

commodity clusters  
allows scalable computing  
to achieve data-parallel  
scalability for big data  
→ cost effective

+

advances in distributed file  
systems to move computation  
to data

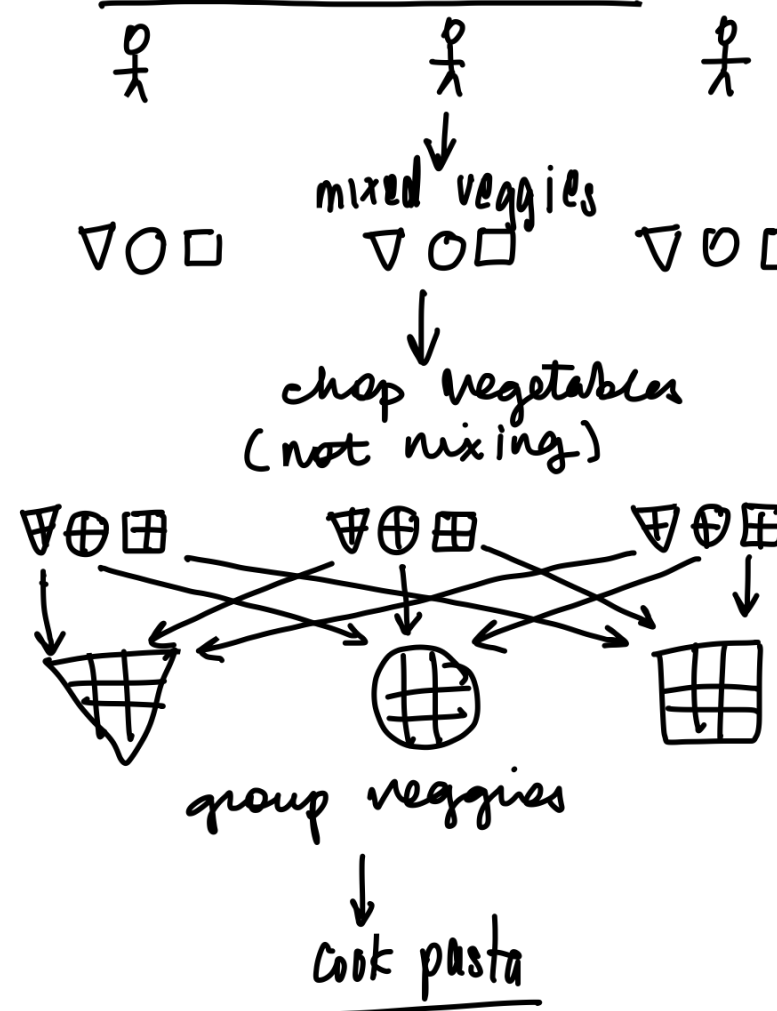
scalable  
big data  
analytics

how to take  
advantage of  
this  
advances

programming  
models

- example (of big data) modeling
- what: cook pasta with tomato sauce
  - time: in an hour (1hr)

USE PARALLELIZATION



it is

abstraction /  
existing machinery /  
infra structure  
set of abstract  
runtime libraries  
and programming languages

can be  
→ low level (machine language)  
→ high level (Java)

form a  
model of  
computation

(summary)  
programming model = abstractions  
runtime libraries + programming languages

FOR  
BIG DATA

programmability  
distributed file systems (DFS)

on top of  
computer programs that  
work efficiently  
cope with issues

requirements

1. support big  
data operations

split volumes of data  
access data fast  
distribute computations  
to nodes — (scheduling many  
parallel tasks at once)

2. handle fault  
tolerance

replicate data partitions  
recover files when  
needed

3. enable adding  
more racks

(new  
resources), — for < more or faster data  
called scaling  
out if needed not losing performance

4. optimized for  
specific data  
types

types  
optimized for  
at least  
one type  
document  
table  
key-value  
graph  
multimedia  
stream

model:  
Map Reduce

big data  
programming  
model

many  
implementations

Hadoop (framework)  
implementation

supports all  
requirements

processing of large data  
split computations into  
different parallel tasks  
make efficient use of  
large commodity clusters

abstracts out

details of parallelization  
fault tolerance  
data distribution  
monitoring  
load balancing