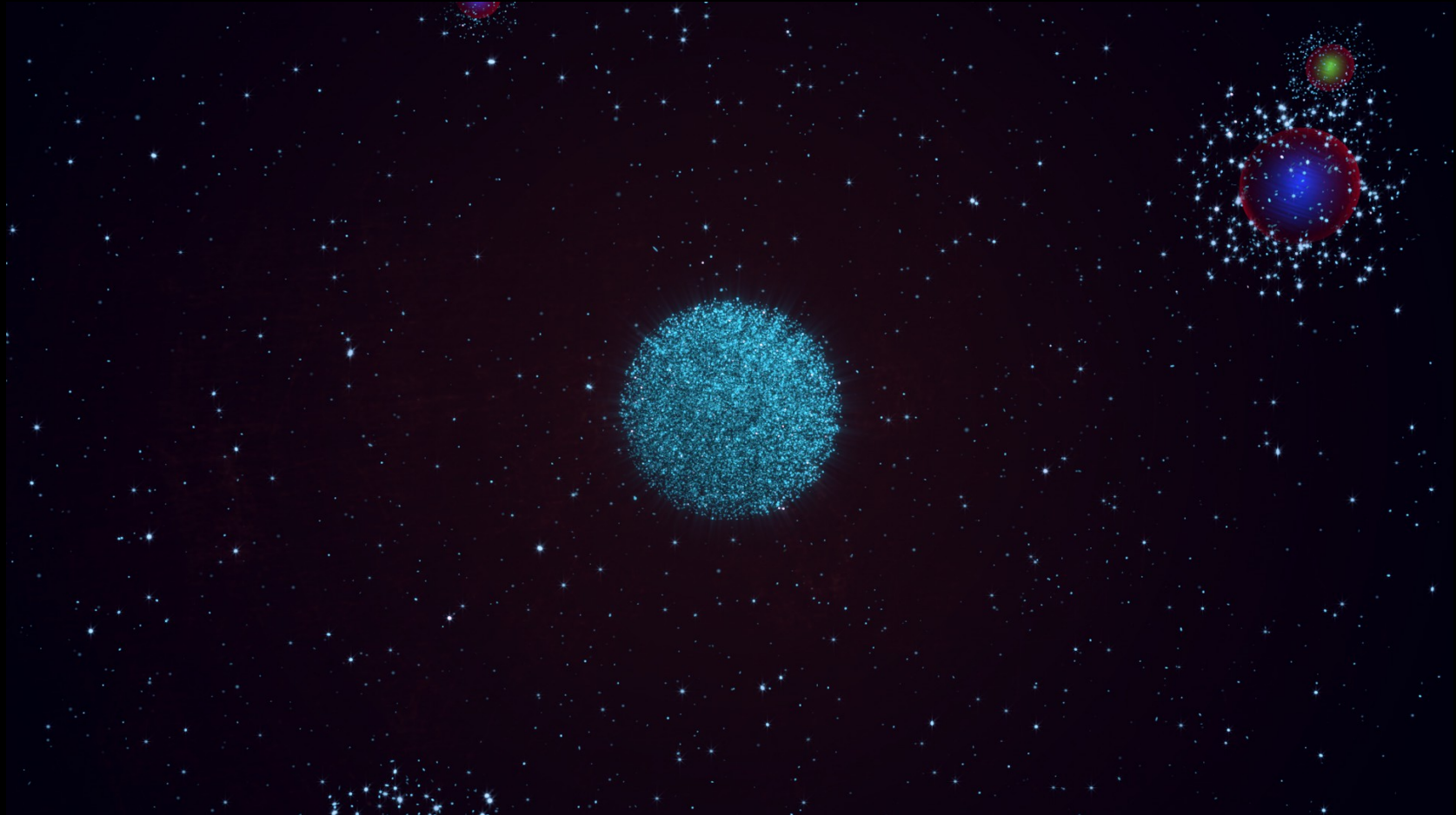


Diagram of internet

2015



Boson de Higgs

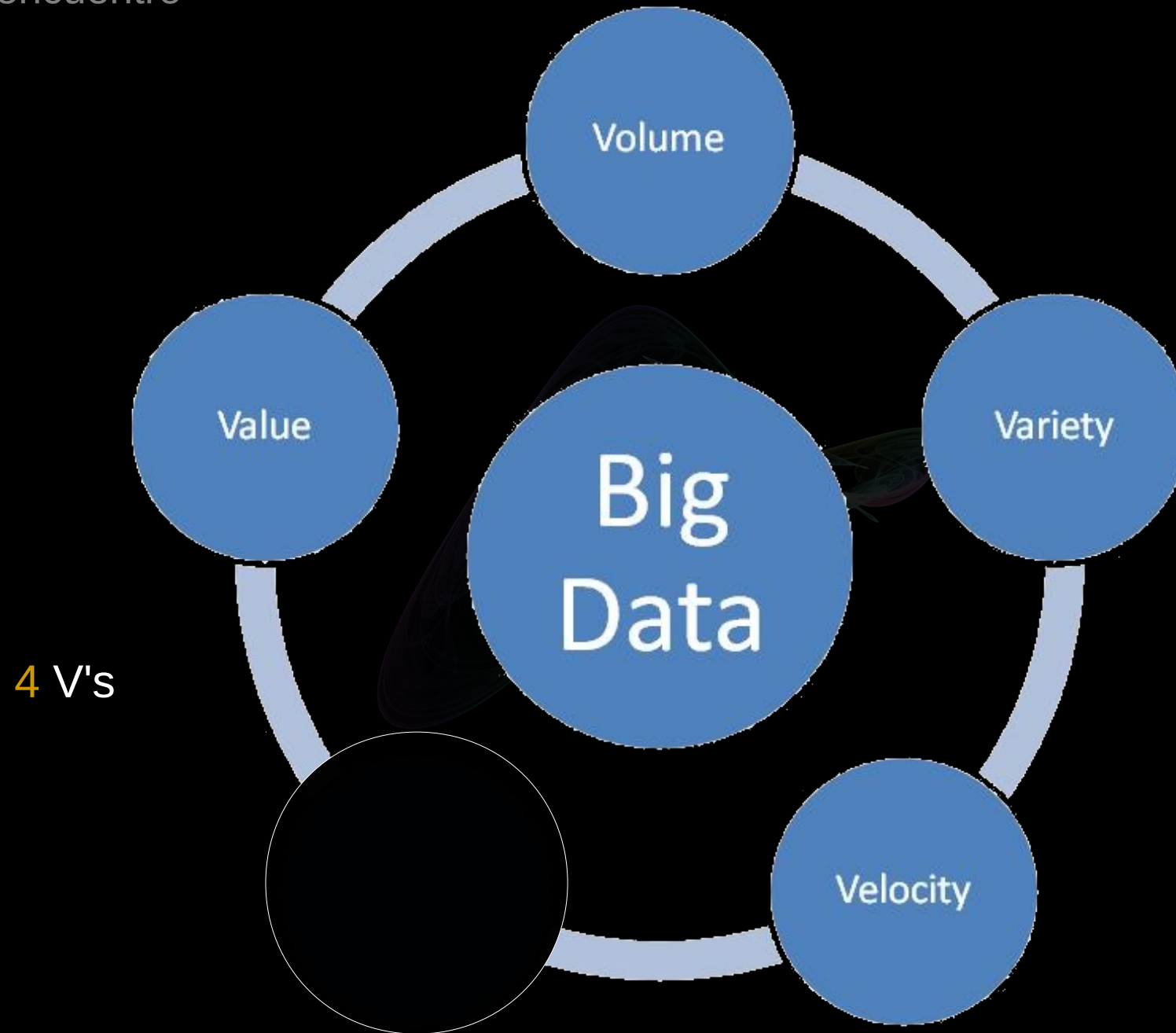
Big Data ?



Big data is a broad term
for data sets so **large** or **complex**
that **traditional** data processing
applications are **inadequate**

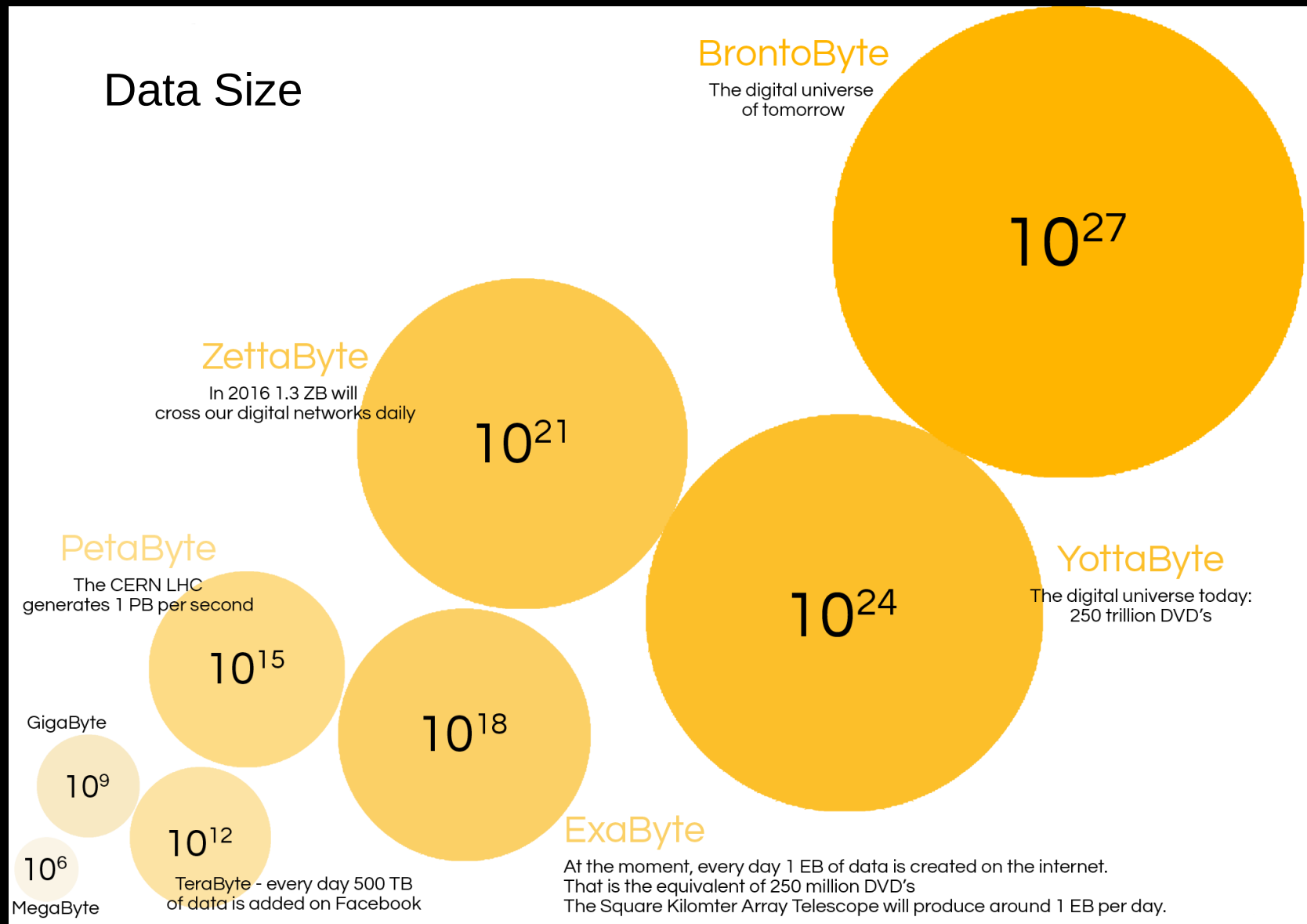
3 V's





5 V's







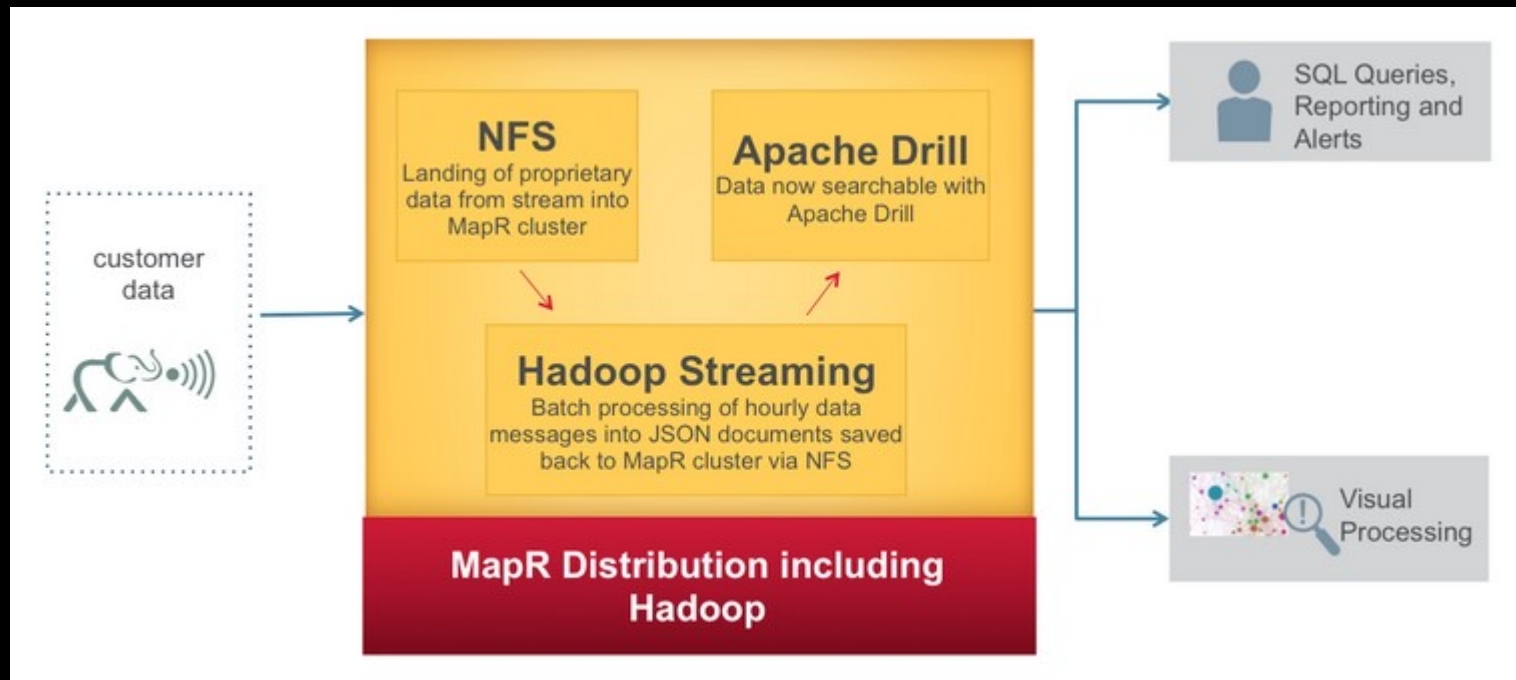
Hadoop ?



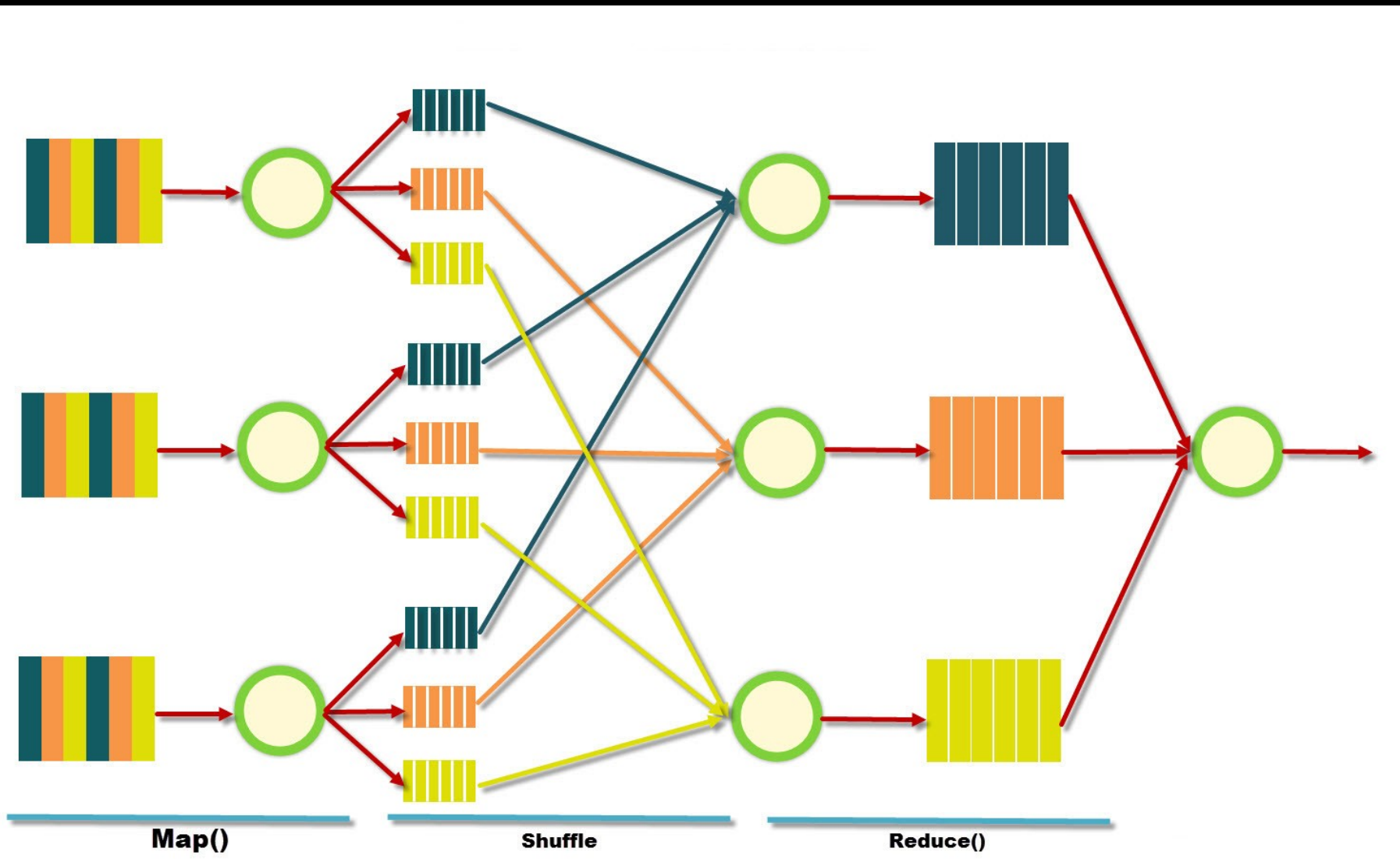
Apache **Hadoop** is an open-source
software framework
written in **Java** for **distributed storage**
and **distributed processing**
of **very large** data sets
on **computer clusters** built from **commodity hardware**



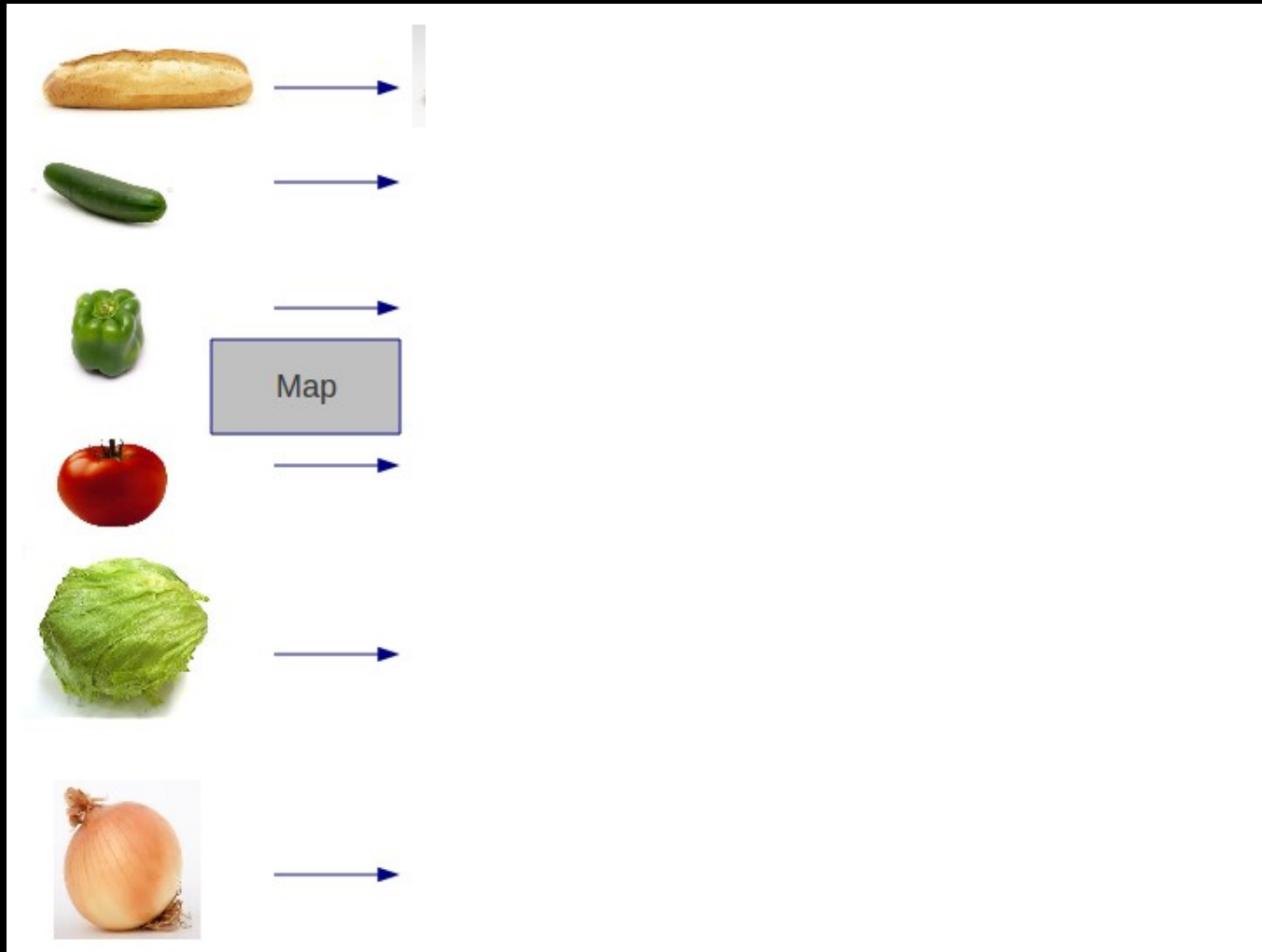
Map Reduce ?



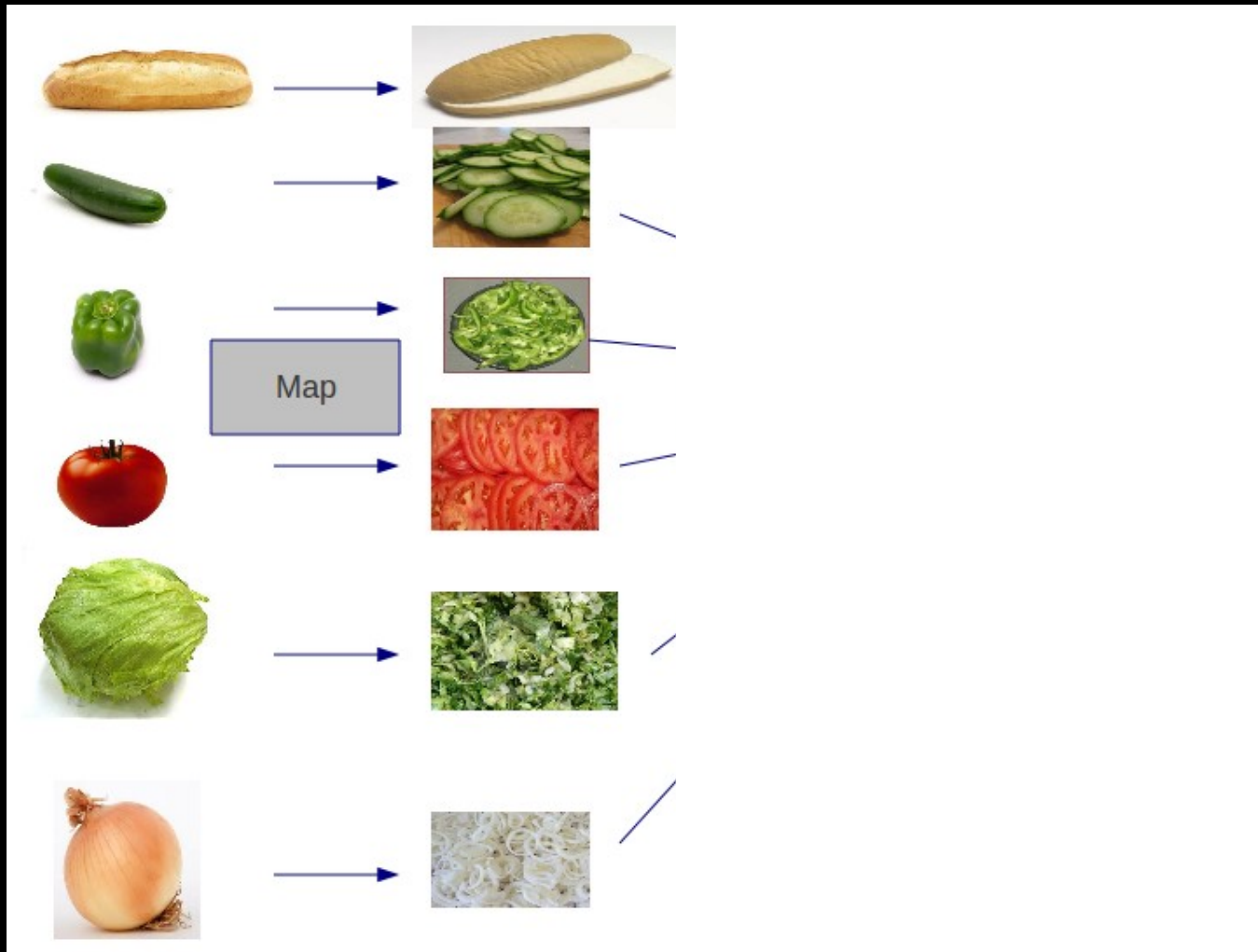
Map Reduce Schema



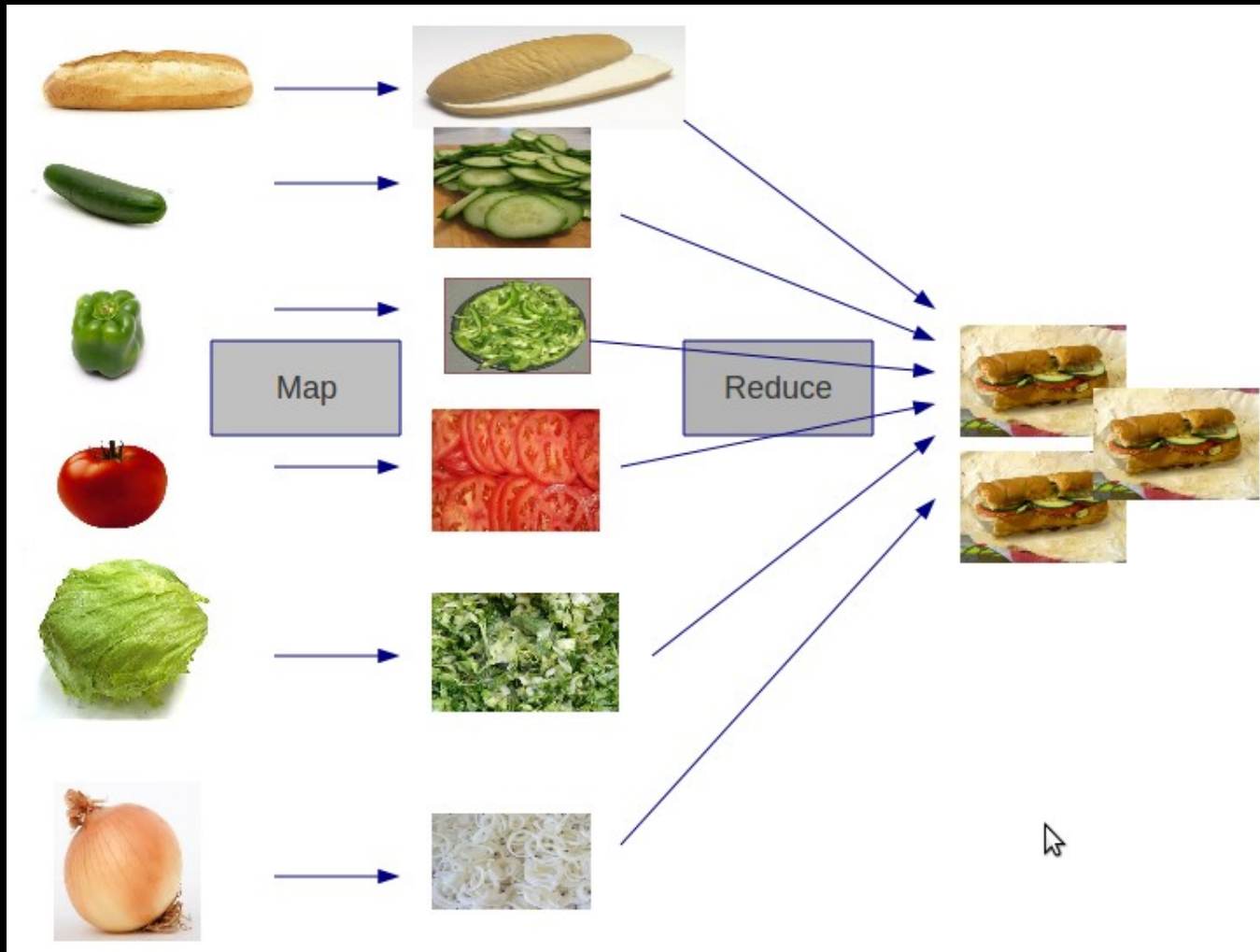
Map Reduce Schema



Map Reduce Schema



Map Reduce Schema



Map Reduce Schema



R Studio ?



library(rhdfs) → Hadoop

library (plyr) → distributed processing

library(rmr2) → Map Reduce



Happy Hacking