# Lessons Learned with Spark at the US Patent & Trademark Office
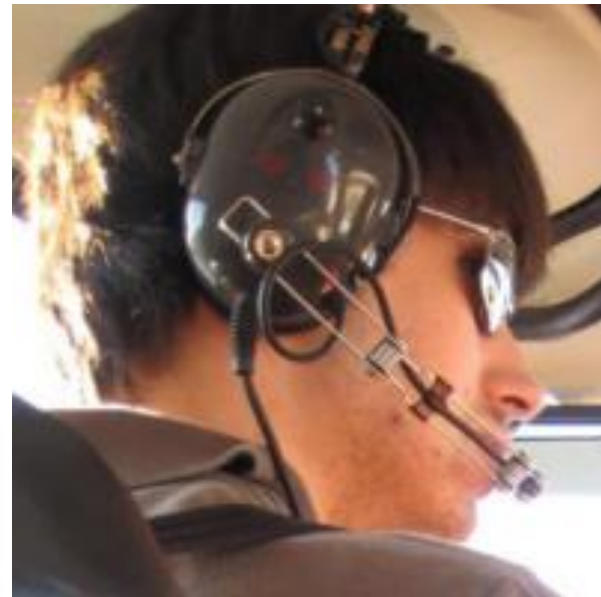
## Christopher Bradford

Big Data Architect at OpenSource Connections
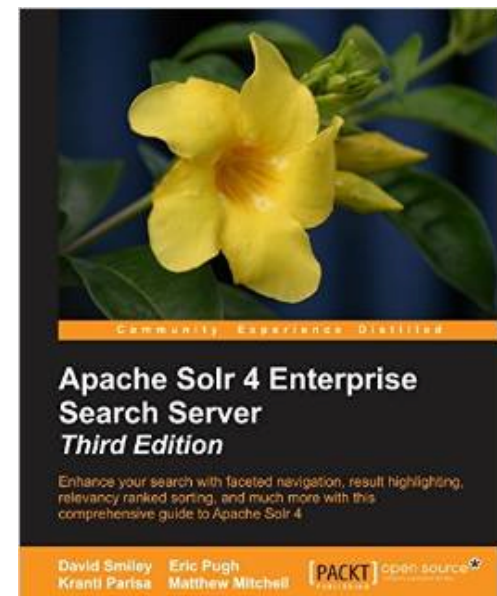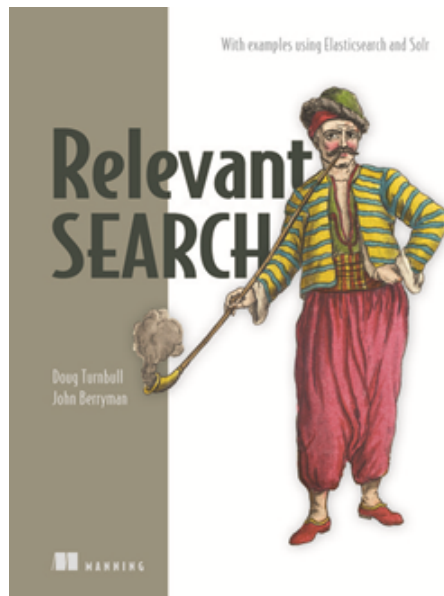
**Christopher Bradford**

Twitter: @bradfordcp
GitHub: bradfordcp

# OpenSource Connections

# Exploring Search Technologies - EST

# EST – Technology Stack

# EST – Data Loading

## CSS Ingestion (CSS2C)



## Solr Ingestion (C2S)

# EST – C2S Process



*Note: some connections are omitted for clarity*

# EST – C2S Process (Scaled Out)



*Note: some connections are omitted for clarity*

## EST – C2S Review



Did it work?

Why change it?

How could we make it better?

# EST – Old C2S Process



*Note: some connections are omitted for clarity*

# EST – Spark C2S Process



Note: some connections are omitted for clarity

# How did this work out?

Poorly

# Poor Performance

```
joinedRDD = …
joinedRDD.foreach()
  document = … // build document
  sc = new SolrConnection()
  sc.push(document)
  sc.disconnect()
// Job is done
```

# Poor Performance

```
sc = new SolrConnection()
sc.push(document)
sc.disconnect()
```

# Optimum Performance

```
joinedRDD = …
sc = new SolrConnection()
joinedRDD.foreach()
  document = … // build document
  sc.push(document)
sc.disconnect()
// Job is done
```

```
joinedRDD = …
joinedRDD.foreachPartition()
  sc = new SolrConnection()
  partition.foreach()
    document = … // build document
    sc.push(document)
  sc.disconnect()
// Job is done
```

*Almost*

# The Solution!

```
joinedRDD = …

joinedRDD.mapPartitions()

  sc = new SolrConnection()

  partition.foreach()

    document = … // build
document

    sc.push(document)

  sc.close()

  return partition.rows

.collect()
```
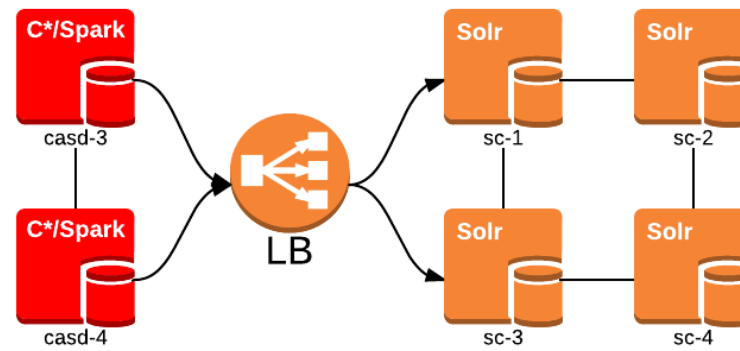
```
joinedRDD = …

joinedRDD.mapPartitions()

  sc = new SolrConnection()

  partition.foreach()

    document = … // build
document

    sc.push(document)

  sc.close()

  return partitions.rows.count

.collect()
```
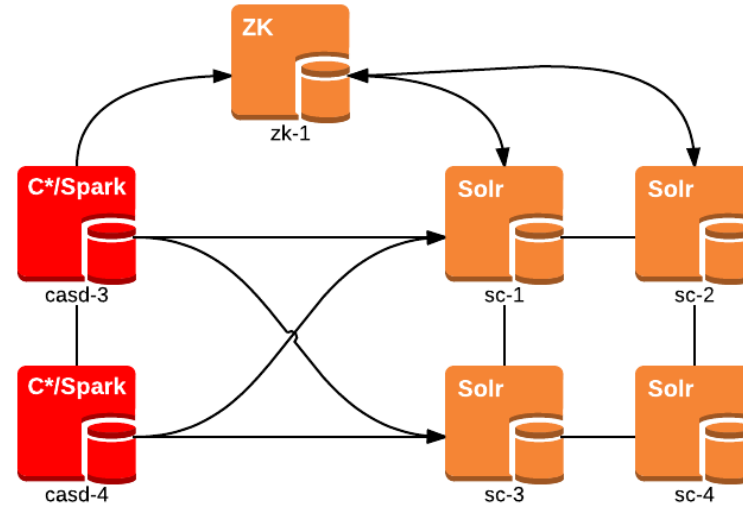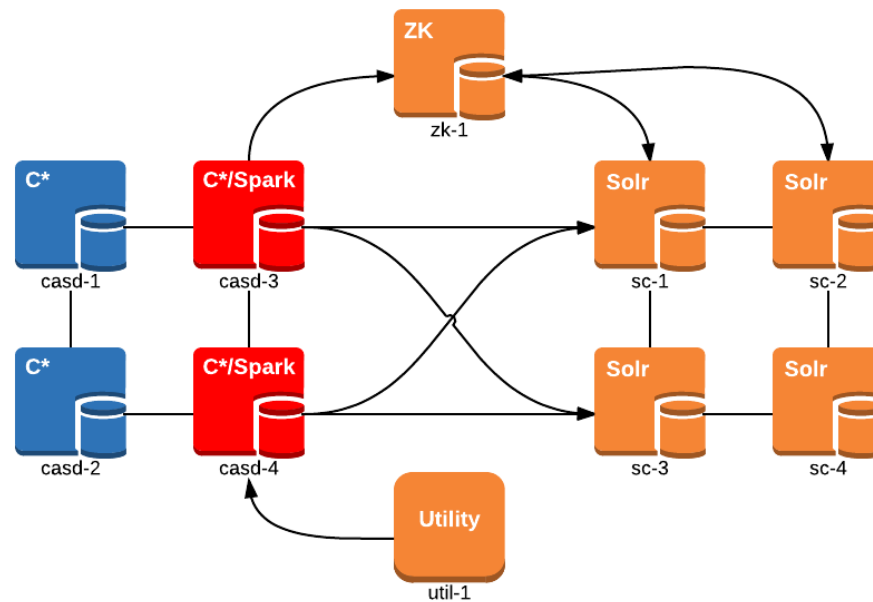
# Results?

# Solr Indexing

# *Better* Solr Indexing

*Note: some connections are omitted for clarity*

# EST – Spark C2S Process v2



Note: some connections are omitted for clarity

# Success?

## YUP

5x faster than the original C2S process (with optimizations)

# What's Next?

- Optimization of the C2S Spark job
- More Spark jobs
- Newer version of Spark & DSE
- Scala Spark jobs instead of Java