# Serverless BM25 Search and BERT Reranking

Mayank Anand, Jiarui Zhang, Shane Ding, Ji Xin, and Jimmy Lin
University of Waterloo

If a server listens in a forest and there was no one there to start it, does it really exist?

# DESIRES

A systems-oriented biennial conference, complementary in its mission to the mainstream Information Access and Retrieval conferences, emphasizing the *innovative technological aspects* of search and retrieval systems.

It gathers researchers and practitioners from both *academia and industry* to discuss the latest innovative and visionary ideas in the field.

# Computing without Servers, V8, Rocket Ships, and Other Batsh*t Crazy Ideas in Data Systems

**Jimmy Lin**
David R. Cheriton School of Computer Science
University of Waterloo

Wednesday, August 29, 2018

# We've done it!
Serverless BM25 Search and BERT Reranking

# Servers

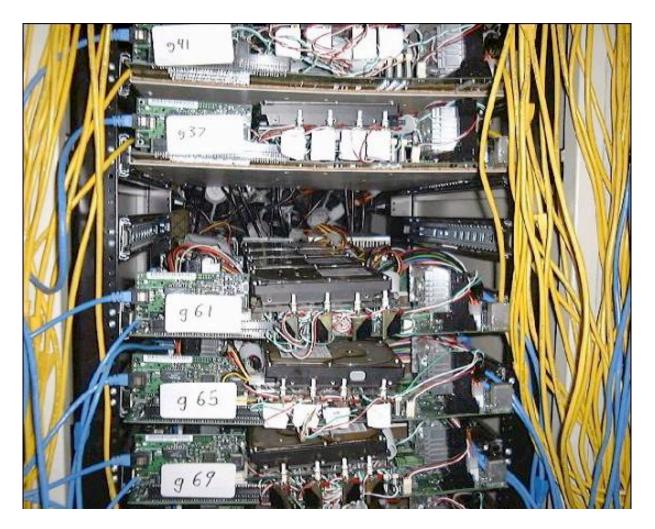The most fundamental building blocks of IR systems
(both software and hardware)

In the beginning…

# "Google" Circa 1997  (google.stanford.edu)



Slide from Jeff Dean, WSDM 2009 keynote

Google

# "Corkboards" (1999)



**Slide from Jeff Dean, WSDM 2009 keynote**

Google

# Google Data Center (2000)



**Slide from Jeff Dean, WSDM 2009 keynote**

# Google Data Center (3 days later)



**Slide from Jeff Dean, WSDM 2009 keynote**

# Servers

… in the cloud

# Early 2001: In-Memory Index



Slide from Jeff Dean, WSDM 2009 keynote
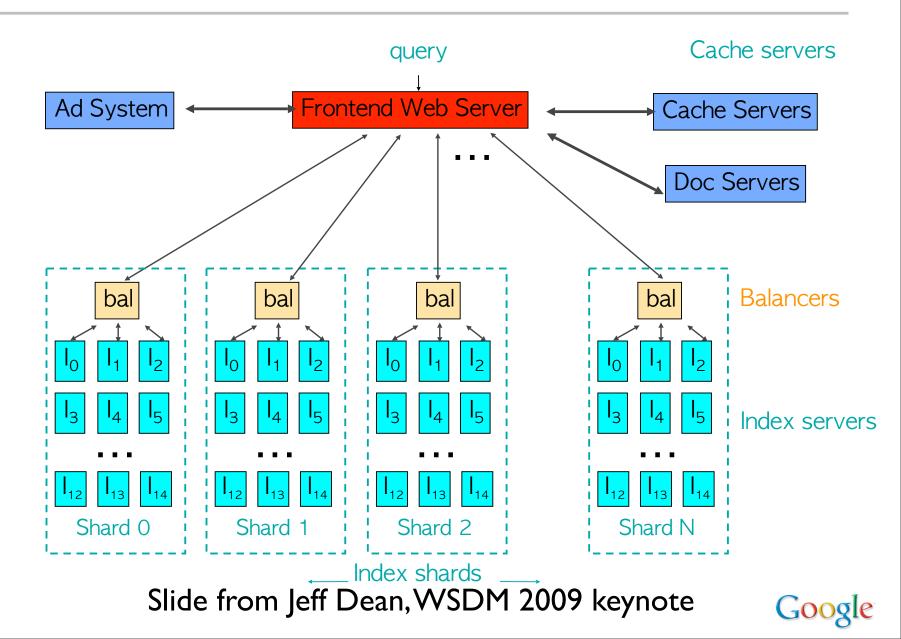
# Challenges Remain
## (Especially if you're not Google)

Always on!
Scaling up… scaling down…
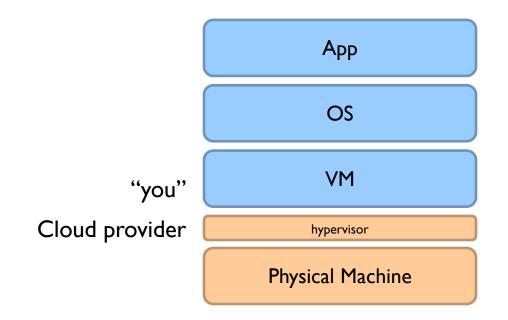Scaling to zero?

# Server**less**

Preliminaries

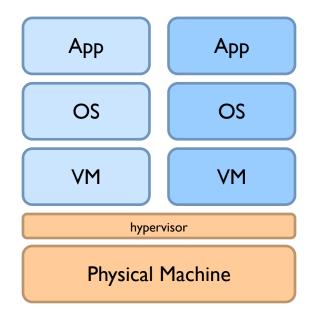Cloud computing allows us to explore different abstractions and organizations of computing

(trend towards disaggregation)

# In the beginning…

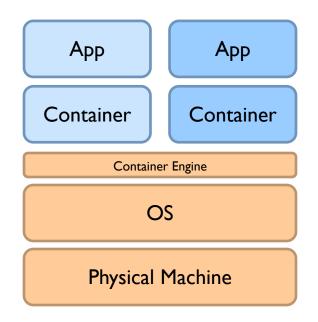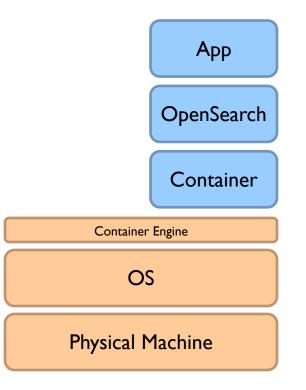# Infrastructure as a Service (IaaS)

App

OS

VM

"you"

Cloud provider

hypervisor

Physical Machine

# Multi-Tenancy

# Containers >> VMs

# Typical Stack

App

OpenSearch

Container

Container Engine

OS

Physical Machine

# Multi-Container Orchestration

| App | App | App | App |
|-----|-----|-----|-----|
| OpenSearch | OpenSearch | OpenSearch | OpenSearch |
| Container | Container | Container | Container |

| Container Engine | Container Engine |
|------------------|------------------|
| OS | OS |
| Physical Machine | Physical Machine |

| App | App | App | App |
|-----|-----|-----|-----|
| OpenSearch | OpenSearch | OpenSearch | OpenSearch |
| Container | Container | Container | Container |

| Container Engine | Container Engine |
|------------------|------------------|
| OS | OS |
| Physical Machine | Physical Machine |

# Platform as a Service

| App | App | App | App |
|-----|-----|-----|-----|

| OpenSearch as as Service |
|---|

| OpenSearch | OpenSearch | OpenSearch | OpenSearch |
|---|---|---|---|

| Container | Container | Container | Container |
|---|---|---|---|

| Container Engine | Container Engine |
|---|---|

| OS | OS |
|---|---|

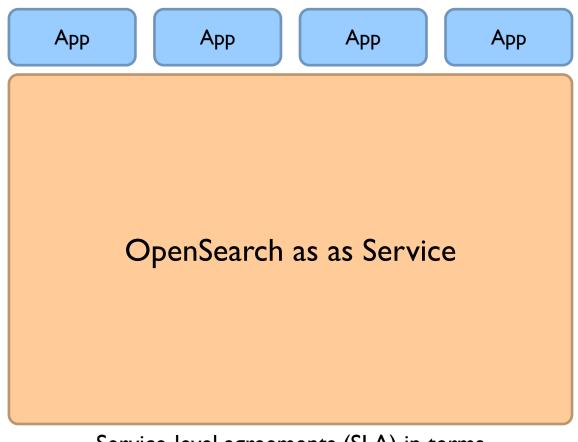| Physical Machine | Physical Machine |
|---|---|

# Platform as a Service

App  App  App  App

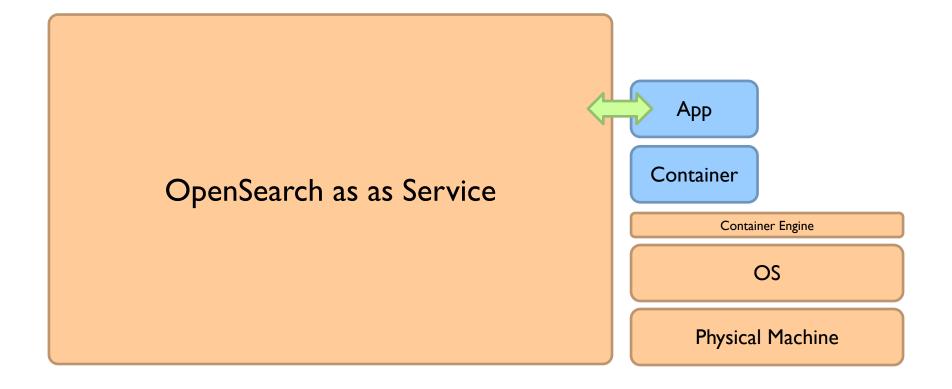OpenSearch as as Service

Service-level agreements (SLA) in terms
of latency, capacity, scalability, etc.

# What about the apps?

OpenSearch as as Service

App

Container

Container Engine

OS

Physical Machine

# Scaling out the apps…

OpenSearch as as Service

App

App

Container **?** Container

Container Engine

OS

Physical Machine

# Operational Semantics of Computing…

$$\frac{\langle E, s \rangle \Rightarrow V}{\langle L := E, s \rangle \longrightarrow (s \uplus (L \mapsto V))}$$

**if** the expression $E$ in state $s$ reduces to value $V$,

**then** the program $L := E$ will update the state $s$ with the assignment $L = V$
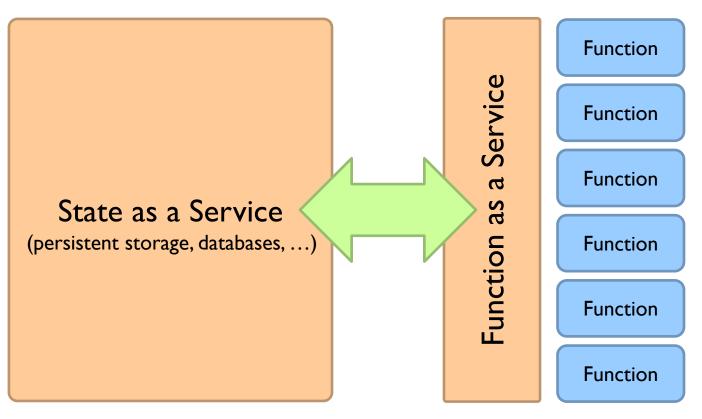
State          State as a Service
               (persistent storage, databases, message queues, …)

Transitions    Function as a Service
               (blocks of code with a well-defined entry and exit points)

# Computing without Servers

Developer: Write a bunch of functions
typically – read state, perform
some computation, update state

State as a Service
(persistent storage, databases, …)

Function as a Service

Function

Function

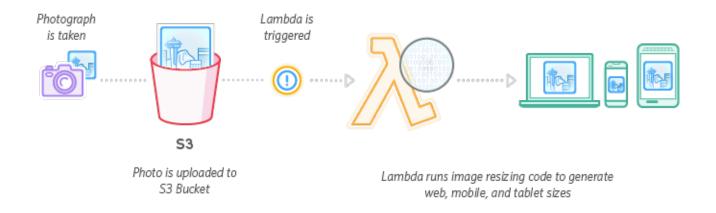Function

Function

Function

Function

Cloud provider handles everything else!
allocation of resources for execution, scaling up and down,
load balancing, cleaning up, etc.
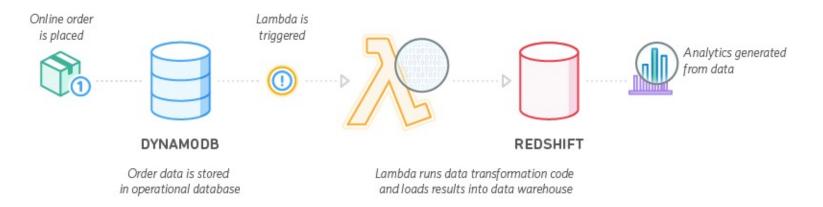
Cost model: pay per function invocation

# Serverless Examples

# Serverless Examples

Serverless computing isn't actually computing without servers…

It's just that servers become someone else's problem!

If a server listens in a forest and there was no one there to start it, does it really exist?

What would a serverless
search engine look like?

# Serverless Search

Client

API Gateway

Query Evaluation

Where do the postings go?

Lambda function as a service

term₁    term₃

term₂

DynamoDB key-value store
(postings lists)

Matt Crane and Jimmy Lin. An Exploration of Serverless Architectures for Information Retrieval. *ICTIR 2017*.

# Serverless Search

**Client**

**API Gateway**

**Query Evaluation**

**Query Evaluation**

**Query Evaluation**

**Lambda function as a service**

term$_1$   term$_3$

term$_2$

**DynamoDB key-value store**
**(postings lists)**

Matt Crane and Jimmy Lin. An Exploration of Serverless Architectures for Information Retrieval. *ICTIR 2017*.

# Serverless Search



Client
Client
Client
Client
Client

API Gateway

Query Evaluation
Query Evaluation
Query Evaluation
Query Evaluation
Query Evaluation
Query Evaluation

. . .

Lambda function as a service

$term_1$   $term_3$

$term_2$

. . .

DynamoDB key-value store
(postings lists)

Matt Crane and Jimmy Lin. An Exploration of Serverless Architectures for Information Retrieval. *ICTIR 2017*.

I got 99 problems but scaling ain't one!

# How well does serverless search work?

tl; dr – *not very well…*



Query Latency

Gov2 collection (25m docs), topic 701-850

Matt Crane and Jimmy Lin. An Exploration of Serverless Architectures for Information Retrieval. *ICTIR 2017.*
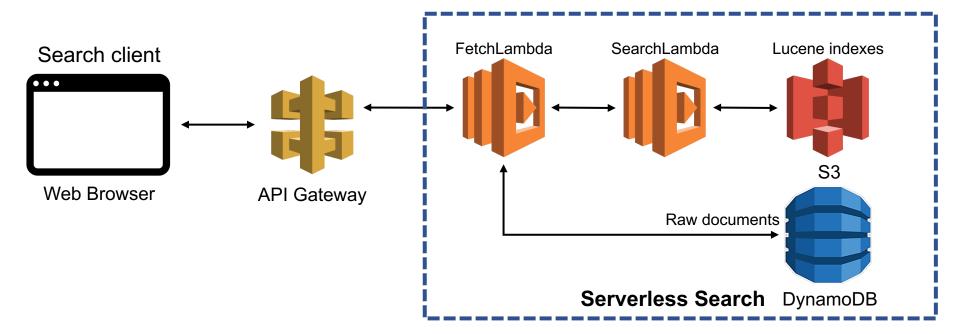
# Take 2
## (This work)

# Technical Highlights

## tl;dr – Serverless Lucene

Index structures stored on S3
Query evaluation in AWS Lambda

*Minimal modifications to "vanilla" Lucene!*

## What about query latency?

"cold instance" startup – loads indexes in memory
"warm instance" execution – indexes already in memory

*No different from "standard" in-memory search!*

Search client

Web Browser

API Gateway

FetchLambda

SearchLambda

Lucene indexes

S3

Raw documents

**Serverless Search**

DynamoDB

Serverless Lucene isn't enough!

# Retrieve + Rerank

Select some promising texts

Rerank selected texts

# Retrieve + Rerank

Select some promising texts

*Serverless!*

Lucene

+

monoBERT
(cross-encoder reranker)

*Serverless?*

# Yea, so what about serverless BERT?

(Good thing it's embarrassingly parallel!)

# Technical Highlights
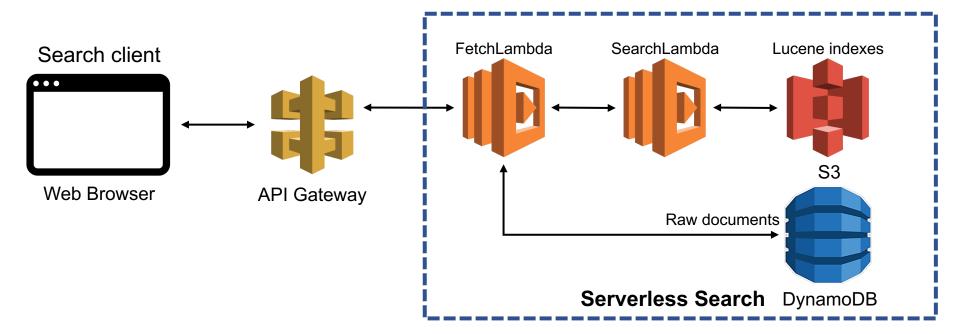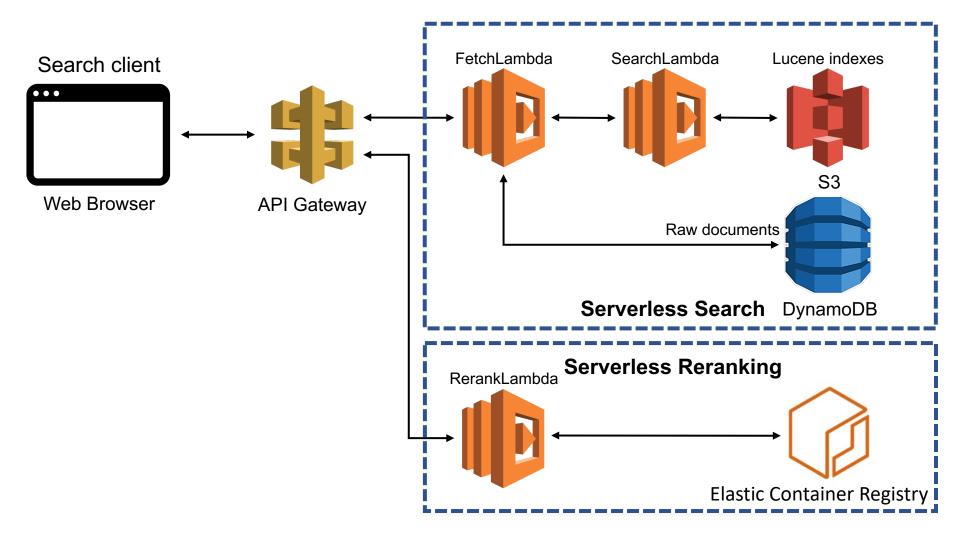
tl;dr – Serverless BERT inference in AWS Lambda
"Early-Exit" optimizations* to reduce inference latency
Main technical challenge is model size – solved by Elastic Container Registry

*It's SMOP!*

The bad: CPU inference
The good: Massive parallelism (100-way parallelism)
The ugly: Image packaging is a pain!

* Xin et al. Early Exiting BERT for Efficient Document Ranking. *SustaiNLP 2020*.

Search client

Web Browser

API Gateway

**Serverless Search**

FetchLambda

SearchLambda

Lucene indexes

S3

Raw documents

DynamoDB

**Serverless Reranking**

RerankLambda

Elastic Container Registry

# Results
## (MS MARCO passage ranking)

Setup: reranking 1000 hits from Lucene with monoBERT

| Stage | Latency (s/Q) | | | Cost |
| --- | --- | --- | --- | --- |
| | Mean | P50 | P99 | (/100Q) |
| BM25 | 0.65 | 0.65 | 0.92 | $0.022 |
| DynamoDB Fetch | 0.95 | 0.96 | 1.06 | - |
| BERT reranking | 11.21 | 10.64 | 17.90 | $15.90 |
| End to end | 12.81 | 12.24 | 19.35 | $16.00 |
| BERT reranking (V100) | 26.21 | 25.52 | 36.64 | $2.20 |

(We confirmed that effectiveness is the same for serverless vs. server-based deployments)

Wait, how can CPU be faster than GPU?
7-8× more expensive = breakeven at 85%-90% idle

# Objections

It's still too slow! (12s end-to-end)
It's still too expensive! ($0.16 per query)

Agreed… but this is only the beginning!
Serverless infrastructure will become more efficient.
Lots of neural inference acceleration techniques to try.

The price of minimal management overhead and infinite scaling?

Serverless search is worth considering?!