

Introduction to Big Data

Assignment Project Exam Help
with Spark

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro



BerkeleyX

This Lecture

The Big Data Problem

Hardware for Big Data **Assignment Project Exam Help**

Distributing Wo **<https://eduassistpro.github.io/>**

Handling Failures and Slow Ma **Add WeChat edu_assist_pro**

Map Reduce and Complex Jobs

Apache Spark

Some Traditional Analysis Tools

- Unix shell commands, Pandas, R

Assignment Project Exam Help

<https://eduassistpro.github.io/>

All r Add WeChat edu_assist_pro
single machine

The Big Data Problem

- Data growing faster than computation speeds
Assignment Project Exam Help
- Growing data sets
 - » Web, mobile, sci <https://eduassistpro.github.io/>
- Storage getting cheaper
Add WeChat edu_assist_pro
 - » Size doubling every 18 months
- But, stalling CPU speeds and storage bottlenecks



Big Data Examples

- Facebook's daily logs: 60 TB
 - 1,000 genome
 - Google web index
 - Cost of 1 TB of disk: ~\$35
 - Time to read 1 TB from disk: 3 hours
(100 MB/s)
- Assignment Project Exam Help
<https://eduassistpro.github.io/>
Add WeChat edu_assist_pro

The Big Data Problem

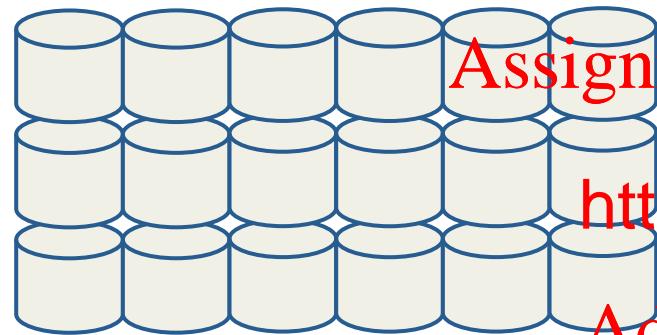
- A single machine can no longer process or even store all the data! **Assignment Project Exam Help**
- Only solution <https://eduassistpro.github.io/>
Add WeChat edu_assist_pro

Assignment Project Exam Help

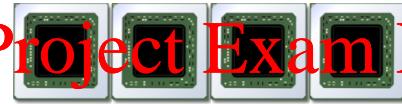
How <https://eduassistpro.github.io/> hing?

Add WeChat edu_assist_pro

Hardware for Big Data



Assignment Project Exam Help



<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Lots of hard drives ... a

Hardware for Big Data

One big box?

(1990's solution)

<https://eduassistpro.github.io/>

But, expensive

- » Low volume
- » All “premium” hardware

Add WeChat edu_assist_pro

And, still not big enough!

Image: Wikimedia Commons / User:Tonusamuel

Hardware for Big Data

Consumer-grade hardware

Not “gold plated”

Many desktop-like

Easy to add capacity

Cheaper per CPU/disk

<https://eduassistpro.github.io/>



Image: Steve Jurvetson/Flickr

Complexity in software

Problems with Cheap Hardware

Failures, Google's numbers:

1-5% hard drives/year
Assignment Project Exam Help

0.2% DIMMs/
<https://eduassistpro.github.io/>

Network speeds versus shared
Add WeChat edu_assist_pro

Much more latency

Network slower than storage

Uneven performance

What's Hard About Cluster Computing?

- How do we split work across machines?

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

How do you count the number of occurrences of each word in a document?

Assignment Project Exam Help

“I am Sam
I am Sam
Sam I am
Do you like
Green eggs and ham?”

<https://eduassistpro.github.io/>

Add  WeChat edu_assist_pro

you: I
like: I

...

One Approach: Use a Hash Table

“I am Sam
Assignment Project
Exam Help
I am Sa
Sam I a
Do you like
<https://eduassistpro.github.io/>
Add WeChat edu_assist_pro
Green eggs and ham?”

One Approach: Use a Hash Table

“I am Sam Assignment Project Exam Help { | : | }

I am Sa <https://eduassistpro.github.io/>

Sam I a Add WeChat edu_assist_pro

Do you like

Green eggs and ham?”

One Approach: Use a Hash Table

“I am Sam Assignment Project 1, Exam Help
I am Sa <https://eduassistpro.github.io/>
Sam I a Add WeChat edu_assist_pro
Do you like
Green eggs and ham?”

One Approach: Use a Hash Table

“I am **Sam** Assignment Project Exam Help
I am Sa <https://eduassistpro.github.io/>
Sam I a Add WeChat `edu_assist_pro`
Do you like Sa
Green eggs and ham?”

One Approach: Use a Hash Table

“I am Sam Assignment Project [1.2, Exam Help
 I am Sa <https://eduassistpro.github.io/>
Sam I a Add WeChat edu_assist_pro
Do you like Sa

Green eggs and ham?”

What if the Document is Really Big?

"I am Sam
I am Sam
Sam I am

Do you like
Green eggs and ham?
I do not like them

Sam I am
I do not like
Green eggs and ham
Would you like them
Here or there?
..."

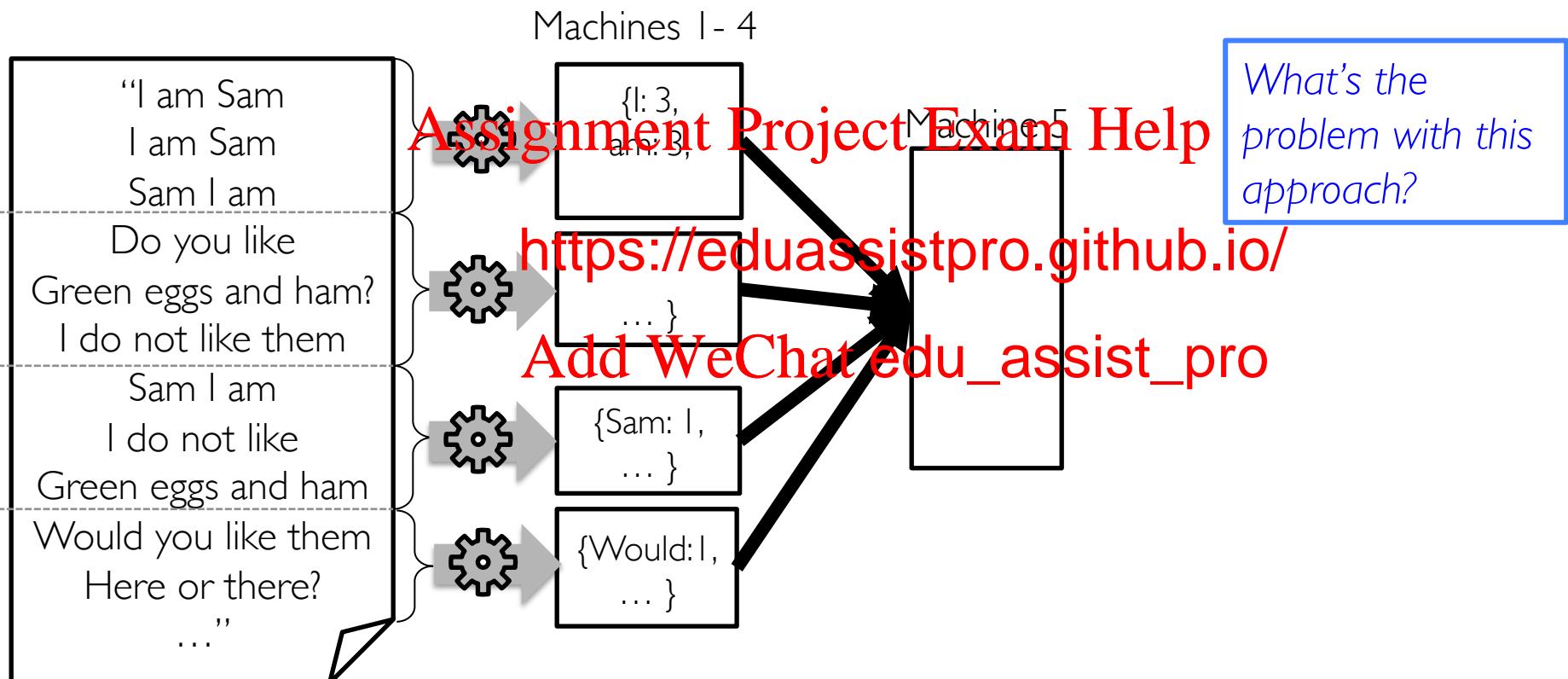


Assignment Project Exam Help

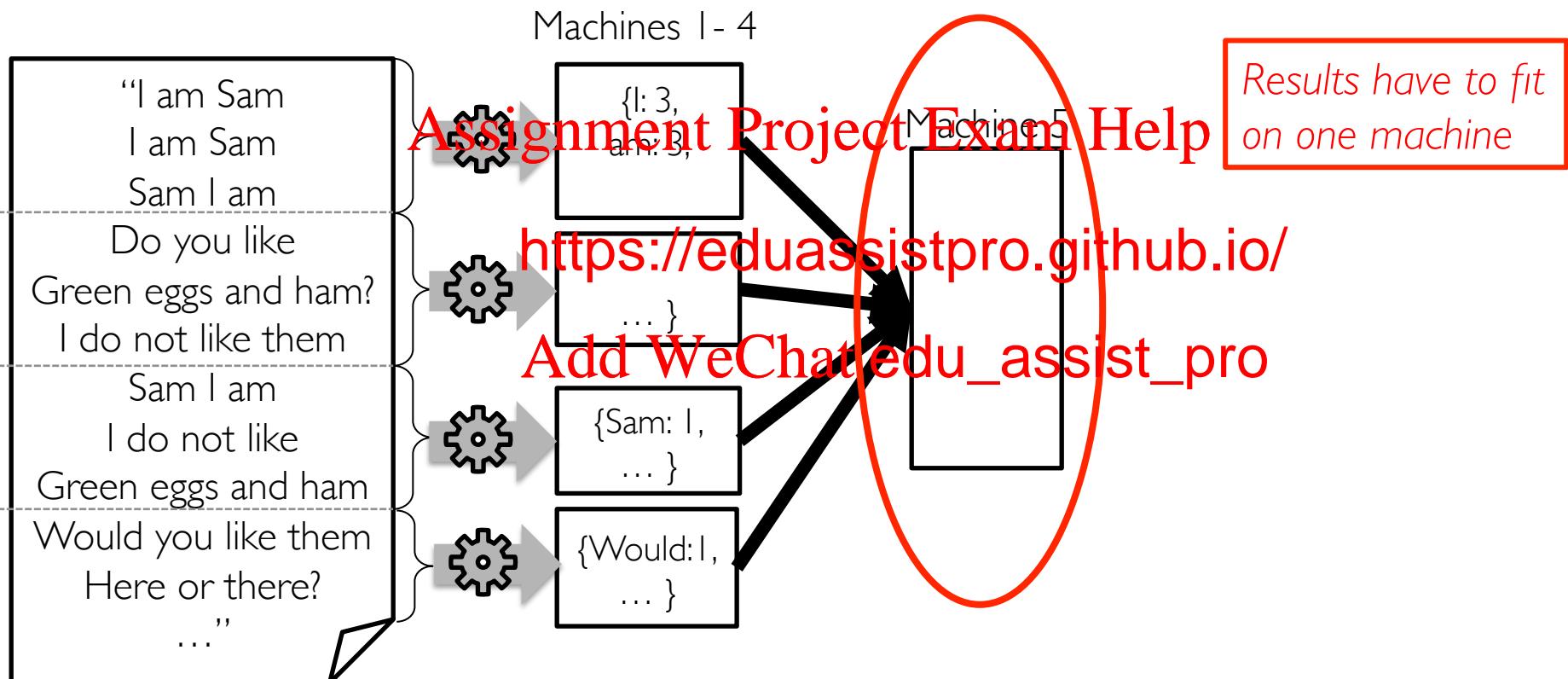
<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

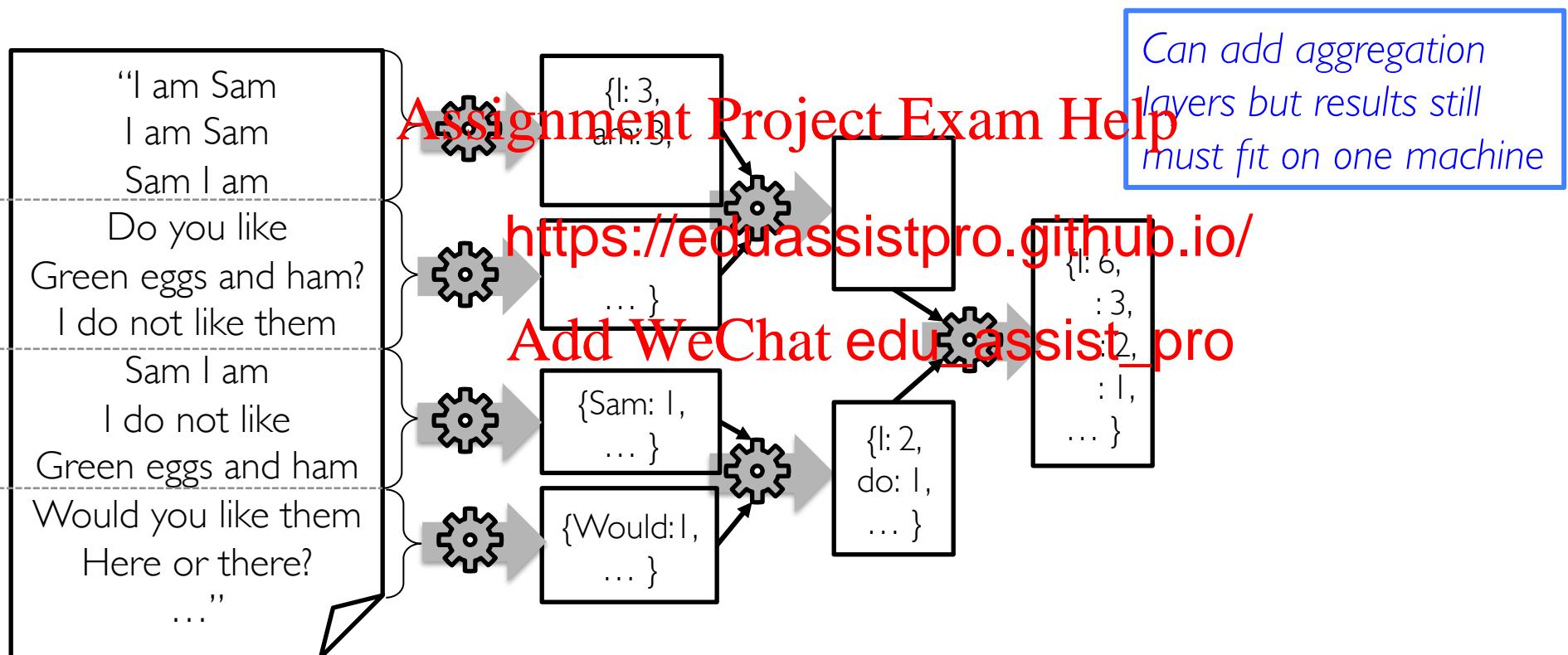
What if the Document is Really Big?



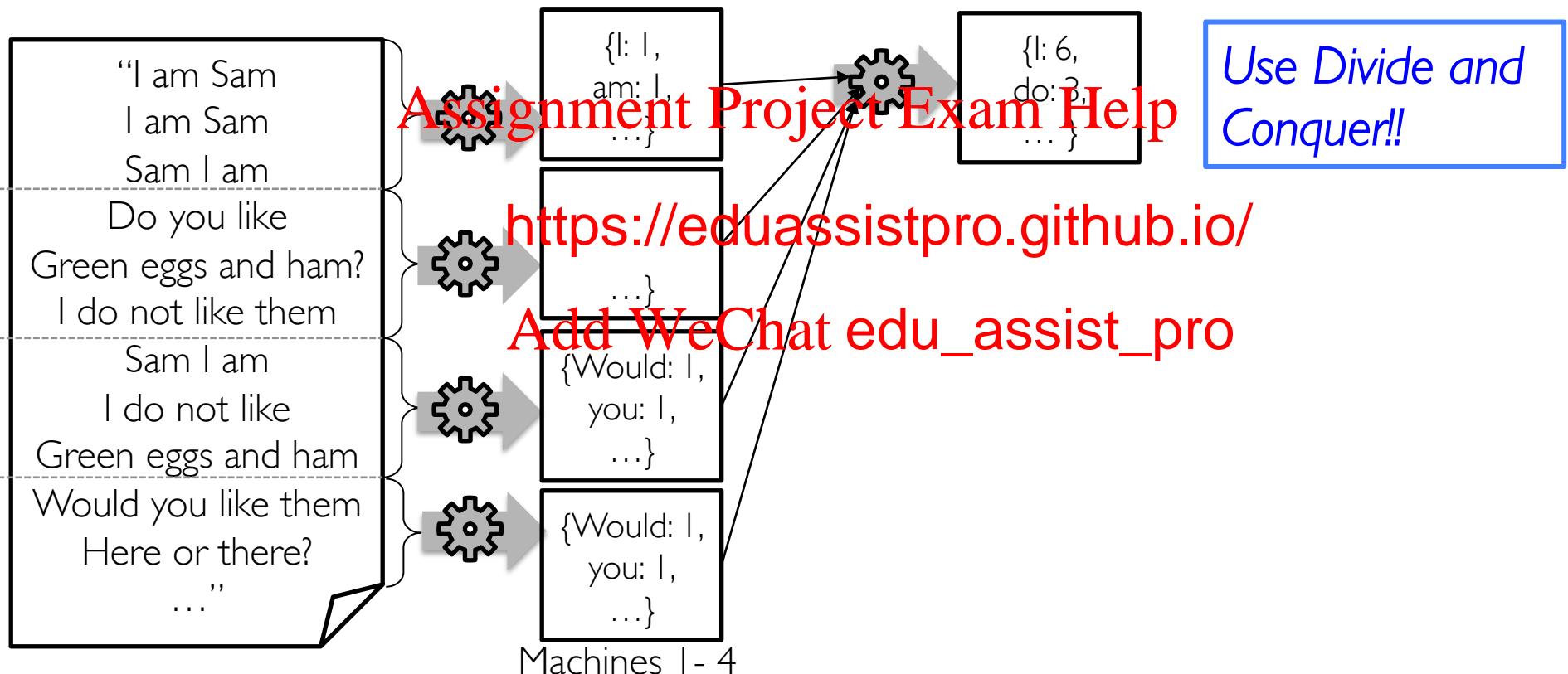
What if the Document is Really Big?



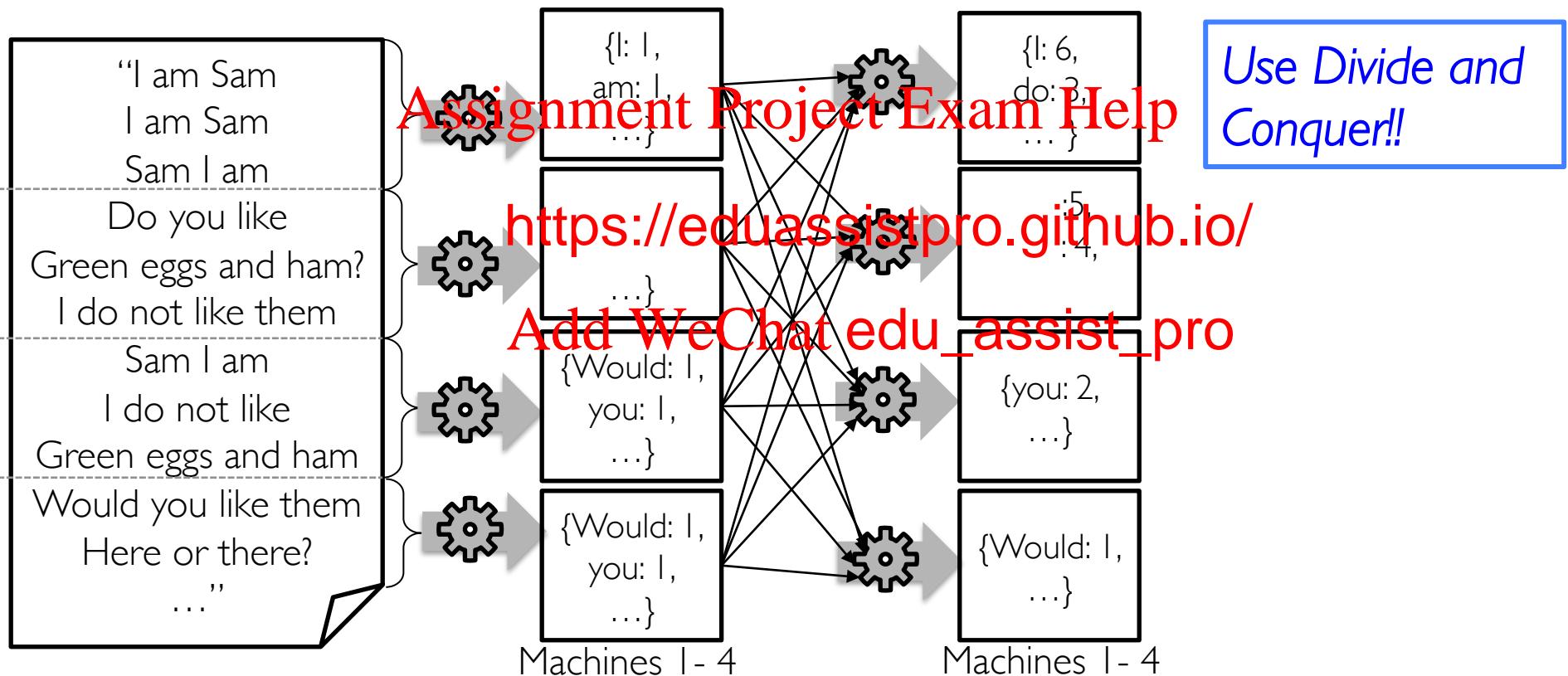
What if the Document is Really Big?



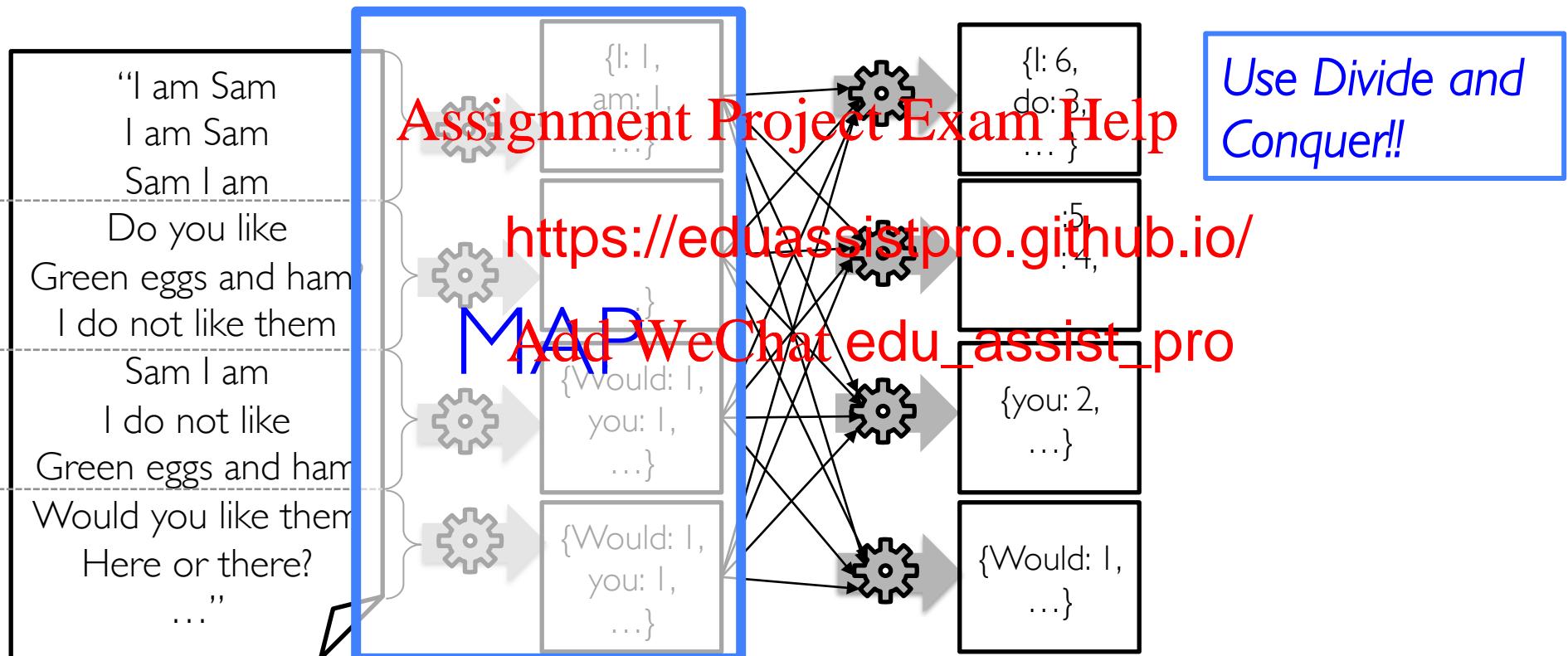
What if the Document is Really Big?



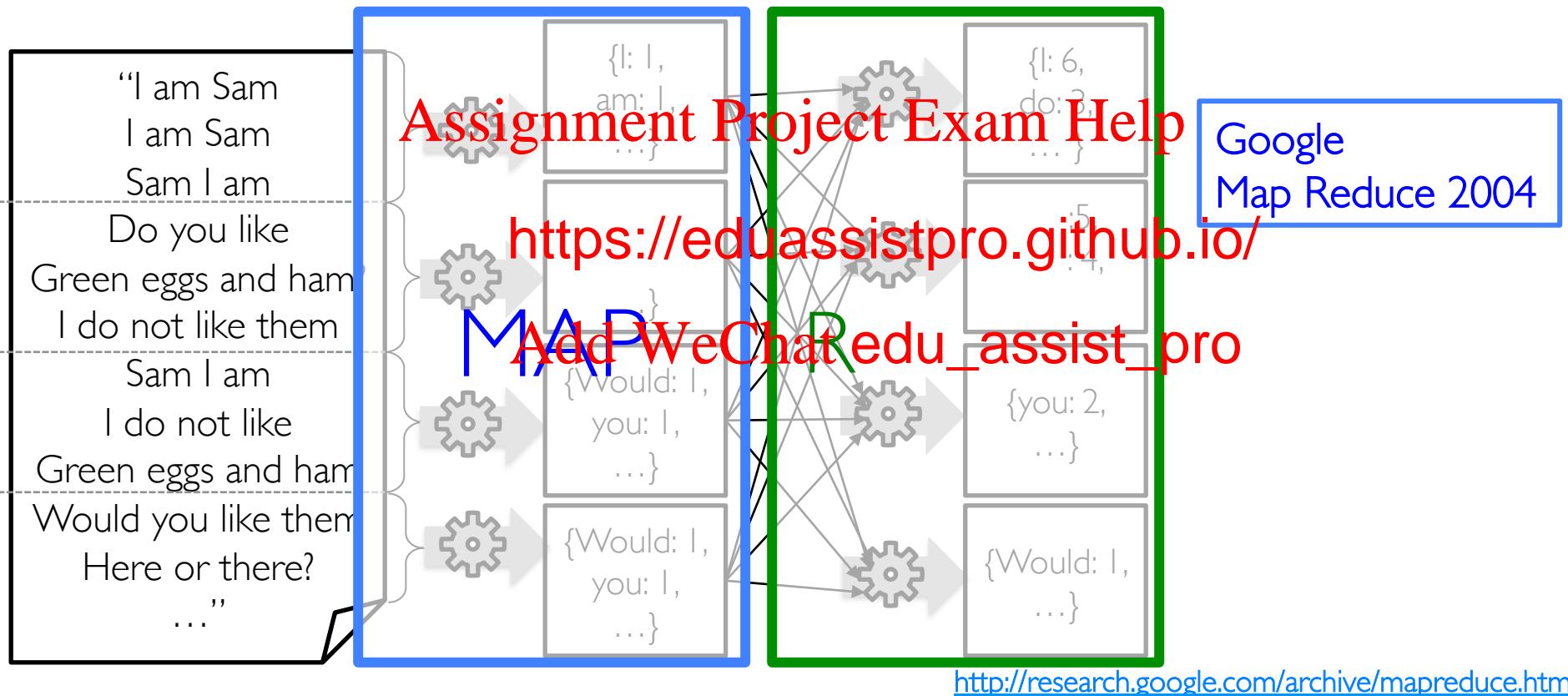
What if the Document is Really Big?



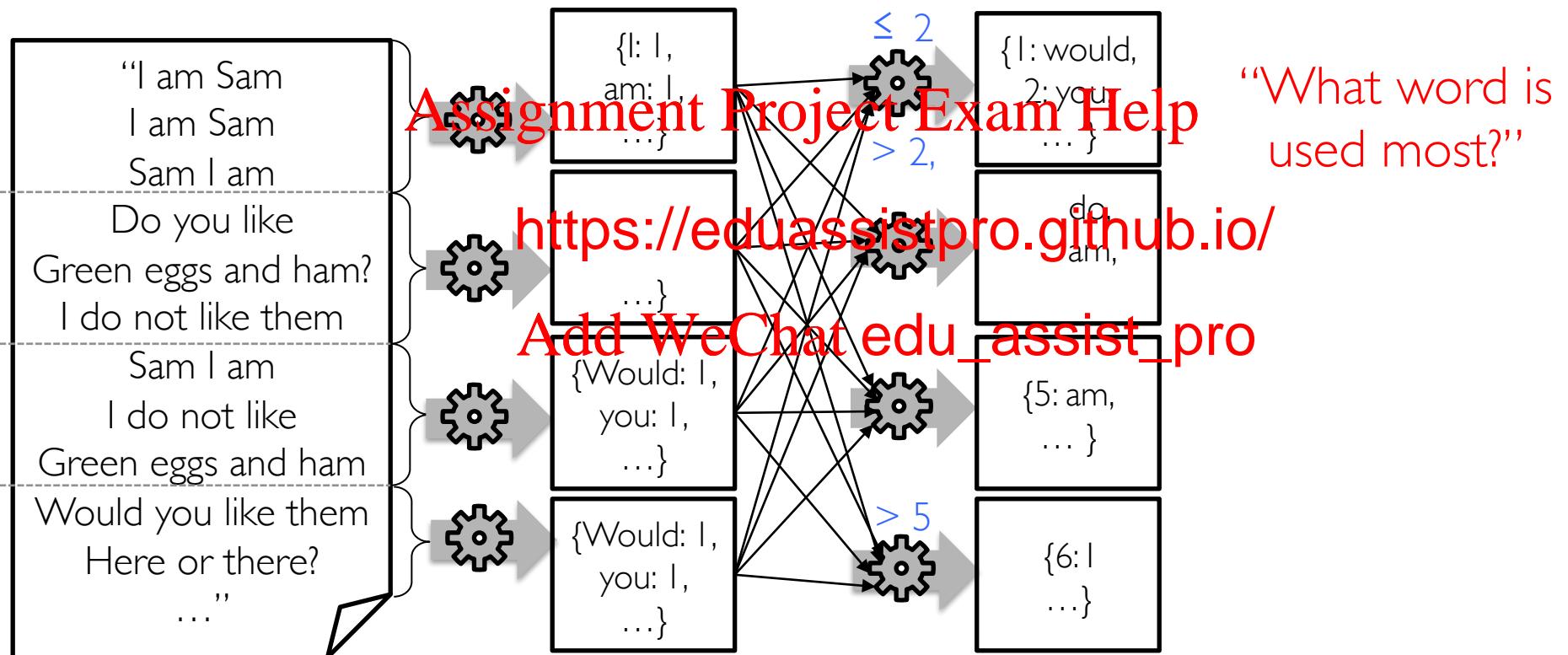
What if the Document is Really Big?



What if the Document is Really Big?



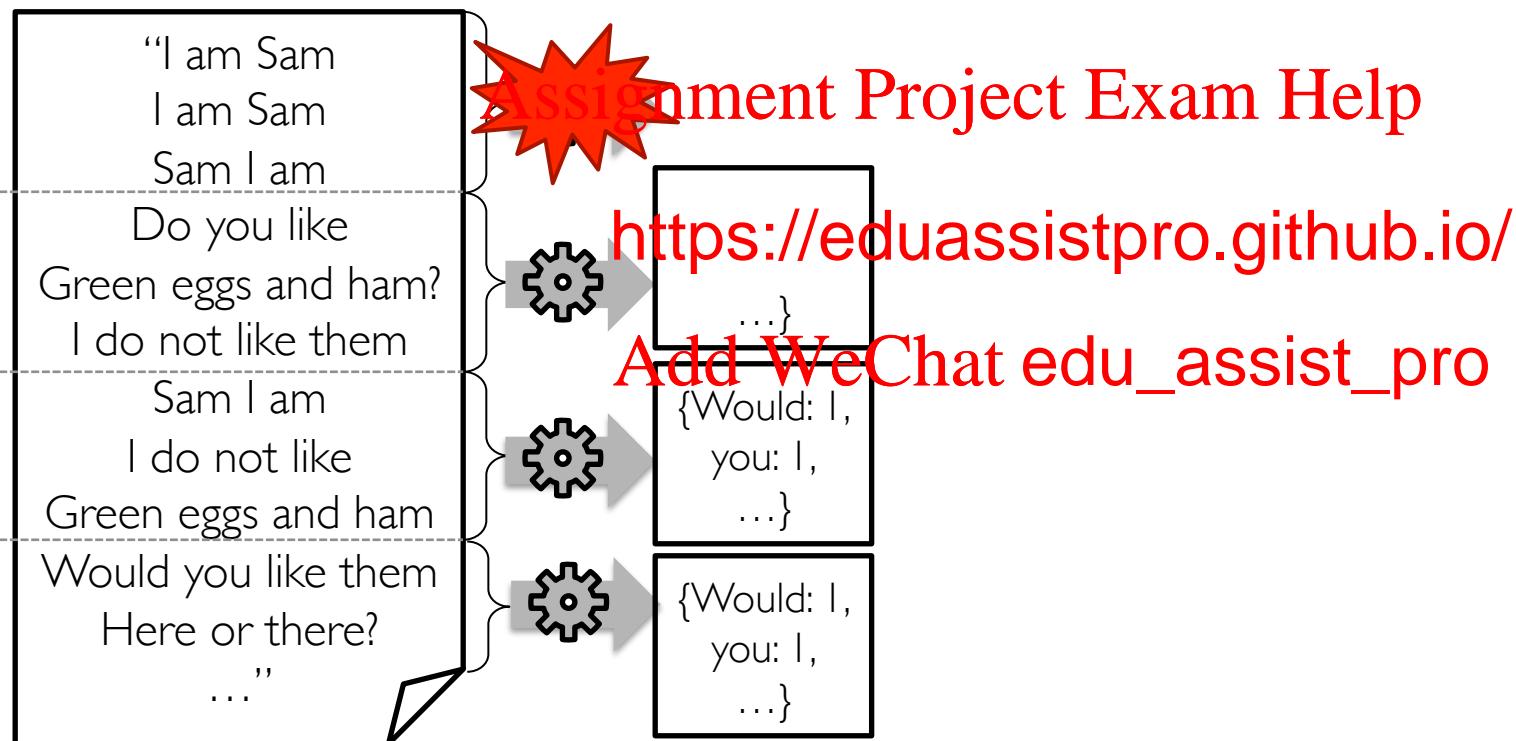
Map Reduce for Sorting



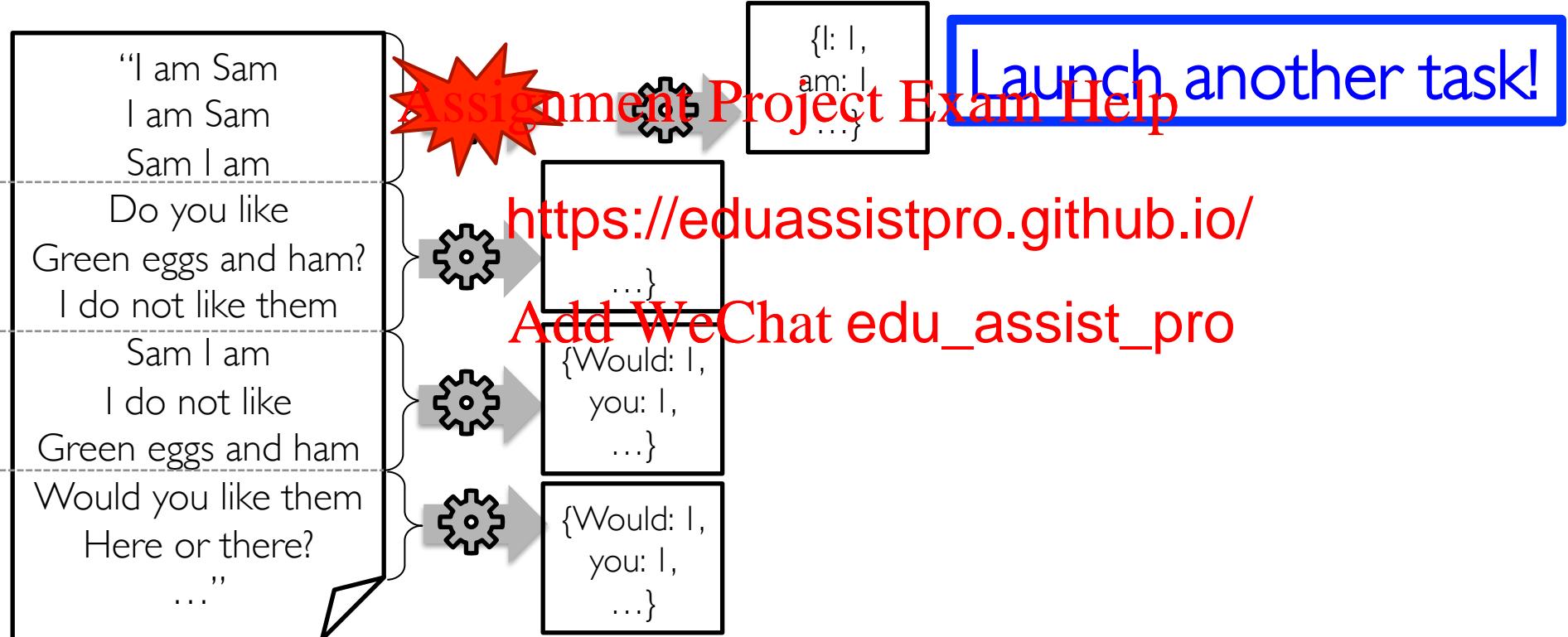
What's Hard About Cluster Computing?

- How to divide work across machines?
 - » Must consider network, data locality
 - » Moving data may be very expensive
- How to deal with <https://eduassistpro.github.io/>
 - » A server fails every 3 years → with 10 faults/day
 - » Even worse: stragglers (not failed, b
es)

How Do We Deal with Failures?



How Do We Deal with Machine Failures?

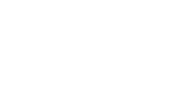


How Do We Deal with Slow Tasks?

"I am Sam
I am Sam
Sam I am

Do you like
Green eggs and ham?
I do not like them

Sam I am
I do not like
Green eggs and ham
Would you like them
Here or there?
..."

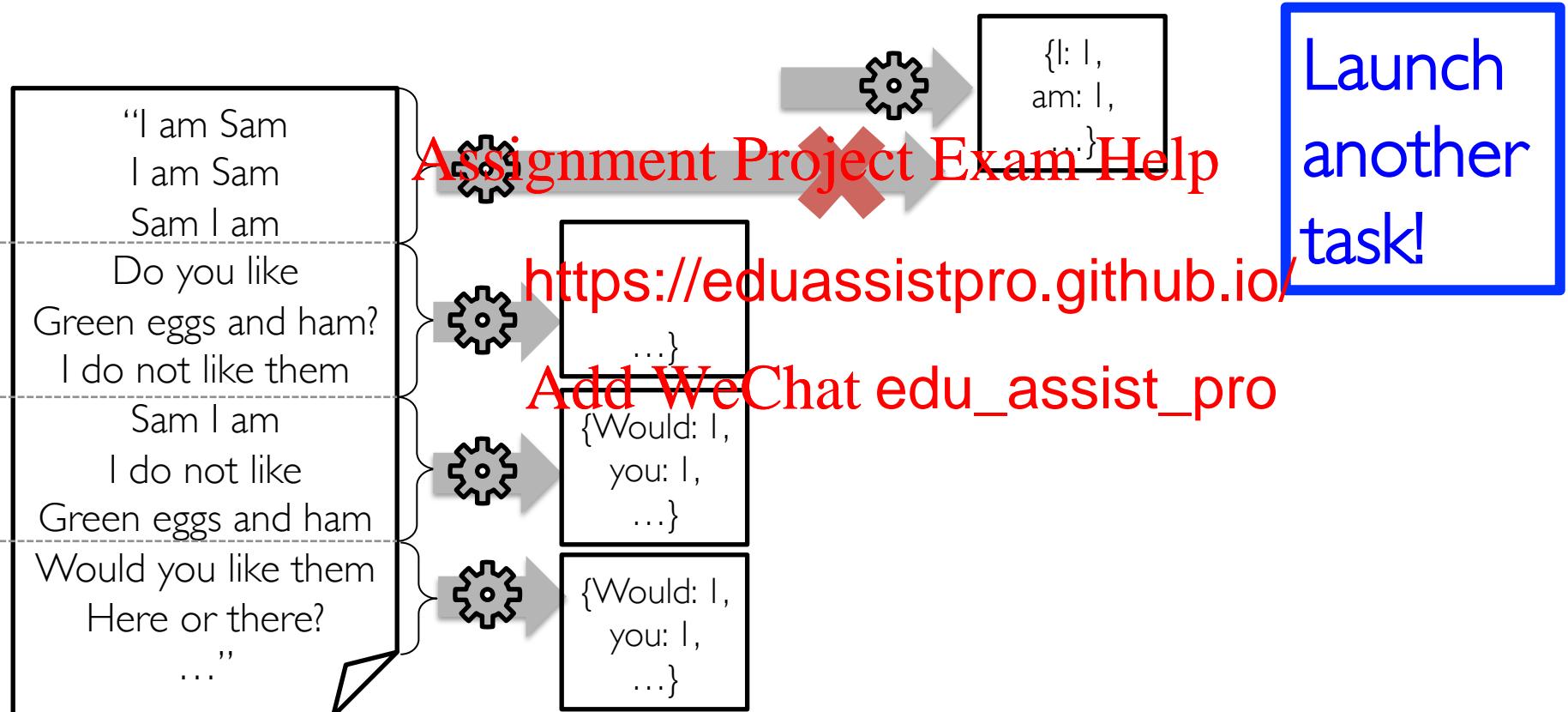


Assignment Project Exam Help

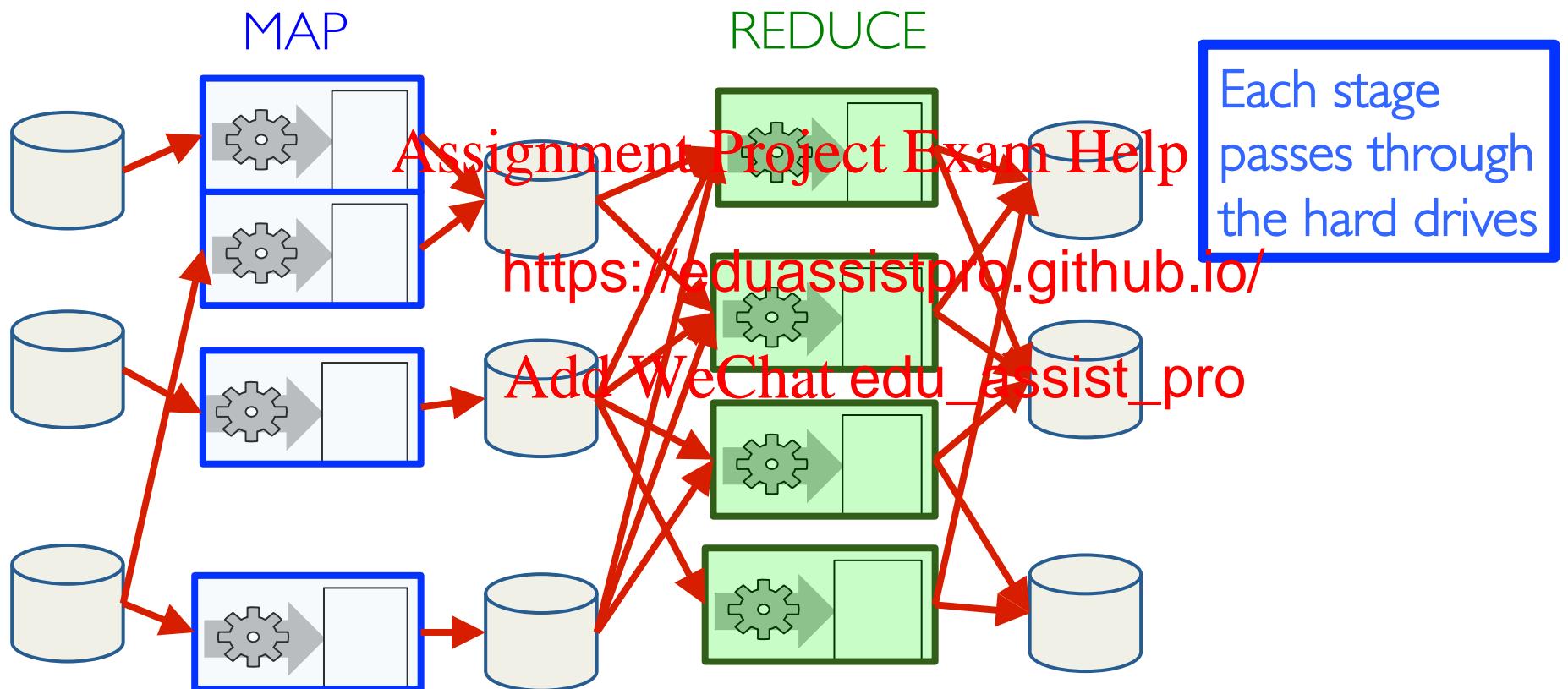
<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

How Do We Deal with Slow Tasks?

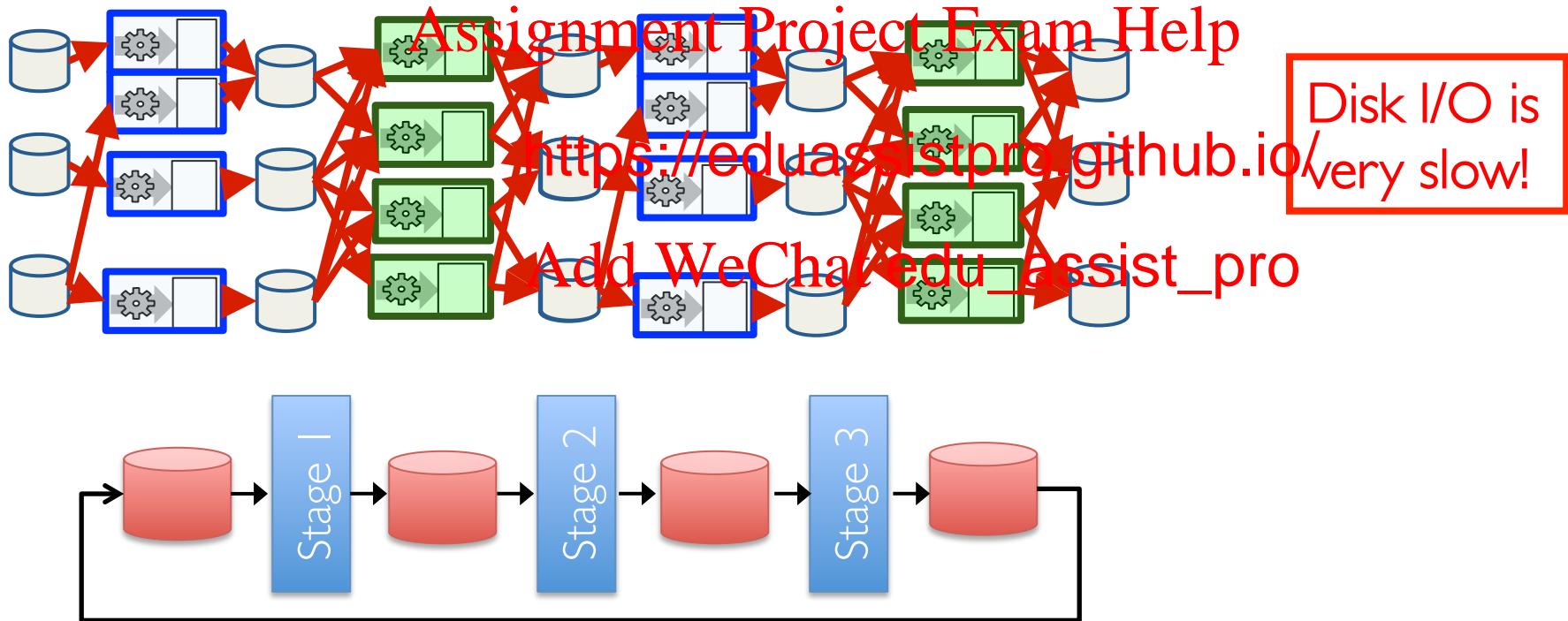


Map Reduce: Distributed Execution



Map Reduce: Iterative Jobs

- Iterative jobs involve a lot of disk I/O for each repetition



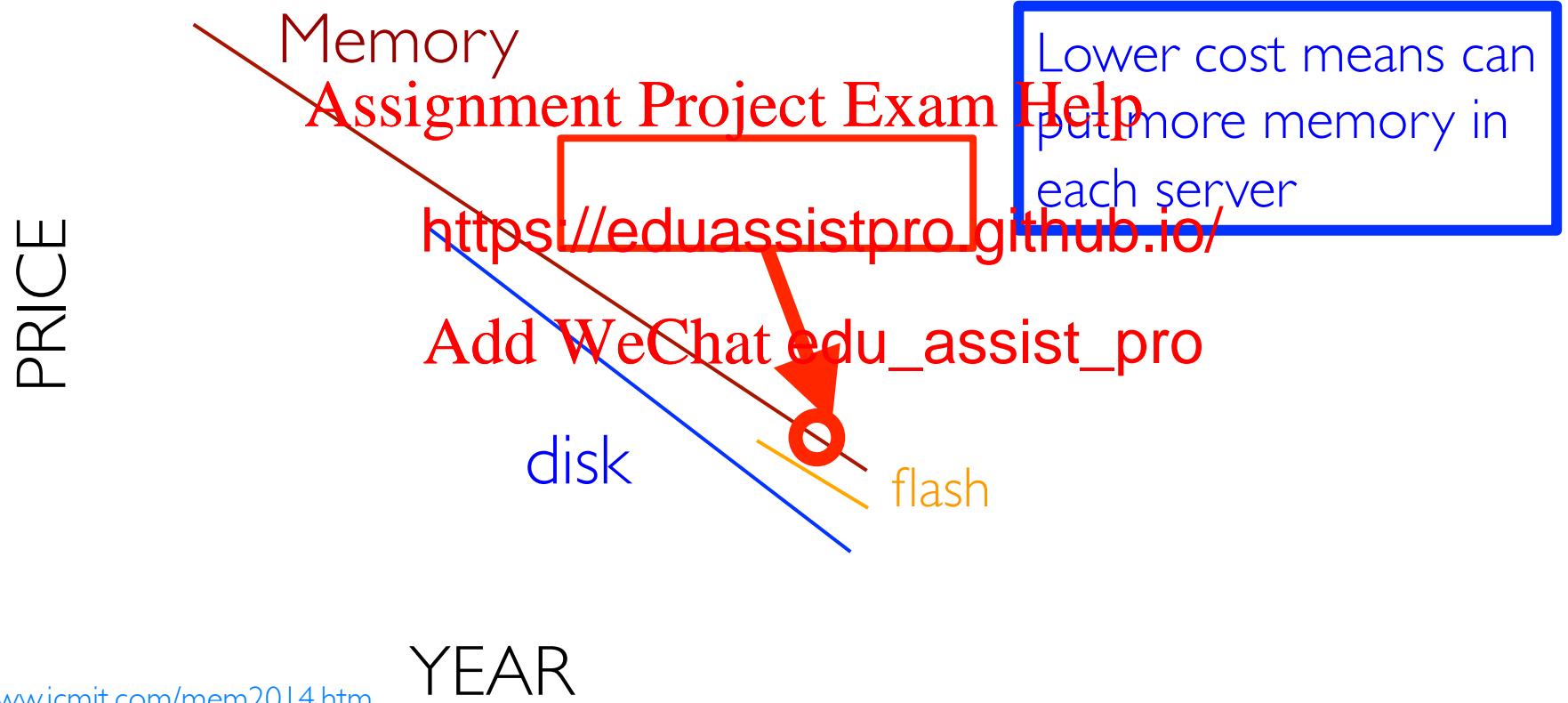
Apache Spark Motivation

- Using Map Reduce for complex jobs, interactive queries and online processing provides ~~Assignment Project~~ ~~Fast~~ ~~Disk~~ ~~I/O~~

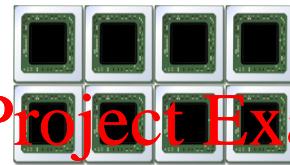
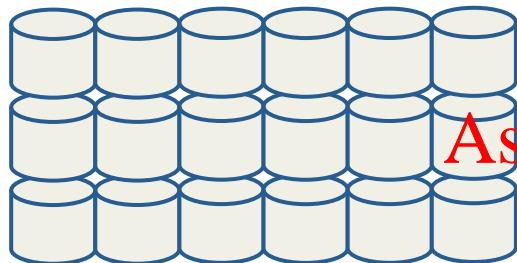


Disk I/O is very slow

Tech Trend: Cost of Memory



Hardware for Big Data



Lots of hard driv

... and memory!



Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

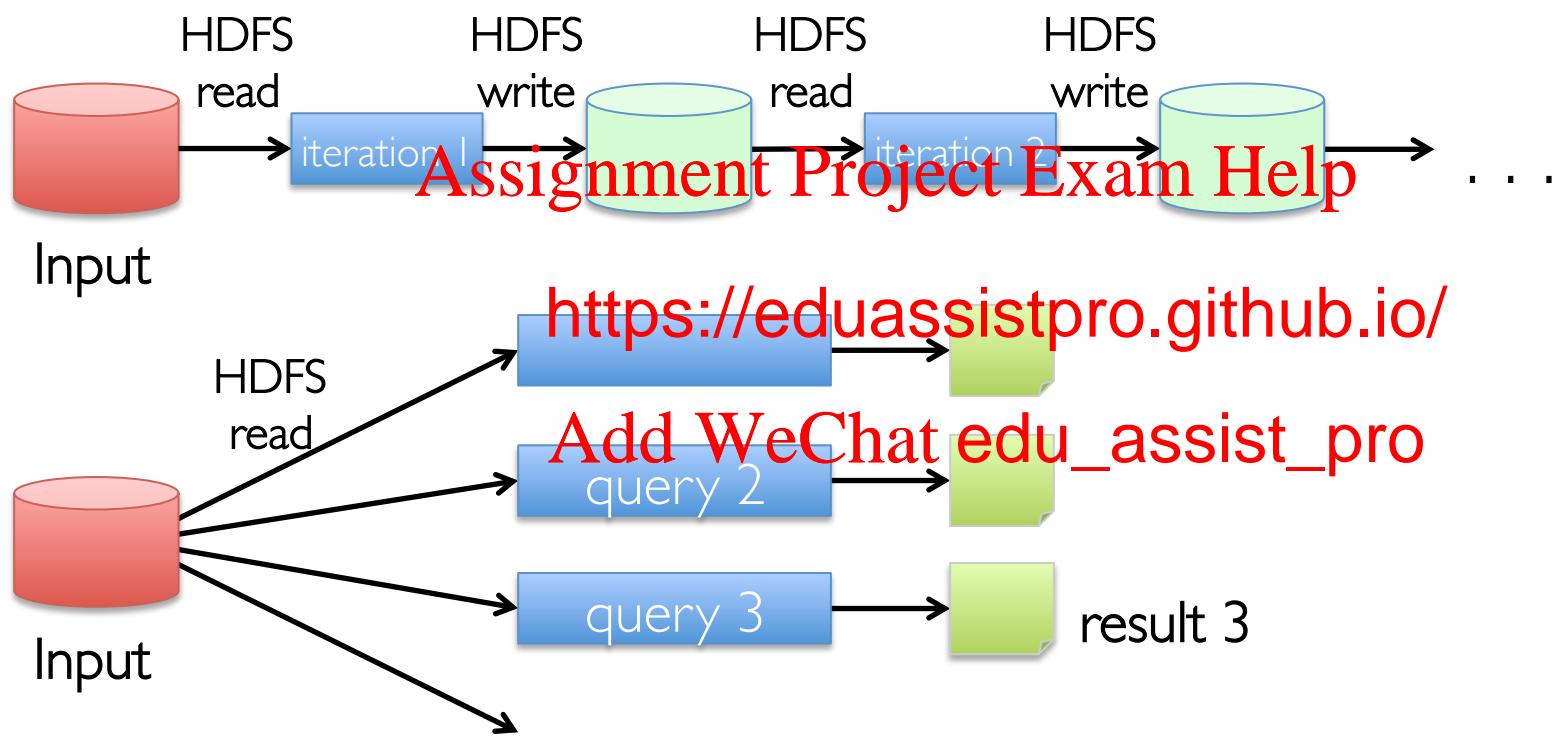
Opportunity

- Keep more data *in-memory*
Assignment Project Exam Help
<https://eduassistpro.github.io/>
- Create new dis
ine:
Add WeChat edu_assist_pro

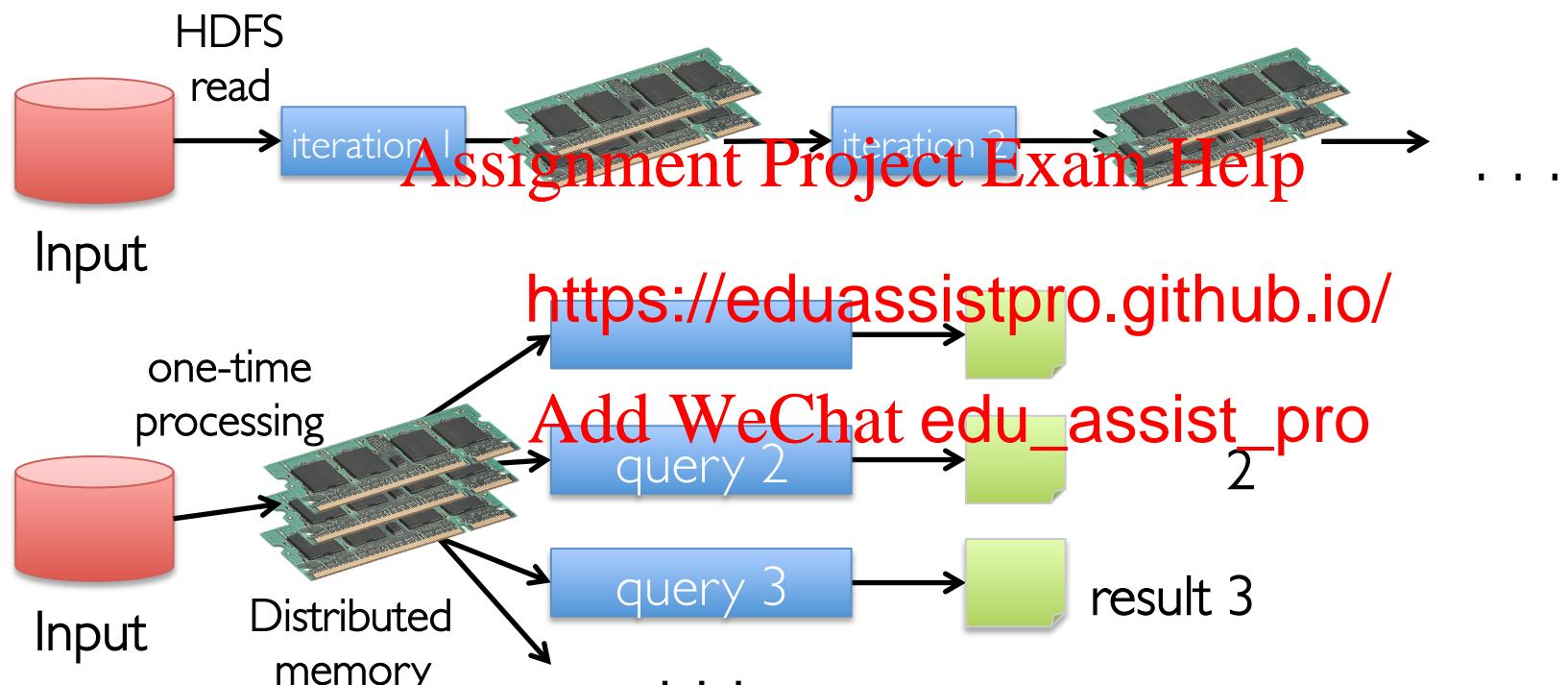


http://people.csail.mit.edu/matei/papers/2010/hotcloud_spark.pdf

Use Memory Instead of Disk



In-Memory Data Sharing



10-100x faster than network and disk

Resilient Distributed Datasets (RDDs)

- Write programs in terms of operations on distributed datasets
- Partitioned collections of objects spread across a cluster, stored in memory or on disk
 - Assignment Project Exam Help
<https://eduassistpro.github.io/>
- RDDs built and manipulated through a set of parallel transformations (map) and actions (count, collect, save)
- RDDs automatically rebuilt on machine failure

The Spark Computing Framework

- Provides programming abstraction and parallel runtime to hide complexities of fault-tolerance and slow machines
<https://eduassistpro.github.io/>
- “Here’s an operation, run it on the data”
 - » I don’t care where it runs (you say)
 - » In fact, feel free to run it twice on different nodes

Spark Tools

Spark
SQL

Assignment Project Exam Help
Spar
Stream
<https://eduassistpro.github.io/>
Add WeChat edu_assist_pro

Apache Spark

Spark and Map Reduce Differences

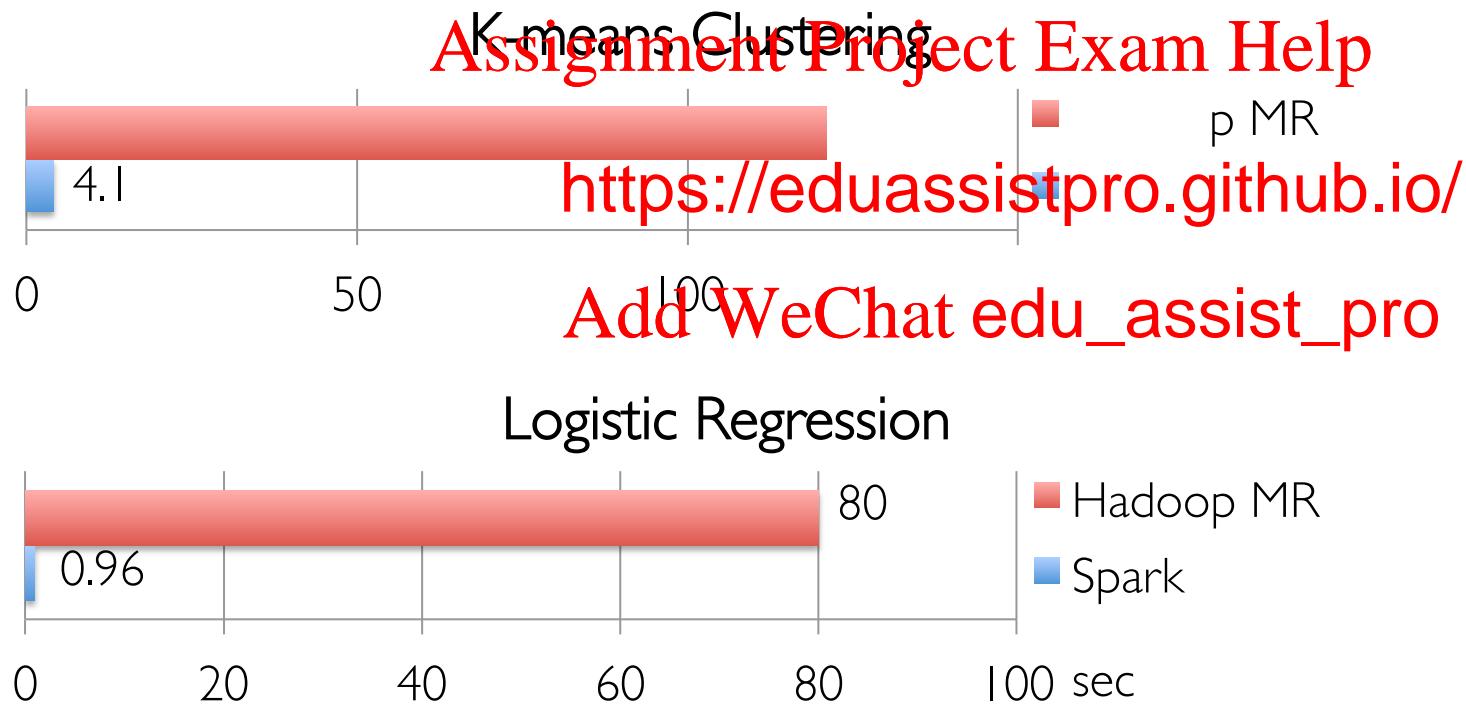
	Hadoop Map Reduce	Spark
	Assignment Project Exam Help	
Storage	Disk only	https://eduassistpro.github.io/
Operations	Map and Reduce	Add WeChat edu_assist_pro
Execution model	Batch	Batch, interactive, streaming
Programming environments	Java	Scala, Java, R, and Python

Other Spark and Map Reduce Differences

- Generalized patterns
 ⇒ unified design for many use cases
- Lazy evaluation <https://eduassistpro.github.io/>
 ⇒ reduces wait states, better
- Lower overhead for starting jobs
- Less expensive shuffles

In-Memory Can Make a Big Difference

- Two iterative Machine Learning algorithms:



First Public Cloud Petabyte Sort

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Daytona Gray 100 TB
sort benchmark record
(tied for 1st place)

<http://databricks.com/blog/2014/11/05/spark-officially-sets-a-new-record-in-large-scale-sorting.html>

Spark Expertise Tops Big Data Median Salaries

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Over 800 respondents across 53 countries and 41 U.S. states

<http://www.oreilly.com/data/free/2014-data-science-salary-survey.csp>

History Review

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Historical References

- circa 1979 – Stanford, MIT, CMU, etc.: set/list operations in LISP, Prolog, etc., for parallel processing
<http://www-formal.stanford.edu/jmc/history/lisp/lisp.htm>
- circa 2004 – Google. MapReduce: Simplified Data Processing on Large Clusters
Jeffrey Dean and Sanjay G
<http://research.google.com/>
- circa 2006 – Apache Hadoop
Doug Cutting
<http://research.yahoo.com/files/cutting.pdf>
- circa 2008 – Yahoo!: web scale search indexing
Hadoop Summit, HUG, etc.
<http://developer.yahoo.com/hadoop/>
- circa 2009 – Amazon AWS: Elastic MapReduce
Hadoop modified for EC2/S3, plus support for Hive, Pig, Cascading, etc.
<http://aws.amazon.com/elasticmapreduce/>

Spark Research Papers

- *Spark: Cluster Computing with Working Sets*

Matei Zaharia, Mosharaf Chowdhury, Michael J. Franklin, Scott Shenker, Ion Stoica
USENIX HotCloud (2010)
people.csail.mit.edu/matei

<https://eduassistpro.github.io/>

- *Resilient Distributed Datasets: A Fault-Tolerant Abstraction for In-Memory Cluster Computing*

Matei Zaharia, Mosharaf Chowdhury, Tathagata Das,
Ankur Dave, Justin Ma, Murphy McCauley, Michael J. Franklin,
Scott Shenker, Ion Stoica
NSDI (2012)

[usenix.org/system/files/conference/nsdi12/nsdi12-final138.pdf](https://www.usenix.org/system/files/conference/nsdi12/nsdi12-final138.pdf)