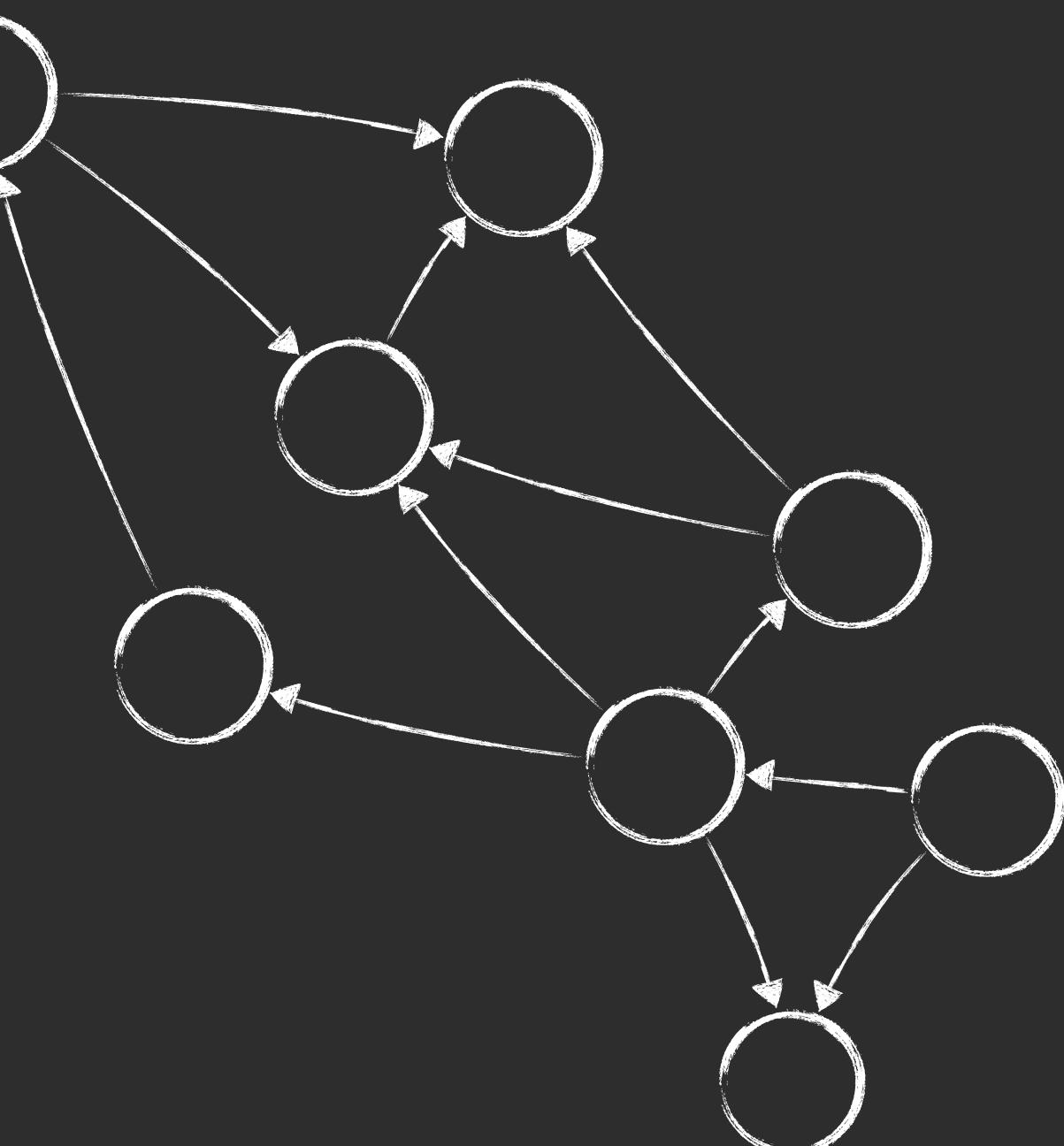


O'REILLY®

oscon

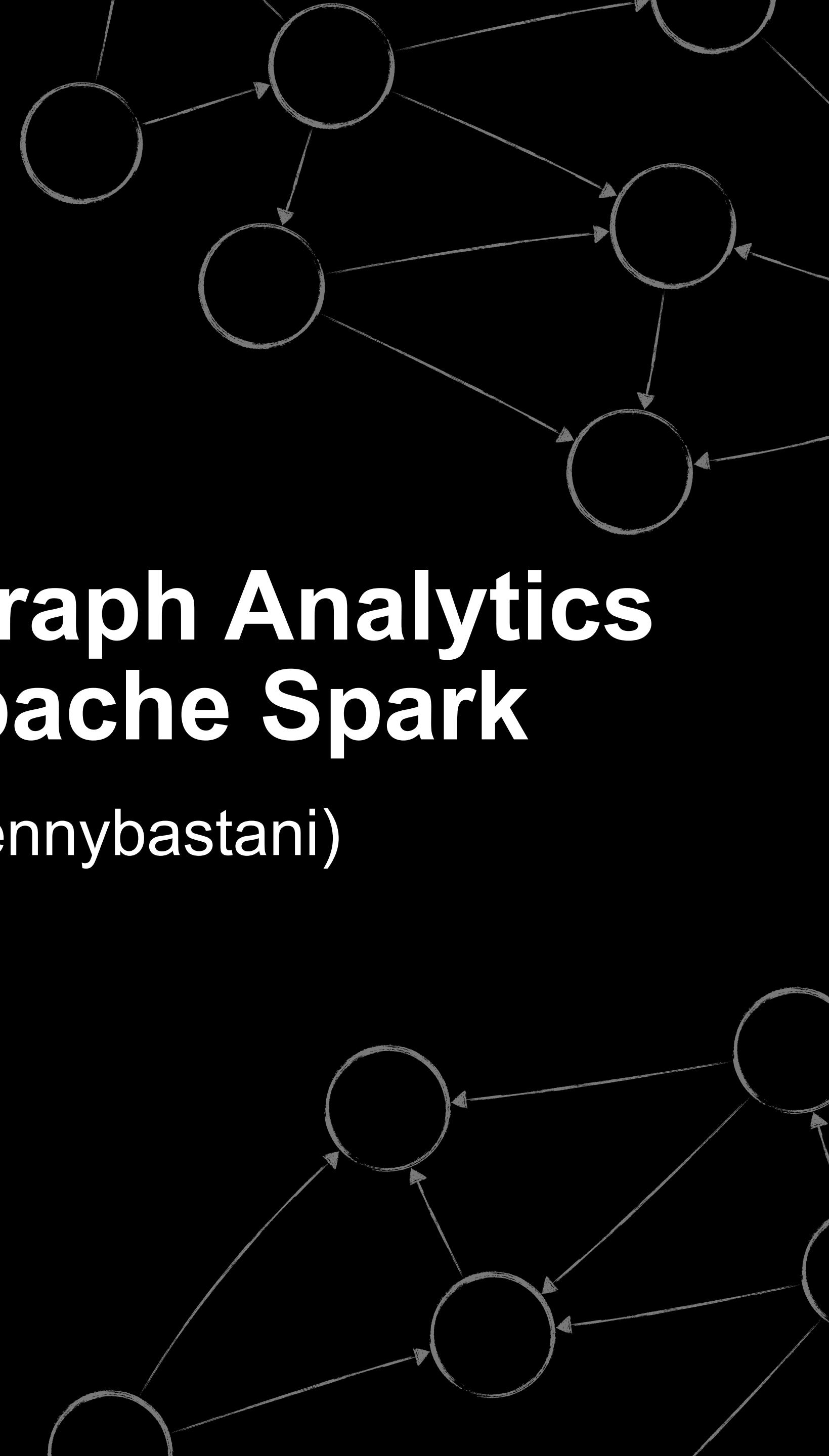


oscon.com

#oscon

Open Source Big Graph Analytics on Neo4j with Apache Spark

Kenny Bastani (@kennybastani)



Speaker Introduction

- Graph database enthusiast
- Building microservice architectures
- Lead developer at Digital Insight
- Co-author of upcoming O'Reilly book — Spring Boot Essentials

#oscon



O'REILLY®
Oscon

The Problem

It's hard to analyze graphs at scale



The importance of graph algorithms

- PageRank gave us Google
- Friend of a friend gave us Facebook
- Collaborative filtering makes Netflix recommendations awesome

#oscon



O'REILLY®
Oscon

Why is it so hard to do this stuff?

#oscon



O'REILLY®
Oscon

Enemy #1:

Relational databases store data in ways that make it difficult to extract graphs for analysis



Enemy #2:

If you still think Big Data is a buzz word
You haven't had to feel the pain of failing at it.

When you hit a wall because your data is too big

You start to see what this big data thing is all about.

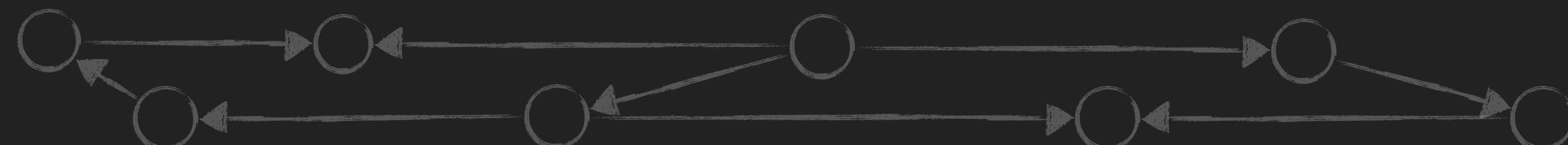


Distributed File Systems

Distributed file systems are a foundational component of big data analytics

Chops things into manageable sized blocks, usually 64MB

Spreads those blocks out across a cluster of VM resources



Hadoop MapReduce

Worth mentioning, Hadoop started this whole distributed MapReduce thing

You could translate the raw data from a CSV and turn it into a map of keys to values

Keys are distributed per node and used to reduce the values into a partitioned analysis



Graph algorithms can be evil at scale

It depends on the complexity of your graph

How many strongly connected components you have

But since some graph algorithms like PageRank are iterative

You have to iterate from one stage and use the results of the previous stage



#oscon

It doesn't matter how many nodes you have in your cluster

For iterative graph algorithms, the complexity of the graph will make you or break you

Graphs with high complexity need a lot of memory to be processed iteratively

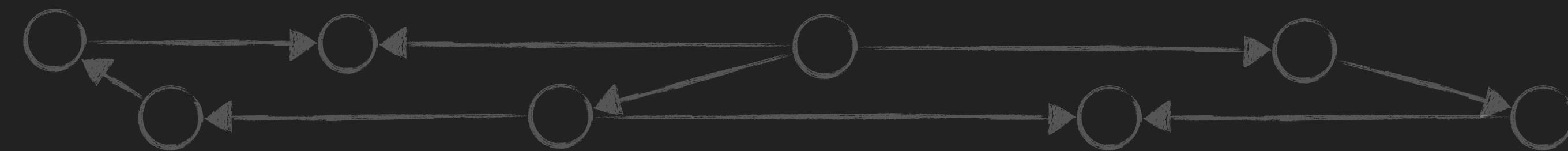


O'REILLY®
Oscon

Neo4j Mazerunner Project

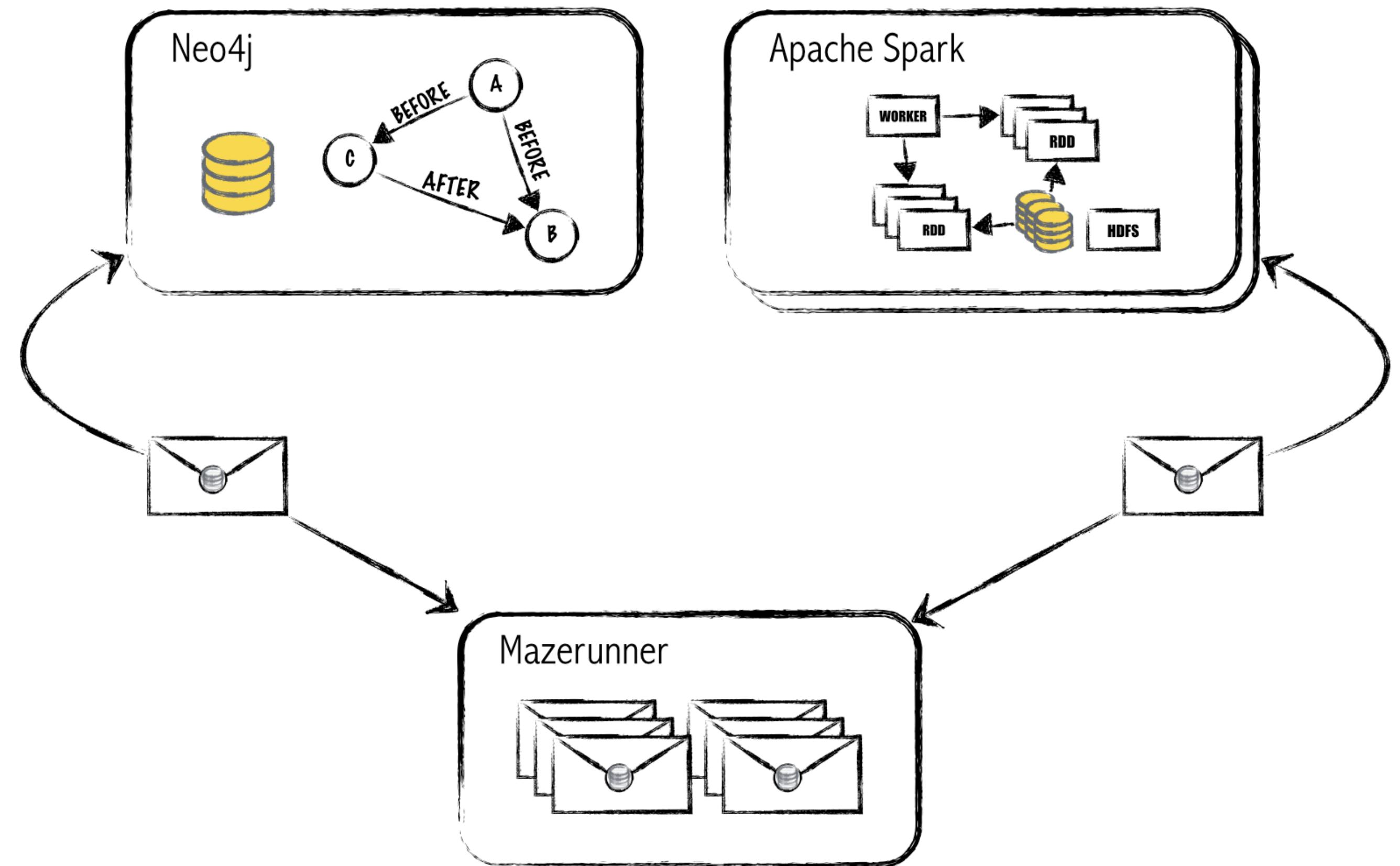
2-way Graph ETL

#oscon



O'REILLY®
Oscon

What is Neo4j Mazerunner?



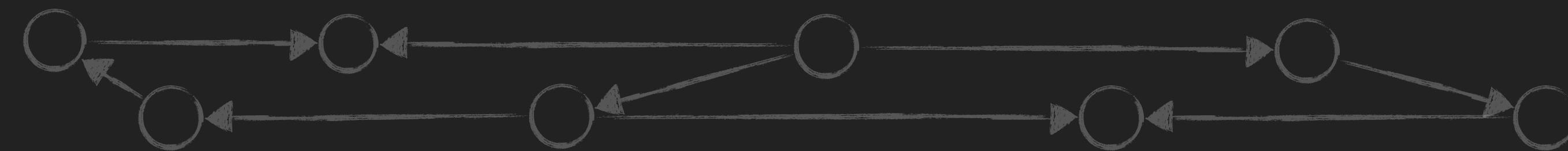
The basic idea is...

Graph databases need ETL so you can analyze your data and look it up later.



Docker

If you're not up on Docker, let me give you a quick intro.



Docker

Docker is a VM framework that lets you easily create a recipe for an image and deploy applications with ease. The idea is that infrastructure and operational complexity makes it hard for agile development of new products.



Why?

If I am an engineer on a product team, I want to choose my own software libraries and languages to solve problems.



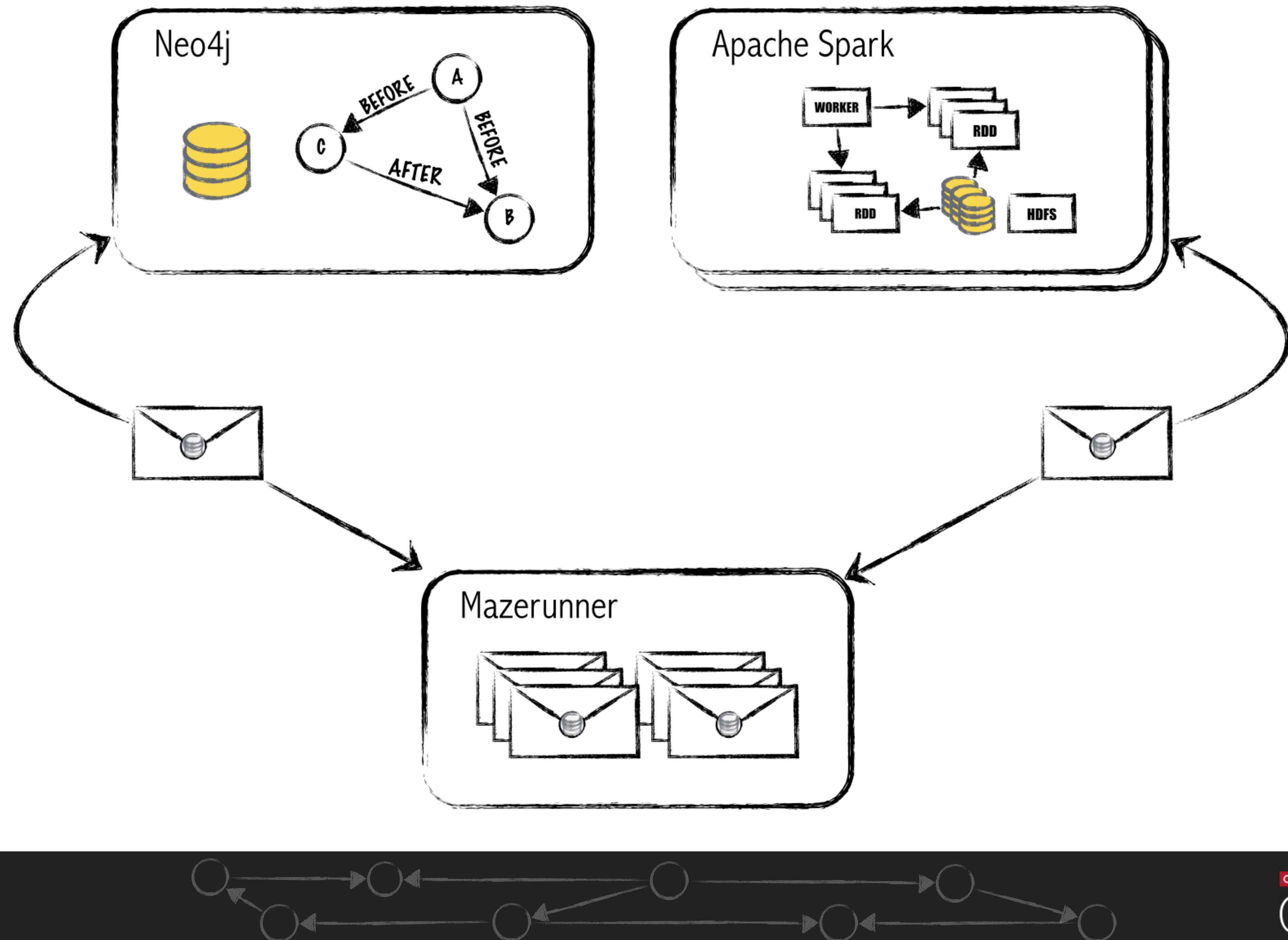
Microservices and Docker

- If want to build a new service, use whatever application framework you want. As long as you communicate over REST.
- Docker gives you the freedom to use Neo4j, or MySQL, or MongoDB or whatever application dependency you want inside your container.

#oscon



O'REILLY®
Oscon



Docker Compose

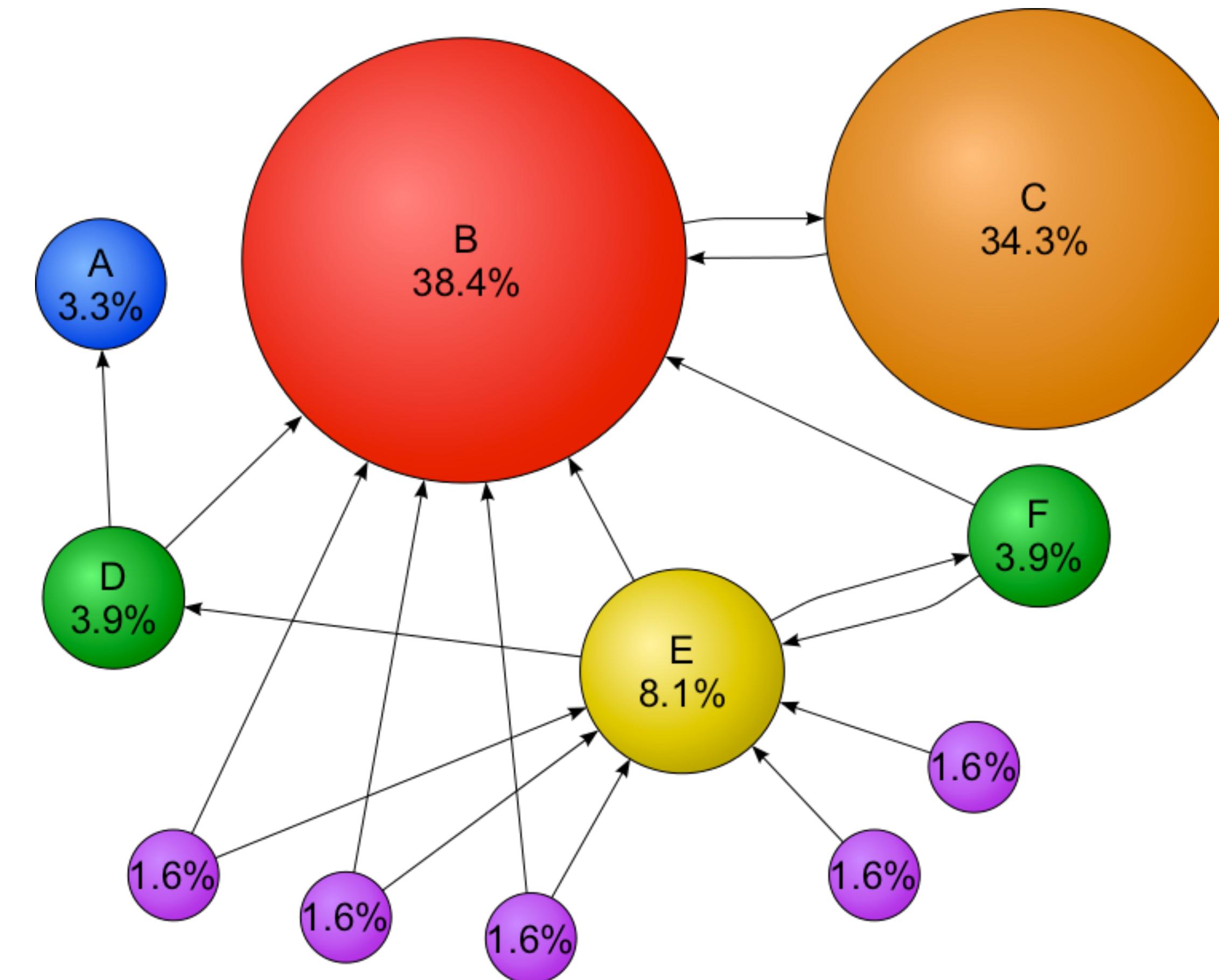
Docker Compose allows you to run multi-container applications
It uses a single YAML file



```
1  hdfs:  
2    image: sequenceiq/hadoop-docker:2.4.1  
3    command: /etc/bootstrap.sh -d -bash  
4  mazerunner:  
5    image: kbastani/neo4j-graph-analytics:latest  
6    links:  
7    - hdfs  
8  graphdb:  
9    image: kbastani/docker-neo4j:latest  
10   ports:  
11   - "7474:7474"  
12   - "1337:1337"  
13   volumes:  
14   - /opt/data  
15   links:  
16   - mazerunner  
17   - hdfs
```



PageRank

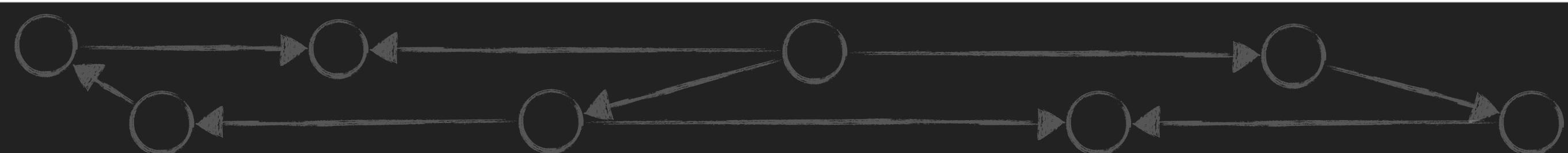
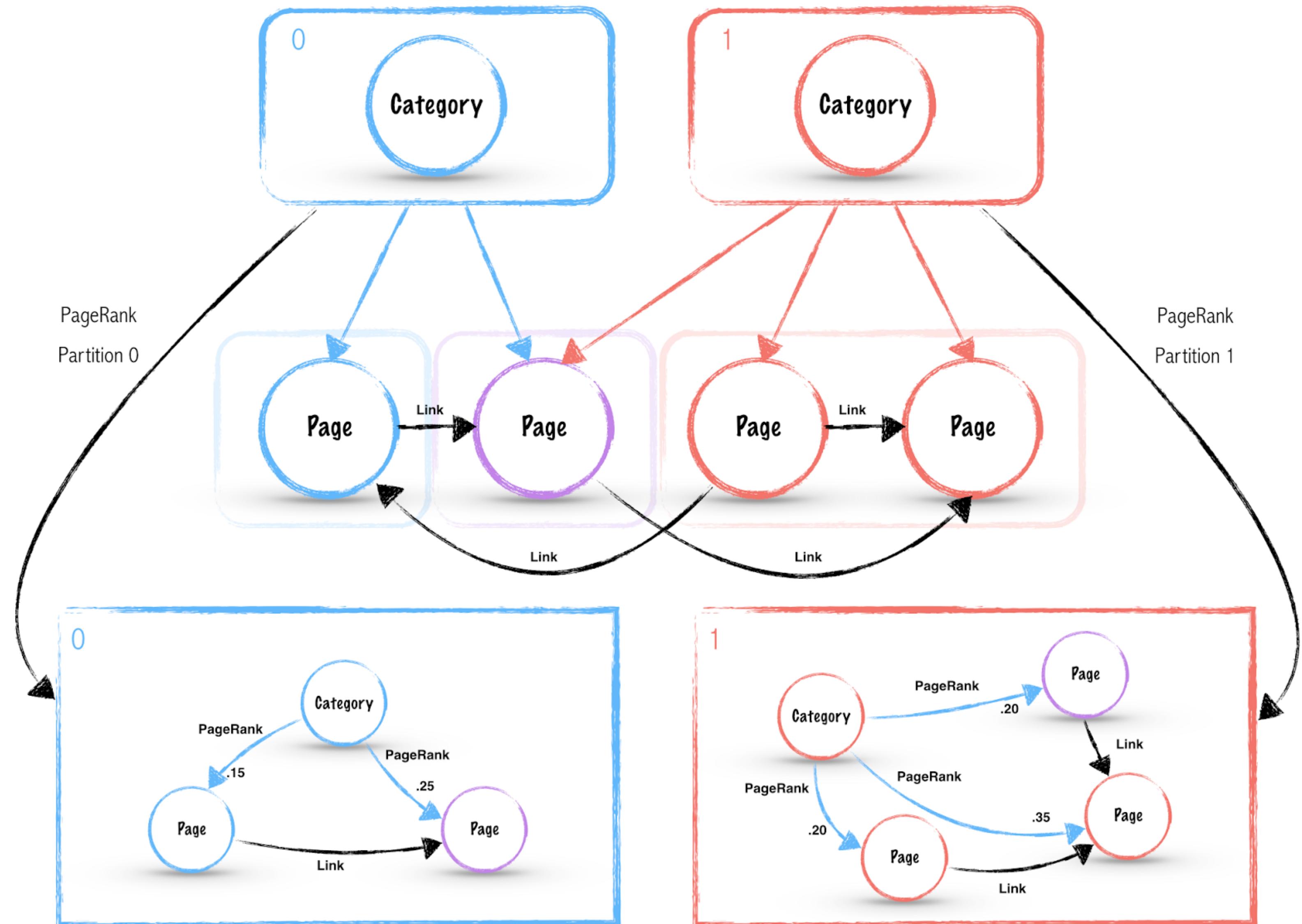


Distributed PageRank

#oscon



O'REILLY®
Oscon



Questions?

kennybastani.com

#oscon



O'REILLY®
Oscon